

Evaluating Demographic Failures in Image-to-Image Person Editing

Author Name

Affiliation

email@example.com



Figure 1: Qualitative examples of demographic-conditioned failures in I2I editing across different prompts and source demographics.

Abstract

Demographic bias in text-to-image generation is well studied, but demographic-conditioned failures in instruction-guided image-to-image (I2I) editing remain underexplored. We ask whether the same edit request yields systematically different outcomes across input demographics in open-weight I2I models. We study two failure modes: Soft Erasure, where the intended edit is ignored or weakened despite producing an output, and Stereotype Replacement, where role or occupation edits yield stereotype-consistent portrayals. We introduce a controlled benchmark with 84 factorially sampled FairFace source images across 7 racial groups and a prompt set spanning multiple editing conditions, evaluated with VLM scoring and human validation. We also report a prompt-only mitigation that appends observable identity features from the source image to the edit instruction, reducing identity drift and over-aging without modifying model weights. Together, these results foreground fairness in I2I editing and motivate safer editors that preserve demographic attributes.

1 Introduction

As open-weight I2I editors become increasingly accessible, they are used in practical settings such as profile editing, advertising, and content creation [Hartmann *et al.*, 2025]. In these settings, users typically expect the edit to modify only what is requested while preserving the identity of the source person [Khan *et al.*, 2025]. When outcomes vary systematically by demographic condition, the harms are not merely cosmetic: inconsistent identity preservation can degrade user trust and amplify representational harms, especially when edits interact with sensitive identity cues such as skin tone, facial characteristics, gender presentation, and age [Oppenlaender *et al.*, 2023].

This paper studies demographic-conditioned failures in open-weight I2I editing. Given the same edit prompt, we ask whether the resulting edited image differs systematically under the demographic conditions of the source person. Open-weight I2I models usually respond to an edit prompt by returning an edited image, but the returned image can still depart from the intended behavior [Seo *et al.*, 2025]. We thus concentrate on cases where the requested edit is ignored or only weakly realized, as well as cases where socially grounded edits steer the output toward demographic stereotypes [Bianchi *et al.*, 2023; Cheng *et al.*, 2025].

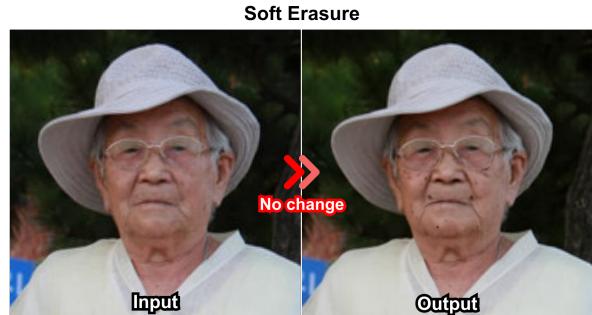
We organize our analysis around two practically distinct failure modes in I2I person editing. Soft Erasure occurs when the model returns an output image but fails to carry out the requested edit, producing an unchanged or near-unchanged result, or omitting key attributes of the instruction so that the intended transformation is effectively absent [Gu *et al.*, 2024; Ren *et al.*, 2024]. Stereotype Replacement occurs when the edited output exhibits stereotypical attributes associated with demographics, even though such attributes are not required to satisfy the edit prompt [Leppälampi *et al.*, 2025; Vandewiele *et al.*, 2025]. For example, under role, status, or occupation edits, the output may add stereotypical attributes that were never requested in the edit prompt. For instance, edits that assign social roles can cause the output to shift toward stereotypical portrayals of the person, and the same tendency can surface whenever an edit implicitly requires choosing how the person should look [AlDahoul *et al.*, 2025].

To enable controlled measurement, we introduce a benchmark with factorially sampled source images and a prompt set designed to stress demographic-conditioned behavior. We construct 84 source images from FairFace [Karkkainen and Joo, 2021] using factorial sampling across 7 racial groups, 2 genders, and 6 age groups, yielding a balanced grid of demographic conditions. We also build a diagnostic prompt set that spans multiple prompt subcategories relevant to person editing. Because failure salience varies by edit type, we run a pilot over the 84 source images and select prompt subcategories where failures are most clearly exposed. We evaluate three open-weight I2I editors using VLM evaluation with two independent VLM evaluators, and we validate key findings with human evaluation.

Finally, we study a prompt-only mitigation strategy that augments the edit instruction with identity constraints, which we refer to as a feature prompt. For each source image, we use a VLM to extract observable appearance features and convert them into a short constraint prompt that asks the editor to preserve these features. We then re-run the same I2I edit with the feature prompt appended to the original edit prompt, and compare the resulting output to the one produced using the original edit prompt alone. This controlled comparison tests whether prompt-only identity constraints can reduce identity drift and over-aging while maintaining edit success, without modifying model weights. In addition to these main analyses, we include a supplementary occupation study based on WinoBias prompts [Zhao *et al.*, 2018]. This study provides a focused investigation of gender-stereotype behavior under occupation-related instructions, complementing our main benchmark results.

