

# Selective Test-Time Debiasing for CLIP via Reward Gating

Anonymous ACL submission

## Abstract

Vision language models (VLMs) demonstrate strong zero-shot performance, but often perpetuate social stereotypes in person-centric queries, yielding skewed demographic distributions. Current debiasing methods apply uniform bias corrections across all input queries regardless of their bias sensitivity, creating a fundamental fairness–utility trade-off. Strong debiasing distorts semantically meaningful information in bias-insensitive queries, while weak debiasing fails to mitigate stereotypes in bias-sensitive ones. This one-size-fits-all approach hampers simultaneously achieving high utility on bias-insensitive queries and fairness on bias-sensitive queries. We introduce **Reward-Gated Test-Time Adaptation (RG-TTA)**, a reinforcement learning-based test-time adaptation framework that selectively applies debiasing based on input sensitivity. RG-TTA adaptively triggers fairness regularization based on the bias sensitivity of each input during test-time policy adaptation, while focusing exclusively on optimizing cross-modal alignment for bias-insensitive inputs. Experiments on fairness benchmarks (e.g., FairFace, UTKFace) demonstrate substantial bias reduction while simultaneously improving zero-shot utility, resolving the trade-off of uniform debiasing. Code will be made publicly available.

## 1 Introduction

Vision Language Models (VLMs) have demonstrated exceptional zero-shot capabilities across a wide range of multimodal tasks (Deng et al., 2009; Plummer et al., 2015), reaching the stage of real-world applications. By learning joint representations from web-scale image-text pairs, these models achieve strong cross-modal alignment without task-specific fine-tuning. However, this same training paradigm causes VLMs to internalize social stereotypes (Birhane et al., 2021) present in

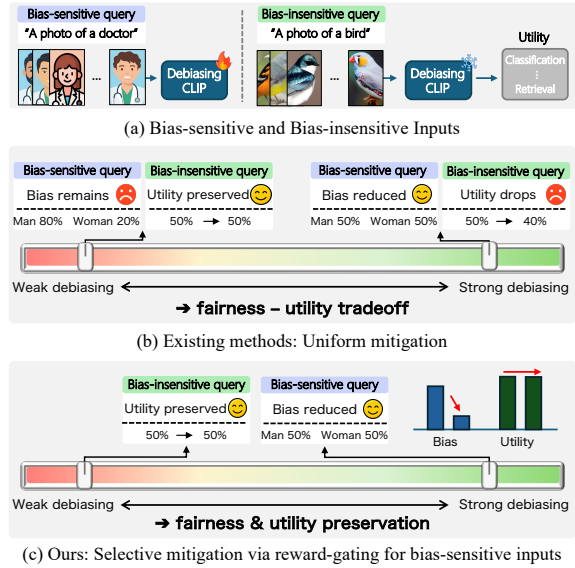


Figure 1: (a) We categorize inputs into **bias-sensitive** and **bias-insensitive**, where only the former requires debiasing intervention. (b) Existing methods apply **uniform mitigation**, creating a structural trade-off: weak debiasing retains bias in sensitive queries (left), while strong debiasing distorts insensitive queries, degrading utility (right). (c) Our approach employs **selective mitigation via reward-gating**, which applies strong debiasing only to bias-sensitive inputs while preserving insensitive ones, ensuring both fairness and utility.

their training data, leading to biased outputs that reflect and potentially amplify social prejudices (Hall et al., 2023; Hamidieh et al., 2024; Janghorbani and De Melo, 2023; Zhao et al., 2021). These biases manifest most critically in person-centric queries, where models produce skewed demographic distributions. For instance, querying “a photo of a doctor” yields disproportionately male images, or certain occupations become strongly associated with specific racial groups. Such behavior poses serious risk of reinforcing discriminatory decision-making.

