

LatentCompass: T2I Diffusion Steering via Orthogonal Attribute Spaces for Debiasing, Concept Erasure, and Red Teaming

Anonymous Author(s)
 Affiliation
 Address
 email

Abstract

1 Text-to-image (T2I) diffusion models suffer from biased results stemming from
 2 entangled generative priors and a lack of accurate control over outputs. Current
 3 mitigation attempts rely on imprecise, adversarially-vulnerable prompt and text
 4 embedding interventions, or they require prohibitive and invasive fine-tuning. Fur-
 5 ther, text-based methods can only control *descriptive* attributes, i.e., what an image
 6 depicts, but not *evaluative* attributes, i.e., how it is perceived by an external judge.
 7 We propose LatentCompass, an exemplar-based approach that enables disentangled
 8 and controllable generation for both descriptive and evaluative concepts in
 9 a training-free manner. LatentCompass steers the generative trajectory by (a)
 10 constructing a nonlinear, low-dimensional, and orthogonal attribute space via a
 11 closed-form solution that explicitly isolates desired concepts, (b) computing an
 12 optimal shift in the constructed space, and (c) reflecting the corresponding shift
 13 in the T2I latent space. Extensive evaluations demonstrate that LatentCompass
 14 effectively (i) mitigates generative stereotypes by 100%, (ii) reduces unsafe con-
 15 cept generation by 58%, (iii) enhances aesthetic quality by 27% on average, (iv)
 16 boosts red-teaming success rates against Deepfake detectors by up to 47%, and
 17 (v) enables high-fidelity style and face attribute editing without attribute leakage.
 18 **CAUTION: This paper includes content that may be inappropriate or offensive.**

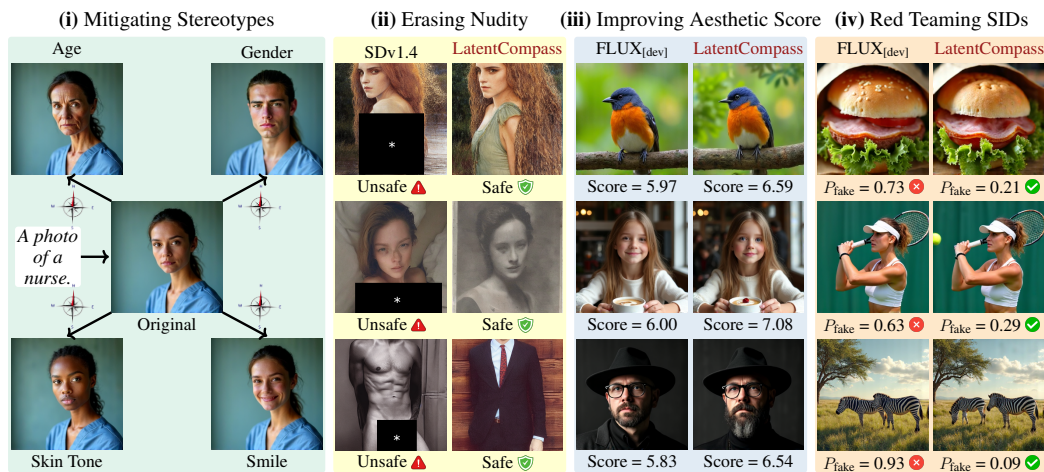


Figure 1: **Independent control of descriptive and evaluative concepts in T2I models:** (i) Facial attribute steering, which can be used to debias profession generation, (ii) improving content safety by nudity erasure, (iii) enhancing human aesthetic preferences of synthesized images, (iv) Red teaming, by finding examples that can deceive synthetic image detectors (SID).

1 Introduction

20 **The case for controllable generation.** While text-to-image (T2I) diffusion models excel at high-
 21 fidelity synthesis, their generative priors are entangled due to the correlations in their training data

22 [12, 21, 37, 59]. Consequently, T2I generation defaults to stereotypical attribute combinations,
23 neglects rare joint conditions, and fails to correctly follow user intentions, perpetuating harmful
24 societal biases [59, 15, 18, 54]. A *controllable* T2I model overcomes these limitations by overriding
25 default priors to synthesize targeted distributions, ensuring prompt faithfulness, and suppressing
26 explicit, harmful, or copyrighted concepts.

27 **A practical control mechanism** must satisfy three key requirements. **(R1) Attribute generality:**
28 it must handle both *descriptive attributes*, which are easily captured by natural language (e.g.,
29 demographic traits or artistic style), and *evaluative attributes*, which are judge-dependent properties
30 that defy text prompting entirely, such as human aesthetic preference or the “fakeness” boundary
31 learned by a deepfake detector. **(R2) Disentanglement:** modulating a target attribute must leave
32 all others unchanged. **(R3) Efficiency and extensibility:** the mechanism must operate directly on
33 a frozen, pre-trained T2I model without any weight updates, so that new control concepts can be
34 introduced incrementally and at minimal computational cost.

35 **Existing approaches for controlling T2I generation** fall into four families. **(1) Training and**
36 **fine-tuning based methods** [26, 30, 32, 35, 22] directly optimize model weights to instill new
37 concepts, but their computational cost and lack of modularity do not support R3. **(2) Prompt-based**
38 **methods** [59, 47, 45, 38, 55] steer generation through text alone, but fail on two fronts: training-data
39 correlations cause models to ignore exact instructions; and evaluative concepts such as perceived
40 aesthetic quality or realness have no natural linguistic expression, making R1 unachievable by
41 construction. **(3) Text-encoder interventions** [58, 1] operate directly on text embeddings, but remain
42 vulnerable to adversarial prompts that implicitly smuggle in unintended concepts, undermining
43 disentanglement and violating R2. **(4) Classifier guidance** [16, 2] steers the generative trajectory
44 using gradients from an external classifier, but requires a separate classifier per concept group and
45 produces attribute leakage due to non-orthogonal guidance directions, violating R2.

46 **We propose LatentCompass**, an exemplar-based method that satisfies R1–3 by decomposing
47 controllable generation into three subproblems. **(P1)** Given a set of generated images labeled
48 with attributes of interest, LatentCompass constructs, in *closed-form* (satisfies R3), a nonlinear
49 attribute space where each axis encodes an independent control attribute such as gender, style,
50 perceived aesthetic, and content safety, handling both descriptive and evaluative attributes (satisfies
51 R1). Crucially, the axes are explicitly constructed to be **mutually orthogonal**, turning a high-
52 dimensional entangled representation into a compact semantic *compass*; movement along one axis
53 isolates that attribute while minimizing interference with the rest (satisfies R2). **(P2)** Utilizing this
54 geometric framework, the intermediate representation of the image is projected onto the attribute
55 space and shifted toward a target semantic profile via likelihood maximization. **(P3)** Finally, this
56 optimized state is integrated back into the T2I sampling trajectory by solving the pre-image, ensuring
57 the final output adheres to the desired control attributes (satisfies R2).

58 **Applications of LatentCompass** includes steering facial attributes to mitigate stereotypical biases
59 concerning perceived gender and skin tone in T2I generation of medical professions (§ 4.1), improving
60 content safety by erasing toxic concepts such as nudity (§ 4.3), steering images to better match human
61 aesthetic preferences (§ 4.5), and red teaming of synthetic image detectors (SID) by shifting images
62 toward their failure modes (§ 4.4). We show examples of these applications in Fig. 1. Due to its exact
63 solution, LatentCompass can efficiently incorporate additional concepts and applications at little cost.

Summary of Contributions.

- **Exemplar-Based Steering Framework.** We propose LatentCompass, an exemplar-based approach for steering T2I diffusion models. By deriving a *nonlinear and orthogonal* attribute space, our method enables precise control over both *descriptive* and *evaluative* semantics (R1) while ensuring *disentangled* attribute shifting (R2) on both categorical and continuous labels.
- **Efficient & Extensible Approach.** Construction of this attribute space is computationally efficient and does not require iterative training, as it has a *closed-form* solution. Further, new concepts can be extended in a modular way using only a minimal labeled exemplar set (R3).
- **Broad Applications.** We demonstrate the efficacy of LatentCompass across diverse tasks, namely, achieving a **100% mitigation** in generative stereotypes, **up to 58%** improvement in unsafe concept erasure, a **27%** enhancement in human aesthetic preference, and **up to 47%** boost in red-teaming success against synthetic image detectors in a black-box manner.

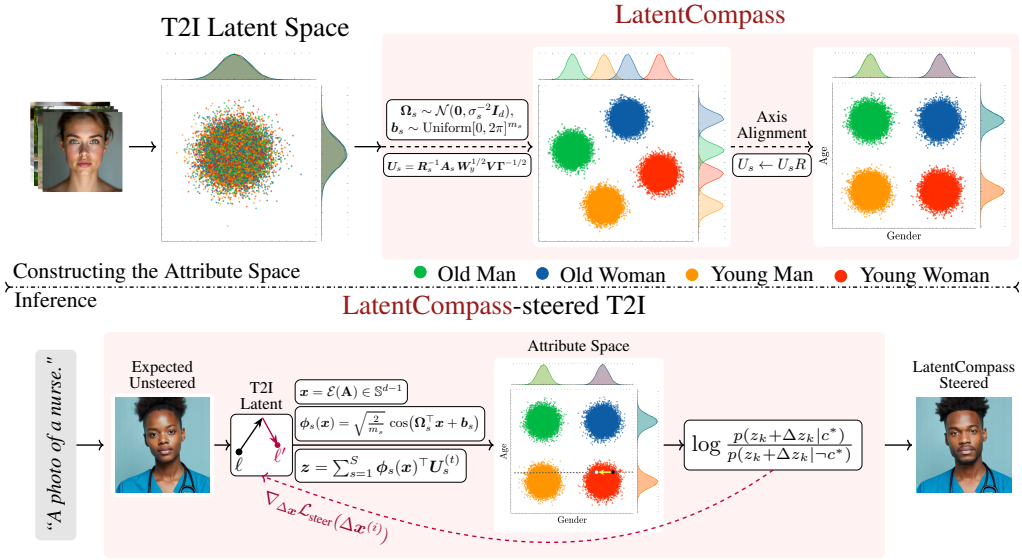


Figure 2: **Overview of LatentCompass:** (Top) LatentCompass uses labeled examples to construct an attribute space where each axis controls a target attribute. (Bottom) At inference, LatentCompass steers the T2I latent space to shift the image into the desired attribute region.

65 2 Problem Setup: Controllable Generation in Text-to-Image Diffusion Models

66 **Notation.** Scalars are denoted by lowercase letters (e.g., λ , σ), vectors by bold lowercase (e.g., x ,
67 z), and matrices by bold uppercase (e.g., X , U). Operators and sets are represented using calligraphic
68 letters (e.g., \mathcal{E} , \mathcal{D} , \mathcal{K}). For a square matrix K , we write $\text{Tr}\{K\}$ for its trace.

69 **Problem Formulation.** Let $G : \mathcal{L} \times \mathcal{C} \rightarrow \mathcal{A}$ denote a text-to-image (T2I) diffusion model that
70 produces an image $\mathbf{A} = G(\ell_T, c)$ from an initial noise latent $\ell_T \in \mathcal{L}$ and a text prompt condition
71 $c \in \mathcal{C}$. We associate each generated image with K concepts of interest, each captured by a scalar
72 attribute $y_k \in \mathbb{R}$, $k = 1, \dots, K$. These attributes span two broad categories: (i) **Descriptive**
73 **attributes**, which are visually interpretable properties that can be easily described by text prompts,
74 such as, facial attributes, nudity, or violence, and, (ii) **Evaluative attributes**, which lack a direct
75 visual correlate but are measurable by an external judge, such as the *realness* score of a synthetic
76 image detector (SID) or a human *aesthetic preference* score.

77 Concept-guided generation is then the task of *steering* a T2I model along a chosen subset of the
78 attributes. Given a target subset $\mathcal{K} \subseteq \{1, \dots, K\}$, target values $\{y_k^*\}_{k \in \mathcal{K}}$, and a measurement function
79 $Y : \mathcal{A} \rightarrow \mathbb{R}^K$, we seek a *steering mechanism* \mathcal{S} that produces a modified image $\hat{\mathbf{A}} = \mathcal{S}(G, c, \mathcal{K})$
80 whose attributes on \mathcal{K} match the targets. Formally:

Problem 1. Concept-Guided Generation

$$\mathcal{S}^* = \arg \min_{\mathcal{S}} \mathbb{E}_{k \sim \mathcal{K}} \left[\|Y_k(\mathcal{S}(G, c, \mathcal{K})) - y_k^*\|^2 \right].$$

81
82 **Prob. 1 alone is insufficient**, as it can drive the target attributes to their desired values while also
83 perturbing the remaining concepts, i.e., attribute leakage. For example, shifting *age* in a face model
84 may inadvertently change the *expression* if the two attributes are correlated in the data. An ideal
85 steering mechanism must satisfy a *preservation constraint*: for every non-target concept $j \notin \mathcal{K}$, its
86 value on the steered image must remain within an ϵ -neighborhood of its value on the unsteered image,

$$|Y_j(\mathcal{S}(G, c, \mathcal{K})) - Y_j(G(\ell_T, c))| \leq \epsilon, \quad \forall j \notin \mathcal{K}. \quad (1)$$

87 3 LatentCompass: An Approach for Concept-Guided Generation

88 **Naive Approach.** The most direct way to control an attribute in a T2I model is through the text
89 prompt itself. However, prompt-based control is fundamentally limited. (L1) Evaluative attributes
90 cannot be expressed in natural language, as these concepts are external to the generative process (e.g.,
91 a T2I model has no inherent notion of what makes an image appear *real* to a particular SID). (L2)
92 T2I models do not reliably follow prompts due to training-data correlations [6, 31]. For instance,

93 prompting FLUX_[dev] with “a photo of an old man with a smooth, clean-shaven face” occasionally
 94 yields bearded faces, as the model has learned a strong co-occurrence between *old*, *male*, and *beard*.

95 **Exemplar-Based Guidance.** An alternative is to learn concept directions from *examples*. Given a
 96 set of images labeled with attributes of interest, one can discover directions in the T2I latent space
 97 along which each attribute varies. This resolves both the above-mentioned limitations. Evaluative
 98 attributes become accessible the moment a labeling oracle (i.e., an external judge) can annotate
 99 examples with the target property (addresses L1). Furthermore, operating directly in the latent space
 100 provides control over the internal representation, enabling calibrated edits that bypass the entangled
 101 text-conditioning pathway entirely (addresses L2).

102 **LatentCompass** steers the image generation process via exemplar-based guidance by decomposing
 103 the task into three constituent subproblems: **(P1)** constructing a nonlinear, disentangled representation
 104 space from a limited set of examples annotated with either categorical or continuous labels, where
 105 individual axes encode independent attributes; **(P2)** computing an optimal shift within this manifold;
 106 and **(P3)** propagating the shift back into the latent trajectory of the underlying T2I diffusion process.

107 3.1 P1: Constructing the Attribute Space

108 Let \mathcal{E} be a frozen pretrained feature extractor that maps an image to a zero-mean, L_2 -normalized
 109 embedding $\mathbf{x} = \mathcal{E}(\mathbf{A}) \in \mathbb{S}^{d-1}$. For each timestep t of the T2I sampling trajectory, we seek an
 110 encoder $g_t : \mathbb{R}^d \rightarrow \mathbb{R}^K$ that maps \mathbf{x} to an attribute coordinate $z = g_t(\mathbf{x})$ while satisfying three
 111 requirements: **(E1) Nonlinear Expressivity:** the mapping captures the highly nonlinear statistical
 112 dependencies between embeddings and attributes, **(E2) Orthogonal Disentanglement:** the axes
 113 are mutually orthogonal under the data distribution to prevent attribute leakage (Eq. (1)), and **(E3)**
 114 **Semantic Alignment:** each axis z_k is maximally dependent on its corresponding attribute y_k .

115 To satisfy requirements **(E1–3)**, we model g_t as a kernelized supervised projection in a reproducing
 116 kernel Hilbert space (RKHS), which captures nonlinear dependencies between the embedding and the
 117 attributes **(E1)** while admitting an orthogonality constraint compatible with the data covariance **(E2)**.
 118 For N exemplars, rather than working with the $N \times N$ kernel matrix directly, we approximate the
 119 RKHS via Random Fourier Features (RFF) [44], which provide an explicit finite-dimensional map
 120 $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ with $\phi(\mathbf{u})^\top \phi(\mathbf{v}) \approx k(\mathbf{u}, \mathbf{v})$, reducing the $O(N^3)$ inversion bottleneck to $O(N)$ and
 121 keeping g_t differentiable in \mathbf{x} . A single kernel bandwidth, however, fixes the spatial frequency at
 122 which similarity is measured, which is too rigid since attributes are most discriminable at different
 123 scales. We therefore extend the kernel construction to S independent RFF banks. In the following,
 124 we describe the encoder fitting at a single time step, dropping the subscript t for clarity.

