037

038

039

040

041

042

043

044

045

046

047

048

049

050

# Generating Synthetic Data via Augmentations for Improved Facial Resemblance in DreamBooth and InstantID

Anonymous CVPR submission

Paper ID \*\*\*\*\*

#### Abstract

001 The personalization of Stable Diffusion for generating pro-002 fessional portraits from amateur photographs is a bur-003 geoning area, with applications in various downstream contexts. This paper investigates the impact of augmen-004 tations on improving facial resemblance when using two 005 prominent personalization techniques: DreamBooth and 006 InstantID. Through a series of experiments with diverse 007 subject datasets, we assessed the effectiveness of various 008 augmentation strategies on the generated headshots' fi-009 010 *delity to the original subject. We introduce* FaceDistance, *a* wrapper around FaceNet, to rank the generations based on 011 facial similarity, which aided in our assessment. Ultimately, 012 013 this research provides insights into the role of augmentations in enhancing facial resemblance in SDXL-generated 014 portraits, informing strategies for their effective deployment 015 016 in downstream applications.

# **017 1. Introduction**

Personalized text-to-image generation has gained traction
with the rise of models like Stable Diffusion (SD). However, training SD on small, user-specific datasets presents
challenges, such as identity retention, overfitting, and artifact generation. Augmentation techniques are widely used
in deep learning to improve generalization, but their role in
personalized text-to-image generation is underexplored.

025In this work, we analyze the effect of augmentations on026personalized SD models trained with few-shot and zero-027shot methods, particularly DreamBooth and InstantID. We028investigate how different augmentations impact model per-029formance and whether they enhance the realism and consis-030tency of generated images.

In particular, we analyze both classical and generative augmentation strategies to bridge the gap between limited real data and high-fidelity synthetic outputs. By refining facial features and preserving identity through targeted GenAI-based augmentations, such as InstantID, we



Figure 1. Pipeline for creating personalized images based on synthetically generated images through classical and GenAI-based augmentations for better downstream resemblance in DreamBooth-generated images.

aim to improve the applicability of personalized generation in scenarios where synthetic data must closely mirror real-world characteristics. We analyze under which conditions we can ensure that "GenAI outputs improve GenAI outputs", avoiding a data quality collapse, providing best practices and heuristics.

Our contributions include:

- Analysis of classical augmentation techniques such as flipping, cropping, color enhancement, and background modifications.
- Using InstantID as a fast way of enhancing the userspecific dataset using the diffusion model itself.
- We conduct a survey to evaluate how white-collar workers perceive personalized generations from DreamBooth and InstantID under various augmentation strategies.

112

115

## **2. Background and Related Work**

In this section, we present the foundational concepts and
prior research relevant to our work on augmentation techniques for few-shot personalization in diffusion models.

# 055 2.1. Text-to-Image Diffusion Models

Text-to-image diffusion models generate high-quality images from natural language descriptions by gradually denoising random Gaussian noise guided by text embeddings [12, 15]. Stable Diffusion [15] employs a latent diffusion approach that operates in a compressed latent space rather than pixel space, reducing computational requirements while maintaining generative quality.

# **063 2.2. Subject-Driven Image Generation**

Subject-driven image generation creates images featuring
specific subjects with high fidelity while maintaining their
identity across contexts [16]. Key approaches include:

DreamBooth [16] fine-tunes the U-Net of Stable Diffusion using 3-5 images of a specific subject. It preserves the semantic prior through class-specific prior preservation loss and uses a rare token with weak prior to refer to the subject with the prompt format "a [V] [class noun]".

072InstantID [20] is a zero-shot method that combines fa-073cial feature extraction with text conditioning. It extracts five074key facial landmarks to condition the position and orienta-075tion of the generated face, providing greater control over the076output.

We standardize our experiments using the same SDXL
model for both techniques to ensure fair methodological
comparison.

# 080 2.3. Image Augmentation Techniques

## 081 2.3.1. Classical Image Augmentations

Classical image augmentation techniques include geometric transformations (flipping, rotation, scaling, cropping),
photometric adjustments (brightness, contrast, saturation,
hue), and noise injections (Gaussian, salt-and-pepper).
These predefined transformations maintain semantic integrity while introducing controlled diversity to expand limited training datasets.

## **089 2.3.2.** Augmentations in Diffusion Models

Data augmentation enhances diffusion model performance
while reducing computational demands [19]. Key approaches include mixing-based augmentations that interpolate between existing samples [9] and consistency regularization techniques that enforce invariance to specific transformations [8, 11]. Our work investigates these techniques
specifically for few-shot personalization applications.