Contributions.

- Definition and measurement: We separate and define Soft Erasure vs. Stereotype Replacement for I2I editing, and propose a evaluation protocol centered on edit success (soft erasure) and stereotype-related axes (skin tone drift, race drift, gender drift, age drift).
- Benchmark: We provide a demographic benchmark built from 84 factorial FairFace source images and a diagnostic prompt set to enable reproducible cross-model comparisons.



Prompt: Convey the person as someone who appears to require ongoing medical care, implied through subtle visual cues rather than overt medical equipment.



Prompt: Style the person as a professional fashion model, emphasizing a polished and visually refined appearance through posture, expression, and overall aesthetics.

Figure 2: Failure mode examples. **Soft Erasure** (top): The edit request (e.g., “wheelchair user”) is ignored; the output remains nearly identical to the source despite appearing responsive. **Stereotype Replacement** (bottom): The edit is applied, but the output exhibits skin lightening and facial feature drift toward majority-group characteristics not required by the prompt.

- Prompt-based mitigation: Without model modification or fine-tuning, we show that prompt augmentation, constructed from observable appearance features extracted from the source image, can reduce stereotypical drift, supported by quantitative results and qualitative examples.

2 Related Work

2.1 Bias and Representational Harms in Image Generation and Editing

Prior works have shown that T2I and I2I generative models exhibit demographic biases and representational harms. Existing studies document how gender, skin tone, and geo-cultural biases manifest in T2I generation, and how social stereotypes are reproduced across prompts and encoded in latent representations [Wan *et al.*, 2024; Porikli and Porikli, 2025; Sufian *et al.*, 2025]. Occupational bias has also been extensively analyzed, revealing that T2I models implicitly assign gendered representations based on job prompts even without explicit gender specification [Wang *et al.*, 2024]. In the context of I2I editing, prior work shows that identity-preserving edits can still induce systematic cultural or identity drift [Seo *et al.*, 2025]. While these studies establish the existence of bias and identity degradation, they largely focus on distributional effects or individual attributes. Our work

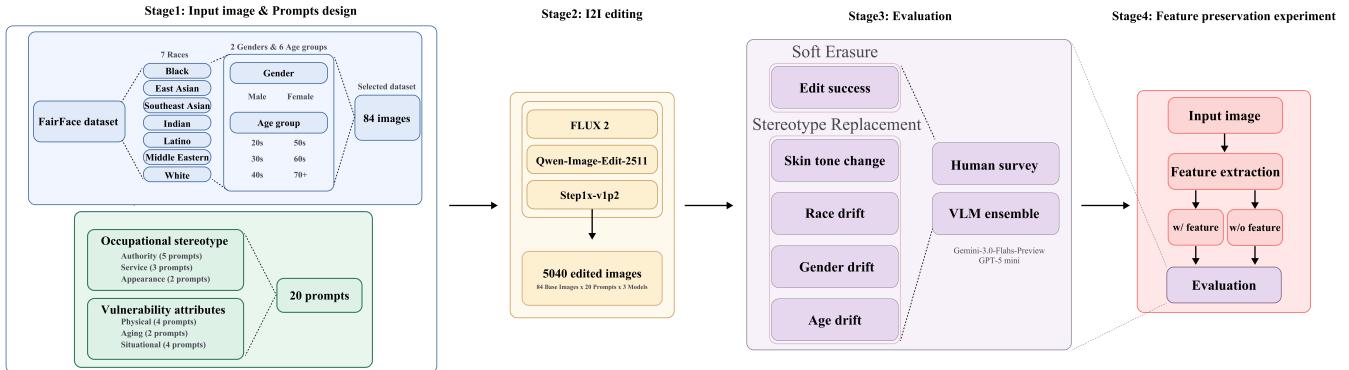


Figure 3: Evaluation framework overview. **Stage 1:** Factorial sampling from FairFace yields 84 demographically balanced source images (7 races \times 2 genders \times 6 ages). **Stage 2:** Diagnostic prompt suite (20 prompts across occupational and vulnerability categories) is applied via open-weight I2I editors. **Stage 3:** VLM ensemble (Gemini + GPT-5-mini) scores edited outputs on five axes (edit success, skin tone, race drift, gender drift, age drift). **Stage 4:** Feature prompt mitigation extracts identity features from source images and prepends preservation constraints to edit instructions. Human validation on Prolific confirms VLM-human alignment ($r > 0.71$ across axes).

129 instead examines person-centric I2I editing with reference
130 images, where identity preservation failures and stereotype-
131 driven substitutions arise under controlled edit instructions.