125 **Per-scale supervised covariances (E1).** For each scale $s \in \{1, \dots, S\}$, the RFF bank is defined as

$$\phi_s(\mathbf{x}) = \sqrt{\frac{2}{m_s}} \cos(\boldsymbol{\Omega}_s^\top \mathbf{x} + \mathbf{b}_s), \quad \boldsymbol{\Omega}_s \sim \mathcal{N}(\mathbf{0}, \sigma_s^{-2} \mathbf{I}_d), \quad \mathbf{b}_s \sim \text{Unif}[0, 2\pi]^{m_s}, \quad (2)$$

126 with $\boldsymbol{\Omega}_s \in \mathbb{R}^{d \times m_s}$ and per-bank RFF dimension m_s .

127 **Joint axis construction across scales (E2).** Rather than concatenating the S banks into a single
 128 high-dimensional feature space and solving one large supervised problem, we treat each bank
 129 as an *independent source of evidence* about the attributes and fuse them through a joint $K \times K$
 130 diagonalization. We first define the per-scale, variance-normalized design matrix $\boldsymbol{\Phi}_s \in \mathbb{R}^{N \times m_s}$,
 131 the per-scale cross-covariance $\mathbf{A}_s := \boldsymbol{\Phi}_s^\top \mathbf{Y} \in \mathbb{R}^{m_s \times K}$, where $\mathbf{Y} \in \mathbb{R}^{N \times K}$ is the mean-centered
 132 label matrix, and the regularized, mean-centered per-scale Gram $\mathbf{R}_s := \boldsymbol{\Phi}_s^\top \mathbf{H} \boldsymbol{\Phi}_s + \lambda \mathbf{I}_{m_s}$, where
 133 $\mathbf{H} = \mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top$ is the centering matrix. Then, we seek K axes, i.e., one per attribute, by solving

$$\begin{aligned} \max_{\{\mathbf{U}_s\}_{s=1}^S} & \text{Tr} \left\{ \mathbf{W}_y^{1/2} \left(\sum_{s=1}^S \mathbf{U}_s^\top \mathbf{A}_s \right)^\top \left(\sum_{s=1}^S \mathbf{U}_s^\top \mathbf{A}_s \right) \mathbf{W}_y^{1/2} \right\} \\ \text{s.t.} & \sum_{s=1}^S \mathbf{U}_s^\top \mathbf{R}_s \mathbf{U}_s = \mathbf{I}_K, \end{aligned} \quad (3)$$

134 where $\mathbf{W}_y = \text{diag}(w_1, \dots, w_K) \succ 0$ is an optional per-attribute weighting and $\lambda > 0$ regularizes
 135 against rank deficiency in $\boldsymbol{\Phi}_s^\top \boldsymbol{\Phi}_s$ when $N < m_s$. The trace objective maximizes the joint statistical
 136 dependence [25, 3] between the fused projection $\sum_s \mathbf{U}_s^\top \mathbf{A}_s$ and the labels **(E3)**; the summed
 137 Mahalanobis constraint enforces orthogonality of the recovered axes under the joint feature geometry
 138 of the S banks **(E2)**. See the derivation of Eq. (3) from the population HSIC in § G.1.

139 **Closed-form solution.** Eq. (3) admits a closed-form solution via a single $K \times K$ eigendecomposition.
 140 Define the joint supervised matrix

$$M = \mathbf{W}_y^{1/2} \left(\sum_{s=1}^S \mathbf{A}_s^\top \mathbf{R}_s^{-1} \mathbf{A}_s \right) \mathbf{W}_y^{1/2} \in \mathbb{R}^{K \times K}, \quad (4)$$

141 and let $M = \mathbf{V}\mathbf{\Gamma}\mathbf{V}^\top$ be its eigendecomposition with $\mathbf{\Gamma} = \text{diag}(\gamma_1, \dots, \gamma_K) \succ 0$. The per-bank
 142 projection matrices are recovered as

$$\mathbf{U}_s = \mathbf{R}_s^{-1} \mathbf{A}_s \mathbf{W}_y^{1/2} \mathbf{V} \mathbf{\Gamma}^{-1/2} \in \mathbb{R}^{m_s \times K}. \quad (5)$$

143 The derivation of this closed-form solution by Lagrangian analysis is given in § G.2.

144 **Semantic Axis Alignment (E3).** Although Eq. (3) recovers the optimal supervised subspace, its
 145 basis is rotationally invariant, lacking a 1-to-1 correspondence with target attributes. To resolve this
 146 identifiability problem and enforce disentanglement, we apply Orthogonal Procrustes alignment [48].
 147 By computing the cross-covariance $\mathbf{C} = \mathbf{Z}^\top \mathbf{Y}$ and its SVD $\mathbf{C} = \mathbf{U}_C \mathbf{\Sigma}_C \mathbf{V}_C^\top$, we derive the optimal
 148 rotation $\mathbf{R} = \mathbf{U}_C \mathbf{V}_C^\top$. Updating the projection matrices via $\mathbf{U}_s \leftarrow \mathbf{U}_s \mathbf{R}$ anchors the representation to
 149 a canonical orientation, establishing the practical coordinate system required for precise downstream
 150 steering. The optimality of this rotation under the Frobenius criterion is established in § G.3.

151 **Encoder at inference.** Given a new embedding \mathbf{x} at timestep t , the final attribute coordinate is
 152 computed by projecting the multi-scale features through the timestep-specific aligned basis:

$$\mathbf{z} = g_t(\mathbf{x}) = \sum_{s=1}^S \phi_s(\mathbf{x})^\top \mathbf{U}_s^{(t)} \in \mathbb{R}^K. \quad (6)$$

153 3.2 P2: Computing the Optimal Coordinate Shift

154 Given the disentangled attribute space, we seek an optimal intervention $\Delta \mathbf{z}$ to achieve a target
 155 attribute profile \mathbf{y}^* . By design (E2, E3), each axis independently controls its corresponding attribute,
 156 reducing the manipulation to a set of independent 1D optimizations.

157 **Gaussian Density Estimation.** For each attribute $k \in \{1, \dots, K\}$ and class $c \in \{-1, +1\}$,
 158 we model the class-conditional density $p(z_k | c) \sim \mathcal{N}(\mu_{k,c}, \sigma_{k,c}^2)$ under the assumption that the
 159 distribution of a low-dimensional projection of high-dimensional data is approximately normal [17].
 160 For continuous labels, we set a threshold at $\mathbb{E}[y_k]$ to obtain binary c . The sufficient statistics are
 161 computed from the projected examples \mathcal{Z} at timestep t : $\mu_{k,c} = \mathbb{E}_{\mathbf{z} \in \mathcal{Z}}[z_k | c]$, $\sigma_{k,c}^2 = \text{Var}_{\mathbf{z} \in \mathcal{Z}}[z_k | c]$.

162 **Log-Likelihood Ratio Objective.** To steer z_k toward a target class c^* , we maximize the log-
 163 likelihood ratio (LLR) between the target class c^* and the opposing class $-c^*$ of shifted coordinate
 164 $\tilde{z}_k = z_k + \Delta z_k$. Under the Gaussian assumption, this yields a closed-form quadratic objective:

$$\mathcal{L}_{\text{LLR}}^{(c^*)}(\tilde{z}_k) = \log \frac{p(\tilde{z}_k | c^*)}{p(\tilde{z}_k | -c^*)} = \frac{(\tilde{z}_k - \mu_{k,-c^*})^2}{2\sigma_{k,-c^*}^2} - \frac{(\tilde{z}_k - \mu_{k,c^*})^2}{2\sigma_{k,c^*}^2} + \log \frac{\sigma_{k,-c^*}}{\sigma_{k,c^*}}. \quad (7)$$

165 **Optimizing the Shift.** The optimal multi-attribute shift $\Delta \mathbf{z}^*$ maximizes this objective across the
 166 targeted attribute subset \mathcal{K} , regularized by $\beta \geq 0$ to penalize extreme deviations:

$$\Delta \mathbf{z}^* = \arg \max_{\Delta \mathbf{z}} \sum_{k \in \mathcal{K}} \mathcal{L}_{\text{LLR}}^{(c_k^*)}(z_k + \Delta z_k) - \beta \|\Delta \mathbf{z}\|_2^2. \quad (8)$$

167 3.3 P3: Propagating the Shift into the Generative Trajectory

168 While the steering objective is formulated in the K -dimensional attribute space, the actual interven-
 169 tion must be applied on the T2I diffusion latent ℓ_t . Since our encoder $g_t(\mathbf{x})$ is constructed from
 170 differentiable operations (RFFs and linear projections), we can propagate the steering signal directly
 171 back to the image latent via gradient-based optimization.

172 **Latent Optimization.** We seek a latent perturbation $\Delta \ell$ such that the edited latent $\ell' = \ell + \Delta \ell$
 173 achieves the desired attribute profile for targets \mathcal{K} , while preserving the original class assignments
 174 for all non-target attributes $j \notin \mathcal{K}$. We project the perturbed latent to the clean data manifold via
 175 the Tweedie estimate $\hat{\ell}_0(\ell_t + \Delta \ell)$, decode it to pixel space via VAE decoder \mathcal{D} , and extract its

176 zero-mean, L_2 -normalized embedding $\mathbf{x}' = \mathcal{E}(\mathcal{D}(\hat{\ell}_0)) \in \mathbb{S}^{d-1}$. On this hypersphere, we optimize
 177 $\Delta\ell$ to maximize target evidence c_k^* while penalizing non-target deviations c_j^{orig} :

$$\mathcal{L}_{\text{steer}}(\Delta\ell) = \sum_{k \in \mathcal{K}} \mathcal{L}_{\text{LLR}}^{(c_k^*)}(\tilde{z}_k(\Delta\ell)) + \frac{1}{K - |\mathcal{K}|} \sum_{j \notin \mathcal{K}} \mathcal{L}_{\text{LLR}}^{(c_j^{\text{orig}})}(\tilde{z}_j(\Delta\ell)), \quad (9)$$

178 where $\tilde{z}_m(\Delta\ell) = [g_t(\mathbf{x}')]_m$ denotes the m -th shifted coordinate.

179 **Latent Update and Early Stopping Criterion.** We initialize the perturbation as $\Delta\ell^{(0)} = \mathbf{0}$ and
 180 refine it using projected gradient ascent. To regularize the update, we maximize the joint evidence
 181 while explicitly projecting the updated state back onto the hypersphere defined by $\|\ell\|_2$:

$$\widetilde{\Delta\ell} = \Delta\ell^{(i)} + \eta \nabla_{\Delta\ell} \mathcal{L}_{\text{steer}}(\Delta\ell^{(i)}), \quad \Delta\ell^{(i+1)} = \left(\frac{\ell + \widetilde{\Delta\ell}}{\|\ell + \widetilde{\Delta\ell}\|_2} \|\ell\|_2 \right) - \ell, \quad (10)$$

182 where η is the step size. Rather than running optimization for a fixed number of steps, which
 183 frequently leads to attribute over-saturation and artifact generation, we terminate the optimization
 184 early once the evidence for all targeted attributes safely exceeds a statistical confidence threshold τ_k :

$$\min_{k \in \mathcal{K}} \mathcal{L}_{\text{LLR}}^{(c_k^*)}(\tilde{z}_k(\Delta\ell^{(i+1)})) \geq \tau_k(\mu_{k,c^*}, \sigma_{k,c^*}). \quad (11)$$

185 Upon reaching the stopping criterion, we resume the image generation using the final edited latent.

186 4 Experimental Evaluation

187 We evaluate LatentCompass across multiple domains for
 188 both *descriptive* and *evaluative* attributes. For descriptive
 189 attributes, we consider the tasks of mitigating bias and
 190 stereotype in medical profession generation (§ 4.1), disen-
 191 tangled face attribute editing (§ 4.2), and improving content
 192 safety by erasing nudity (§ 4.3). For evaluative attributes,
 193 we perform two experiments: we red team synthetic image
 194 detectors (SID) by steering T2I models to produce false
 195 negatives (§ 4.4), and we try to synthesize images that bet-
 196 ter match the aesthetic taste of observers (§ 4.5). We use
 197 FLUX_[dev][7] and SDv1.4 [47] as the T2I generator.

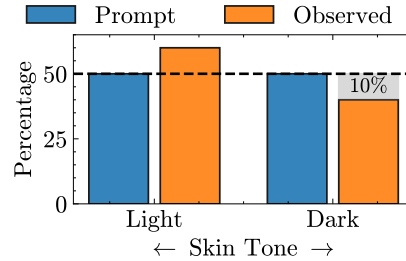


Figure 3: Prompt non-adherence.

198 4.1 Bias and Stereotype Mitigation via LatentCompass

199 **Motivation.** T2I models inherit demographic correlations from their
 200 training data: prompts for occupations consistently default to a single
 201 race-gender stereotype, even when the prompt itself is demographically
 202 neutral [4, 15, 37]. Crucially, explicit prompting often fails to override
 203 these priors, occasionally yielding images that contradict the requested
 204 demographic attributes (see Fig. 3). We evaluate LatentCompass’s ability
 205 to mitigate these demographic predispositions *without* altering the prompt
 206 that elicits the concept.

207 **Setup & Metrics.** For each concept $C \in \{\text{Doctor}, \text{Nurse}\}$
 208 we generate 300 images and measure the two binary attributes
 209 $\{\text{perceived gender}, \text{skin tone}\}$. LatentCompass is applied to the minority
 210 classes for each concept. We report the empirical conditional density $P_{D,C}(y)$ in the generated set
 211 D alongside two metrics: **1) Stereotype Score** [15], $\Psi = \max(0, P_{D,C}(y) - P_{D,C}^*(y))$, where
 212 $P_C^*(y)$ is the real-world occupational density obtained from U.S. census data. By construction, the
 213 Stereotype Score penalizes only *directional* over-representation of the stereotypical attribute, separ-
 214 ating genuine stereotype reinforcement from harmless under-representation. **2) Demographic Parity**
 215 **(DP)** as $\text{DP} = |P_{D,C}(y=0) - P_{D,C}(y=1)|$, which measures imbalance between the two classes
 216 regardless of direction. The two metrics are complementary: stereotype score targets sociologically
 217 meaningful violations of the real-world distribution, while DP captures raw distributional imbalance.

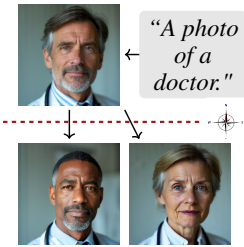


Figure 4: (Top) Unsteered (Bottom) LatentCompass.

Table 1: Evaluation of gender and skin tone bias and stereotypes in Doctor and Nurse images generated by unsteered versus LatentCompass-steered FLUX_[dev].

Method	Gender Stereotype			Gender Bias		Skin Tone Stereotype			Skin Tone Bias		
	$P_C(y = \text{male})$	$P_C^*(y = \text{male})$	$\Psi \downarrow$	$P_C(y = \text{male})$	$DP \downarrow$	$P_C(y = \text{light})$	$P_C^*(y = \text{light})$	$\Psi \downarrow$	$P_C(y = \text{light})$	$DP \downarrow$	
Doctor	Unsteered	82.6	61.0	21.6	82.6	65.2	99.0	94.0	5.0	99.0	98.0
	LatentCompass	61.0	61.0	0.0	49.1	1.8	94.0	94.0	0.0	50.3	0.6
Nurse	Unsteered	89.3	88.0	1.3	89.3	68.6	98.0	89.0	9.0	98.0	96.0
	LatentCompass	88.0	88.0	0.0	50.8	1.6	89.0	89.0	0.0	50.0	0.0

218 **Results.** Tab. 1 summarizes the impact of LatentCompass on neutralizing systemic biases regarding
 219 perceived gender and skin tone. Our framework significantly attenuates gender and skin tone
 220 stereotypes, yielding average reductions of Ψ by **100%**. Further, we reduce DP by **98.57%**
 221 on average, demonstrating robust performance in mitigating historical biases [5] inherent in T2I priors.
 222 Fig. 4 visualizes this effect: the steered samples preserve the visual cues of the prompted profession
 223 (uniform, stethoscope, clinical setting) while diversifying gender, confirming that the orthogonality
 224 of the attribute axes prevents leakage into the occupation concept itself.

225 **Remark.** LatentCompass aligns the demographics of the generated data with any user-specified reference distribution, mitigating both stereotypes and biases.

226 4.2 Attribute Disentanglement in LatentCompass Steering

227 **Motivation.** When editing one attribute, others should remain unaffected (R2). To evaluate this disentanglement capability of LatentCompass, we edit four facial attributes, $\{\text{age, skin tone, perceived gender, expression}\}$, and measure the corresponding attribute leakage.