## 2.4. Face Processing Approaches

FaceNet [17] maps facial images to a 128-dimensional em-<br/>bedding space where similar faces are positioned closely<br/>together. The standard pipeline uses MTCNN [22] for face<br/>detection before embedding generation, with cosine dis-<br/>tance metrics for similarity assessment [18].098<br/>099

Alternative approaches include faceswapping methods103[6, 10] and augmented reality techniques for virtual try-on104applications [7]. While these provide real-time capabilities,105they often lack the flexibility and integration capabilities of106diffusion-based approaches.107

Our research builds on these foundations to investi-<br/>gate how strategic data augmentation can improve few-shot<br/>personalization in diffusion models, focusing on identity<br/>preservation and recontextualization.108<br/>109111

3. Methodology

We use augmentations across various Subject Datasets to 113 see if there is an overall improvement in generated pictures. 114

#### **3.1. Subject Datasets.**

Our dataset consists of 3 to 15 images per participant, with 116 n = 10 participants. To maintain a naturalistic data col-117 lection process, we instructed them: "Can you send me 118 portrait/selfie-style photos of your face in different places? 119 The more different places, the better." By avoiding rigid 120 guidelines, we ensured that the collected images reflect real-121 istic user behavior. As a result, our findings are well-aligned 122 with real-world data distributions, enhancing the transfer-123 ability and applicability of our results. 124

Our dataset exhibits a diverse range of environmental 125 conditions, facial orientations, and image qualities, ensur-126 ing variability that mirrors real-world scenarios. The im-127 ages encompass different backgrounds, lighting conditions, 128 and subject behaviors, contributing to a dataset that is both 129 representative and robust. For instance, some images fea-130 ture cluttered or irregular backgrounds (e.g., Baker-Zoe, 131 Bottle-Hugo), while others are captured in controlled en-132 vironments (Biometric-Kora). Variation in gaze direction 133 is also present, with Doctor-Nina not looking at the cam-134 era, while 3D-Gary includes dynamic head movements ex-135 tracted from a video. Additionally, differences in personal 136 appearance and accessories are observed, such as Farmer-137 Lisa wearing a helmet and Staircase-Judy wearing makeup. 138 Lighting conditions range from well-lit (Vacation-Anna) to 139 suboptimal (2024-Kora), further enhancing the dataset's re-140 alism. These characteristics make our dataset a valuable 141 resource for evaluating model performance under uncon-142 strained, real-world conditions. 143



Figure 2. Pipeline for creating Gen-AI augmented personalized data via InstantID. Based on one or more input images of a person, we run it through the InstantID pipeline but with augmented landmarks and prompts. The landmarks are taken from the input image and slightly perturbed for good resemblance. The collected synthetic dataset is then further used for downstream DreamBooth training. Figure modified from [20].

#### 144 **3.2. Dataset Augmentations**

We apply augmentations individually to evaluate each tech-nique's performance improvement independently.

147Classical AugmentationsStandard techniques in-148clude: (i) Random Horizontal Flip with  $p \in \{0, \frac{1}{2}, 1\}$ ,149and (ii) Color Jitter varying brightness, contrast, saturation150 $(\pm 5, \pm 15)$  and hue  $(\pm 5^{\circ})$ .

**Background Augmentation** We use U<sup>2</sup>-Net [14] for
subject isolation, testing both base and human segmentation models. Backgrounds include *flat colors, patterns* from
Wikimedia [5], and Flickr Places.

Blending Techniques We separately evaluate *Alpha Blending* and *Poison Blending* through both automated and
manual techniques.

**Resizing Methods** We compare: (i) downsampling
then upsampling, (ii) upsampling only, and (iii) original dimensions. Methods include ESRGAN [21], Lanczos, and
Bilinear.

