

Ridge Scale Aligned Diffusion for Identity Preserving and Style Controllable Fingerprint Synthesis

Anonymous CVPR submission

Paper ID *****

Abstract

001 *Training robust fingerprint recognizers is constrained by*
 002 *privacy regulations and limited intra-class diversity in public*
 003 *datasets. We propose a controllable fingerprint synthesis*
 004 *framework built on Stable Diffusion with ControlNet*
 005 *and a Multi-IP-Adapter derived from IP-Adapter to gener-*
 006 *ate structure-consistent, style-diverse fingerprints condi-*
 007 *tioned on a content image. Our design supports two modes:*
 008 *(i) augmenting a real identity with diverse sensor styles,*
 009 *and (ii) privacy-motivated synthesis using virtual identi-*
 010 *ties sampled from a DDPM prior to reduce direct depen-*
 011 *dence on real identity data. To enable fine control, we*
 012 *introduce ridge-scale normalization and dual-mask spa-*
 013 *tial injection to better separate ridge regions from back-*
 014 *ground during generation. Experiments demonstrate high*
 015 *visual fidelity, consistent retention of content-guided struc-*
 016 *tural cues under style transfer, and improved downstream*
 017 *recognition. In particular, joint training with real and our*
 018 *synthetic data boosts ViT TAR@FAR=0.1% from 80.08%*
 019 *to 91.28%, exceeding joint-training results obtained with*
 020 *FPGAN-Control and PrintsGAN.*

021 1. Introduction

022 Deep fingerprint recognition models require large-scale
 023 data with sufficient *intra-class* variation (multiple impres-
 024 sions per identity under different sensors, pressure, and cap-
 025 ture conditions). However, collecting and sharing real finger-
 026 prints is heavily restricted due to privacy constraints:
 027 biometric leakage is irreversible, and regulations often pre-
 028 vent large-scale acquisition and release. As a result, public
 029 datasets are limited in both scale and diversity, leading to
 030 overfitting and poor cross-sensor generalization.

031 Generative models offer a promising alternative by syn-
 032 thesizing realistic fingerprints to expand diversity without
 033 additional data collection, including GANs [1] and diffu-
 034 sion models such as DDPM [2]. Prior GAN-based finger-
 035 print synthesis methods can generate high-quality im-

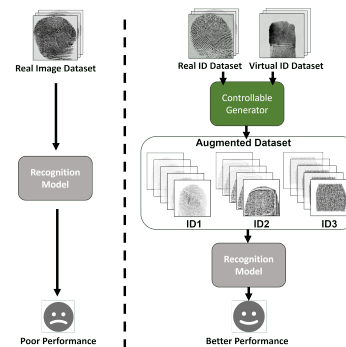


Figure 1. Illustration of structure-consistent augmentation: synthetic fingerprints preserve identity while introducing diverse styles, addressing limited intra-class diversity.

036 ages [3, 4], but structure-consistent augmentation remains
 037 challenging due to limited controllability and identity du-
 038 plication at scale [3]. Diffusion-based fingerprint synthesis
 039 has recently shown strong ridge-detail modeling [5–9], yet
 040 existing approaches still lack fine-grained control over (i)
 041 cross-sensor style transfer, (ii) global shape variation (e.g.,
 042 plain vs. rolled), and (iii) explicit separation of ridge fore-
 043 ground from background.

044 We propose a controllable fingerprint synthesis frame-
 045 work based on Stable Diffusion [10], combining Control-
 046 Net [11] for structure-consistent generation and a *Multi-*
 047 *IP-Adapter* derived from IP-Adapter [12] for disentangled
 048 style/shape/text control. The framework supports (1) **aug-**
 049 **mentation mode** that generates diverse impressions from
 050 a real identity, and (2) **privacy-motivated mode** that uses
 051 DDPM-sampled virtual identities as control inputs [2]. To
 052 stabilize cross-style fusion, we further introduce ridge-scale
 053 normalization (ridge frequency estimation follows standard
 054 orientation-field based formulations [13]) and dual-mask
 055 spatial injection.

056 As illustrated in Fig. 1, our framework produces multiple
 057 diverse impressions for each identity, improving intra-class
 058 diversity for robust recognition training.

059 **Contributions.** (1) A unified diffusion-based pipeline
 060

060 for **structure-consistent**, style-diverse fingerprint synthesis
 061 with real-identity augmentation and **privacy-motivated**
 062 virtual-identity generation modes. (2) Multi-domain conditioning
 063 (style/acquisition-type/text) with decoupled cross-attention
 064 and dual-mask injection for spatial control. (3) Ridge-scale
 065 normalization to align ridge frequency between the content
 066 and the **target** style. (4) Extensive evaluation showing
 067 **consistent retention of content-guided structural cues**
 068 and **substantial** recognition gains (e.g., ViT TAR@FAR=0.1%:
 069 80.08% \rightarrow 91.28% with joint training).
 070

071 2. Related Work

072 2.1. Fingerprint Synthesis

073 Classical simulators (e.g., SFinGe) generate fingerprints via
 074 handcrafted ridge models and noise, but are limited in realism
 075 and controllability [14]. GAN-based methods improve visual
 076 quality for fingerprint synthesis [15–18], yet often suffer
 077 from unstable training and limited coverage at scale; structure-
 078 consistent augmentation remains challenging due to issues such
 079 as identity duplication and constrained controllability [3, 4].
 080 Recent diffusion-based approaches show superior ridge detail
 081 modeling and have been explored for fingerprint synthesis and
 082 augmentation [5–9], but many lack explicit mechanisms for
 083 cross-sensor style control, global shape variation, and
 084 foreground/background disentanglement.
 085

086 2.2. Controllable Generation and Style Transfer

087 Diffusion models enable high-fidelity synthesis [2, 10], but
 088 text-only prompting is often insufficient for preserving precise
 089 fingerprint structure. ControlNet injects structural conditions
 090 (e.g., edges) to anchor spatial layout [11], while IP-Adapter
 091 transfers reference-image priors via decoupled attention [12].
 092 Related controllable diffusion plugins such as T2I-Adapter
 093 further demonstrate that lightweight adapters can provide
 094 strong controllability while keeping the diffusion backbone
 095 frozen [19]. Fingerprint generation imposes stricter
 096 constraints than generic style transfer: small geometric
 097 distortions can corrupt ridge topology and minutiae. Our
 098 method combines ControlNet-based structural anchoring [11]
 099 with multi-domain IP-Adapter conditioning [12], augmented
 100 by ridge-scale normalization (based on standard ridge/
 101 orientation analysis [13]) and dual-mask injection to
 102 preserve ridge integrity under strong style changes.