232 **Setup & Results.** We steer a set of 1000 images for $K = 4$ attributes. We measure the success rate of the edits, while also measuring the preservation of the other $K - 1 = 3$ attributes. In Fig. 5 we show the success rate along the diagonal, while the off-diagonal entries denote the leakage rate. On average, LatentCompass achieves a steering success rate of **99.5%**, while having a preservation rate of **99.7%**.

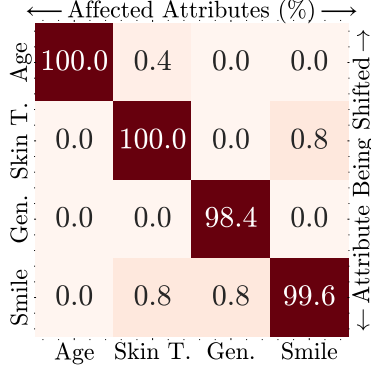


Figure 5: Edit Confusion Matrix.

239 **Remark.** LatentCompass precisely edits the target attribute without affecting non-target ones.

240 4.3 Ensuring Content Safety by Erasing Nudity via LatentCompass

241 **Motivation.** T2I models can proliferate nudity, violence, or other harmful and not-safe-for-work (NSFW) content. Existing training-free concept erasure approaches [58] suppress NSFW concepts mostly through prompt embedding manipulation, but these methods remain susceptible to adversarial and implicit prompts as they are *input-level* interventions [58]. Instead, we explore an *output-level* erasure approach where the T2I latents containing nudity are explicitly steered toward safe subspaces.

246 **Setup & Benchmarks.** We use LatentCompass to shift a continuous attribute, $Nudity \in [0, 1]$ strictly towards zero. Following the settings in [58], we evaluate LatentCompass on removing nudity in SDv1.4 on three adversarial prompt benchmarks: P4D [11], Ring-A-Bell [53], and MMA-Diffusion [57]. To assess erasure effectiveness, we report the **Attack Success Rate (ASR)** of these benchmarks against each erasure method. To measure generation quality, we employ **Fréchet Inception Distance (FID)** [27] and **CLIP score** on COCO [33].

252 **Tab. 2** compares the effectiveness of LatentCompass and baselines in erasing nudity. LatentCompass is able to achieve state-of-the-art erasure by nearly eliminating adversarial vulnerabilities across all benchmarks. By explicitly steering latents at the output level, our method yields substantial absolute reductions in Attack Success Rate (ASR) over the strongest fine-tuning-free baselines, and even surpasses fine-tuning approaches. Crucially, this robust safety mechanism avoids the severe utility trade-offs typical of concept erasure. LatentCompass attains the highest image fidelity (FID) among all evaluated methods while maintaining text-alignment (CLIP) on par with the foundational model. Ultimately, these results demonstrate that our output-level intervention provides a fundamentally stronger, more resilient defense against implicit and adversarial prompts than existing input-level manipulations, successfully decoupling NSFW concept erasure from generation degradation.

Table 2: Comparison of Attack Success Rate (ASR) (\downarrow) for nudity erasure on SD-v1.4, across adversarial prompt benchmarks.

Method	Fine-Tuning Free	P4D \downarrow	Ring-A-Bell \downarrow	MMA-Diffusion \downarrow	COCO	
					FID \downarrow	CLIP \uparrow
SD-v1.4 [46]	-	98.7	83.1	95.7	-	31.3
ESD [22]	\times	75.0	52.8	87.3	-	30.7
SA [26]	\times	62.3	32.9	20.5	54.98	30.6
CA [30]	\times	92.7	77.3	85.5	40.99	31.2
MACE [35]	\times	14.6	7.6	18.3	52.24	29.4
SDID [32]	\times	93.3	69.6	90.7	22.99	30.5
UCE [23]	\checkmark	66.7	33.1	86.7	31.25	31.3
RECE [24]	\checkmark	38.1	13.4	67.5	37.60	30.9
SLD-Medium [49]	\checkmark	93.4	64.6	94.2	<u>31.47</u>	31.0
SLD-Strong [49]	\checkmark	86.1	62.0	92.0	40.88	29.6
SLD-Max [49]	\checkmark	74.2	57.0	83.7	50.51	28.5
SAFREE [58]	\checkmark	38.4	11.4	58.5	36.35	31.1
LatentCompass (ours)	\checkmark	2.7 (-35)	1.2 (-10)	0.0 (-58)	22.08	<u>31.2</u>

Remark. LatentCompass eliminates nearly all adversarial nudity vulnerabilities by steering latents at the output level, preserving higher image fidelity than both input-level and fine-tuning baselines.

262

263

4.4 LatentCompass as a Tool for Red Teaming Synthetic Image Detectors

264

Motivation. As T2I models close the gap between synthetic and real imagery, synthetic image detectors (SIDs) have emerged as the primary defense against malicious uses such as

265

266

267

268

disinformation and identity fraud [51, 39, 29, 14]. Red teaming SIDs unveils their failure modes and is essential for hardening these detectors before deployment. As SIDs are mostly used in proprietary settings, it is important to red team them in a black-box manner. We use LatentCompass to steer an *evaluative* attribute, i.e., detector *realness*, which cannot be described through text.

269

270

271

272

273

274

275

Setup. Following [14], we use FLUX_[dev] as the T2I model and prompts from COCO, and report Attack Success Rate (ASR) against three *black-box* SIDs: UFD [39], RINE [29], and FatFormer [34]. We use LatentCompass to flip a binary label *Fakeness* $\in \{0, 1\}$ to zero.

276

277

278

279

280

281

282

283

284

285

Results. Tab. 3 reports ASR across the three SIDs. Although PolyJuice [14] boosts ASR over the unsteered baseline, it plateaus at 86% on average. LatentCompass closes this gap and tightens the variance across detectors. This improvement is a result of the non-linear geometry of our attribute space, as PolyJuice only uses linear shifts in the T2I latent space. Further, LatentCompass adds a negligible overhead ($\sim 3s$) for generating a false negative compared to a PolyJuice ($\sim 193s$), therefore showing the practicality of our approach in dataset generation. In Tab. 4, we also show that LatentCompass also outperforms transfer-based attacks, such as DiffPGD [56], that performs white-box attacks on a proxy detector trained on RINE labels.

Table 3: LatentCompass against T2I-specific SIDs at 256×256 .

Method	UFD [39]	RINE [29]	FatFormer [34]	AVG \pm STD	Time (s) \downarrow
No Steering	67.6	52.4	55.1	58.4 ± 6.6	~ 5.0
PolyJuice [14]	96.3	81.2	83.2	86.8 ± 6.6	~ 198.6
LatentCompass	100.0 (+32)	99.5 (+47)	96.4 (+41)	98.6 ± 1.6	~ 8.2

Table 4: LatentCompassASR (%) against transfer-based baseline.

Method	ASR (\uparrow)
Unsteered	52.4
DiffPGD ^t (x^n) [56]	57.2
DiffPGD ^t (x_0^n) [56]	70.1
LatentCompass (ours)	99.5 (+46)

286

287

Remark. LatentCompass on average achieves 98.6% black-box red-teaming success against SIDs.

4.5 Aligning T2I Generation with Human Aesthetic Preferences using LatentCompass

288

289

290

291

292

293

294

295

296

297

298

Motivation. A prime example of an evaluative attribute is the aesthetic quality of an image. Since human preference is intrinsically tied to visual appeal, aligning T2I models to maximize aesthetic scores enhances their practical utility and drives broader adoption across real-world applications.

Setup. We use LAION-Aesthetics_Predictor V2 [50] as a proxy for human preference. LatentCompass is fit on the examples labeled by this predictor and then used to enhance a continuous attribute *Aes* $\in [0, 10]$ for 1000 COCO prompts.

Results. From Fig. 6, we observe that an unsteered FLUX_[dev] obtains *Aes* of ~ 5.5 on average, while LatentCompass-

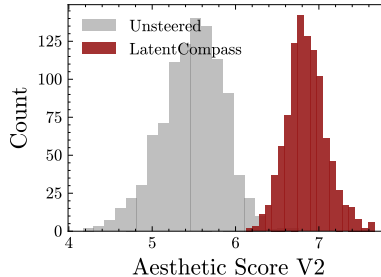


Figure 6: Improving ASV2 of COCO-prompted images via LatentCompass.

steered images have an average score of ~ 7 . Interestingly, though LatentCompass is fitted using FLUX_[dev]-generated examples with a maximum *Aes* of 7.4, we find that LatentCompass-steered samples have a maximum *Aes* > 9. This suggests that the nonlinear direction found by LatentCompass is highly generalizable and enables extrapolation.

Remark. LatentCompass significantly improves visual appeal, shifting average aesthetic scores from 5.5 to 7.0 and successfully extrapolating beyond its fitting exemplars.

5 Related Work

Conventional steering approaches typically rely on prompt engineering [1, 55, 38] or CLIP-based text guidance [59, 47, 45]. However, these approaches struggle to override entrenched training-data correlations, and are inherently incapable of steering non-linguistic, evaluative attributes. To bypass text bottlenecks, classifier-guidance [1, 16, 59] steers latent trajectories using gradients from external classifiers. While operating directly on latents, this approach struggles with attribute leakage [59].

Generative stereotypes and bias. T2I models inherently default to skewed demographic distributions, an issue often exacerbated rather than resolved by prompt engineering [49, 59]. Although open-set auditing has improved [18, 15, 4, 37, 36, 8, 9], current surface-level text mitigations fail to explicitly disentangle correlated attributes within generative priors, leaving the visual latent space structurally biased and resistant to prompt-agnostic correction. The ability of LatentCompass in disentangled attribute editing allows us to mitigate visual bias and stereotypes in T2I models.

Red-teaming of synthetic image detectors. Although diffusion-based adversarial attacks [56, 13, 10] are effective on SIDs, they demand white-box access, which limits their utility for commercial black-box detectors. Alternatively, while black-box methods like PolyJuice [14] red team SIDs via linear latent shifts, unlike LatentCompass, these linear approximations cannot capture the highly non-linear geometry of *fakeness*, necessitating expensive hyperparameter search.

Concept erasure. Training-based unlearning methods permanently alter model weights to suppress concepts, but suffer from prohibitive computational costs and catastrophic forgetting [26, 30, 32, 35, 22]. Alternatively, training-free inference interventions bypass retraining by manipulating classifier-free guidance (SLD [49]) or steering text embeddings away from toxic subspaces (SAFREE [58]). However, since these models are vision-unaware, they struggle with safeguarding against subversive prompts that induce harmful results. In contrast, LatentCompass presents a complementary approach that intervenes directly on the image latents, addressing the shortcomings of text-only intervention.

6 Concluding Remarks

In this paper, we present LatentCompass, an exemplar-based, training-free method for controllable text-to-image generation that operates directly on the latent space of frozen T2I diffusion models. To overcome entrenched generative stereotypes and predispositions, LatentCompass constructs a closed-form, nonlinear, and explicitly orthogonal attribute subspace from a small set of generated labeled examples, then steers the T2I by propagating the shift back into the diffusion trajectory. Owing to its orthogonal-by-construction design, a single attribute compass supports semantic steering, stereotype mitigation, concept erasure, and black-box red-teaming of synthetic image detectors, with new concepts added by simply refitting the closed-form encoder on a new labeled set.

Limitations & Future Directions. Similar to any training-free steering approach, LatentCompass is constrained to the T2I model’s existing generative data manifold and cannot synthesize out-of-distribution concepts. Moreover, it requires a limited search for the steering magnitude, which is standard for all the guidance methods. Finally, because the base model cannot faithfully render continuous demographic spectra, our attribute evaluations for *perceived gender* and *skin tone* (§ 4.1, § 4.2) are restricted to binary categories, presenting an avenue for future work.

Ethical Considerations. While LatentCompass is designed to mitigate harmful biases and stress-test safety-critical detectors, its red-teaming capability is dual-use and could be misused to evade synthetic image detection in the wild. We strongly oppose such use and intend LatentCompass solely for responsible auditing, debiasing, and detector hardening; we discuss potential defenses against malicious deployment in Appendix § F.

348 **References**

- 349 [1] Yuanhao Ban, Ruochen Wang, Tianyi Zhou, Minhao Cheng, Boqing Gong, and Cho-Jui Hsieh. Un-
350 derstanding the impact of negative prompts: When and how do they take effect? In *ECCV*. Springer,
351 2024.
- 352 [2] Arpit Bansal, Hong-Min Chu, Avi Schwarzschild, Soumyadip Sengupta, Micah Goldblum, Jonas Geiping,
353 and Tom Goldstein. Universal guidance for diffusion models. In *Proceedings of the IEEE/CVF Conference*
354 *on Computer Vision and Pattern Recognition*, pages 843–852, 2023.
- 355 [3] Elnaz Barshan, Ali Ghodsi, Zohreh Azimifar, and Mansoor Zolghadri Jahromi. Supervised principal
356 component analysis: Visualization, classification and regression on subspaces and submanifolds. *Pattern*
357 *Recognition*, 2011.
- 358 [4] Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori
359 Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. Easily accessible text-to-image generation
360 amplifies demographic stereotypes at large scale. In *ACM Conference on Fairness, Accountability, and*
361 *Transparency*, 2023.
- 362 [5] Abeba Birhane, Sepehr Dehdashtian, Vinay Prabhu, and Vishnu Boddeti. The dark side of dataset scaling:
363 Evaluating racial classification in multimodal models. *FAccT '24*, page 1229–1244. Association for
364 Computing Machinery, 2024.
- 365 [6] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models
366 with reinforcement learning. In *The Twelfth International Conference on Learning Representations*, 2024.
- 367 [7] BlackForestLabs. FLUX. [https://blackforestlabs.ai/
368 announcing-black-forest-labs/](https://blackforestlabs.ai/announcing-black-forest-labs/), 2024.
- 369 [8] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to
370 computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in Neural*
371 *Information Processing Systems*, 2016.
- 372 [9] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial
373 gender classification. In *ACM Conference on Fairness, Accountability, and Transparency*, 2018.
- 374 [10] Jianqi Chen, Hao Chen, Keyan Chen, Yilan Zhang, Zhengxia Zou, and Zhenwei Shi. Diffusion models for
375 imperceptible and transferable adversarial attack. *IEEE Transactions on Pattern Analysis and Machine*
376 *Intelligence*, 2024.
- 377 [11] Zhi-Yi Chin, Chieh Ming Jiang, Ching-Chun Huang, Pin-Yu Chen, and Wei-Chen Chiu. Prompt-
378 ing4Debugging: Red-teaming text-to-image diffusion models by finding problematic prompts. In *Proceed-*
379 *ings of the 41st International Conference on Machine Learning*, pages 8468–8486, 2024.
- 380 [12] Jaemin Cho, Abhay Zala, and Mohit Bansal. Dall-eval: Probing the reasoning skills and social biases of
381 text-to-image generation models. In *Proceedings of the IEEE/CVF International Conference on Computer*
382 *Vision (ICCV)*, October 2023.
- 383 [13] Xuelong Dai, Kaisheng Liang, and Bin Xiao. Advdiff: Generating unrestricted adversarial examples using
384 diffusion models. In *European Conference on Computer Vision*, pages 93–109. Springer, 2024.
- 385 [14] Sepehr Dehdashtian, Mashrur Morshed, Jacob Seidman, Gaurav Bharaj, and Vishnu Naresh Boddeti.
386 PolyJuice Makes It Real: Black-Box, Universal Red-Teaming for Synthetic Image Detectors. In *Neural*
387 *Information Processing Systems*, 2025.
- 388 [15] Sepehr Dehdashtian, Gautam Sreekumar, and Vishnu Naresh Boddeti. OASIS Uncovers: High-Quality
389 T2I Models, Same Old Stereotypes. In *International Conference on Learning Representations*, 2025.
- 390 [16] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in*
391 *neural information processing systems*, 34:8780–8794, 2021.
- 392 [17] Persi Diaconis and David Freedman. Asymptotics of graphical projection pursuit. *The annals of statistics*,
393 pages 793–815, 1984.
- 394 [18] Moreno D’Incà, Elia Peruzzo, Massimiliano Mancini, DeJia Xu, Vidit Goel, Xingqian Xu, Zhangyang
395 Wang, Humphrey Shi, and Nicu Sebe. OpenBias: Open-set Bias Detection in Text-to-Image Generative
396 Models. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024.
- 397 [19] Bradley Efron. Tweedie’s formula and selection bias. *Journal of the American Statistical Association*, 106
398 (496):1602–1614, 2011.