162 Cropping Strategies Five approaches: (i) SDXL dimensions [13], (ii) automated center cropping to 1MP, (iii) downsample-then-crop to 1MP at various aspect ratios, (iv)

manual eight-variation cropping, and (v) MTCNN facebased cropping. 165

ColorAdjustmentAdobeLightroomauto-adjustment enhances visual quality.168

Generative AugmentationUsing InstantID, we gen-erate new subject images with prompts from dolphin1702.2.1 - Mistral 7B [3] and varied facial landmarks171(Figure 2).172

#### 3.3. Hardware, Software, and Hyperparamters

experiments were conducted 174 A11 on а single NVIDIA GeForce RTX 3090 with 24GB VRAM. 175 We use sd-scripts[4] for DreamBooth and 176 ComfyUI\_InstantID[1] for InstantID experiments, 177 inheriting all bias in their pipeline, if any. 178

Results of DreamBooth finetuning a diffusion model179(DM) greatly depends on the DMs ability of generating images. We use RealVisXL\_V4.0 [2], which is a community finetune of SDXL for realistic image generation.181182183

232

233

234

235

236

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

183 Default prompt we use "A professional headshot of a
184 subject wearing a suit standing in a well-lit studio, DSLR"
185 as the default prompt. Empirical evidence suggests includ186 ing the gender as *man* and *woman* gives generated images
187 gender characteristics based on western culture, which we
188 preferred.

For DreamBooth, we use "a [V] [man|woman]" where [V] denotes the rare token. InstantID doesn't have a special prompt and work with any text. LLM generated prompts are useful in both.

## **193 3.4. FaceDistance Metric**

For a subject image dataset, we calculate their embedding using FaceNet [17], which maps similar faces to similar locations on a hypersphere. Then, we calculate the mean face vector  $v_{real}$ . For a given generated face, we measure its embedding vectors cosine distance to  $v_{real}$ .

FaceDistance is a useful technique for distinguishing between "good" and "bad" generations. This can be used to
rank generated images based on their similarity (lower is
better). It can be used to discard the largest n% of distances
to improve personalization pipelines.

For our subject datasets, the mean cosine distance of  $v_{real}$  to real images is  $\bar{v}_{within \ real} \approx 0.13$ . We notice  $\max(v_{within \ real}) = 0.22$  and  $\min(v_{within \ real}) = 0.05$ .

## **4. Experiment Results**

We try to achieve higher facial similarity via DreamBoothand InstantID using highlighted augmentations.

210 Despite selecting a realistic image generation model, achieving photorealistic generation of an individual's face 211 remains challenging without imposing strict constraints on 212 213 the subject images. We have relaxed many of these constraints to enhance usability, as expecting an average user 214 215 to compile a dataset of themselves without understanding 216 the underlying image generation techniques presents a sig-217 nificant challenge. Ensuring high subject fidelity is crucial 218 for these methods to be effective in downstream applications, as humans are highly sensitive to variations in facial 219 220 features compared to textures.

One major issue with datasets without great constraints 221 222 is that the images is not a good representation of the person. It can be compared to having difficulty recognizing a 223 224 person in real life whom you only saw in photographs. We 225 observe this phenomenon for small datasets with size  $\leq 3$ . In these cases, the generated images is a good reflection of 226 the dataset (if someone doesn't know them in real life, they 227 are likely to claim that these pictures are good. Otherwise 228 the generated images are not a good representation of the 229 230 real person).

#### 4.1. DreamBooth

We configure our hyperparameters such that recontextualization capabilities can be sacrificed for high facial fidelity. Identity preservation is hard in DreamBooth. so we rather overfit to achieve high subject fidelity and have limited freedom in generations.

The common theme in augmentations is that if the aug-237 mented image has any kind of artifact/anomaly, then the 238 rare token will be associated with it. The supporting ob-239 servations are (i) When background is replaced with a ge-240 ometric pattern (from wikimedia patterns), the model will 241 focus on learning the pattern than the subject (ii) When im-242 age is upscaled with ESRGAN, the texture ESRGAN in-243 troduces say present in generations (iii) the masks gener-244 ated with U<sup>2</sup>-Net is not pixel-perfect. and results in a mix 245 around hair/air boundary. This mix becomes associated 246 with the subject. The human segmentation models train-247 ing data was not highly accurate around hairs but was better 248 in identifying body parts. The base model is performs bet-249 ter around hair and was overall better. The robustness of 250 human segmentation model is not needed. (iii) any kind 251 of color jitter is visible in generated images. For exam-252 ple the saturation change of 0.1 is clearly present in gen-253 erations. (iv) using Adobe Lightroom as a preprocessing 254 step resulted in better color graded generations compared 255 to non-preprocessed datasets. (v) datasets with low contrast 256 (e.g. exclusively Polaroid pictures) resulted in copying the 257 photography style/lighting from the pictures — though this 258 can be attributed to our hyperparameter configuration. 259

Because of the low recontextualization capabilities, backgrounds becomes highly associated with the rare token. Replacing the background with **Pastel Colors** and **Rainbow Colors** led to eccentric and often unrealistic images, with the latter occasionally generating pictures without subjects. **Gray** offered the highest resemblance to the subject, while **Dark Gray** caused the model to disassociate the subject from its context. Because of problems with U<sup>2</sup>-Net, **Light Gray** background outperformed **Dark Gray**, especially in bright environments.**Wikimedia Patterns** slowed down learning and degraded the image quality across all generations. Lastly, **Studio Backdrops** introduced irregularities that reduced the quality of the generated images which can be thought as similar to Wikimedia Patterns because backdrops has patterns.