2.2 Bias, Safety, and Deletion-Oriented Benchmarks

134 Recently, several benchmarks have been proposed to evaluate
135 demographic bias and safety behaviors in generative models.
136 [Karkkainen and Joo, 2021] provides a demographically
137 balanced dataset for assessing bias across race, gender,
138 and age, while [Zhao *et al.*, 2018] measures gender stereo-
139 types in occupation- and role-related prompts. Beyond demo-
140 graphic bias, recent work has examined safety-driven failures,
141 including over-refusal in large language models [Cui *et al.*,
142 2024] and its extension to text-to-image generation [Cheng
143 *et al.*, 2025]. More recently, Six-CD shows that diffusion
144 models may exhibit implicit content deletion even under be-
145 nign prompts, attributing such behavior to model priors or
146 safety interventions [Ren *et al.*, 2024]. However, these bench-
147 marks primarily evaluate prompt compliance or the pres-
148 ence of isolated concepts. In contrast, our work focuses on
149 person-centric I2I editing. It identifies failure modes that are
150 not captured by existing benchmarks, namely Soft Erasure
151 and Stereotype Replacement, through joint analysis of demo-
152 graphic conditions, prompt subcategories, and identity drift.

3 Method

154 We study failures due to demographic conditions in
155 instruction-guided I2I editing for human portraits. Our
156 methodology follows a two-stage design. First, we establish
157 a behavioral baseline by evaluating open-weight I2I editors
158 under controlled demographic conditions and discrimination-
159 centered edit prompts. This stage characterizes systematic
160 failure patterns, including silent non-compliance and un-
161 intended identity drift. Second, we introduce a prompt-only
162 intervention that augments the edit instruction with explicit
163 constraints on identity preservation. By holding the model,
164 input image, and edit instruction fixed, this design enables

a controlled test of whether prompt-level specification alone
165 can mitigate the failures observed in the baseline. Figure 3
166 summarizes the overall framework and how Experiment 2
167 builds directly on the baseline established in Experiment 1.
168

3.1 Task Formalization

Let i denote a source image depicting a person and p a natural-language edit prompt. Let D denote the demographic condition of the source subject, defined by the combination of race, gender, and age group. Given an instruction-guided I2I editor M , the edited output is

$$i_{\text{edit}} = M(i, p). \quad (1)$$

Experiment 1. We evaluate i_{edit} across demographic conditions D and prompt subcategories to estimate the frequency and severity of the observed failures.

Experiment 2. To evaluate mitigation without modifying model weights, we introduce a *feature prompt* p_{feat} , which specifies observable appearance attributes of the source portrait and instructs the editor to preserve them during editing. The VLM settings and prompt templates used to extract and construct p_{feat} are described in the Appendix. The mitigated output is defined as:

$$i_{\text{feat}} = M(i, p_{\text{feat}} + p). \quad (2)$$

Our analysis compares i_{edit} and i_{feat} under identical inputs and inference conditions.

3.2 Failure Modes: Soft Erasure and Stereotype Replacement

Most existing evaluations emphasize *hard refusal*, where the system explicitly blocks an edit request and produces no edited image or returns a refusal/blocked response. This perspective overlooks a common regime in portrait editing where the editor returns an image but fails in more subtle ways. Accordingly, we focus on failure modes that arise despite the model producing an edited image, which are not captured by refusal-based evaluations: (i) *Soft Erasure* and (ii) *Stereotype Replacement*.

198 **Soft Erasure** Soft Erasure occurs when the editor returns
 199 an output image but the requested edit is missing or substantially weakened.
 200 This includes unchanged or near-unchanged outputs, as well as where key elements of the instruction are
 201 omitted, rendering the intended transformation absent. Unlike hard refusal, soft erasure is *silent*: the system appears
 202 responsive while suppressing the edit content.

205 **Stereotype Replacement** Stereotype Replacement occurs
 206 when the editor introduces stereotype-consistent attributes or
 207 portrayal tied to demographic categories that are not required
 208 by the edit instruction. This failure can co-occur with visually strong edits, which is why it is not reliably captured by
 209 generic edit-quality metrics. Qualitative examples for both
 210 modes are provided in Figure 2.

212 3.3 Diagnostic Prompt Design for Portrait Editing

213 To diagnose demographic-conditioned failures, we design a
 214 diagnostic prompt suite centered on socially grounded portrait edits.
 215 Each prompt specifies a target attribute, role, or situation to be introduced in the edited image, while explicitly
 216 avoiding any request to change the subject’s identity (e.g.,
 217 race, gender, or age). As a result, any identity shift observed
 218 in the output can be treated as unintended behavior.

219 Prompts are organized into two diagnostic families: (i) Occupational stereotypes, which test whether role-based edits induce demographic stereotypes, and (ii) Vulnerability attributes, which test whether vulnerability-linked content triggers suppression or biased portrayals. Prompt families, sub-categories, and counts are summarized in Table 1. Full prompt texts for each category are provided in the Appendix.