- 399 [20] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi,
400 Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution
401 image synthesis. In *Forty-first international conference on machine learning*, 2024.
- 402 [21] Felix Friedrich, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Patrick Schramowski, Sasha
403 Luccioni, and Kristian Kersting. Fair diffusion: Instructing text-to-image generation models on fairness.
404 *arXiv preprint arXiv:2302.10893*, 2023.
- 405 [22] Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts from
406 diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2023.
- 407 [23] Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified concept
408 editing in diffusion models. In *Proceedings of the IEEE/CVF winter conference on applications of computer
409 vision*, pages 5111–5120, 2024.
- 410 [24] Chao Gong, Kai Chen, Zhipeng Wei, Jingjing Chen, and Yu-Gang Jiang. Reliable and efficient concept
411 erasure of text-to-image diffusion models. In *European Conference on Computer Vision*, pages 73–88.
412 Springer, 2024.
- 413 [25] Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical dependence
414 with hilbert-schmidt norms. In *International conference on algorithmic learning theory*. Springer, 2005.
- 415 [26] Alvin Heng and Harold Soh. Selective amnesia: A continual learning approach to forgetting in deep
416 generative models. *Advances in Neural Information Processing Systems*, 36:17170–17194, 2023.
- 417 [27] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans
418 trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information
419 processing systems*, 30, 2017.
- 420 [28] E. Kokiopoulou, J. Chen, and Y. Saad. Trace optimization and eigenproblems in dimension reduction
421 methods. *Numerical Linear Algebra with Applications*, 2011.
- 422 [29] Christos Koutlis and Symeon Papadopoulos. Leveraging representations from intermediate encoder-blocks
423 for synthetic image detection. In *European Conference on Computer Vision*, pages 394–411. Springer,
424 2024.
- 425 [30] Nupur Kumari, Bingliang Zhang, Sheng-Yu Wang, Eli Shechtman, Richard Zhang, and Jun-Yan Zhu.
426 Ablating concepts in text-to-image diffusion models. In *Proceedings of the IEEE/CVF international
427 conference on computer vision*, pages 22691–22702, 2023.
- 428 [31] Tony Lee, Michihiro Yasunaga, Chenlin Meng, Yifan Mai, Joon Sung Park, Agrim Gupta, Yunzhi Zhang,
429 Deepak Narayanan, Hannah Teufel, Marco Bellagente, et al. Holistic evaluation of text-to-image models.
430 *Advances in Neural Information Processing Systems*, 36:69981–70011, 2023.
- 431 [32] Hang Li, Chengzhi Shen, Philip Torr, Volker Tresp, and Jindong Gu. Self-discovering interpretable
432 diffusion latent directions for responsible text-to-image generation. In *Proceedings of the IEEE/CVF
433 Conference on Computer Vision and Pattern Recognition*, pages 12006–12016, 2024.
- 434 [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár,
435 and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision—ECCV 2014:
436 13th European conference, zurich, Switzerland, September 6-12, 2014, proceedings, part v 13*, pages
437 740–755. Springer, 2014.
- 438 [34] Huan Liu, Zichang Tan, Chuangchuang Tan, Yunchao Wei, Jingdong Wang, and Yao Zhao. Forgery-
439 aware adaptive transformer for generalizable synthetic image detection. In *Proceedings of the IEEE/CVF
440 Conference on Computer Vision and Pattern Recognition*, pages 10770–10780, 2024.
- 441 [35] Shilin Lu, Zilan Wang, Leyang Li, Yanzhu Liu, and Adams Wai-Kin Kong. Mace: Mass concept erasure
442 in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
443 Recognition*, 2024.
- 444 [36] Sasha Luccioni, Christopher Akiki, Margaret Mitchell, and Yacine Jernite. Stable bias: Evaluating societal
445 representations in diffusion models. *Advances in Neural Information Processing Systems*, 2024.
- 446 [37] Ranjita Naik and Besmira Nushi. Social biases through the text-to-image generation lens. In *AAAI/ACM
447 Conference on AI, Ethics, and Society*, 2023.
- 448 [38] Viet Nguyen, Anh Nguyen, Trung Dao, Khoi Nguyen, Cuong Pham, Toan Tran, and Anh Tran. Super-
449 charged one-step text-to-image diffusion models with negative prompts. In *Proceedings of the IEEE/CVF
450 International Conference on Computer Vision*, pages 18004–18013, 2025.

- 451 [39] Utkarsh Ojha, Yuheng Li, and Yong Jae Lee. Towards universal fake image detectors that generalize
452 across generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
453 Recognition*, pages 24480–24489, 2023.
- 454 [40] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre
455 Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual
456 features without supervision. *Transactions on Machine Learning Research*, 2023.
- 457 [41] Maitreya Patel, Song Wen, Dimitris N. Metaxas, and Yezhou Yang. Flowchef: Steering of rectified flow
458 models for controlled generations. In *Proceedings of the IEEE/CVF International Conference on Computer
459 Vision (ICCV)*, pages 15308–15318, October 2025.
- 460 [42] Qwen Team. Qwen3.5: Towards native multimodal agents, February 2026. URL [https://qwen.ai/
461 blog?id=qwen3.5](https://qwen.ai/blog?id=qwen3.5).
- 462 [43] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
463 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from
464 natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmlR,
465 2021.
- 466 [44] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in neural
467 information processing systems*, 20, 2007.
- 468 [45] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional
469 image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 470 [46] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution
471 image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer
472 Vision and Pattern Recognition (CVPR)*, June 2022.
- 473 [47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution
474 image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer
475 vision and pattern recognition*, pages 10684–10695, 2022.
- 476 [48] Peter H Schönemann. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31(1):
477 1–10, 1966.
- 478 [49] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion:
479 Mitigating inappropriate degeneration in diffusion models. In *IEEE/CVF Conference on Computer Vision
480 and Pattern Recognition*, 2023.
- 481 [50] Christoph Schuhmann and LAION. Improved aesthetic predictor. [https://github.com/
482 christophschuhmann/improved-aesthetic-predictor](https://github.com/christophschuhmann/improved-aesthetic-predictor), 2023.
- 483 [51] Kaede Shiohara and Toshihiko Yamasaki. Detecting deepfakes with self-blended images. In *Proceedings
484 of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18720–18729, 2022.
- 485 [52] Not AI Tech. Nudenet. <https://github.com/notAI-tech/NudeNet>, 2023.
- 486 [53] Yu-Lin Tsai, Chia-Yi Hsu, Chulin Xie, Chih-Hsun Lin, Jia-You Chen, Bo Li, Pin-Yu Chen, Chia-Mu Yu,
487 and Chun-Ying Huang. Ring-a-bell! how reliable are concept removal methods for diffusion models?
488 *arXiv preprint arXiv:2310.10012*, 2023.
- 489 [54] Adriana Fernández de Caleyá Vázquez and Eduardo C Garrido-Merchán. A Taxonomy of the Biases of
490 the Images created by Generative Artificial Intelligence. *arXiv preprint arXiv:2407.01556*, 2024.
- 491 [55] Max Woolf. Stable diffusion 2.0 and the importance of negative prompts for good results, 2022, 2022.
- 492 [56] Haotian Xue, Alexandre Araujo, Bin Hu, and Yongxin Chen. Diffusion-based adversarial sample generation
493 for improved stealthiness and controllability. *Advances in Neural Information Processing Systems*, 36:
494 2894–2921, 2023.
- 495 [57] Yijun Yang, Ruiyuan Gao, Xiaosen Wang, Tsung-Yi Ho, Nan Xu, and Qiang Xu. Mma-diffusion:
496 Multimodal attack on diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision
497 and Pattern Recognition*, pages 7737–7746, 2024.
- 498 [58] Jaehong Yoon, Shoubin Yu, Vaidehi Patil, Huaxiu Yao, and Mohit Bansal. Safree: Training-free and
499 adaptive guard for safe text-to-image and video generation. In *International Conference on Learning
500 Representations*, 2025.

- 501 [59] Cheng Zhang, Xuanbai Chen, Siqi Chai, Chen Henry Wu, Dmitry Lagun, Thabo Beeler, and Fernando
502 De la Torre. Iti-gen: Inclusive text-to-image generation. In *Proceedings of the IEEE/CVF International*
503 *Conference on Computer Vision*, pages 3969–3980, 2023.
- 504 [60] Yinglin Zheng, Hao Yang, Ting Zhang, Jianmin Bao, Dongdong Chen, Yangyu Huang, Lu Yuan, Dong
505 Chen, Ming Zeng, and Fang Wen. General facial representation learning in a visual-linguistic manner. In
506 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.

507	Outline of Appendix	
508	A Implementation Details	14
509	A.1 Generating Examples.	14
510	A.1.1 Prompt templates.	14
511	A.1.2 Attribute Labels.	15
512	A.1.3 T2I Generation Settings.	15
513	A.1.4 Feature Extractors	15
514	A.2 Steering	15
515	A.2.1 Estimating The Denoised Latents	15
516	A.2.2 Gradient Calculation	16
517	A.3 Generation Hyperparameters	16
518	A.4 Experiments Compute Resources	16
519	A.5 Algorithmic Summary of LatentCompass	16
520	A.6 Threat Model for Red Teaming Experiment	18
521	B Extended Qualitative Examples	20
522	C Summary of Key Notations	23
523	D Ablation on the Number of Exemplars	25
524	E LatentCompass Attribute Space Visualization	25
525	F Potential Defense Mechanisms Against LatentCompass	26
526	G Proofs	26
527	G.1 Derivation of the Multi-Scale Supervised Axis Objective	26
528	G.2 Derivation of the Closed-Form Solution to Eq. (3)	28
529	G.3 Optimality of the Procrustes Axis Alignment	29
530	A Implementation Details	
531	A.1 Generating Examples.	
532	For all our experiments in § 4, we first construct an example set by generating 40K images with the	
533	base T2I model. In Fig. 11, we perform an ablation on the effect of the number of examples used to	
534	construct the LatentCompass attribute space.	
535	A.1.1 Prompt templates.	
536	(i) Face Attribute Steering: For § 4.1 and § 4.2, we generate a common set of images following the	
537	template "A realistic photo of a <age> <skin tone> <perceived gender>	
538	with a <smile> expression.", where, $\text{age} \in \{\text{young}, \text{old}\}$, $\text{skin tone} \in \{\text{light}, \text{dark}\}$,	
539	$\text{perceived gender} \in \{\text{male}, \text{female}\}$, and $\text{smile} \in \{\text{neutral}, \text{smiling}\}$. We also append some	
540	additional phrases at the end of the prompt to improve the generation quality: "RAW color	
541	photograph, photorealistic, shot on DSLR, 85mm lens, shallow depth	

542 of field, natural skin texture, pores visible, realistic lighting,
543 no stylization."

544 **(ii) Red Teaming SID and Steering Aesthetic Preference:** For both § 4.4 and § 4.5, we generate a
545 common set of examples using a subset of text prompts from the COCO dataset [33]. No additional
546 quality prompts are appended.

547 **(iii) Content Safety Improvement:** We use 1K prompts from the MMA-Diffusion [57] benchmark,
548 specifically the ‘target_prompt’ column that contains non-adversarial toxic prompts detectable by
549 safety filters. We also use a large language model (LLM) to generate a set of 1K ‘safe’ prompt
550 candidates, and then manually check/correct each prompt with human supervision.

551 A.1.2 Attribute Labels.

552 **(i) Face Attribute Steering:** The four binary attributes (age, skin tone, perceived gender, smile) are
553 independently and identically distributed (i.i.d.) according to a Bernoulli distribution with $p = 0.5$.
554 Since we construct the prompts ourselves, we have access to the attributes associated with every
555 text prompt. After image generation, we also classify each face with an external multimodal LLM,
556 Qwen3.5-9B [42]. We refer to the attributes specified before generation, and the result attributes
557 labeled by Qwen, as ‘prompt’ and ‘observed’ labels respectively in Fig. 3. For fitting the encoder, we
558 solely use the ‘prompt’ labels; Qwen3.5-9B is used for evaluation only.

559 **(ii) Red Teaming SIDs:** We label the set of generated COCO images with hard binary labels (Real,
560 Fake) with three SIDs, UFD [39], RINE [29], FatFormer [34].

561 **(iii) Steering Aesthetic Preference:** We label our set of generated COCO images with the Aesthetic
562 Predictor V2 model [50].

563 **(iv) Content Safety Improvement:** We use NudeNet [52] to label each image with an unsafe score
564 $\in [0, 1]$.

565 A.1.3 T2I Generation Settings.

566 **(i) Face Attribute Steering:** We use FLUX_[dev] [7] with 30 sampling steps, a Flow Matching Euler
567 Discrete scheduler, and a guidance of 3.5.

568 **(ii) Red Teaming SID and Steering Aesthetic Preference:** We use FLUX_[dev] with 50 sampling
569 steps, a Flow Matching Euler Discrete scheduler, and a guidance of 3.5.

570 **(iii) Content Safety Improvement:** We use SDv1.4 [47] with 50 sampling steps, a DDIM scheduler,
571 and a guidance of 7.

572 For each timestep t in the generation process, we compute and store the denoised latent $\hat{\ell}_0$ from the
573 current latent ℓ_t .

574 A.1.4 Feature Extractors

575 For face attribute steering, we use features from FaRL [60], which shares the same architecture
576 as CLIP [43]. For red teaming SIDs and steering aesthetic preference, we used CLIP, specifically
577 clip-vit-large-patch-14. For nudity erasure, instead of visio-lingual features, we used
578 purely visual Dino V2 [40] features.

579 A.2 Steering

580 A.2.1 Estimating The Denoised Latents

581 Before feature extraction with FaRL, CLIP, or Dino-V2, it is necessary to approximate the clean
582 datapoint $\hat{\ell}_0$ at every step. For diffusion models, we can do this by Tweedie’s formula [19],

$$\hat{\ell}_0 = \frac{\ell_t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(\ell_t, t, c)}{\sqrt{\bar{\alpha}_t}} \quad (12)$$

583 where $\bar{\alpha}_t$ is the diffusion noise schedule, and $\epsilon_\theta(\cdot)$ denotes the predicted noise. However, we note
 584 that FLUX_[dev] is a rectified flow model [20] (which are closely related to diffusion models), where
 585 the denoising can be simplified to a single Euler step with stepsize t ,

$$\hat{\ell}_0 = \ell_t - \mathbf{v}_\theta(\ell_t, t, c) \cdot t \quad (13)$$

586 where $\mathbf{v}_\theta(\cdot)$ is the predicted velocity.

587 A.2.2 Gradient Calculation

588 For a latent ℓ_t at timestep t , we first compute the denoised estimate $\hat{\ell}_0$ (Eq. (12), Eq. (13)), and
 589 then decode $\hat{\ell}_0$ into an image $\hat{\mathbf{A}}_0 = \mathcal{D}(\hat{\ell}_0)$. We then use the feature extractor \mathcal{E} to extract features
 590 $\mathbf{x} = \mathcal{E}(\hat{\mathbf{A}}_0)$, and then project it on to the K-axis attribute space as $\mathbf{z} = g_t(\mathbf{x})$. For brevity, we can
 591 summarize these sequence of actions as a single differentiable transformation, $\mathbf{z} := \text{ToCoord}(\ell_t)$.

592 It is apparent that computing the gradient of $\Delta\ell_t$ (Eq. (10)) w.r.t. the optimization objective (Eq. (9))
 593 would require backpropagating through the diffusion model itself, due to Eq. (12), Eq. (13). While
 594 this may be tractable for a smaller model like SDv1.4, it is infeasible for a larger model like FLUX_[dev]
 595 with 12B parameters.

596 For FLUX_[dev], we treat the velocity prediction $\mathbf{v}_\theta(\ell_t, t, c)$ as a constant term \mathbf{v}_t instead of a function
 597 of ℓ_t , which simplifies Eq. (13) to,

$$\hat{\ell}_0 = \ell_t - \mathbf{v}_t \cdot t. \quad (14)$$

598 Under Eq. (14), computing the gradient of ℓ_t no longer requires backpropagating through the actual
 599 T2I model. This approach is originally used by Bansal et al. [2], and referred to as *backward universal*
 600 *guidance*. However, backward guidance requires injecting noise back into ℓ_t to ensure they remain
 601 on the trajectory manifold. More recently, Patel et al. [41] theoretically show that for rectified flows,
 602 due to smoothness properties, $\nabla_{\hat{\ell}_0} \mathcal{L}$ has a negligible error compared to the exact gradient $\nabla_{\ell_t} \mathcal{L}$. As
 603 such, we use the approximation in Eq. (14), and **we do not compute the gradient through the T2I**
 604 **model for FLUX_[dev].**

605 A.3 Generation Hyperparameters

606 We present the hyperparameters per experiment in Tab. 5. D_1 , D_2 , and D_3 represent a set of face
 607 image examples, a set of nude/safe image examples, and a set of COCO-prompt images respectively,
 608 previously described in § A.1.