**Random Horizontal Flip** slowed learning due to face asymmetry, which confused facial features. **Random Rotation** caused distorted images and introduced black padding bars, which also can be seen in augmented subject images. **Color Jitter** led to undesirable results, as brightness, contrast, saturation, and hue changes were linked to rare tokens, causing erratic generations.

Both Alpha Blending and Poison Blending are discour-<br/>aged, as they require careful manual processing to achieve282283

321

322

323

324

325

326

327

328



(a) Real Images

(b) DreamBooth results with classical augmentations (crop, resize, and color)

(c) DreamBooth results with GenAI augmentations & without classical augmentations

Figure 3. Example improvement of including Instantid generated images in the subject dataset *Vacation-Anna*. The model is prompted with *default prompt* with batchsize 4. The results are **not** cherry-picked to resemble the downstream application use. Although (b) is visually more interesting, the method in (c) is more consistent across many subject datasets.

284 good results. These techniques are not straightforward to285 apply and can lead to undesirable artifacts if not handled286 properly.

287 Images around 1 Megapixel performed best, providing a balanced resolution for high-quality generation. Upscal-288 ing with ESRGAN introduced visible artifacts, especially 289 290 around facial features. Upscaling with Lanczos was effective, particularly when starting from larger images. How-291 ever, if the initial dataset contained low-resolution images, 292 the generated images exhibited facial blurring due to the na-293 294 ture of the Lanczos algorithm. The difference between bicubic and Lanczos was negligible. Downscaling resulted in 295 lower-quality generations compared to using original-sized 296 images. It should be noted that our testing output resolution 297 was  $1024 \times 1024$ . 298

299 InstantID Augmentation Datasets augmented with InstantID yield clearly superior performance. The added im-300 301 ages need to be diverse (i.e., generated with various text 302 conditioning and different keypoint images). Since we trade 303 recontextualization abilities for increased facial similarity, 304 generating the same person in similar contexts is beneficial. 305 DreamBooth achieves similar facial similarity compared to InstantID but allows for greater control. The rigidity caused 306 by the keypoint images is eliminated. However, this method 307 is more computationally expensive than raw InstantID. Ad-308 309 ditionally, achieving proper prompt diversity can be challenging. I prefer InstantID over DreamBooth. 310

The ratio of real to InstantID-generated images dependsentirely on the diversity of the generated images. One rule

of thumb is that no single concept should comprise more313than 25% of the dataset. For example, if images labeled as314"a [V] man in a library" exceed 25%, DreamBooth training315will associate the rare token with the concept. This results316in a final DreamBooth model that is unusable due to a complete loss of recontextualization ability caused by overfitting.319

Since InstantID generations are highly realistic, one can generate additional images with it to better represent the subject during DreamBooth training. We use the same diffusion model for both InstantID and DreamBooth to integrate the subject more effectively into the model without altering the subject's context. This ensures that the dataset distribution remains closer to the diffusion model's generation space.

#### 4.1.1. FaceDistance

We tried to select the "best" DreamBooth checkpoint by 329 generating images of "a [V] man" in different contexts for 330 all checkpoints and ranking them using FaceDistance. This 331 method was able to discard obviously bad checkpoints (e.g. 332 anomalies in generations, unable to generate the subject, di-333 vergence) but is not able to rank "good enough" checkpoints 334 335 within themselves. (Figure 4 The FaceNet manifold isn't sensitive to very similar looking people. For a given hyper-336 parameter configuration, a few tests show when the model 337 will be converged to its best state (usually between 3k and 338 6k steps) and since FaceDistance isn't able to differentiate 339 betweent them, FaceDistance isn't a useful tool for this pur-340 pose. 341

Despite these challenges, FaceDistance appears to be 342

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

414

415

416

417

418

419

420

421



Figure 4. FaceDistance Distributions of 2000 Samples from Different Saved Dreambooth finetunes of SDXL Real. *Closeup-Kora* is used. The KDE for each looks like a normal distribution.

functioning for loosely ranking generated images. This im-proves the user experience.