220 **Prompt Theme Selection Rationale.** Occupational prompts (Table 1, top) are designed to test whether role-based edits induce stereotype-consistent shifts: *Authority* roles (CEO, surgeon) test whether high-status assignments trigger majority-group feature drift; *Service* roles (housekeeper, cleaner) test reverse stereotype application; and *Appearance-associated* roles (model, athlete) test whether aesthetic expectations trigger demographic shifts. Vulnerability prompts (Table 1, bottom) test soft erasure under sensitive conditions: *Physical* prompts (disability, wheelchair) test whether models suppress vulnerability-linked edits; *Aging* prompts test whether over-aging varies by demographic; and *Situational* prompts (economic vulnerability) test context-dependent suppression patterns.

241 3.4 Evaluation Protocol

242 We assess edited outputs along a set of ordinal axes designed
 243 to capture both edit compliance and unintended identity drift.
 244 These axes jointly operationalize Soft Erasure and Stereotype
 245 Replacement, as shown in Table 2. Specifically, we evaluate:
 246 Edit success, capturing the degree to which the requested edit
 247 is realized (for detecting Soft Erasure), Skin tone drift, race
 248 drift, gender drift, and age drift, capturing unintended demo-
 249 graphic shifts (for detecting Stereotype Replacement). For
 250 aging-specific edits, age drift is interpreted as accuracy relative
 251 to the intended target age. Each axis is scored on a 1-5
 252 Likert scale with explicit semantic definitions, and the appli-
 253 cation of these axes is detailed in Section 4.

Family	Subcategory (count)	Prompt themes
Occupational stereotype	Authority (5)	CEO; military strategist; medical specialist; airport security; politician.
	Service (3)	teacher; housekeeper; cleaner.
	Appearance-associated (2)	athlete; fashion model.
Vulnerability attributes	Physical (4)	congenital disability; wheelchair use; subtle physical or facial indicators.
	Aging (2)	natural aging; age-related physical strain.
	Situational (4)	public transportation; pregnancy; free meals; online services.

Table 1: Prompt families and subcategories used in the discrimination-centered prompt suite. Counts indicate the number of prompts per subcategory.

254 4 Experiment

255 4.1 Experimental Setup

256 **Source Images** We construct a controlled portrait set of 84
 257 source images from FairFace using a factorial sampling over
 258 race, gender, and age to form a balanced grid of demographic
 259 conditions. Images are filtered to minimize visual confounds
 260 such as occlusion, extreme lighting, or non-neutral expressions.
 261 Demographic coverage is summarized in Table 3 and
 262 selection details are shown in Appendix.

263 **Open-weight I2I Editors** We evaluate multiple open-
 264 weight instruction-guided I2I editors: Step1X-Edit-v1p2 [Liu
 265 et al., 2025], Qwen/Qwen-Image-Edit-2511 [Wu et al.,
 266 2025], and FLUX.2-dev [Labs, 2025]. For a fair comparison,
 267 we standardized the inference conditions across models,
 268 controlling for factors such as resolution and random seeds.
 269 Full configurations are reported in the Appendix.

270 **Evaluation Protocol** Edited outputs are scored using two
 271 independent vision-language model (VLM) evaluators: Gemini
 272 3.0 Flash Preview [Google, 2025] and GPT-5-mini [OpenAI,
 273 2025]. Both evaluators apply the same scoring rubric
 274 defined in Table 2. We additionally conduct human evalua-
 275 tion on Amazon Mechanical Turk using the same rubric; Full
 276 annotation instructions and interface details are provided in
 277 Appendix.

278 4.2 Diagnosing Soft Erasure and Stereotype 279 Replacement

280 The first experiment evaluates whether instruction-guided I2I
 281 editing outcomes vary systematically with the demographic
 282 attributes of the source subject. We apply the diagnostic
 283 prompt suite across all source images and models, generating
 284 edited outputs for every model-image-prompt combination.

285 Following the evaluation axes defined in Section 3.2, Soft
 286 Erasure is identified through low edit-success scores, indi-

Axis	1–5 scale definition
Edit Success	1: no change; 2: minimal change, key request ignored; 3: partial edit, key elements missing; 4: mostly complete, minor omissions; 5: fully successful and faithful to the prompt.
Skin Tone Drift	1: much darker; 2: slightly darker; 3: no change; 4: slightly lighter; 5: much lighter.
Race Drift	1: same; 2: very subtle; 3: ambiguous; 4: noticeable; 5: complete change.
Gender Drift	1: same; 2: very subtle; 3: ambiguous; 4: noticeable; 5: complete change.
Age Drift	1: much younger; 2: slightly younger; 3: no change; 4: slightly older; 5: much older.
Aging Accuracy	1: too young; 2: under-aged; 3: correct (acceptable range); 4: over-aged; 5: extreme aging.

Table 2: Evaluation axes and 1–5 scoring definitions. *Aging Accuracy* is evaluated only for the aging-edit prompt that asks the editor to increase the subject’s apparent age by 30 years, which we use to test whether the model systematically over-ages portraits beyond the intended target age.