Table 5: Hyperparameter configurations across different experiments.

Hyperparameter	§ 4.1	§ 4.2	§ 4.3	§ 4.4	§ 4.5
Example Set	D_1	D_1	D_2	D_3	D_3
Sampling Steps, T	30	30	50	50	50
Shift Strength, η	{2, 3, 4}	{2, 3, 4}	{1, 2.5, 5}	{1, 2, 3}	{1, 2, 3}
Max PGD Steps, i_{\max}	30	30	20	20	20
Stopping Criterion, τ	$\mu_{k,c^*} \pm 1\sigma_{k,c^*}$	$\mu_{k,c^*} \pm 1\sigma_{k,c^*}$	$\mu_{k,c^*} - 1\sigma_{k,c^*}$	μ_{k,c^*}	$\mu_{k,c^*} + 1\sigma_{k,c^*}$

609 A.4 Experiments Compute Resources

610 For all steps of experiments, including dataset generation, attribute space construction, and inference
 611 steering, we used eight NVIDIA Quadro RTX 8000 GPUs, each with 48 GB of memory. The primary
 612 computational bottleneck arises from the memory requirements of the T2I models during image
 613 generation; LatentCompass itself adds negligible overhead.

614 A.5 Algorithmic Summary of LatentCompass

615 We summarize LatentCompass as two procedures: a one-time *fitting* stage that constructs the attribute
 616 encoder from labeled exemplars, and an *inference* stage that steers the T2I latent toward a target

Algorithm 1 LatentCompass: Fitting the Attribute Encoder

Input: Labeled exemplars $\{(x_i, y_i)\}_{i=1}^N$ with zero-mean, L_2 -normalized embeddings $x_i \in \mathbb{S}^{d-1}$ and labels $y_i \in \mathbb{R}^K$; diffusion timestep t .

Output: Per-bank projections $\{U_s^{(t)}\}_{s=1}^S$, RFF parameters $\{\Omega_s, \mathbf{b}_s\}_{s=1}^S$, alignment rotation \mathbf{R} , class statistics $\{(\mu_{k,c}, \sigma_{k,c}^2)\}_{k,c}$.

```
1:  $\mathbf{H} \leftarrow \mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top$  {Centering matrix}
2:  $\mathbf{Y} \leftarrow \mathbf{H}\mathbf{Y}$  {Mean-center the labels ( $N \times K$ )}
3: for  $s = 1, \dots, S$  do
4:   Sample  $\Omega_s \sim \mathcal{N}(\mathbf{0}, \sigma_s^{-2} \mathbf{I}_d)$  and  $\mathbf{b}_s \sim \text{Unif}[0, 2\pi]^{m_s}$ 
5:    $\Phi_s \leftarrow$  matrix with rows  $\sqrt{2/m_s} \cos(\Omega_s^\top x_i + \mathbf{b}_s)^\top$ 
6:    $\mathbf{A}_s \leftarrow \Phi_s^\top \mathbf{Y} \in \mathbb{R}^{m_s \times K}$ 
7:    $\mathbf{R}_s \leftarrow \Phi_s^\top \mathbf{H} \Phi_s + \lambda \mathbf{I}_{m_s}$  {Centered Gram matrix}
8: end for
9:  $\mathbf{M} \leftarrow \mathbf{W}_y^{1/2} (\sum_{s=1}^S \mathbf{A}_s^\top \mathbf{R}_s^{-1} \mathbf{A}_s) \mathbf{W}_y^{1/2} \in \mathbb{R}^{K \times K}$ 
10: Eigendecompose  $\mathbf{M} = \mathbf{V}\mathbf{\Gamma}\mathbf{V}^\top$  with  $\mathbf{\Gamma} = \text{diag}(\gamma_1, \dots, \gamma_K) \succ 0$ 
11: for  $s = 1, \dots, S$  do
12:    $\mathbf{U}_s \leftarrow \mathbf{R}_s^{-1} \mathbf{A}_s \mathbf{W}_y^{1/2} \mathbf{V}\mathbf{\Gamma}^{-1/2} \in \mathbb{R}^{m_s \times K}$ 
13: end for
14:  $\mathbf{Z} \leftarrow \sum_{s=1}^S \Phi_s \mathbf{U}_s \in \mathbb{R}^{N \times K}$  {Training projections}
15:  $\mathbf{C} \leftarrow \mathbf{Z}^\top \mathbf{Y}$ ; SVD:  $\mathbf{C} = \mathbf{U}_C \mathbf{\Sigma}_C \mathbf{V}_C^\top$ 
16:  $\mathbf{R} \leftarrow \mathbf{U}_C \mathbf{V}_C^\top$  {Optimal Procrustes rotation}
17: for  $s = 1, \dots, S$  do
18:    $U_s^{(t)} \leftarrow U_s \mathbf{R}$ 
19: end for
20:  $\mathbf{Z} \leftarrow \mathbf{Z} \mathbf{R}$ 
21: for  $k = 1, \dots, K$  and  $c \in \{-1, 1\}$  do
22:    $\mu_{k,c} \leftarrow \mathbb{E}_{i: y_{i,k}=c} [Z_{i,k}]$ ;  $\sigma_{k,c}^2 \leftarrow \text{Var}_{i: y_{i,k}=c} [Z_{i,k}]$ 
23: end for
24: return  $\{U_s^{(t)}, \Omega_s, \mathbf{b}_s\}_{s=1}^S, \mathbf{R}, \{(\mu_{k,c}, \sigma_{k,c}^2)\}_{k,c}$ 
```

617 attribute profile. The two stages share no per-prompt state: once Alg. 1 has been run for a given
618 concept set, Alg. 2 can be invoked on any prompt without revisiting the labeled exemplars, realizing
619 the *build once, steer anywhere* property discussed in § 3.

620 **Fitting (Alg. 1).** Given N labeled exemplars $\{(x_i, y_i)\}_{i=1}^N$, the fitting procedure samples S inde-
621 pendent RFF banks at bandwidths $\{\sigma_s\}$ (Eq. (2)), assembles the per-bank cross-covariances \mathbf{A}_s
622 and Gram matrices \mathbf{R}_s , and solves the joint multi-scale objective in Eq. (3) via the closed-form
623 expressions in Eqs. (4) and (5). The recovered axes are then aligned to the attribute basis through
624 Orthogonal Procrustes (§ 3.1), and per-class Gaussian sufficient statistics $(\mu_{k,c}, \sigma_{k,c}^2)$ are estimated
625 on the aligned coordinates for use at inference. Computationally, the dominant step is a single $K \times K$
626 eigendecomposition of \mathbf{M} in Eq. (4); everything else is linear in N .

627 **Steering (Alg. 2).** At inference, a T2I latent \mathbf{x} at timestep t is projected into the attribute space
628 via Eq. (6). The optimal per-axis shift is obtained in closed form by maximizing the regularized
629 log-likelihood ratio of the target class against its complement (P2). This shift is then propagated back
630 to the d -dimensional latent space through the differentiable encoder g_t by gradient ascent on Eq. (9),
631 with the iteration in Eq. (10) terminated as soon as the per-axis LLR confidence exceeds the threshold
632 τ . The resulting steered latent $\mathbf{x} + \Delta \mathbf{x}^*$ replaces the original embedding, and the diffusion sampler
633 resumes its trajectory unchanged.

634 **Modularity.** New control concepts are added by re-running Alg. 1 on a fresh labeled exemplar set,
635 without modifying the underlying T2I model or any previously fitted encoders. Alg. 2 is agnostic
636 to which concepts are present in \mathcal{K} , so categorical attributes (e.g., perceived gender), continuous
637 evaluative attributes (e.g., aesthetic score), and binary safety labels (e.g., nudity) are all handled by
638 the same inference procedure.

Algorithm 2 LatentCompass: Steering at Inference

Input: Frozen T2I diffusion latent ℓ_t at timestep t ; target subset $\mathcal{K} \subseteq \{1, \dots, K\}$ and target class profiles c_k^* for $k \in \mathcal{K}$; encoder components $\{U_s^{(t)}, \Omega_s, \mathbf{b}_s\}_{s=1}^S$ and class statistics $\{(\mu_{k,c}, \sigma_{k,c}^2)\}_{k,c}$ from Alg. 1.

Output: Steered diffusion latent $\ell_t + \Delta\ell^*$.

```
1:  $\Delta\ell^{(0)} \leftarrow \mathbf{0}$ 
2: for  $i = 0, 1, \dots, I_{max} - 1$  do
3:   {(i) Project noisy latent to clean hypersphere embedding}
4:    $\hat{\ell}_0 \leftarrow \text{Tweedie}(\ell_t + \Delta\ell^{(i)})$ 
5:    $\mathbf{x}' \leftarrow \text{Norm}(\mathcal{E}(\mathcal{D}(\hat{\ell}_0))) \in \mathbb{S}^{d-1}$  {Zero-mean,  $L_2$ -normalized}
6:   {(ii) Map embedding to orthogonal attribute space}
7:   for  $s = 1, \dots, S$  do
8:      $\phi_s(\mathbf{x}') \leftarrow \sqrt{2/m_s} \cos(\Omega_s^\top \mathbf{x}' + \mathbf{b}_s)$ 
9:   end for
10:   $\tilde{\mathbf{z}}^{(i)} \leftarrow \sum_{s=1}^S \phi_s(\mathbf{x}')^\top U_s^{(t)} \in \mathbb{R}^K$ 
11:  {(iii) Compute steering objective}
12:   $\mathcal{L}_{\text{steer}} \leftarrow \sum_{k \in \mathcal{K}} \mathcal{L}_{\text{LLR}}^{(c_k^*)}(\tilde{z}_k^{(i)}) + \lambda_{\text{pres}} \sum_{j \notin \mathcal{K}} \mathcal{L}_{\text{LLR}}^{(c_j^{\text{orig}})}(\tilde{z}_j^{(i)})$ 
13:  {(iv) Hyperspherical projected gradient update}
14:   $\Delta\ell \leftarrow \Delta\ell^{(i)} + \eta \nabla_{\Delta\ell} \mathcal{L}_{\text{steer}}$ 
15:   $\Delta\ell^{(i+1)} \leftarrow \left( \frac{\ell_t + \Delta\ell}{\|\ell_t + \Delta\ell\|_2} \|\ell_t\|_2 \right) - \ell_t$ 
16:  {(v) Early stopping on target confidence}
17:  if  $\min_{k \in \mathcal{K}} \mathcal{L}_{\text{LLR}}^{(c_k^*)}(\tilde{z}_k^{(i)}) \geq \tau_k$  then
18:     $\Delta\ell^* \leftarrow \Delta\ell^{(i+1)}$ ; break
19:  end if
20: end for
21: if stopping criterion not met then
22:    $\Delta\ell^* \leftarrow \Delta\ell^{(I_{max})}$ 
23: end if
24: return  $\ell_t + \Delta\ell^*$  {Resume T2I sampling}
```

639 A.6 Threat Model for Red Teaming Experiment

640 Following the framework established in PolyJuice [14], we formalize the red-teaming of Synthetic
641 Image Detectors (SIDs) through the lens of an adversarial attacker. The threat model is defined across
642 three dimensions: the attacker’s objective, knowledge, and capabilities.

643 **Attacker’s Objective.** The primary goal of the attacker is evasion. Given a target SID deployed
644 to distinguish between real photographs and synthetically generated media, the attacker aims to
645 synthesize an image $A = G(l_T, c)$ that successfully deceives the SID into classifying it as “real.”
646 Formally, the attacker seeks to shift the evaluative attribute (detector fakeness) to minimize the
647 predicted synthetic score, thereby maximizing the False Negative Rate (FNR) of the target detector
648 without compromising the perceptual quality or prompt-alignment of the generated image.

649 **Attacker’s Knowledge (Black-Box Access).** We operate under a strict black-box threat model
650 regarding the target SID. The attacker has no prior knowledge of the detector’s internal architecture,
651 network weights, defense mechanisms, or training data distributions. The attacker’s visibility is
652 entirely limited to API-level query access, meaning they can only submit generated images to the
653 SID and observe the resulting output predictions (either a continuous confidence score or a discrete
654 binary label) used to fit the LatentCompass attribute space.

655 **Attacker’s Capabilities (Unrestricted Latent Intervention).** While the attacker lacks white-box
656 access to the SID, they possess full white-box access to the underlying Text-to-Image (T2I) diffusion
657 model used for generation. Unlike conventional adversarial attacks that rely on imperceptible, L_p -
658 norm bounded noise added directly to the pixel space post-generation, our attacker performs an

659 *unrestricted* attack. The attacker is permitted to intervene directly in the intermediate latent space
660 of the T2I model during the sampling trajectory. By employing LatentCompass to shift the latents
661 toward the “realness” manifold, the attacker synthesizes images that are inherently adversarial by
662 construction, entirely bypassing standard pixel-level perturbation defenses.

663 **B Extended Qualitative Examples**

664 We provide qualitative examples of our applications in figures Fig. 7 Fig. 8, Fig. 9, and Fig. 10.



Figure 7: Stereotype and bias mitigation for nurse and doctor generation.

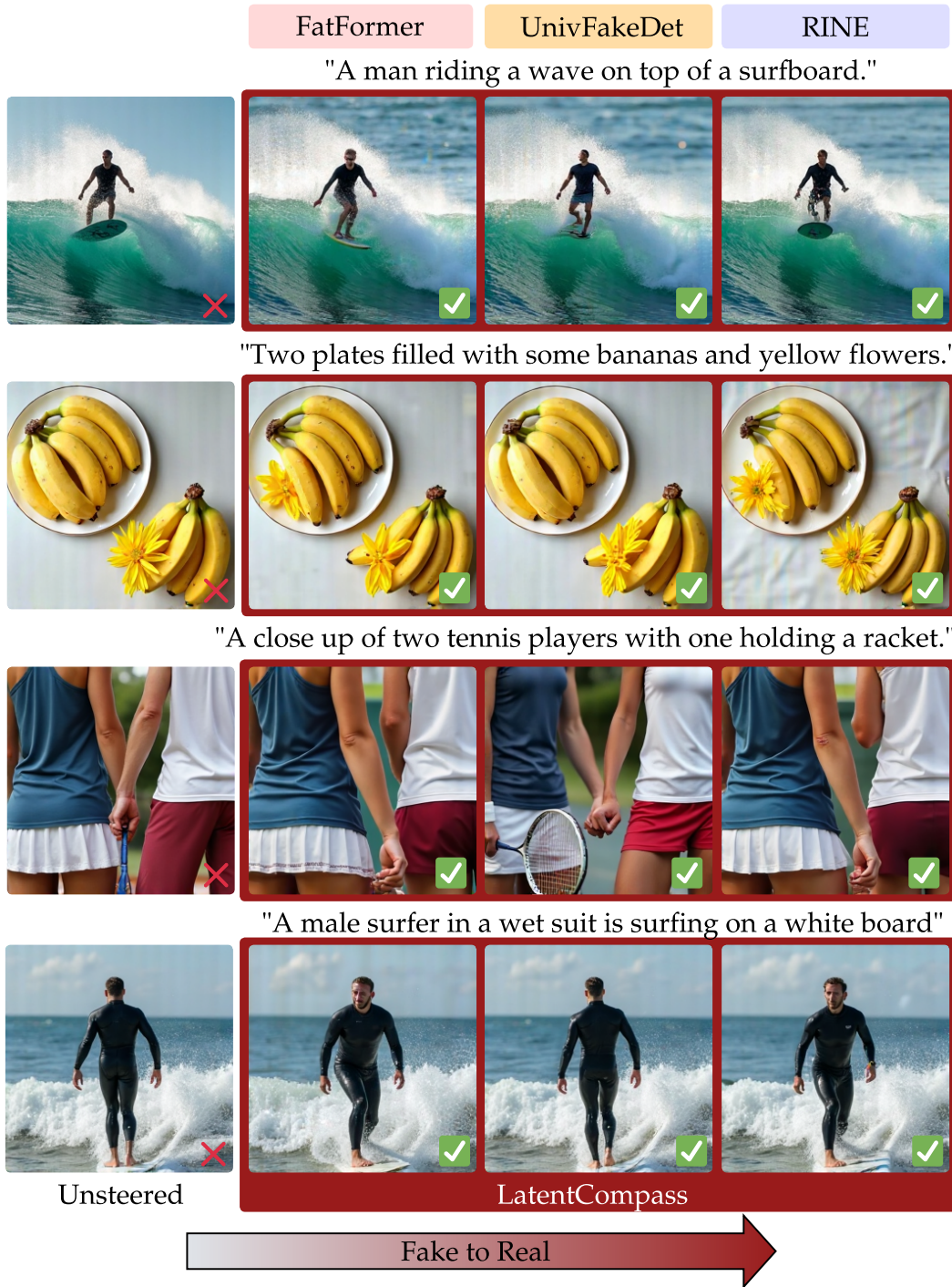


Figure 8: **(Left)** A synthetic image detected as fake by FatFormer, UFD, and RINE. **(Right)** LatentCompass-steered images that deceive each detector.