## **345 4.2. InstantID**

The effectiveness of InstantID is highly dependent on thequality and characteristics of the provided reference images.

#### 348 4.2.1. Face Embedding

We conducted experiments to determine the optimal num-349 350 ber of reference images that balances usability and facial similarity. Our findings confirm those of [20], demonstrat-351 352 ing that using multiple reference images results in increased facial similarity. When only one reference image is pro-353 vided, the generated face is heavily influenced by the spe-354 cific appearance captured in that single image. We attribute 355 this limitation to insufficient information being extracted by 356 357 the Face Encoder from a single perspective. Our analysis 358 indicates that four reference images provide satisfactory results in most cases, with diminishing returns observed be-359 360 yond this number. Since reference images are cropped and aligned before being processed by the face encoder, users 361 362 have considerable flexibility in selecting images without 363 compromising model performance.

#### **364 4.2.2. Landmarks Image**

We observe that facial landmarks exert strong conditioning
influence, often rendering text prompts ineffective for controlling the subject's position. The generated image consistently replicates the face placement, orientation, and size
specified by the provided keypoints, due to the five-point
landmark system employed.

For practical applications, users frequently struggle to understand how face positioning in the landmark image transfers to the generated output. This communication challenge often results in user dissatisfaction with generated images, despite the issue stemming from suboptimal conditioning input. To address this limitation, we propose two solutions: 2shot generation and face replacement.

In **2-shot generation**, we collect subject reference images  $(s_1, \ldots, s_n)$  and a separate image representing the desired pose and composition  $s_{kpts}$ . These are used as reference images and the keypoints image, respectively. While the resulting output  $s_{out}$  is generally satisfactory, using facial landmarks from one person to generate another reduces facial similarity due to structural differences in the five keypoints (eyes, nose, mouth). We hypothesize this stems from imbalanced conditioning weights. Performance improves when replacing  $s_{kpts}$  with a previously generated image of the subject, yielding better facial similarity while maintaining compositional control.

In **face replacement**, users interact with a simple tool to manipulate (move/rotate/resize) their cropped face on a canvas matching the diffusion model's output dimensions. This approach eliminates the similarity issues caused by using another person's facial landmarks. However, the method performs poorly when none of the reference images show the subject facing the camera (deviations > 30 degrees). User satisfaction was higher with this approach compared to 2-shot generation, which we attribute to increased interactivity and faster generation times.

#### 4.2.3. Augmentations

Due to InstantID's architectural design, rotational and 402 shape-altering augmentations proved ineffective. Back-403 ground replacement and similar context-modifying aug-404 mentations degraded similarity because the resulting arti-405 facts fall outside the distribution of images encountered dur-406 ing training by the model provided in [20]. The trained 407 model demonstrates robustness to meaningful color ad-408 justments, rendering color modifications unnecessary for 409 well-lit scenes. For low-resolution images, traditional up-410 scaling methods (Lanczos/bicubic) performed adequately, 411 while neural network-based upscaling introduced novel ar-412 tifacts unseen during training, resulting in reduced quality. 413

# 5. Survey

We conducted the survey to evaluate the viability of AIgenerated portraits for professional use and to compare the performance of DreamBooth and InstantID in generating realistic headshots. 97 white-collar workers from diverse professional backgrounds participated in the online survey. Numerical data can be found in Suppl. 11 and questionary can be found in Suppl. 12.

Overall Performance of Generated Portraits422generated by DreamBooth and InstantID performed simi-<br/>larly across multiple aspects, including overall quality, fa-<br/>cial detail clarity, identity preservation, perceived level of<br/>editing, and background quality. Using high-quality subject422424424425425

516

datasets led to slightly better results in most categories, except for "Editing," where participants indicated familiarity
and acceptance of traditional Photoshop-enhanced portraits.

430 Method Preferences A slightly higher percentage of par-431 ticipants (4%) preferred the standardized portraits from InstantID over the more flexible outputs of DreamBooth. 432 InstantID was often perceived as more professional, likely 433 due to its consistent "Photoshopped look," which resonated 434 435 with a broader audience. Open-ended responses highlighted diverse preferences, with participants emphasizing factors 436 such as lighting, pose, angle, expression, detail, color, and 437 background. 438

Facial Similarity DreamBooth demonstrated superior facial similarity between real images individuals and their
generated portraits compared to InstantID. More participants identified InstantID images as depicting a different
person than the reference. DreamBooth consistently maintained a higher level of facial similarity across both highand low-quality subject datasets.