103 3. Proposed Method

104 Our proposed framework is built on Stable Diffusion [10]
 105 and is designed to achieve structure-consistent, style diversity,
 106 and spatially controllable synthesis (Fig. 2). The diffusion
 107 process is guided by ControlNet [11], which enforces
 108 structural consistency using a Canny edge representation

Algorithm 1: Controllable Structure-Consistent Fingerprint Synthesis Pipeline

Input: Content fingerprint image C (real or virtual), style image S , text prompt T , fingerprint mask M_{fg} , background mask M_{bg}

Output: Synthetic fingerprint image \hat{I}

Step 1: Ridge-Scale Normalization;

Compute average ridge widths w_C and w_S ; resize and pad C with ratio $\gamma = \frac{w_S}{w_C}$.

Step 2: Structure-Consistent Conditioning;

Extract a structural representation E_C from C using a Canny edge detector; condition Stable Diffusion with ControlNet on E_C .

Step 3: Multi-Domain Feature Extraction;

Extract style features F_{sty} from S (style encoder); extract shape features F_{sh} (shape encoder); encode text prompt to F_{txt} (text encoder).

Step 4: Decoupled Cross-Attention Fusion;

Fuse features via Eq. (1).

Step 5: Spatially-Guided Injection;

Apply dual-mask injection via Eq. (2).

Step 6: Image Synthesis;

Run the denoising U-Net conditioned on $\{F_{sty}, F_{sh}, F_{txt}\}$ and ControlNet guidance to output \hat{I} .

tation of the input fingerprint. To flexibly control visual appearance, a *Multi-IP-Adapter* derived from IP-Adapter [12] injects style, shape, and text features through decoupled cross-attention, ensuring style diversity while maintaining content-guided structural cues. Spatial controllability is further enhanced with a dual-mask strategy and ridge-scale normalization, which separately condition the fingerprint and background regions and align ridge widths for seamless fusion. We estimate ridge width by patch-wise ridge spacing using an orientation-field based approach [13]; details are provided in the supplementary material. All input and style images are preprocessed by isotropic resizing and white padding to 1024×1024 , ensuring consistent spatial alignment.

The overall generation algorithm is presented in Algorithm 1.

3.1. Structure-Consistent Conditioning

To ensure Structure-Consistent Conditioning, we employ ControlNet [11] with Canny edge detection as a strong structural constraint during the diffusion process. The Canny-processed control image captures ridge flow and topology with thin boundaries that align well with ControlNet’s conditioning. In our implementation, Canny thresholds are set to 50/200 with an aperture size of 3.

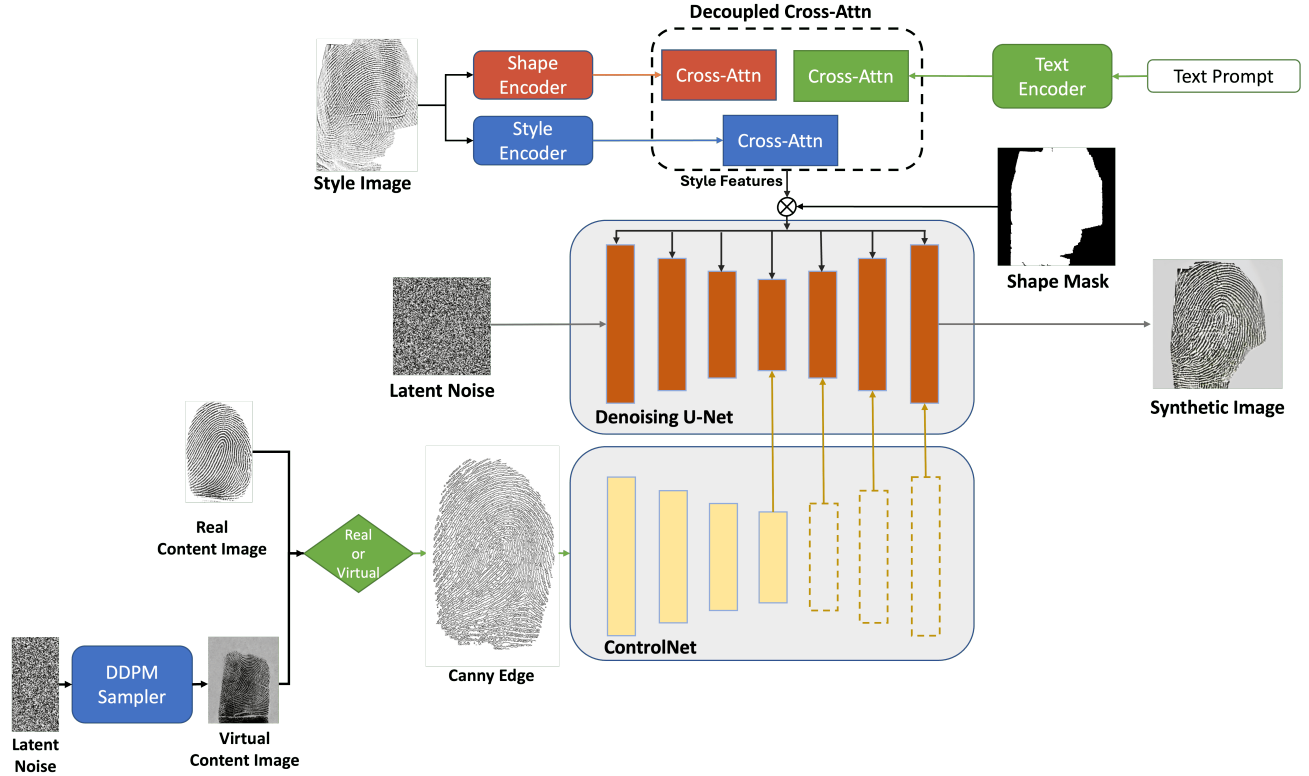


Figure 2. Overview of our proposed method. The lower part preserves content-guided ridge structure via ControlNet and U-Net conditioned on real or virtual content images. The upper part injects style, shape, and text cues through decoupled cross-attention, guided by a shape mask, to generate structurally faithful and style-controllable fingerprints.