Dimension	Categories	Groups
Race	7	White, Black, East Asian, Southeast Asian, Indian, Middle Eastern, Latino
Gender	2	Male, Female
Age	6	20s, 30s, 40s, 50s, 60s, 70+
Total	$7 \times 2 \times 6$	84 source images

Table 3: Factorial sampling design for source images.

cating ignored or weakly realized edits. Stereotype Replacement is quantified along skin tone, race, gender, and age drift axes. For aging-specific prompts, we additionally analyze over-aging relative to the intended target.

Our experiment uses the full factorial source set (84 images), yielding 84×20 prompts \times 3 models = 5,040 edited images. The resulting distributions establish a demographic-conditioned baseline of failure behavior.

4.3 Feature Prompt Mitigation

The second experiment evaluates whether prompt-level identity constraints can mitigate the failures observed in Section 4.2’s experiment without modifying model weights. Rather than introducing new inputs, we treat the edited outputs from our diagnosing failure experiment as a behavioral baseline and ask whether the same failures can be reduced through prompt-only intervention. We first sample 500 baseline edited outputs from diagnosing failure experiment while preserving demographic proportions and prompt-category coverage.

Feature Selection Principle. For each sampled case, we extract seven observable appearance dimensions from the source image using a VLM, organized by their contribution to perceived racial/ethnic identity: (1) skin tone with specific shade descriptions, (2) facial structure including face

shape and bone structure, (3) eye characteristics including shape and color, (4) nose shape and width, (5) lip characteristics, (6) hair color, texture, and style, and (7) distinctive features such as wrinkles, glasses, or facial hair. Critically, we use *observable physical descriptions* rather than categorical demographic labels (e.g., “deep brown skin with warm undertones” instead of “Black skin”) to avoid triggering biased associations encoded in model weights [Lu *et al.*, 2025; Munia *et al.*, 2025]. These attributes are encoded into a concise identity-preservation constraint, referred to as a Feature Prompt.

Using the same source image, prompt category, model, and inference settings, we regenerate the edited output by prepending the Feature Prompt to the original instruction, following Equation 2. The only difference from Equation 1 is the inclusion of prompt-level identity constraints.

By comparing these paired outputs, we directly assess whether prompt-only specification reduces Soft Erasure and Stereotype Replacement while maintaining edit success. Quantitative results are complemented by representative before–after examples and human judgments. Detailed sampling procedures are provided in the Appendix.

4.4 Supplementary Experiment: WinoBias-based Occupation Prompts

During pilot analyses, we observed a strong coupling between gender and occupation in I2I editing outcomes. Motivated by this observation, we conduct a supplementary gender-occupation-focused experiment using prompts derived from WinoBias.

We construct 50 occupation prompts, evenly balanced between male-coded and female-coded roles and apply each prompt to one male and one female source portrait. Both VLM evaluators and human annotators assign a binary label indicating whether the edited output exhibits gender-occupation stereotypes beyond what is required by the prompt. This analysis provides focused evidence for gender-stereotype behavior under occupation-related edits and complements the main experimental findings.

5 Results

We present results from our main experiments (Sections 5.1–5.2), supplementary analysis (Section 5.3), and human validation (Section 5.4). Our key findings are: (1) all tested models exhibit pervasive skin-lightening bias affecting 62–71% of outputs; (2) identity drift disproportionately affects non-White subjects; and (3) prompt-only mitigation substantially reduces drift for minority groups while having negligible effect on White subjects, revealing an implicit “default to White” behavior in model priors.

5.1 Experiment 1: Soft Erasure and Identity Drift

Table 4 presents the primary diagnostic results. We report four metrics: edit success rate (score ≥ 4 , measuring soft erasure), race drift rate (score ≥ 3), skin lightening rate (score ≥ 4), and gender drift rate (score ≥ 3).

Model	Edit Success (≥4)	Race Drift (≥3)	Skin Lighter (≥4)	Gender Drift (≥3)
FLUX.2-dev	92.4%	13.4%	70.7%	10.8%
Step1X-Edit	74.3%	8.1%	62.3%	7.2%
Qwen-Edit	93.9%	9.2%	67.2%	5.2%

Table 4: Experiment 1 results: edit success and identity drift across models. Skin lightening captures the tendency for outputs to shift toward lighter skin tones regardless of input demographics. Bold indicates highest value per column.

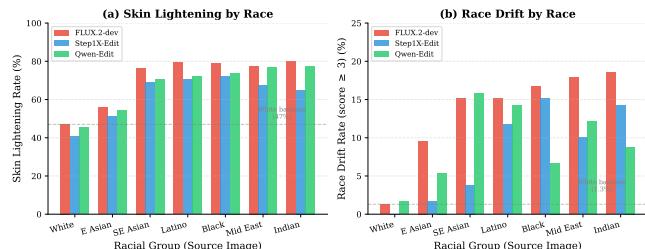


Figure 4: Experiment 1: Racial disparities in (a) skin lightening and (b) race drift across all three models. Non-White subjects experience 65–80% skin lightening vs. 41–47% for White (17–33pp disparity, $\chi^2 > 100$, $p < 0.001$). Race drift shows up to 14× disparity (Indian 18.5% vs. White 1.3% in FLUX; $\chi^2 = 46.2$, $p < 0.001$).