Noticing AI Generations Most white-collar workers 446 struggled to identify AI-generated headshots when not ex-447 plicitly prompted, often focusing on well-known but absent 448 449 flaws commonly associated with AI generation. Among a subset of participants (n = 77) who regularly notice 450 AI-generated images in daily life, the generated portraits 451 blended well with conventional studio photographs. How-452 ever, participants who actively use AI for image creation 453 454 (n = 29) demonstrated better identification skills. This 455 group was more likely to recognize DreamBooth images as AI-generated, possibly due to DreamBooth's popularity, 456 while InstantID generations, being more niche, had a near 457 50/50 chance of being identified as AI. 458

## 459 6. Discussion

Our experiments offer insights into augmentation strategies 460 for improving facial resemblance in personalized text-to-461 image generation using DreamBooth and InstantID. While 462 463 classical augmentations are common in deep learning, ap-464 plying them to few-shot personalization can yield undesirable results. Geometric transformations like flipping and 465 rotation disrupted learning due to face asymmetry and ar-466 tifacts. Color jittering caused erratic generations by asso-467 ciating color shifts with DreamBooth's rare token. Back-468 ground augmentations with U<sup>2</sup>-Net introduced segmenta-469 470 tion imperfections, especially around hair, which the model learned. Replacing backgrounds with patterns or studio 471 472 backdrops also degraded image quality. However, auto color adjustment with Adobe Lightroom improved color 473 474 grading.

Generative augmentation via InstantID proved more ef-475 fective for enhancing facial similarity in DreamBooth train-476 ing. By generating diverse synthetic images with varied 477 prompts and facial landmarks, we enriched the dataset with 478 realistic examples, aligning it with the diffusion model's 479 space. However, maintaining a balance between real and 480 InstantID-generated images is crucial to avoid overfitting 481 and loss of recontextualization. 482

FaceDistance provided a quantitative measure of facial 483 similarity but became less useful for hyperparameter tun-484 ing once a certain fidelity level was reached. A user survey 485 among white-collar workers showed that both DreamBooth 486 and InstantID performed similarly in quality, clarity, iden-487 tity preservation, editing, and background. A slight pref-488 erence emerged for the "Photoshopped look" of InstantID 489 portraits. While DreamBooth achieved better facial similar-490 ity, many participants struggled to distinguish AI-generated 491 images from real ones, particularly those unfamiliar with 492 AI tools. Users actively engaged in AI image creation were 493 more likely to identify DreamBooth images as synthetic, 494 possibly due to its higher popularity. 495

InstantID's effectiveness depends on reference image 496 quality and diversity. Using multiple references (around 497 four) improved similarity by enriching information for the 498 Face Encoder. Facial landmarks strongly influenced pose 499 and composition, sometimes overriding text prompts. We 500 explored 2-shot generation and interactive face replace-501 ment to enhance control, with the latter showing higher 502 user satisfaction. Rotational and shape-altering augmen-503 tations were ineffective, and background modifications re-504 duced similarity. Traditional upscaling worked well for 505 low-resolution images, whereas neural network-based up-506 scaling introduced artifacts. 507

# 7. Limitations

A key limitation is that InstantID-based augmentation509reduces realism in generated images.510Booth remains more flexible for personalized generation,511InstantID-enhanced datasets still outperform unaugmented512ones. Given the baseline model's photorealism constraints,513using generative augmentation to refine its training data is a514practical approach.515

## 8. Conclusion

This study examined augmentation strategies for improv-<br/>ing facial resemblance in personalized image generation<br/>using DreamBooth and InstantID. Classical augmentations<br/>can introduce artifacts that degrade facial fidelity, requiring<br/>careful application.517519520521

We found generative augmentation with InstantID to be highly effective for improving DreamBooth training. Creating diverse, realistic synthetic images while maintaining a 524 525 balanced ratio with real data prevents overfitting.

526 User surveys confirmed that both DreamBooth and In527 stantID produce high-quality, professional-looking head528 shots, often indistinguishable from real photos. While
529 DreamBooth excels in facial similarity, InstantID's consis530 tent output appears more polished.

For practical use, employing multiple reference images
enhances facial information capture. Improving control
over pose and composition through landmarks is crucial,
with interactive face replacement showing promise.

535Overall, our findings provide insights into augmenta-536tion strategies for personalized image generation, guiding537their application in tasks requiring high facial fidelity. Fu-538ture work should explore advanced generative augmentation539techniques and better user control over InstantID outputs.