Our framework supports two complementary input modes. In the *augmentation mode*, real fingerprint images are directly used as ControlNet inputs. In the *privacy-motivated mode*, we employ a DDPM-based sampler [2] to synthesize virtual identities, which are then used as control inputs for ControlNet. The DDPM sampler is trained from scratch on 5,000 samples from NIST SD 302 [20].

3.2. Multi-IP-Adapter for Fingerprint Style Generation

The IP-Adapter [12] injects visual priors from reference images or text into diffusion models. We extend it to disentangle *style*, *shape*, and *text* control for fingerprint synthesis (Fig. 2).

3.2.1. Style Control

We adopt a CLIP-based image encoder (e.g., OpenCLIP-ViT-H-14 [21]) as the style encoder without additional fine-tuning.

3.2.2. Text-Guided Style Control

Text prompts are encoded by a CLIP text encoder [21]. Among all components, only the text-guided style adapter is fine-tuned on NIST SD 302 [20] (5,000 samples, 100

epochs, $\text{lr}=1 \times 10^{-4}$, AdamW, batch size=8 on a single NVIDIA A6000 GPU).

3.2.3. Shape Control

To explicitly control global geometry (e.g., plain vs. rolled), we extract *shape* features from the style image using a CLIP image encoder and inject them as a separate conditioning stream, preventing the ControlNet structure from dominating shape and enabling type-specific synthesis.

3.2.4. Decoupled Cross-Attention

We fuse text, style, and shape conditions through independent cross-attention branches:

$$Z_{\text{new}} = \text{Attn}(Q, K^t, V^t) + \lambda_1 \text{Attn}(Q, K^{\text{sty}}, V^{\text{sty}}) + \lambda_2 \text{Attn}(Q, K^{\text{sh}}, V^{\text{sh}}). \quad (1)$$

We set $\lambda_1 = \lambda_2 = 1.0$ for all experiments.

3.3. Spatially-Guided Feature Injection via Foreground/Background Masking

To enhance spatial controllability, we use a fingerprint mask M_{fg} and a background mask $M_{\text{bg}} = 1 - M_{\text{fg}}$ (mask generation details are provided in the supplementary material).

172 Given style features F_{sty} and a background embedding F_{bg} ,
173 we inject:

$$174 F_{\text{out}} = M_{\text{fg}} \odot F_{\text{sty}} + M_{\text{bg}} \odot F_{\text{bg}}. \quad (2)$$

175 4. Experimental Results

176 4.1. Visual Evaluation

177 We conducted a qualitative visual evaluation to demon-
178 strate the style control and generalization ability of our fin-
179 gerprint synthesis framework. In total, we generated fin-
180 gerprints under 11 distinct style conditions, covering three
181 acquisition types: rolled, slap, and plain. Eight styles
182 correspond to sensors seen during training (Fig. 3), all
183 from the NIST SD 302 dataset: MorphoWave_Desktop
184 (roll), ANDLN2N (roll), LIVETOUCH_QUATTRO (roll),
185 S120 (roll), Undisclosed (roll; as labeled in the origi-
186 nal dataset), Crossmatch_L_SCAN_1000PX (slap), Cross-
187 match_L_SCAN_1000PX (roll), and ANDLN2N (plain).
188 The remaining three styles represent unseen sensor domains
189 from the FVC 2004 dataset (Fig. 3)—Crossmatch V300
190 (plain), U.are.U 4000 (plain), and FingerChip FCD4B14CB
191 (plain)—which are not included in training, enabling evalu-
192 ation of cross-domain generalization.

193 For each style, we display three images: the first two
194 are synthetic fingerprints generated by our method, while
195 the third is a real fingerprint sample acquired from the cor-
196 responding sensor. This side-by-side presentation enables
197 direct visual assessment of realism, style fidelity, and struc-
198 tural consistency.

199 The generated images exhibit diverse visual characteris-
200 tics and faithfully reproduce differences in ridge contrast,
201 background noise, and image tone across sensors and ac-
202 quisition types. For styles seen during training, our model
203 closely matches the appearance of real samples and mim-
204 ics sensor-specific traits. For unseen styles from FVC
205 2004 [22], our model transfers stylistic attributes (e.g.,
206 ridge sharpness and background artifacts) without exposure
207 during training, demonstrating encouraging cross-domain
208 generalization. Notably, synthesized fingerprints preserve
209 ridge-flow from the control image while adapting to target
210 style cues.

211 4.2. Matching Analysis between Synthetic and Real 212 Content Images

213 To assess structural consistency after style transfer, we per-
214 form matching between each generated fingerprint and its
215 corresponding input content image. As shown in Fig. 4 and
216 Fig. 5, the left image in each subfigure is the input content
217 image, while the right image is the synthetic output. The
218 top row shows styles seen during training, and the bottom
219 row shows unseen styles.

220 For matching, we adopt LightGlue [23] and XFeat [24],
221 two deep learning-based feature matching methods that

Table 1. NFIQ 2.0 score (\uparrow) [25] comparison between synthetic and real images.

Sensor (Type)	Synthetic	Real	$\Delta\mu$	p -value
	$\mu \pm \sigma$	$\mu \pm \sigma$		
MorphoWave_Desktop (roll)	30.79 \pm 9.96	30.70 \pm 12.65	+0.09	0.945
ANDLN2N (roll)	40.79 \pm 7.45	39.65 \pm 5.96	+1.14	0.292
LIVETOUCH_QUATTRO (roll)	26.89 \pm 15.11	31.67 \pm 20.29	-4.78	0.236
S120 (roll)	31.51 \pm 10.87	43.38 \pm 8.61	-11.87	3.44×10^{-11}
Undisclosed (roll)	27.40 \pm 14.55	39.09 \pm 14.38	-11.69	3.40×10^{-11}
Crossmatch_L_SCAN_1000PX (roll)	37.71 \pm 12.34	48.11 \pm 20.43	-10.40	4.59×10^{-4}
Crossmatch_L_SCAN_1000PX (slap)	37.02 \pm 9.27	40.30 \pm 6.37	-3.28	5.20×10^{-13}
ANDLN2N (plain)	41.61 \pm 8.57	40.44 \pm 5.17	+1.17	0.154

Note: $\Delta\mu = \mu_{\text{syn}} - \mu_{\text{real}}$. Higher NFIQ 2.0 scores indicate better image quality.