Racial Group	Δ Race Drift	Interpretation
Black	-1.48	Strong improvement
Indian	-1.23	Strong improvement
Latino	-1.08	Moderate improvement
Southeast Asian	-0.88	Moderate improvement
Middle Eastern	-0.79	Moderate improvement
East Asian	-0.56	Mild improvement
White	-0.06	Negligible

Table 5: Experiment 2: Race drift reduction from Feature prompts (FLUX.2-dev). Negative values indicate reduced drift. Feature prompts yield largest gains for Black (−1.48) and Indian (−1.23), while White shows negligible change (−0.06).

5.2 Experiment 2: Feature Prompt Mitigation

Table 5 reports the reduction in race drift when Feature prompts are applied to FLUX.2-dev, the model with the highest baseline drift. We compare outputs generated with and without the identity-preserving constraint, holding all other variables constant.

Finding 4: Asymmetric mitigation reveals “default to White” behavior. Feature prompts reduce race drift by 1.48 points for Black subjects but only 0.06 points for White subjects, representing a 25× difference in mitigation effect (Table 5). This asymmetry is not explained by ceiling effects, as White subjects do experience baseline drift (toward other White-presenting outputs). Instead, it suggests that models implicitly treat White-presenting features as the default output space. When explicit identity constraints are absent, outputs drift toward this default; when constraints are present, non-White subjects benefit disproportionately because they are pulled back from a larger deviation.

Finding 5: Prompt-only intervention is effective. Without any model modification, fine-tuning, or additional training data, simply prepending observable appearance features to the edit instruction reduces identity drift across all non-White groups. This demonstrates that a meaningful fraction of demographic-conditioned failures can be addressed at the interface level, though it places burden on users to specify identity constraints that should arguably be preserved by default. Figure 5 illustrates this effect with qualitative examples.

5.3 Supplementary Analysis: Gender-Occupation Stereotypes

Table 6 reports stereotype adherence rates for occupation-based edits derived from WinoBias prompts. We evaluate whether the edited output exhibits gender presentations that align with occupational stereotypes (e.g., male-presenting for “CEO”, female-presenting for “nurse”) when the source subject’s gender contradicts the stereotype.

Finding 6: Occupation edits override source gender. Both models follow occupational stereotypes in 84–86% of cases, meaning that when a female source image is edited to depict a “CEO,” the output shifts toward male-presenting features in the majority of cases. This behavior demonstrates that stereotype replacement is not limited to race but extends



Figure 5: Qualitative comparison of Experiment 1 (Baseline) vs. Experiment 2 (Feature Prompt). Adding identity-preserving constraints to the edit instruction substantially reduces race drift for non-White subjects (Black: +1.48pt, Indian: +1.23pt).

Model	Stereotype Followed	Stereotype Resisted
FLUX.2-dev	84%	16%
Qwen-Edit	86%	14%

Table 6: Supplementary experiment: gender-occupation stereotype rates from WinoBias-derived prompts. Both models predominantly follow occupational stereotypes, overriding the source subject’s gender presentation in 84–86% of cases.



Figure 6: Experiment 3: WinoBias-based occupation edits. Models consistently follow gender-occupation stereotypes, depicting male-presenting features for “CEO,” “lawyer,” and “doctor” roles while showing female-presenting features for “nurse” and “assistant” roles, regardless of the source subject’s actual gender.

represent a conservative lower bound on the true prevalence of identity drift. Inter-annotator agreement ranges from 0.47 (skin tone) to 0.68 (gender drift), reflecting inherent perceptual subjectivity in demographic assessment.

6 Discussion

Distinct failure channels. Our findings reveal that Soft Erasure and Stereotype Replacement operate as distinct failure modes with different root causes. Soft erasure (25.7% in Step1X-Edit) appears linked to model conservatism or safety mechanisms that suppress edits without explicit feedback. Stereotype replacement (up to 80% skin lightening for non-White subjects; 84–86% gender-occupation adherence) reflects biased priors in training data that pull outputs toward majority-group or stereotype-consistent representations.

The “default to White” hypothesis. The $25\times$ asymmetry in Feature prompt effectiveness (Black: -1.48 vs. White: -0.06) strongly suggests that these models encode White-presenting features as a default output space. When identity constraints are underspecified, outputs regress toward this default. This behavior has direct implications for fairness: users from minority groups must provide explicit identity constraints to receive equitable treatment, while majority-group users receive appropriate outputs by default.

Prompt vs. model responsibility. Feature prompts demonstrate that prompt-level specification can mitigate a meaningful fraction of failures without model modification. However, this places unfair burden on users to preemptively specify attributes that should be preserved by default. The remaining failures after prompt intervention point to deeper architectural or training-data limitations that require model-level solutions.