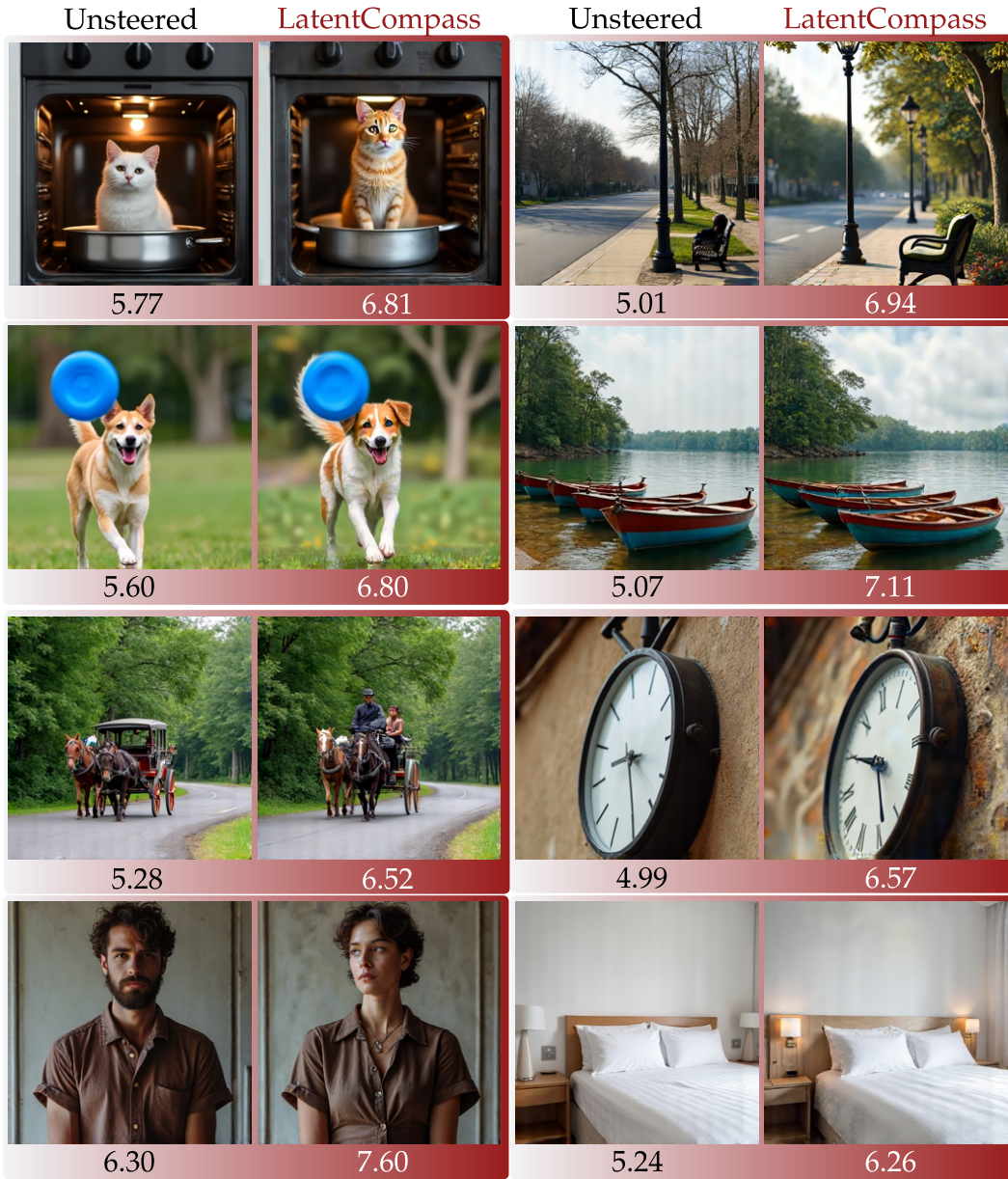


Figure 9: Improving the aesthetic score of images generated from COCO prompts.

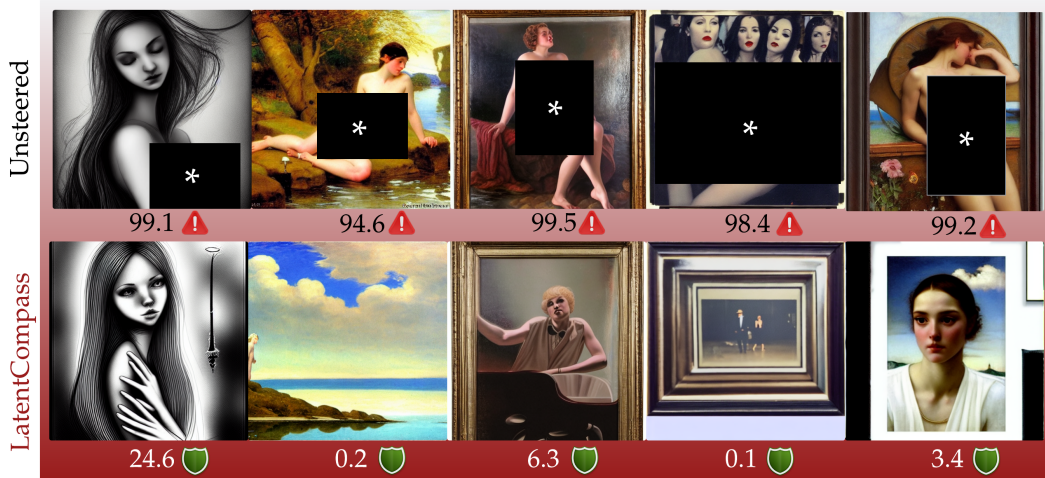


Figure 10: Erasing nudity for P4D [11] benchmark, which provides the exact random seed and guidance level associated with each toxic text prompt. The images generated from SDv1.4, and scores (0-100) given by NudeNet [52], where a score > 60 is counted as a nude image.

665 **C Summary of Key Notations**

666 To facilitate reading and provide a quick reference for the mathematical symbols, variables, and
 667 operators used throughout this work, we summarize our key notations in Tab. 6. The table is organized
 668 chronologically by the main components of our problem formulation and proposed approach.

Table 6: Summary of Key Notations used in the paper.

Symbol	Meaning	Reference
Problem Setup & T2I Model		
\mathbf{A}	Generated image	§ 3.1
K	Total number of concepts/attributes of interest	§ 3.1
y_k, \mathbf{y}^*	Scalar attribute representing concept k , and target attribute profile	§§ 3 and 3.1
\mathcal{K}	Target subset of attributes to be steered	§ 3
c^*, c_j^{orig}	Target class for steering, and original class assignment for non-targets	§ 3 and Eq. (9)
ℓ_t, ℓ	T2I diffusion latent at timestep t , and current latent state	§§ 3.1 and 3.3
Attribute Space Construction (P1)		
\mathcal{E}	Frozen pretrained feature extractor	§ 3.1
\mathbf{x}	d -dimensional zero-mean, L_2 -normalized image embedding ($\mathbf{x} = \mathcal{E}(\mathbf{A})$)	§ 3.1
g_t	Encoder mapping embeddings to attribute space at timestep t	§ 3.1
\mathbf{z}, z_k	K -dimensional attribute coordinate vector and its k -th component	§ 3.1 and Eq. (6)
S	Number of independent Random Fourier Features (RFF) scales/banks	§ 3.1
ϕ_s	RFF explicit finite-dimensional map for scale s	Eq. (2)
m_s, σ_s	Dimension and RBF kernel bandwidth for the RFF bank at scale s	Eq. (2)
Ω_s, \mathbf{b}_s	Random weights and biases for the RFF bank at scale s	Eq. (2)
\mathbf{Y}, \mathbf{H}	Mean-centered label matrix and centering matrix	§ 3.1
Φ_s, \mathbf{A}_s	Per-scale variance-normalized design matrix and cross-covariance matrix	§ 3.1
\mathbf{R}_s	Regularized, mean-centered per-scale Gram matrix	Eq. (3)
λ	Ridge regularization parameter for the Gram matrix	Eq. (3)
\mathbf{U}_s	Per-bank projection matrix	Eqs. (3) and (5)
\mathbf{W}_y	Optional per-attribute weighting matrix	Eq. (3)
$\mathbf{M}, \mathbf{V}, \mathbf{\Gamma}$	Joint supervised matrix and its eigendecomposition components	Eqs. (4) and (5)
\mathbf{C}, \mathbf{R}	Cross-covariance matrix and optimal Procrustes rotation matrix	§ 3.1
Latent Shift & Optimization (P2 & P3)		
$\mu_{k,c}, \sigma_{k,c}^2$	Mean and variance of class-conditional density for attr. k , class c	§ 3
$\mathcal{L}_{\text{LLR}}^{(c^*)}$	Log-likelihood ratio objective targeting class c^*	§ 3 and Eq. (9)
$\Delta \mathbf{z}, \tilde{z}_k$	Intervention shift in attribute space, and shifted coordinate	§ 3 and Eq. (9)
β	Regularization coefficient penalizing extreme attribute deviations	§ 3
$\Delta \ell$	Latent perturbation applied to the T2I diffusion trajectory	Eqs. (9) and (10)
$\hat{\ell}_0, \mathcal{D}$	Tweedie estimate of the clean data manifold and VAE decoder	§ 3.3
$\mathcal{L}_{\text{steer}}$	Objective function for gradient-based latent optimization	Eq. (9)
η	Projected gradient ascent step size	Eq. (10)
τ_k	Statistical confidence threshold for the early stopping criterion	§ 3.3

669 D Ablation on the Number of Exemplars

670 A fundamental advantage of LatentCompass is its modular, exemplar-based approach to constructing
 671 the attribute space. To systematically evaluate the data efficiency and robustness of our closed-form
 672 fitting procedure, we conduct an ablation study varying the number of labeled exemplars used to fit
 673 the encoder.

674 **Setup.** We fit the LatentCompass encoder using varying
 675 sizes of labeled data subsets, specifically $N \in$
 676 $\{300, 600, 768, 1200, 2400, 4800, 9600, 19200, 38400\}$.
 677 For each subset, we evaluate the quality of the constructed
 678 orthogonal attribute space by measuring the separability
 679 of the projected attribute clusters on a held-out validation
 680 set, quantified via the average Validation AUROC.

681 **Results and Analysis.** As illustrated in Figure 11, the
 682 cluster separability remains remarkably stable across the
 683 entire evaluated range. Crucially, the average AUROC
 684 curve plateaus very early; even when the encoder is fitted
 685 with a severely restricted subset of only $N = 300$ exem-
 686 plars, the separability metric shows negligible degradation
 687 compared to the attribute space constructed using the full
 688 dataset of $N = 38,400$ samples.

689 **Conclusion.** These results empirically validate that Lat-
 690 entCompass is highly sample-efficient and fundamentally
 691 insensitive to the size of the fitting dataset. Since the multi-scale supervised axis objective optimally
 692 extracts the principal directions of attribute variance without iterative training, the non-linear, or-
 693 thogonal axes can be reliably anchored with a minimal number of examples. This data efficiency
 694 practically eliminates the annotation bottleneck, allowing users to rapidly introduce and steer novel,
 695 highly specific concepts using only a few hundred exemplars.

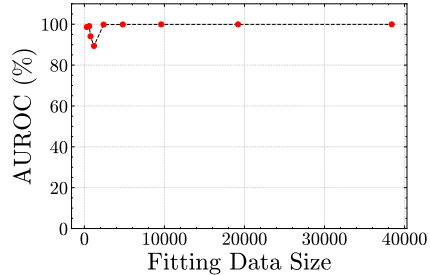


Figure 11: Average Validation AUROC of projected attribute clusters across varying numbers of fitting exemplars. LatentCompass maintains high cluster separability even in extremely low-data regimes.

696 E LatentCompass Attribute Space Visualization

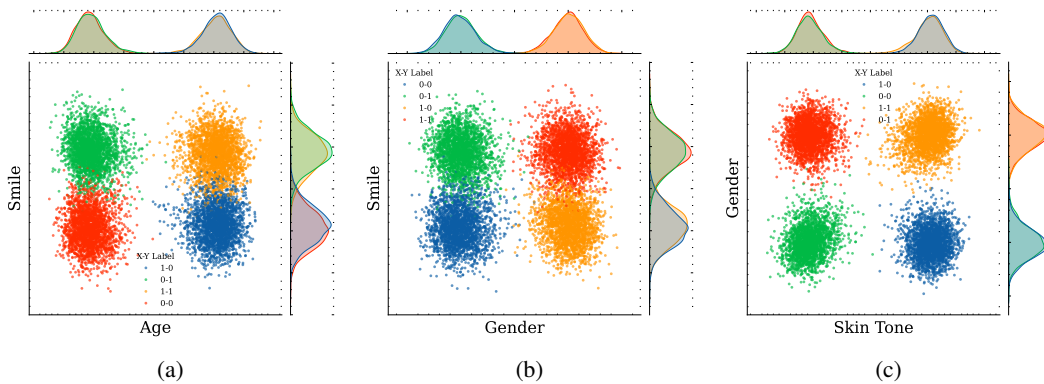


Figure 12: Visualization of attribute spaces constructed by LatentCompass.

697 To qualitatively validate the geometric properties of the representation space learned by LatentCom-
 698 pass, we visualize the projected attribute coordinates $z = g_t(x)$ for a diverse set of generated images.
 699 Fig. 12 presents the pairwise joint distributions for selected descriptive attributes: Age versus Smile
 700 (a), Gender versus Smile (b), and Skin Tone versus Gender (c).

701 The marginal distributions along each axis, visualized via Gaussian kernel density estimation, demon-
 702 strate a clear and well-calibrated separation between the binary states of each target attribute. Further-
 703 more, the two-dimensional scatter plots reveal distinct, tightly clustered quadrants corresponding to
 704 the four possible joint label combinations (e.g., 00, 01, 10, 11).

705 Crucially, these visualizations confirm the efficacy of the orthogonal disentanglement constraint
 706 formulated in Eq. (3). The axes exhibit strict geometric orthogonality; variance along one attribute
 707 axis does not systematically correlate with variance along the other. Paired with the semantic
 708 alignment achieved via Orthogonal Procrustes rotation, this geometric structure guarantees that
 709 traversing the latent space to shift a target attribute (e.g., Age) leaves non-target attributes (e.g., Smile)
 710 fundamentally undisturbed. Consequently, this empirical evidence strongly supports our core claim
 711 that LatentCompass successfully isolates entangled generative priors, facilitating the leakage-free
 712 semantic steering demonstrated throughout our main evaluations.

713 F Potential Defense Mechanisms Against LatentCompass

714 While LatentCompass is introduced as a constructive tool for auditing, debiasing, and concept erasure,
 715 its capacity to manipulate evaluative attributes—such as detector "fakeness"—presents a dual-use
 716 risk. Specifically, malicious actors could employ LatentCompass to systematically evade Synthetic
 717 Image Detectors (SIDs) by steering generated images toward the decision boundaries of real data. To
 718 mitigate these risks, we propose several complementary defense mechanisms, centering primarily on
 719 Adversarial Representation Learning (ARL) and proactive system hardening.

720 **Adversarial Training and Representation Learning (ARL).** The most direct countermeasure
 721 against latent-space evasion is proactive adversarial training. Because LatentCompass is computa-
 722 tionally efficient, training-free, and operates effectively in a black-box manner relative to the SID, it
 723 serves as an ideal generator for adversarial training examples. Defenders can deploy LatentCompass
 724 to continuously red-team their own SIDs, discovering localized failure modes and false negatives.
 725 By incorporating these high-quality, explicitly steered synthetic images back into the SID’s training
 726 curriculum, the detector is forced to learn more robust, invariant representations of synthetic artifacts.
 727 This closed-loop ARL approach ensures that the detector’s decision boundaries co-evolve with
 728 emerging latent-space attack vectors, rather than overfitting to the static, unsteered distributions of
 729 base T2I models.