222 leverage dense keypoints and descriptors across the full im-
223 age. Unlike conventional fingerprint matchers relying on
224 explicit minutiae detection, these methods capture more
225 global and nuanced structural information, making them
226 suitable for evaluating similarity when local textures may
227 shift under style transfer while global ridge-flow patterns
228 remain consistent

229 The results in Fig. 4 and Fig. 5 show substantial cor-
230 respondences between content and synthetic images for
231 both seen and unseen styles. Although LightGlue typically
232 yields more matches due to denser representations, both
233 matchers consistently identify meaningful structural align-
234 ment, suggesting our method retains structure-level ridge-
235 flow cues under significant style variation.

236 4.3. Assessment of Synthetic Fingerprint Quality 237 across Sensors

238 To quantitatively evaluate synthetic fingerprint quality, we
239 compare NFIQ 2.0 [25] and NIQE [26] between real and
240 synthetic images across eight sensor styles. For each sen-
241 sor, we generate synthetic images using real fingerprints as
242 style references under identical style conditions. We evalu-
243 ate 2,000 real and 2,000 synthetic images per sensor. Ta-
244 bles 1 and 2 report the mean and standard deviation of both
245 metrics. We additionally perform Welch’s two-sample t -test
246 per sensor to assess whether differences between real and
247 synthetic scores are statistically significant. Figure 6 visu-
248 alizes NFIQ 2.0 score distributions (blue: real; red: syn-
249 thetic).

250 As shown in Table 1, NFIQ 2.0 scores of synthetic fin-
251 gerprints are generally comparable to real images. For most
252 sensors (e.g., MorphoWave_Desktop (roll) and ANDLN2N
253 (plain)), mean scores are nearly identical, indicating pre-
254 served biometric quality and structural fidelity. How-
255 ever, for certain sensors including S120 (roll) and Cross-

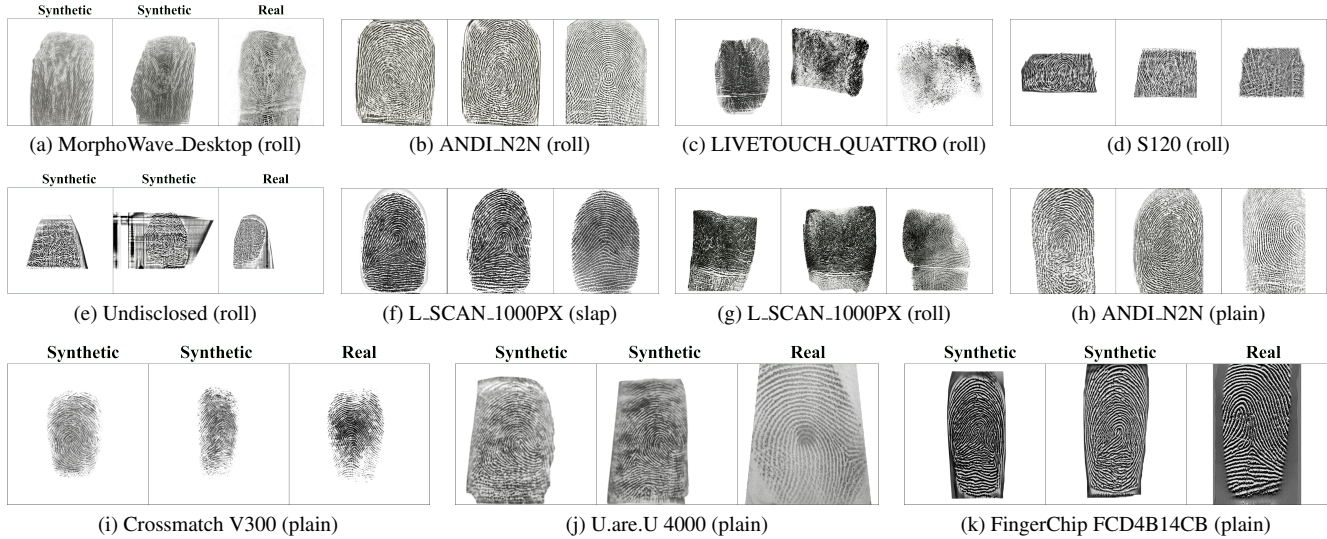


Figure 3. Visual evaluation of generated and real fingerprints. The top two rows show sensors seen during training (NIST SD 302), and the bottom row shows unseen sensors (FVC 2004). For each sensor panel, images (left→right) are two synthetic results and one real image from the same sensor style.

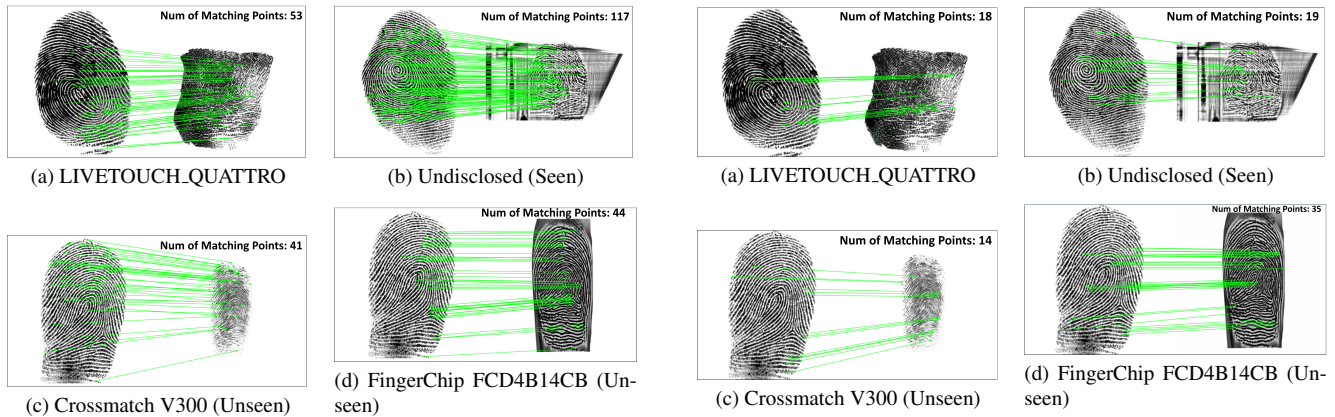


Figure 4. Matching results between the input content image (left) and the synthetic fingerprint (right) using LightGlue [23]. Green lines indicate matched feature points. Averaged over 50 test images, the number of matched points is 241.4 for seen styles and 138.1 for unseen styles.