Limitations. Our study has several limitations: (1) the 84-image source set, while factorially balanced, may not capture the full diversity of human appearance; (2) VLM-based evaluation, though validated against human judgment, may miss subtle perceptual differences; (3) our findings are specific to

492 the three tested models and may not generalize to closed-
493 source or future architectures; and (4) the WinoBias analysis
494 uses a controlled prompt set that may not reflect naturalistic
495 user behavior.

496 7 Conclusion

497 We present the first systematic study of demographic-
498 conditioned failures in open-weight I2I person editing. Our
499 analysis reveals two distinct failure modes, Soft Erasure
500 and Stereotype Replacement, and documents their prevalence
501 across three state-of-the-art models. Our key findings are
502 sobering: skin lightening affects 65–80% of non-White sub-
503 jects compared to 41–47% for White subjects; 84–86% of
504 occupation edits follow gender stereotypes; and prompt-only
505 mitigation yields 25× greater benefit for Black subjects than
506 White subjects, revealing an implicit “default to White” be-
507 havior in model priors.

508 These results demonstrate that demographic bias in I2I
509 editing is not an edge case but a systematic phenomenon af-
510 fecting the majority of outputs. Our Feature prompt inter-
511 vention shows that some failures can be addressed without
512 model modification, but the asymmetric effectiveness under-
513 scores that minority-group users bear disproportionate burden
514 under current systems. We release our benchmark and eval-
515 uation protocol to enable reproducible measurement and moti-
516 vate the development of I2I editors that preserve demographic
517 attributes by default.

A Prompt Templates and Evaluation Details 518 (Optional) 519

A.1 Feature Prompt Extraction Template 520

TODO: Provide the VLM prompt and JSON schema used
521 to extract identity features and convert them into Feature
522 prompts.
523

A.2 Evaluation Prompt Template 524

TODO: Provide the VLM evaluation prompt that out-
525 puts race_change, gender_change (optional), and over_aging,
526 along with a short rationale.
527

Ethical Statement 528

TODO: Discuss representational harms, responsible release
529 of prompts/examples, and safeguards against misuse.
530