730 **Robust Watermarking and Provenance.** Because LatentCompass manipulates the intermediate
 731 diffusion latents x_t to suppress the visual artifacts that SIDs typically rely upon, purely post-hoc
 732 pixel-level detection will inevitably face an arms race. A structural defense is the integration of
 733 robust, imperceptible watermarking into the base T2I model’s decoder. Assuming the watermarking
 734 mechanism operates orthogonally to the semantic attribute spaces constructed by LatentCompass, the
 735 latent shift Δx^* will not destroy the embedded signature. Combined with cryptographic provenance
 736 tracking (e.g., C2PA standards), this ensures the image remains definitively identifiable as synthetic,
 737 regardless of its steered evaluative "realness."

738 **Latent Trajectory Auditing.** In controlled deployment environments (e.g., closed APIs), providers
 739 can defend against malicious steering by auditing the generative process itself. LatentCompass
 740 achieves its target state by refining a perturbation Δx^* on a log-likelihood ratio objective. System
 741 providers can implement anomaly detection on the diffusion trajectory to monitor for statistically
 742 significant, out-of-distribution deviations from the expected prompt-conditioned score function.
 743 Flagging these anomalous latent shifts allows providers to intercept and block adversarial generations
 744 before the final image is decoded.

745 G Proofs

746 G.1 Derivation of the Multi-Scale Supervised Axis Objective

747 **Proposition.** *The multi-scale supervised axis objective in Eq. (3) is the empirical counterpart of the*
 748 *squared Hilbert–Schmidt norm of the supervised cross-covariance operator between the projected*
 749 *embedding and the labels, maximized subject to an empirical variance-preserving constraint under a*
 750 *linear label kernel and an RFF approximation of the embedding kernel.*

751 *Proof.* Let (x, y) be a random pair with $x \in \mathbb{R}^d$ and $y \in \mathbb{R}^K$, and let $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ denote
 752 the feature map of the embedding kernel taken in the RFF approximation; the operator-theoretic

753 statement extends to the underlying RKHS by completion. For a projection matrix $\mathbf{U} \in \mathbb{R}^{m \times K}$
 754 define the projected coordinates $\mathbf{z} = \mathbf{U}^\top \phi(\mathbf{x}) \in \mathbb{R}^K$. Under the linear label kernel $l(\mathbf{y}, \mathbf{y}') = \mathbf{y}^\top \mathbf{y}'$,
 755 the multivariate Hilbert–Schmidt Independence Criterion between \mathbf{z} and \mathbf{y} is the squared Hilbert–
 756 Schmidt norm of the cross-covariance operator $C_{\mathbf{z}\mathbf{y}} : \mathbb{R}^K \rightarrow \mathbb{R}^K$ [25], which acts as $\mathbf{v} \mapsto \mathbb{E}[\mathbf{z}(\mathbf{y}^\top \mathbf{v})]$
 757 and admits the matrix representation

$$C_{\mathbf{z}\mathbf{y}} = \mathbb{E}[\mathbf{z}\mathbf{y}^\top] = \mathbf{U}^\top \underbrace{\mathbb{E}[\phi(\mathbf{x})\mathbf{y}^\top]}_{=: \mathbf{M}_\mu} \in \mathbb{R}^{K \times K}, \quad (15)$$

758 with $\mathbf{M}_\mu \in \mathbb{R}^{m \times K}$. Since $C_{\mathbf{z}\mathbf{y}}$ is a finite-rank operator, its Hilbert–Schmidt norm coincides with the
 759 Frobenius norm of its matrix representation, giving

$$\text{HSIC}(\mathbf{z}, \mathbf{y}) = \|C_{\mathbf{z}\mathbf{y}}\|_{\text{HS}}^2 = \|\mathbf{U}^\top \mathbf{M}_\mu\|_F^2 = \text{Tr}\{(\mathbf{U}^\top \mathbf{M}_\mu)^\top (\mathbf{U}^\top \mathbf{M}_\mu)\}. \quad (16)$$

760 To assign different importance to different attributes we introduce per-attribute weights $w_k > 0$
 761 collected in $\mathbf{W}_y = \text{diag}(w_1, \dots, w_K) \succ 0$ and replace \mathbf{M}_μ by $\mathbf{M}_\mu \mathbf{W}_y^{1/2}$, which scales the k -th
 762 column of \mathbf{M}_μ by $w_k^{1/2}$. Distributing the transpose, the weighted dependence measure takes the form

$$\|\mathbf{U}^\top \mathbf{M}_\mu \mathbf{W}_y^{1/2}\|_F^2 = \text{Tr}\{\mathbf{W}_y^{1/2} (\mathbf{U}^\top \mathbf{M}_\mu)^\top (\mathbf{U}^\top \mathbf{M}_\mu) \mathbf{W}_y^{1/2}\}. \quad (17)$$

763 Given N i.i.d. training pairs $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ with feature matrix $\Phi \in \mathbb{R}^{N \times m}$ (rows $\phi(\mathbf{x}_i)^\top$) and label
 764 matrix $\mathbf{Y} \in \mathbb{R}^{N \times K}$, the population cross-covariance $\mathbf{M}_\mu = \mathbb{E}[\phi(\mathbf{x})\mathbf{y}^\top]$ is replaced by its empirical
 765 counterpart $\widehat{\mathbf{M}}_\mu = \frac{1}{N} \Phi^\top \mathbf{Y}$. Substituting into Eq. (17) and absorbing the constant $\frac{1}{N^2}$ into the
 766 maximization yields the empirical unconstrained supervised covariance objective

$$\text{Tr}\{\mathbf{W}_y^{1/2} (\mathbf{U}^\top \Phi^\top \mathbf{Y})^\top (\mathbf{U}^\top \Phi^\top \mathbf{Y}) \mathbf{W}_y^{1/2}\}. \quad (18)$$

767 To prevent infinite scaling of the projection matrix \mathbf{U} and to enforce that the discovered latent axes
 768 are uncorrelated under the empirical data distribution, the maximization must be constrained by
 769 the empirical feature covariance. The standard centered HSIC [25] replaces $\phi(\mathbf{x})$ and \mathbf{y} by their
 770 centered counterparts, which empirically corresponds to multiplying Φ and \mathbf{Y} by the centering matrix
 771 $\mathbf{H} = \mathbf{I}_N - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top$. Assuming the labels are mean-centered such that $\mathbf{1}_N^\top \mathbf{Y} = \mathbf{0}^\top$, it follows
 772 that $\mathbf{H}\mathbf{Y} = \mathbf{Y}$. Consequently, the cross-covariance is inherently invariant to the mean shift of the
 773 uncentered RFFs: $\Phi^\top \mathbf{H}\mathbf{Y} = \Phi^\top \mathbf{Y}$.

774 However, to enforce true decorrelation, the constraint space must utilize the fully centered feature
 775 covariance. The uncentered Gram matrix $\Phi^\top \Phi$ contains a rank-1 spike in the direction of the
 776 empirical RFF mean $\boldsymbol{\mu}_\phi$, which structurally distorts the generalized eigenvalue problem. To recover
 777 the true centered subspace, we constrain the projection using the centered, regularized Gram matrix
 778 $\mathbf{R} = \Phi^\top \mathbf{H}\Phi + \lambda \mathbf{I}_m$.

779 Crucially, this formulation allows us to bypass explicit feature centering at inference. For a new
 780 embedding \mathbf{x} , the uncentered projection evaluates to $z_k = \phi(\mathbf{x})^\top \mathbf{U}$. The theoretically centered
 781 projection is $\tilde{z}_k = (\phi(\mathbf{x}) - \boldsymbol{\mu}_\phi)^\top \mathbf{U} = z_k - c_k$, where $c_k = \boldsymbol{\mu}_\phi^\top \mathbf{U}$ is a fixed scalar offset. Because
 782 this offset applies uniformly across the entire manifold, the downstream class-conditional Gaussian
 783 means $\mu_{k,c}$ estimated in P2 absorb this shift identically. Since the Log-Likelihood Ratio evaluates the
 784 difference $(z_k - \mu_{k,c}) = (\tilde{z}_k + c_k - (\tilde{\mu}_{k,c} + c_k)) = (\tilde{z}_k - \tilde{\mu}_{k,c})$, the uncentered inference projection
 785 is mathematically equivalent to the centered projection under the LLR objective.

786 Now, extending this to the multi-scale setting, we substitute the RFF approximation from Eq. (2) for
 787 each bank s with design matrix $\Phi_s \in \mathbb{R}^{N \times m_s}$ and define $\mathbf{A}_s := \Phi_s^\top \mathbf{Y} \in \mathbb{R}^{m_s \times K}$. For the multi-
 788 scale encoder $\mathbf{z}(\mathbf{x}) = \sum_s \mathbf{U}_s^\top \phi_s(\mathbf{x})$, the empirical cross-covariance with \mathbf{Y} is $\frac{1}{N} \sum_s \mathbf{U}_s^\top \Phi_s^\top \mathbf{Y} =$
 789 $\frac{1}{N} \sum_s \mathbf{U}_s^\top \mathbf{A}_s$. Substituting into Eq. (17) and absorbing the $\frac{1}{N^2}$ constant gives the constrained
 790 optimization problem:

$$\max_{\{\mathbf{U}_s\}} \text{Tr}\left\{\mathbf{W}_y^{1/2} \left(\sum_{s=1}^S \mathbf{U}_s^\top \mathbf{A}_s\right)^\top \left(\sum_{s=1}^S \mathbf{U}_s^\top \mathbf{A}_s\right) \mathbf{W}_y^{1/2}\right\} \quad \text{s.t.} \quad \sum_{s=1}^S \mathbf{U}_s^\top \mathbf{R}_s \mathbf{U}_s = \mathbf{I}_K, \quad (19)$$

791 which yields the exact objective and constraints presented in Eq. (3). The summed Mahalanobis
 792 constraint with $\mathbf{R}_s = \Phi_s^\top \mathbf{H}\Phi_s + \lambda \mathbf{I}_{m_s}$ enforces orthogonality of the recovered axes under the
 793 empirical Mahalanobis geometry of each bank, which is the kernelized analogue of the orthonormality
 794 constraint in standard PCA [3], and the ridge term $\lambda > 0$ regularizes against rank deficiency when
 795 $N < m_s$. \square

796 **G.2 Derivation of the Closed-Form Solution to Eq. (3)**

797 **Proposition.** *The solution to the block-diagonal constrained optimization problem in Eq. (3) is given*
 798 *by the per-bank projection matrices in Eq. (5), recovered from the $K \times K$ eigendecomposition of M*
 799 *in Eq. (4).*

800 *Proof.* Define $\mathbf{A}_s := \Phi_s^\top \mathbf{Y} \in \mathbb{R}^{m_s \times K}$ and the centered, regularized Gram matrix $\mathbf{R}_s := \Phi_s^\top \mathbf{H} \Phi_s +$
 801 $\lambda \mathbf{I}_{m_s} \succ 0$ for each bank $s = 1, \dots, S$, and stack them into

$$\mathbf{U} := \begin{bmatrix} \mathbf{U}_1 \\ \vdots \\ \mathbf{U}_S \end{bmatrix} \in \mathbb{R}^{M \times K}, \quad \mathbf{A} := \begin{bmatrix} \mathbf{A}_1 \\ \vdots \\ \mathbf{A}_S \end{bmatrix} \in \mathbb{R}^{M \times K}, \quad \mathbf{R} := \text{blkdiag}(\mathbf{R}_1, \dots, \mathbf{R}_S) \in \mathbb{R}^{M \times M}, \quad (20)$$

802 where $M := \sum_s m_s$. Under this stacking, $\sum_s \mathbf{U}_s^\top \mathbf{A}_s = \mathbf{U}^\top \mathbf{A}$ and $\sum_s \mathbf{U}_s^\top \mathbf{R}_s \mathbf{U}_s = \mathbf{U}^\top \mathbf{R} \mathbf{U}$, so
 803 the optimization problem in Eq. (3) can be written as

$$\max_{\mathbf{U}} \text{Tr} \left\{ \mathbf{W}_y^{1/2} \mathbf{A}^\top \mathbf{U} \mathbf{U}^\top \mathbf{A} \mathbf{W}_y^{1/2} \right\} \quad \text{s.t.} \quad \mathbf{U}^\top \mathbf{R} \mathbf{U} = \mathbf{I}_K. \quad (21)$$

804 The Lagrangian of Eq. (21), with symmetric multiplier $\Lambda \in \mathbb{R}^{K \times K}$ for the constraint, yields the
 805 stationarity condition

$$\mathbf{A} \mathbf{W}_y \mathbf{A}^\top \mathbf{U} = \mathbf{R} \mathbf{U} \Lambda, \quad (22)$$

806 which is a generalized eigenvalue problem (GEVP). Since $\mathbf{R} \succ 0$, the Cholesky factor $\mathbf{R} = \mathbf{L} \mathbf{L}^\top$
 807 exists with $\mathbf{L} = \text{blkdiag}(\mathbf{L}_1, \dots, \mathbf{L}_S)$ and $\mathbf{R}_s = \mathbf{L}_s \mathbf{L}_s^\top$. Define the change of variable $\mathbf{V} := \mathbf{L}^\top \mathbf{U}$,
 808 so that $\mathbf{U} = \mathbf{L}^{-\top} \mathbf{V}$ and the constraint becomes $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_K$. Substituting into Eq. (22) gives

$$\mathbf{L}^{-1} \mathbf{A} \mathbf{W}_y \mathbf{A}^\top \mathbf{L}^{-\top} \mathbf{V} = \mathbf{V} \Lambda. \quad (23)$$

809 Define $\mathbf{M}' := \mathbf{L}^{-1} \mathbf{A} \mathbf{W}_y \mathbf{A}^\top \mathbf{L}^{-\top} \in \mathbb{R}^{M \times M}$. Note that \mathbf{M}' has rank at most K , since $\mathbf{A} \mathbf{W}_y \mathbf{A}^\top$
 810 has rank at most K . Its non-zero spectrum coincides with that of the $K \times K$ matrix

$$\mathbf{M} = \mathbf{W}_y^{1/2} \mathbf{A}^\top \mathbf{L}^{-\top} \mathbf{L}^{-1} \mathbf{A} \mathbf{W}_y^{1/2} = \mathbf{W}_y^{1/2} \mathbf{A}^\top \mathbf{R}^{-1} \mathbf{A} \mathbf{W}_y^{1/2}, \quad (24)$$

811 which is the matrix in Eq. (4). Indeed, since \mathbf{R} is block-diagonal, $\mathbf{R}^{-1} = \text{blkdiag}(\mathbf{R}_s^{-1})$ and
 812 $\mathbf{A}^\top \mathbf{R}^{-1} \mathbf{A} = \sum_s \mathbf{A}_s^\top \mathbf{R}_s^{-1} \mathbf{A}_s$.