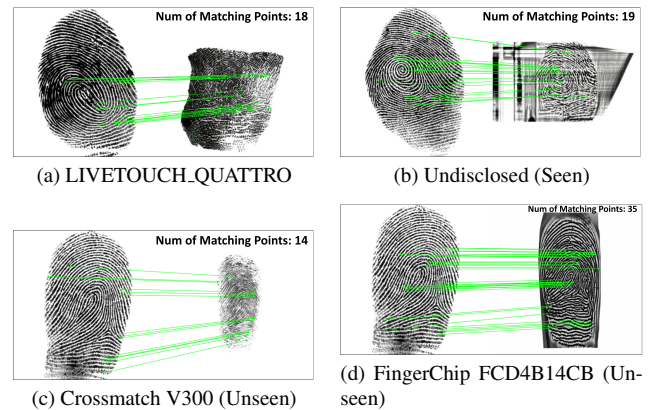


Figure 5. Matching results between the input content image (left) and the synthetic fingerprint (right) using XFeat [24]. Green lines indicate matched feature points. Averaged over 50 test images, the number of matched points is 81.5 for seen styles and 77.1 for unseen styles.

256 match_L_SCAN_1000PX (roll/slap), the mean difference
 257 ($\Delta\mu < 0$) is statistically significant, suggesting slightly
 258 lower NFIQ 2.0 for synthetic images. This gap likely re-
 259 flects sensor-specific statistics (e.g., ridge frequency, micro-
 260 contrast, or device noise) that are not fully captured, rather
 261 than a general degradation in generative quality.

262 In contrast, Table 2 shows that synthetic images often
 263 achieve lower NIQE (a proxy for perceptual cleanliness un-
 264 der natural-image statistics) than real fingerprints for most
 265 sensors. For example, MorphoWave_Desktop (roll) and
 266 Undisclosed (roll) show cleaner ridge textures and reduced

background noise. Even for Crossmatch_L_SCAN_1000PX
 267 (slap) and ANDI_N2N (plain), NIQE remains comparable,
 268 suggesting that style transfer does not introduce noticeable
 269 perceptual degradation.
 270

271 While both NFIQ 2.0 and NIQE are useful quality met-
 272 rics, they capture different aspects: NFIQ 2.0 assesses bio-
 273 metric usability (e.g., ridge clarity and minutiae extractabil-
 274 ity), whereas NIQE measures perceptual naturalness via de-
 275 viations from natural-image statistics. Consequently, some
 276 sensors can exhibit divergent trends (e.g., S120 (roll) hav-
 277 ing lower NFIQ 2.0 but better NIQE), indicating that the

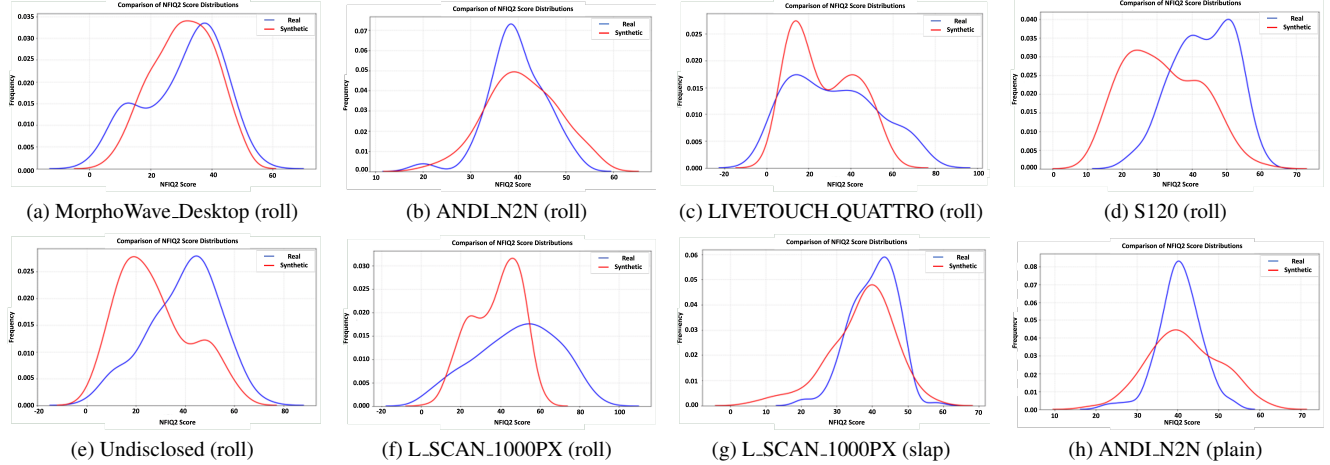


Figure 6. Distributions of NFIQ 2.0 scores across eight sensor styles (real: blue; synthetic: red).

Table 2. NIQE score (\downarrow) [26] comparison between synthetic and real images.

Sensor (Type)	Synthetic	Real	$\Delta\mu$	p -value
	$\mu \pm \sigma$	$\mu \pm \sigma$		
MorphoWave_Desktop (roll)	5.02 ± 0.67	6.47 ± 0.35	-1.45	1.48×10^{-61}
ANDLN2N (roll)	7.00 ± 0.96	7.24 ± 0.59	-0.24	0.0579
LIVETOUCH_QUATTRO (roll)	6.16 ± 1.01	12.45 ± 2.24	-6.29	1.00×10^{-75}
S120 (roll)	8.08 ± 0.91	9.57 ± 0.88	-1.49	4.64×10^{-35}
Undisclosed (roll)	8.33 ± 1.51	7.22 ± 2.54	+1.11	1.15×10^{-5}
Crossmatch_L_SCAN_1000PX (roll)	6.68 ± 1.15	12.12 ± 0.79	-5.44	1.30×10^{-3}
Crossmatch_L_SCAN_1000PX (slap)	7.06 ± 1.17	7.43 ± 3.05	-0.37	3.50×10^{-94}
ANDLN2N (plain)	7.57 ± 1.24	7.54 ± 0.53	+0.03	0.7566

Note: $\Delta\mu = \mu_{\text{syn}} - \mu_{\text{real}}$. Lower NIQE values indicate better perceptual quality.

278 image may look perceptually clean while being slightly less
279 optimal for biometric quality.