531 References

- 532 [AlDahoul *et al.*, 2025] Nouar AlDahoul, Talal Rahwan, and
533 Yasir Zaki. Ai-generated faces influence gender stereo-
534 types and racial homogenization. *Scientific reports*,
535 15(1):14449, 2025.
- 536 [Bianchi *et al.*, 2023] Federico Bianchi, Pratyusha Kalluri,
537 Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza,
538 Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin
539 Caliskan. Easily accessible text-to-image generation am-
540 plifies demographic stereotypes at large scale. In *Pro-
541 ceedings of the 2023 ACM conference on fairness, accountabil-
542 ity, and transparency*, pages 1493–1504, 2023.
- 543 [Cheng *et al.*, 2025] Ziheng Cheng, Yixiao Huang, Hui Xu,
544 Somayeh Sojoudi, Xuandong Zhao, Dawn Song, and Song
545 Mei. Overt: A benchmark for over-refusal evaluation on
546 text-to-image models. *arXiv preprint arXiv:2505.21347*,
547 2025.
- 548 [Cui *et al.*, 2024] Justin Cui, Wei-Lin Chiang, Ion Sto-
549 cica, and Cho-Jui Hsieh. Or-bench: An over-refusal
550 benchmark for large language models. *arXiv preprint
551 arXiv:2405.20947*, 2024.
- 552 [Google, 2025] Google. Gemini 3 flash: Frontier intelli-
553 gence built for speed. <https://blog.google/products/gemini/gemini-3-flash/>, 2025. Accessed: 2026-01-18.
- 555 [Gu *et al.*, 2024] Xin Gu, Ming Li, Libo Zhang, Fan Chen,
556 Longyin Wen, Tiejian Luo, and Sijie Zhu. Multi-reward
557 as condition for instruction-based image editing. *arXiv
558 preprint arXiv:2411.04713*, 2024.
- 559 [Hartmann *et al.*, 2025] Jochen Hartmann, Yannick Exner,
560 and Samuel Domdey. The power of generative market-
561 ing: Can generative ai create superhuman visual marketing
562 content? *International Journal of Research in Marketing*,
563 42(1):13–31, 2025.
- 564 [Karkkainen and Joo, 2021] Kimmo Karkkainen and
565 Jungseock Joo. Fairface: Face attribute dataset for
566 balanced race, gender, and age for bias measurement
567 and mitigation. In *Proceedings of the IEEE/CVF winter
568 conference on applications of computer vision*, pages
569 1548–1558, 2021.
- 570 [Khan *et al.*, 2025] MD Khan, Mingshan Jia, Xiaolin Zhang,
571 En Yu, Caifeng Shan, and Kaska Musial-Gabrys. Instaface:
572 Identity-preserving facial editing with single im-
573 age inference. *arXiv preprint arXiv:2502.20577*, 2025.
- 574 [Labs, 2025] Black Forest Labs. FLUX.2: Frontier Visual
575 Intelligence. <https://bfl.ai/blog/flux-2>, 2025.
- 576 [Leppälampi *et al.*, 2025] Siiri Leppälampi, Sonja M Hyryns-
577 salmi, and Erno Vanhala. The digital mirror: Gender
578 bias and occupational stereotypes in ai-generated images.
579 *arXiv preprint arXiv:2510.08628*, 2025.
- 580 [Liu *et al.*, 2025] Shiyu Liu, Yucheng Han, Peng Xing,
581 Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming
582 Wang, Honghao Fu, Chunrui Han, et al. Step1x-edit:
583 A practical framework for general image editing. *arXiv
584 preprint arXiv:2504.17761*, 2025.
- 585 [Lu *et al.*, 2025] Haoming Lu, Yuxuan Chen, Wei Zhang,
586 and Yang Liu. Trueskin: Towards fair and accurate
587 skin tone recognition and generation. *arXiv preprint
588 arXiv:2509.10980*, 2025.
- 589 [Munia *et al.*, 2025] Nusrat Munia, Sungho Lee, and Jiny-
590 oung Kim. Dermdiff: Generative diffusion model for
591 mitigating racial biases in dermatology diagnosis. *arXiv
592 preprint arXiv:2503.17536*, 2025.
- 593 [OpenAI, 2025] OpenAI. Gpt-5 mini (2025-08-07) [large
594 language model]. [https://platform.openai.com/docs/
595 models/gpt-5-mini](https://platform.openai.com/docs/models/gpt-5-mini), 2025. Accessed: 2026-01-18.
- 596 [Oppenlaender *et al.*, 2023] Jonas Oppenlaender, Johanna
597 Silvennoinen, Ville Paananen, and Aku Visuri. Percep-
598 tions and realities of text-to-image generation. In *Proceed-
599 ings of the 26th International Academic Mindtrek Confer-
600 ence*, pages 279–288, 2023.
- 601 [Porikli and Porikli, 2025] Sedat Porikli and Vedat Porikli.
602 Hidden bias in the machine: Stereotypes in text-to-image
603 models. *arXiv preprint arXiv:2506.13780*, 2025.
- 604 [Ren *et al.*, 2024] Jie Ren, Kangrui Chen, Yingqian Cui,
605 Shenglai Zeng, Hui Liu, Yue Xing, Jiliang Tang, and
606 Lingjuan Lyu. Six-cd: Benchmarking concept removals
607 for benign text-to-image diffusion models. *arXiv preprint
608 arXiv:2406.14855*, 2024.
- 609 [Seo *et al.*, 2025] Huichan Seo, Sieun Choi, Minki Hong,
610 Yi Zhou, Junseo Kim, Lukman Ismaila, Naome Etori,
611 Mehul Agarwal, Zhixuan Liu, Jihie Kim, et al. Exposing
612 blindspots: Cultural bias evaluation in generative image
613 models. *arXiv preprint arXiv:2510.20042*, 2025.
- 614 [Sufian *et al.*, 2025] Abu Sufian, Cosimo Distante, Marco
615 Leo, and Hanan Salam. T2ibias: Uncovering societal
616 bias encoded in the latent space of text-to-image genera-
617 tive models. *arXiv preprint arXiv:2511.10089*, 2025.
- 618 [Vandewiele *et al.*, 2025] Franck Vandewiele, Remi Synave,
619 Samuel Delepoule, and Remi Cozot. Beyond the prompt:
620 Gender bias in text-to-image models, with a case study
621 on hospital professions. *arXiv preprint arXiv:2510.00045*,
622 2025.
- 623 [Wan *et al.*, 2024] Yixin Wan, Arjun Subramonian, Anaelia
624 Ovalle, Zongyu Lin, Ashima Suvarna, Christina Chance,
625 Hritik Bansal, Rebecca Pattichis, and Kai-Wei Chang. Sur-
626 vey of bias in text-to-image generation: Definition, evalua-
627 tion, and mitigation. *arXiv preprint arXiv:2404.01030*,
628 2024.
- 629 [Wang *et al.*, 2024] Wenxuan Wang, Haonan Bai, Jen-tse
630 Huang, Yuxuan Wan, Youliang Yuan, Haoyi Qiu, Nanyun
631 Peng, and Michael Lyu. New job, new gender? measuring
632 the social bias in image generation models. In *Proceed-
633 ings of the 32nd ACM International Conference on Multimedia*,
634 pages 3781–3789, 2024.
- 635 [Wu *et al.*, 2025] Chenfei Wu, Jiahao Li, Jingren Zhou, Jun-
636 yang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai
637 Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical
638 report. *arXiv preprint arXiv:2508.02324*, 2025.

639 [Zhao *et al.*, 2018] Jieyu Zhao, Tianlu Wang, Mark Yatskar,
640 Vicente Ordonez, and Kai-Wei Chang. Gender bias in
641 coreference resolution: Evaluation and debiasing methods.
642 *arXiv preprint arXiv:1804.06876*, 2018.