813 Let $\mathbf{M} = \mathbf{V}_M \mathbf{\Gamma} \mathbf{V}_M^\top$ be the eigendecomposition of \mathbf{M} with $\mathbf{\Gamma} = \text{diag}(\gamma_1, \dots, \gamma_K)$. (Note that \mathbf{V}_M
 814 here corresponds to the eigenvector matrix denoted simply as \mathbf{V} in the main text, as the variable
 815 \mathbf{V} is used locally in this proof for the GEVP change-of-variable). We assume the attributes in the
 816 labeled exemplar set are not perfectly collinear, ensuring \mathbf{M} has full rank K and $\mathbf{\Gamma} \succ 0$. (In cases
 817 of degenerate label correlation, the inverse is replaced by the Moore-Penrose pseudo-inverse $\mathbf{\Gamma}^{\dagger/2}$).
 818 Setting

$$\mathbf{V} = \mathbf{L}^{-1} \mathbf{A} \mathbf{W}_y^{1/2} \mathbf{V}_M \mathbf{\Gamma}^{-1/2}, \quad (25)$$

819 one verifies directly that \mathbf{V} satisfies Eq. (23) with $\Lambda = \mathbf{\Gamma}$ and $\mathbf{V}^\top \mathbf{V} =$
 820 $\mathbf{\Gamma}^{-1/2} \mathbf{V}_M^\top \mathbf{W}_y^{1/2} \mathbf{A}^\top \mathbf{R}^{-1} \mathbf{A} \mathbf{W}_y^{1/2} \mathbf{V}_M \mathbf{\Gamma}^{-1/2} = \mathbf{\Gamma}^{-1/2} \mathbf{V}_M^\top \mathbf{M} \mathbf{V}_M \mathbf{\Gamma}^{-1/2} = \mathbf{\Gamma}^{-1/2} \mathbf{\Gamma} \mathbf{\Gamma}^{-1/2} = \mathbf{I}_K$,
 821 confirming feasibility. The objective at this \mathbf{V} equals $\text{Tr}\{\mathbf{V}^\top \mathbf{M}' \mathbf{V}\} = \text{Tr}\{\Lambda\} = \sum_{k=1}^K \gamma_k$, which
 822 is the maximum attained by selecting the K largest eigenvalues of \mathbf{M} [28]. Reverting the change of
 823 variable gives

$$\mathbf{U} = \mathbf{L}^{-\top} \mathbf{V} = \mathbf{L}^{-\top} \mathbf{L}^{-1} \mathbf{A} \mathbf{W}_y^{1/2} \mathbf{V}_M \mathbf{\Gamma}^{-1/2} = \mathbf{R}^{-1} \mathbf{A} \mathbf{W}_y^{1/2} \mathbf{V}_M \mathbf{\Gamma}^{-1/2}. \quad (26)$$

824 Since $\mathbf{R}^{-1} = \text{blkdiag}(\mathbf{R}_s^{-1})$ and \mathbf{A} is the row-wise stacking of $\{\mathbf{A}_s\}_{s=1}^S$, the s -th block of Eq. (26)
 825 is

$$\mathbf{U}_s = \mathbf{R}_s^{-1} \mathbf{A}_s \mathbf{W}_y^{1/2} \mathbf{V}_M \mathbf{\Gamma}^{-1/2} \in \mathbb{R}^{m_s \times K}, \quad (27)$$

826 which is the expression in Eq. (5). \square

827 **G.3 Optimality of the Procrustes Axis Alignment**

828 **Proposition.** Let $\mathbf{Z} \in \mathbb{R}^{N \times K}$ be the projected training data and $\mathbf{Y} \in \mathbb{R}^{N \times K}$ the label matrix. The
 829 rotation $\mathbf{R}^* = \mathbf{U}_C \mathbf{V}_C^\top$, where $\mathbf{U}_C \mathbf{\Sigma}_C \mathbf{V}_C^\top$ is the SVD of the cross-covariance $\mathbf{C} := \mathbf{Z}^\top \mathbf{Y} \in \mathbb{R}^{K \times K}$,
 830 is the unique orthogonal solution to

$$\min_{\mathbf{R}: \mathbf{R}^\top \mathbf{R} = \mathbf{I}_K} \|\mathbf{Y} - \mathbf{Z}\mathbf{R}\|_F^2. \quad (28)$$

831 *Proof.* Expanding the Frobenius norm in Eq. (28) gives:

$$\begin{aligned} \|\mathbf{Y} - \mathbf{Z}\mathbf{R}\|_F^2 &= \text{Tr}\{(\mathbf{Y} - \mathbf{Z}\mathbf{R})^\top (\mathbf{Y} - \mathbf{Z}\mathbf{R})\} \\ &= \text{Tr}\{\mathbf{Y}^\top \mathbf{Y}\} - 2 \text{Tr}\{\mathbf{Y}^\top \mathbf{Z}\mathbf{R}\} + \text{Tr}\{\mathbf{R}^\top \mathbf{Z}^\top \mathbf{Z}\mathbf{R}\}. \end{aligned} \quad (29)$$

832 Note that because \mathbf{Y} is strictly mean-centered by construction, $\mathbf{Y} = \mathbf{H}\mathbf{Y}$, which implies the cross-
 833 covariance $\mathbf{C} = \mathbf{Z}^\top \mathbf{Y} = \mathbf{Z}^\top \mathbf{H}\mathbf{Y} = (\mathbf{H}\mathbf{Z})^\top \mathbf{Y}$. Thus, the cross-covariance implicitly evaluates the
 834 theoretically correct centered coordinates without requiring explicit centering of \mathbf{Z} .

835 Using the cyclic property of the trace, the middle term becomes $\text{Tr}\{\mathbf{Y}^\top \mathbf{Z}\mathbf{R}\} = \text{Tr}\{\mathbf{R}\mathbf{Y}^\top \mathbf{Z}\} =$
 836 $\text{Tr}\{\mathbf{R}\mathbf{C}^\top\} = \text{Tr}\{\mathbf{C}^\top \mathbf{R}\}$. Since \mathbf{R} is a square orthogonal matrix, $\mathbf{R}\mathbf{R}^\top = \mathbf{I}_K$. The final term
 837 simplifies to $\text{Tr}\{\mathbf{R}^\top \mathbf{Z}^\top \mathbf{Z}\mathbf{R}\} = \text{Tr}\{\mathbf{Z}\mathbf{R}\mathbf{R}^\top \mathbf{Z}^\top\} = \text{Tr}\{\mathbf{Z}\mathbf{Z}^\top\} = \text{Tr}\{\mathbf{Z}^\top \mathbf{Z}\}$, which is independent
 838 of \mathbf{R} . It follows that minimizing the Frobenius distance is exactly equivalent to maximizing the
 839 cross-trace:

$$\max_{\mathbf{R}: \mathbf{R}^\top \mathbf{R} = \mathbf{I}_K} \text{Tr}\{\mathbf{C}^\top \mathbf{R}\}. \quad (30)$$

840 Let $\mathbf{C} = \mathbf{U}_C \mathbf{\Sigma}_C \mathbf{V}_C^\top$ be the Singular Value Decomposition (SVD) of \mathbf{C} , with $\mathbf{\Sigma}_C =$
 841 $\text{diag}(\sigma_1, \dots, \sigma_K)$ and $\sigma_k \geq 0$. Substituting the SVD into the objective and applying the cyclic
 842 property of the trace yields:

$$\text{Tr}\{\mathbf{C}^\top \mathbf{R}\} = \text{Tr}\{\mathbf{V}_C \mathbf{\Sigma}_C \mathbf{U}_C^\top \mathbf{R}\} = \text{Tr}\{\mathbf{U}_C^\top \mathbf{R} \mathbf{V}_C \mathbf{\Sigma}_C\}. \quad (31)$$

843 For any orthogonal \mathbf{R} , define $\mathbf{T} := \mathbf{U}_C^\top \mathbf{R} \mathbf{V}_C$. Since it is the product of orthogonal matrices, \mathbf{T} is
 844 also exactly orthogonal. It follows that:

$$\text{Tr}\{\mathbf{C}^\top \mathbf{R}\} = \text{Tr}\{\mathbf{T} \mathbf{\Sigma}_C\} = \sum_{k=1}^K t_{kk} \sigma_k, \quad (32)$$

845 where t_{kk} denotes the k -th diagonal entry of \mathbf{T} . Since \mathbf{T} is orthogonal, its column vectors have unit
 846 norm, which strictly enforces $|t_{kk}| \leq 1$ for all k [48]. Thus:

$$\sum_{k=1}^K t_{kk} \sigma_k \leq \sum_{k=1}^K \sigma_k = \text{Tr}\{\mathbf{\Sigma}_C\}. \quad (33)$$

847 The theoretical upper bound is attained if and only if $t_{kk} = 1$ for all k , which implies $\mathbf{T} = \mathbf{I}_K$. From
 848 our definition of \mathbf{T} , this requires:

$$\mathbf{U}_C^\top \mathbf{R} \mathbf{V}_C = \mathbf{I}_K \implies \mathbf{R}^* = \mathbf{U}_C \mathbf{V}_C^\top. \quad (34)$$

849 Uniqueness follows from the strict inequality whenever $\sigma_k > 0$ for all k , which holds when \mathbf{C} has
 850 full rank. In practice, if \mathbf{C} is rank-deficient, the solution is unique up to a rotation in the null space of
 851 \mathbf{C} ; since the null directions carry no label-alignment signal, any such arbitrary rotation leaves the
 852 downstream steering unaffected.

853 Updating the per-bank projection matrices via $\mathbf{U}_s \leftarrow \mathbf{U}_s \mathbf{R}^*$ rigidly rotates the semantic axes into
 854 the canonical orientation that minimizes the Frobenius distance to the label matrix \mathbf{Y} . (Note that
 855 while Procrustes alignment preserves scale and does not inherently match the variance of \mathbf{Y} , the
 856 downstream Log-Likelihood Ratio steering objective utilized in LatentCompass is fundamentally
 857 scale-invariant, making an orthogonal rotation strictly sufficient for semantic disentanglement). \square

858 **NeurIPS Paper Checklist**

859 **1. Claims**

860 Question: Do the main claims made in the abstract and introduction accurately reflect the
861 paper’s contributions and scope?

862 Answer: [Yes]

863 Justification: We support our claims with our extensive experimental results in § 4

864 Guidelines:

- 865 • The answer [N/A] means that the abstract and introduction do not include the claims
866 made in the paper.
- 867 • The abstract and/or introduction should clearly state the claims made, including the
868 contributions made in the paper and important assumptions and limitations. A [No] or
869 [N/A] answer to this question will not be perceived well by the reviewers.
- 870 • The claims made should match theoretical and experimental results, and reflect how
871 much the results can be expected to generalize to other settings.
- 872 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
873 are not attained by the paper.

874 **2. Limitations**

875 Question: Does the paper discuss the limitations of the work performed by the authors?

876 Answer: [Yes]

877 Justification: We discuss the limitations of our method in § 6

878 Guidelines:

- 879 • The answer [N/A] means that the paper has no limitation while the answer [No] means
880 that the paper has limitations, but those are not discussed in the paper.
- 881 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 882 • The paper should point out any strong assumptions and how robust the results are to
883 violations of these assumptions (e.g., independence assumptions, noiseless settings,
884 model well-specification, asymptotic approximations only holding locally). The authors
885 should reflect on how these assumptions might be violated in practice and what the
886 implications would be.
- 887 • The authors should reflect on the scope of the claims made, e.g., if the approach was
888 only tested on a few datasets or with a few runs. In general, empirical results often
889 depend on implicit assumptions, which should be articulated.
- 890 • The authors should reflect on the factors that influence the performance of the approach.
891 For example, a facial recognition algorithm may perform poorly when image resolution
892 is low or images are taken in low lighting. Or a speech-to-text system might not be
893 used reliably to provide closed captions for online lectures because it fails to handle
894 technical jargon.
- 895 • The authors should discuss the computational efficiency of the proposed algorithms
896 and how they scale with dataset size.
- 897 • If applicable, the authors should discuss possible limitations of their approach to
898 address problems of privacy and fairness.
- 899 • While the authors might fear that complete honesty about limitations might be used by
900 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
901 limitations that aren’t acknowledged in the paper. The authors should use their best
902 judgment and recognize that individual actions in favor of transparency play an impor-
903 tant role in developing norms that preserve the integrity of the community. Reviewers
904 will be specifically instructed to not penalize honesty concerning limitations.

905 **3. Theory assumptions and proofs**

906 Question: For each theoretical result, does the paper provide the full set of assumptions and
907 a complete (and correct) proof?

908 Answer: [Yes]

909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963

Justification: We provide proofs for our used theories in § G

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Setup of each experiment in § 4 is mentioned in the corresponding subsections. Further, implementation details are provided in § A.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016

Answer: [Yes]

Justification: Although the question is ambiguous to us, we interpret the question as “whether the paper will fully release all code and data in the future”, as opposed to “whether the paper is submitting all code and data as part of supplemental material, at the time of submission”. All the text-to-image generative models we use in this paper (FLUX_[dev] and SDv1.4) are publicly available, and so are the text captions we use in image generation. We also include pseudo-code and detailed instructions to reproduce the main results in the appendix.

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer) necessary to understand the results?

Answer: [Yes]

Justification: Yes, please refer to "Setup" paragraph of the corresponding experiments and § A.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We provide the average and standard deviation of our results in Tab. 3.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- 1017 • The method for calculating the error bars should be explained (closed form formula,
1018 call to a library function, bootstrap, etc.)
- 1019 • The assumptions made should be given (e.g., Normally distributed errors).
- 1020 • It should be clear whether the error bar is the standard deviation or the standard error
1021 of the mean.
- 1022 • It is OK to report 1-sigma error bars, but one should state it. The authors should
1023 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
1024 of Normality of errors is not verified.
- 1025 • For asymmetric distributions, the authors should be careful not to show in tables or
1026 figures symmetric error bars that would yield results that are out of range (e.g., negative
1027 error rates).
- 1028 • If error bars are reported in tables or plots, the authors should explain in the text how
1029 they were calculated and reference the corresponding figures or tables in the text.

1030 8. Experiments compute resources

1031 Question: For each experiment, does the paper provide sufficient information on the com-
1032 puter resources (type of compute workers, memory, time of execution) needed to reproduce
1033 the experiments?

1034 Answer: [Yes]

1035 Justification: We provide compute information in § A.4.

1036 Guidelines:

- 1037 • The answer [N/A] means that the paper does not include experiments.
- 1038 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
1039 or cloud provider, including relevant memory and storage.
- 1040 • The paper should provide the amount of compute required for each of the individual
1041 experimental runs as well as estimate the total compute.
- 1042 • The paper should disclose whether the full research project required more compute
1043 than the experiments reported in the paper (e.g., preliminary or failed experiments that
1044 didn't make it into the paper).

1045 9. Code of ethics

1046 Question: Does the research conducted in the paper conform, in every respect, with the
1047 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1048 Answer: [Yes]

1049 Justification: To our knowledge, we abide by all the guidelines presented in the NeurIPS
1050 Code of Ethics.

1051 Guidelines:

- 1052 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
1053 Ethics.
- 1054 • If the authors answer [No], they should explain the special circumstances that require a
1055 deviation from the Code of Ethics.
- 1056 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
1057 eration due to laws or regulations in their jurisdiction).

1058 10. Broader impacts

1059 Question: Does the paper discuss both potential positive societal impacts and negative
1060 societal impacts of the work performed?

1061 Answer: [Yes]

1062 Justification: We discuss positive societal impacts throughout the paper, and also discuss
1063 potential negative societal impacts in the concluding remarks. Moreover, we provided some
1064 potential defenses against the malicious use of the LatentCompass in § F.

1065 Guidelines:

- 1066 • The answer [N/A] means that there is no societal impact of the work performed.

- 1067
- 1068
- 1069
- 1070
- 1071
- 1072
- 1073
- 1074
- 1075
- 1076
- 1077
- 1078
- 1079
- 1080
- 1081
- 1082
- 1083
- 1084
- 1085
- 1086
- 1087
- If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

1088 **11. Safeguards**

1089 Question: Does the paper describe safeguards that have been put in place for responsible
1090 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
1091 image generators, or scraped datasets)?

1092 Answer: [Yes]

1093 Justification: We do not scrape any data from the internet, but use a well-established and
1094 audited dataset (COCO). All materials related to LatentCompass (model, data, code) will be
1095 extensively audited for safety prior to release. In particular, the LatentCompass-encoders for
1096 nudity erasure and SID red teaming will be gated models, and shared only upon request and
1097 verification.

1098 Guidelines:

- 1099
- 1100
- 1101
- 1102
- 1103
- 1104
- 1105
- 1106
- 1107
- 1108
- The answer [N/A] means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

1109 **12. Licenses for existing assets**

1110 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1111 the paper, properly credited and are the license and terms of use explicitly mentioned and
1112 properly respected?

1113 Answer: [Yes]

1114 Justification: We attribute original owners and respect their license, wherever applicable.

1115 Guidelines:

- 1116
- 1117
- 1118
- 1119
- The answer [N/A] means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.

- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

1131 13. **New assets**

1132 Question: Are new assets introduced in the paper well documented and is the documentation
1133 provided alongside the assets?

1134 Answer: [N/A]

1135 Justification: At the moment of submission we are not releasing any new asset.

1136 Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

1145 14. **Crowdsourcing and research with human subjects**

1146 Question: For crowdsourcing experiments and research with human subjects, does the paper
1147 include the full text of instructions given to participants and screenshots, if applicable, as
1148 well as details about compensation (if any)?

1149 Answer: [N/A]

1150 Justification: No crowdsourcing and research with human subjects.

1151 Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

1160 15. **Institutional review board (IRB) approvals or equivalent for research with human 1161 subjects**

1162 Question: Does the paper describe potential risks incurred by study participants, whether
1163 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1164 approvals (or an equivalent approval/review based on the requirements of your country or
1165 institution) were obtained?

1166 Answer: [N/A]

1167 Justification: No research with human subjects.

1168 Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.

- 1171 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1172 may be required for any human subjects research. If you obtained IRB approval, you
1173 should clearly state this in the paper.
- 1174 • We recognize that the procedures for this may vary significantly between institutions
1175 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1176 guidelines for their institution.
- 1177 • For initial submissions, do not include any information that would break anonymity (if
1178 applicable), such as the institution conducting the review.

1179 **16. Declaration of LLM usage**

1180 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1181 non-standard component of the core methods in this research? Note that if the LLM is used
1182 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
1183 scientific rigor, or originality of the research, declaration is not required.

1184 Answer: [N/A]

1185 Justification: the core method development in this research does not involve LLMs as any
1186 important, original, or non-standard components.

1187 Guidelines:

- 1188 • The answer [N/A] means that the core method development in this research does not
1189 involve LLMs as any important, original, or non-standard components.
- 1190 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
1191 be described.

1192