280 4.4. t-SNE Analysis of Synthetic Sensor Styles

281 To examine whether our model captures sensor-specific
282 style information, we visualize style embeddings of syn-
283 thetic fingerprints using t-SNE. Embeddings are extracted
284 from the style encoder of our model (not an external recog-
285 nition backbone), so the visualization reflects the model’s
286 internal representation of sensor style. Figure 7 shows seen
287 and unseen sensor styles.

288 For Fig. 7a, we select synthetic images from six
289 sensors (MorphoWave_Desktop, ANDLN2N, LIVE-
290 TOUCH_QUATTRO, S120, Undisclosed, Cross-
291 match_L_SCAN_1000PX), removing samples of dif-
292 ferent fingerprint types from the same sensor to isolate

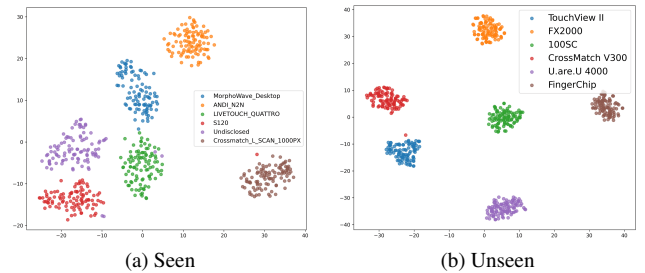


Figure 7. t-SNE visualization of style-encoder embeddings for synthetic fingerprints. (a) Seen sensor styles form compact, well-separated clusters. (b) Unseen sensor styles also form distinct clusters, suggesting the learned style representation extends to novel domains.

sensor-style effects. The resulting clusters are compact
and well separated, indicating that our method captures
sensor-specific style characteristics.

For Fig. 7b, we visualize sensors not present in training (TouchView II, FX2000, 100SC, Crossmatch V300, U.are.U 4000, FingerChip). Unseen styles also form distinct clusters, suggesting encouraging cross-domain generalization to novel sensor domains.

4.5. Training Recognition Model with Synthetic Data

To evaluate the effectiveness of our generative pipeline for augmenting fingerprint recognition, we conduct experiments using ResNet18/34/50/101 [27] and ViT-B/16 [28] backbones, all pretrained on ImageNet-1K [29]. We consider multiple training scenarios with different dataset compositions.

For real-data training, we use NIST SD 302 (Nail-to-Nail, N2N) [20], containing 1,600 real identities with 10

Table 3. TAR@FAR=0.1% and TAR@FAR=0.01% on different training datasets and architectures.

Training Dataset	TAR@FAR=0.1%					TAR@FAR=0.01%				
	ResNet18	ResNet34	ResNet50	ResNet101	ViT	ResNet18	ResNet34	ResNet50	ResNet101	ViT
NIST SD 302 (real) [20]	0.6946	0.7035	0.7132	0.7457	0.8008	0.4988	0.4741	0.5003	0.5039	0.5698
FPGAN-Control [4]	0.7181	0.6682	0.6626	0.6193	0.7875	0.4523	0.4725	0.4833	0.4927	0.5438
PrintsGAN [3]	0.5421	0.5031	0.4561	0.5122	0.5523	0.3951	0.3635	0.3485	0.4201	0.4524
Ours	0.7251	0.7298	0.7251	0.7271	0.7531	0.5320	0.6059	0.5477	0.5303	0.5639
NIST SD 302 + FPGAN-Control	0.8093	0.8514	0.8957	0.8719	0.8935	0.7251	0.7112	0.7478	0.7251	0.7369
NIST SD 302 + PrintsGAN	0.7299	0.7215	0.7337	0.7793	0.7688	0.5474	0.5747	0.6228	0.6329	0.6701
NIST SD 302 + Ours	0.8575	0.8930	0.8938	0.9009	0.9128	0.7850	0.7857	0.8056	0.8169	0.8211

311 impressions each. For synthetic-only training, FPGAN-
312 Control [4] and our method each provide 2,000 synthetic
313 identities with 10 impressions per identity, while Prints-
314 GAN [3] provides 1,500 identities with 10 impressions per
315 identity.

316 In joint training (last three rows of Table 3), the real por-
317 tion consists of 16,000 images from NIST SD 302 (1,600
318 identities \times 10 impressions). The synthetic portion includes
319 20,000 images from FPGAN-Control (2,000 identities \times
320 10), 15,000 images from PrintsGAN (1,500 identities \times
321 10), and 20,000 images from our method (2,000 identities
322 \times 10). For our method, we generate 1,000 identities by aug-
323 menting real identities (from N2N) and 1,000 identities by
324 sampling virtual identities from a pretrained DDPM; each
325 identity is generated with 10 style variations.

326 All models are evaluated on a held-out set of 400 real
327 identities from NIST SD 302. We report TAR at FAR=0.1%
328 and FAR=0.01%.

329 Table 3 shows that our method achieves strong perfor-
330 mance across architectures. Models trained on our synthetic
331 data outperform those trained on other synthetic datasets
332 (FPGAN-Control and PrintsGAN). The gains are more pro-
333 nounced in joint training: NIST SD 302 + Ours consis-
334 tently yields the best results. For example, with ResNet101,
335 NIST SD 302 + Ours reaches TAR=0.9009 at FAR=0.1%,
336 and with ViT it achieves TAR=0.8211 at FAR=0.01%. We
337 attribute these improvements to the diversity and fidelity
338 of our synthetic data. Unlike methods that solely sample
339 new identities from latent space, our approach combines (i)
340 real identity augmentation to model intra-class variation and
341 (ii) virtual identity generation to expand identity coverage,
342 achieving strong robustness and generalization.

343 4.6. Computational Complexity Analysis

344 We compare the computational cost of our diffusion-based
345 pipeline with the GAN-based FPGAN-Control [4]. On an
346 NVIDIA A6000 GPU, our method requires 12.8 seconds
347 per fingerprint image on average, whereas FPGAN-Control
348 requires 0.05 seconds. This difference is expected: GANs
349 generate images in a single forward pass, while diffusion
350 models rely on iterative denoising. Nevertheless, our ap-

proach produces finer ridge structures and more complete
351 details, and we observe improved downstream recognition
352 in our experiments. Since fingerprint generation is primar-
353 ily used for offline data augmentation rather than real-time
354 deployment, the additional cost is acceptable given the im-
355 proved fidelity and downstream utility. 356

357 5. Ablation Study

358 We perform a progressive ablation study from S0 to S5
359 to quantify the effect of each design choice. As shown
360 in Table 4, the text-only baseline (S0) provides limited
361 structural consistency, with LightGlue mean matches of
362 20.01. Adding ControlNet (S1) yields a clear improvement
363 in structural consistency, increasing mean matches to 34.05
364 while also affecting the perceptual proxies (IS/IL).

365 We report three complementary metrics: (i) LightGlue
366 mean matches as a proxy for structure-level alignment be-
367 tween the generated image and its paired content input;
368 (ii) Inception Score (IS) as a proxy for perceptual diver-
369 sity/quality; and (iii) intra-cluster LPIPS distance (IL) to
370 quantify perceptual variation among samples within the
371 same experimental setting (higher IL indicates larger per-
372 ceptual diversity). Since IS/LPIPS are derived from natural-
373 image models, we use them as complementary proxies
374 rather than fingerprint-specific quality metrics.

375 Enabling decoupled cross-attention (S2) further
376 strengthens structure-level alignment (40.74 matches),
377 indicating that separating the conditioning streams helps
378 preserve ridge topology under style transfer. With mask
379 injection (S3), mean matches increases to 42.41 and IL
380 decreases slightly (0.5440), suggesting reduced perceptual
381 variation under explicit foreground/background separation.
382 When introducing the style condition (S4), IS remains
383 comparable (3.170) and IL increases to 0.5996 (greater
384 perceptual variation), while mean matches decreases to
385 40.02, implying stronger style cues can trade off against
386 strict structure-level alignment. Finally, the full model (S5)
387 recovers structure-level alignment (42.52 matches) while
388 keeping IS/IL within a comparable range, showing that
389 combining style and shape conditions can restore structural

Table 4. Progressive ablation from S0 to S5. Results are reported using LightGlue. IS: Inception Score; IL: intra-cluster LPIPS distance (higher indicates larger perceptual diversity).

Metric	S0	S1	S2	S3	S4	S5
IS	2.841	3.678	2.931	3.151	3.170	3.144
IL	0.6795	0.6866	0.5374	0.5440	0.5996	0.5502
LightGlue mean_matches	20.01	34.05	40.735	42.405	40.015	42.52

Table 5. Sensitivity analysis under the full setting S5. IS: Inception Score; IL: intra-cluster LPIPS distance (higher indicates larger perceptual diversity).

Metric	exp11	exp12	exp13	exp14	exp15	exp17	exp18
IS	3.231	2.824	3.285	3.200	3.562	2.7338	2.7109
IL	0.5969	0.5376	0.5962	0.5617	0.5932	0.5409	0.5173
LightGlue mean_matches	38.045	44.275	41.6	40.035	40.305	46.885	59.365

390 fidelity without sacrificing controllability.

391 We additionally conduct sensitivity analysis under the
392 full setting S5 (Table 5). shows that stronger guidance
393 increases mean matches, but tends to reduce IS and IL
394 (i.e., lower perceptual diversity), suggesting that overly
395 strong structural constraints may limit style expressiveness.
396 For decoupled cross-attention weights, symmetric settings
397 ($\lambda_1 = \lambda_2$) produce similar IS/IL with moderate differ-
398 ences in matching (exp13–exp14). The asymmetric setting
399 (exp15, $\lambda_1 = 1.5$, $\lambda_2 = 0.5$) achieves the highest IS (3.562)
400 while keeping IL comparable to other settings, indicating
401 that emphasizing style over shape can improve perceptual
402 diversity without a notable loss in structure consistency.

403 6. Limitation & Future Work

404 Although our method achieves high-quality fingerprint syn-
405 thesis and strong structure-consistent, several limitations re-
406 main. The diffusion-based generation process is computa-
407 tionally intensive due to iterative denoising steps. More-
408 over, the framework currently depends on predefined spatial
409 masks to guide region-specific feature injection, which may
410 limit its adaptability under varying sensor conditions. As
411 potential future improvements, we plan to incorporate more
412 efficient diffusion samplers such as DPM-Solver or Latent
413 Consistency Models (LCM) to reduce sampling steps and
414 inference time. Additionally, exploring automatic or data-
415 driven mask estimation could enhance flexibility and scala-
416 bility for practical deployment.

417 7. Conclusion

418 In this paper, we present a fingerprint image synthesis
419 framework that addresses data scarcity under privacy con-
420 straints in fingerprint recognition. By integrating Stable
421 Diffusion, ControlNet, and Multi-IP-Adapter mod-
422 ules, our method generates **structure-consistent** and style-

diverse fingerprint images, supporting data augmentation
and **privacy-motivated** training by reducing direct depen-
dence on real identity data. Experiments show that the pro-
posed synthetic fingerprints improve downstream recogni-
tion performance, achieving up to an **11 percentage-point**
gain in TAR@FAR=0.1% under joint training compared to
the real-only setting.

Overall, this work provides an effective approach to ex-
pand fingerprint datasets while reducing direct reliance on
real identity data during augmentation. In future work, we
plan to improve generation efficiency and further strengthen
cross-sensor generalization, extending the framework to
other biometric modalities such as palmprint and vein
recognition. We believe this research supports diffusion-
based synthesis toward practical, **privacy-conscious** bio-
metric applications.

An anonymized code repository is available at: https://anonymous.4open.science/r/Privacy-and-Anonymity_Privacy-metrics-1F31/

References

- [1] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 1
- [2] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 2, 3
- [3] Joshua James Engelsma, Steven Grosz, and Anil K Jain. Printsgan: Synthetic fingerprint generator. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):6111–6124, 2022. 1, 2, 7
- [4] Alon Shoshan, Nadav Bhonker, Emanuel Ben Baruch, Ori Nizan, Igor Kviatkovsky, Joshua Engelsma, Manoj Aggarwal, and Gérard Medioni. Fpgan-control: A controllable fingerprint generator for training with synthetic data. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 6067–6076, 2024. 1, 2, 7
- [5] Kejian Li and Xiao Yang. Diffusion probabilistic model based end-to-end latent fingerprint synthesis. In *2023 IEEE 4th International Conference on Pattern Recognition and Machine Learning (PRML)*, pages 343–349. IEEE, 2023. 1, 2
- [6] Freddie Grabovski, Lior Yasur, Yaniv Hacmon, Lior Nisimov, and Stav Nimrod. Difffinger: Advancing synthetic fingerprint generation through denoising diffusion probabilistic models. *arXiv preprint arXiv:2405.04538*, 2024.
- [7] Jingqiao Wang, Zicheng Zhang, and Congying Han. Exploring latent fingerprint synthesis with diffusion probabilistic models. In *Proceedings of the Future Technologies Conference*, pages 78–89. Springer, 2024.
- [8] Mao-Hsiu Hsu, Yung-Ching Hsu, and Ching-Te Chiu. In-painting diffusion synthetic and data augment with feature keypoints for tiny partial fingerprints. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 7(3):396–409, 2024.

- 478 [9] Weizhong Tang, Diego Figuerola, Donglin Liu, Kerstin
479 Johnsson, and Alexandros Sotasakis. Enhancing fingerprint
480 image synthesis with gans, diffusion models, and style trans-
481 fer techniques. *arXiv preprint arXiv:2403.13916*, 2024. 1,
482 2
- 483 [10] Robin Rombach, Andreas Blattmann, Dominik Lorenz,
484 Patrick Esser, and Björn Ommer. High-resolution image
485 synthesis with latent diffusion models. In *Proceedings of
486 the IEEE/CVF conference on computer vision and pattern
487 recognition*, pages 10684–10695, 2022. 1, 2
- 488 [11] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding
489 conditional control to text-to-image diffusion models. In
490 *Proceedings of the IEEE/CVF international conference on
491 computer vision*, pages 3836–3847, 2023. 1, 2
- 492 [12] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-
493 adapter: Text compatible image prompt adapter for text-to-
494 image diffusion models. *arXiv preprint arXiv:2308.06721*,
495 2023. 1, 2, 3
- 496 [13] Lin Hong, Yifei Wan, and Anil Jain. Fingerprint image en-
497 hancement: algorithm and performance evaluation. *IEEE
498 transactions on pattern analysis and machine intelligence*,
499 20(8):777–789, 1998. 1, 2
- 500 [14] Raffaele Cappelli, A Erol, D Maio, and D Maltoni. Syn-
501 thetic fingerprint-image generation. In *Proceedings 15th In-
502 ternational Conference on Pattern Recognition. ICPR-2000*,
503 volume 3, pages 471–474. IEEE, 2000. 2
- 504 [15] Yanming Zhu, Xuefei Yin, and Jiankun Hu. Fingergan: A
505 constrained fingerprint generation scheme for latent finger-
506 print enhancement. *IEEE Transactions on Pattern Analysis
507 and Machine Intelligence*, 45(7):8358–8371, 2023. 2
- 508 [16] Ataher Sams, Homaira Huda Shomee, and SM Mahbubur
509 Rahman. Hq-fingan: High-quality synthetic fingerprint gen-
510 eration using gans. *Circuits, Systems, and Signal Processing*,
511 41(11):6354–6369, 2022.
- 512 [17] Rafael Bouzaglo and Yosi Keller. Synthesis and reconstruc-
513 tion of fingerprints using generative adversarial networks.
514 *arXiv preprint arXiv:2201.06164*, 2022.
- 515 [18] Ogban-Asuquo Ugot, Chika Yinka-Banjo, and Sanjay Misra.
516 Biometric fingerprint generation using generative adversarial
517 networks. In *Artificial Intelligence for Cyber Security: Meth-
518 ods, Issues and Possible Horizons or Opportunities*, pages
519 51–83. Springer, 2021. 2
- 520 [19] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian
521 Zhang, Zhongang Qi, and Ying Shan. T2i-adapter: Learning
522 adapters to dig out more controllable ability for text-to-image
523 diffusion models. In *Proceedings of the AAAI conference on
524 artificial intelligence*, volume 38, pages 4296–4304, 2024. 2
- 525 [20] Gregory P Fiumara, Patricia A Flanagan, John D Grantham,
526 Kenneth Ko, Karen Marshall, Matthew Schwarz, Elham
527 Tabassi, Bryan Woodgate, and Christopher Boehnen. *Nist
528 special database 302: Nail to nail fingerprint chal-
529 lenge*. Gregory P. Fiumara, Patricia A. Flanagan, John D.
530 Grantham, Kenneth Ko . . . , 2019. 3, 6, 7
- 531 [21] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya
532 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,
533 Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning
534 transferable visual models from natural language supervi-
535 sion. In *International conference on machine learning*, pages
536 8748–8763. PmlR, 2021. 3
- 537 [22] Dario Maio, Davide Maltoni, Raffaele Cappelli, Jim L Way-
538 man, and Anil K Jain. Fvc2004: Third fingerprint verifica-
539 tion competition. In *International conference on biometric
540 authentication*, pages 1–7. Springer, 2004. 4
- 541 [23] Philipp Lindenberger, Paul-Edouard Sarlin, and Marc Polle-
542 feys. Lightglue: Local feature matching at light speed. In
543 *Proceedings of the IEEE/CVF international conference on
544 computer vision*, pages 17627–17638, 2023. 4, 5
- 545 [24] Guilherme Potje, Felipe Cadar, André Araujo, Renato Mar-
546 tins, and Erickson R Nascimento. Xfeat: Accelerated fea-
547 tures for lightweight image matching. In *Proceedings of
548 the IEEE/CVF Conference on Computer Vision and Pattern
549 Recognition*, pages 2682–2691, 2024. 4, 5
- 550 [25] Elham Tabassi, Martin Olsen, Oliver Bausinger, Christoph
551 Busch, Andrew Figlarz, Gregory Fiumara, Olaf Henniger,
552 Johannes Merkle, Timo Ruhlend, Christopher Schiel, et al.
553 Nfiq 2 nist fingerprint image quality. 2021. 4
- 554 [26] Anish Mittal, Rajiv Soundararajan, and Alan C Bovik. Mak-
555 ing a “completely blind” image quality analyzer. *IEEE Sig-
556 nal processing letters*, 20(3):209–212, 2012. 4, 6
- 557 [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
558 Deep residual learning for image recognition. In *Proceed-
559 ings of the IEEE conference on computer vision and pattern
560 recognition*, pages 770–778, 2016. 6
- 561 [28] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov,
562 Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,
563 Mostafa Dehghani, Matthias Minderer, Georg Heigold, Syl-
564 vain Gelly, et al. An image is worth 16x16 words: Trans-
565 formers for image recognition at scale. *arXiv preprint
566 arXiv:2010.11929*, 2020. 6
- 567 [29] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li,
568 and Li Fei-Fei. Imagenet: A large-scale hierarchical image
569 database. In *2009 IEEE conference on computer vision and
570 pattern recognition*, pages 248–255. Ieee, 2009. 6