Existing debiasing approaches for VLMs (Wang et al., 2021b; Chuang et al., 2023; Zhang et al., 2025) share common design philosophy: they ap-

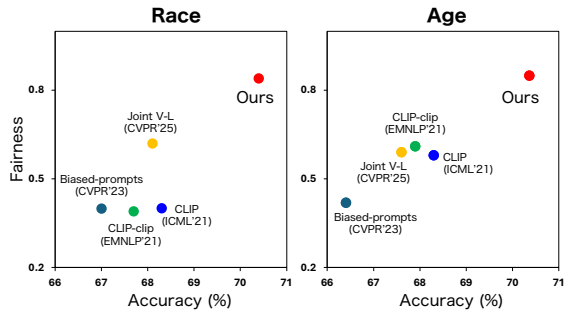


Figure 2: Fairness versus utility for Race and Age. **Accuracy** is measured as ImageNet zero-shot top-1 accuracy (%), and **Fairness** is measured as  $1 - \text{MaxSkew}@1000$ <sup>1</sup> (higher is better). Existing query-independent debiasing baselines (Chuang et al., 2023; Wang et al., 2021b; Zhang et al., 2025) exhibit a fairness–utility trade-off, whereas our method improves fairness while achieving higher accuracy.

ply fixed bias correction uniformly across all language queries, regardless of whether individual queries are sensitive to demographic biases. Although conceptually simple, this uniform mitigation strategy causes fundamental fairness–utility trade-off that we illustrated in Figure 1(b). Here, we use utility to refer to general-purpose zero-shot performance on downstream tasks (e.g., image classification and cross-modal retrieval).

For bias-sensitive queries, model’s predictions are often entangled with demographic attributes, and strong debiasing is necessary for fair outcomes. In contrast for bias-insensitive queries, ground-truth semantics are largely orthogonal to demographic attributes, and the model’s original predictions already reflect accurate cross-modal alignment. When a uniform debiasing framework is applied to both types of queries, one of two failure modes inevitably occurs. As we demonstrate in Figure 2, this inflexibility creates a trade-off where existing methods must compromise between fairness in sensitive queries and utility in general tasks, unable to excel at both simultaneously.

We believe that the aforementioned limitation is a consequence of the query-independent design paradigm. Since uniform debiasing methods cannot distinguish between inputs that require mitigation and those that do not, they are constrained to operate in a compromise regime. This observation motivates a paradigm shift toward adaptive debiasing that selectively activates mitigation based on

<sup>1</sup>MaxSkew@1000 is computed from the protected-attribute distribution within the top-1000 retrieved samples for neutral queries; see Sec. 3 for details.

the bias sensitivity of each input. To this end, we propose **Reward-Gated Test-Time Adaptation for CLIP (RG-TTA)**, a reinforcement learning (RL)-based framework designed for selective debiasing. Our key insight is that debiasing should be treated as a per-query decision rather than global transformation, enabling the model to adapt its behavior dynamically based on input characteristics. As illustrated in Figure 3, RG-TTA operates through an episodic test-time adaptation protocol. For each incoming query, we first assess its bias sensitivity by quantifying the alignment discrepancy between the query semantics and a set of demographic attributes. Based on this assessment, an adaptive reward-gating strategy dynamically triggers the fairness-regularized objective only for bias-sensitive queries, ensuring the preservation of the original cross-modal alignment for neutral inputs.

Empirical evaluations on multiple fairness benchmarks—including FairFace (Kärkkäinen and Joo, 2021), UTKFace (Zhang et al., 2017), and the challenging FACET (Gustafson et al., 2023) dataset—demonstrate that RG-TTA significantly reduces social bias across various demographic attributes. Notably, our framework effectively resolves the fairness–utility trade-off by achieving substantial bias reduction alongside higher accuracy on tasks such as ImageNet-1K (Deng et al., 2009) compared to existing query-independent baselines. By performing episodic optimization with this gated reward, RG-TTA provides a practical design principle for mitigating bias without broadly disrupting the alignment of vision-language models.

## 2 Related Work

**Social debiasing in vision-language models.** Social debiasing for CLIP-style VLMs is broadly categorized into (i) *training-based* approaches (Alabdulmohsin et al., 2024; Hirota et al., 2025; Zhang et al., 2025), which suppress sensitive-attribute signals by introducing additional objectives or modules during learning, and (ii) *training-free* approaches (Chuang et al., 2023; Gerych et al., 2024), which keep the foundation model fixed and apply post-hoc adjustments to embeddings or outputs. The former includes joint debiasing methods that align and remove biases across both modalities, while the latter estimates bias directions from attribute prompts and removes them via lightweight projections. However, many existing

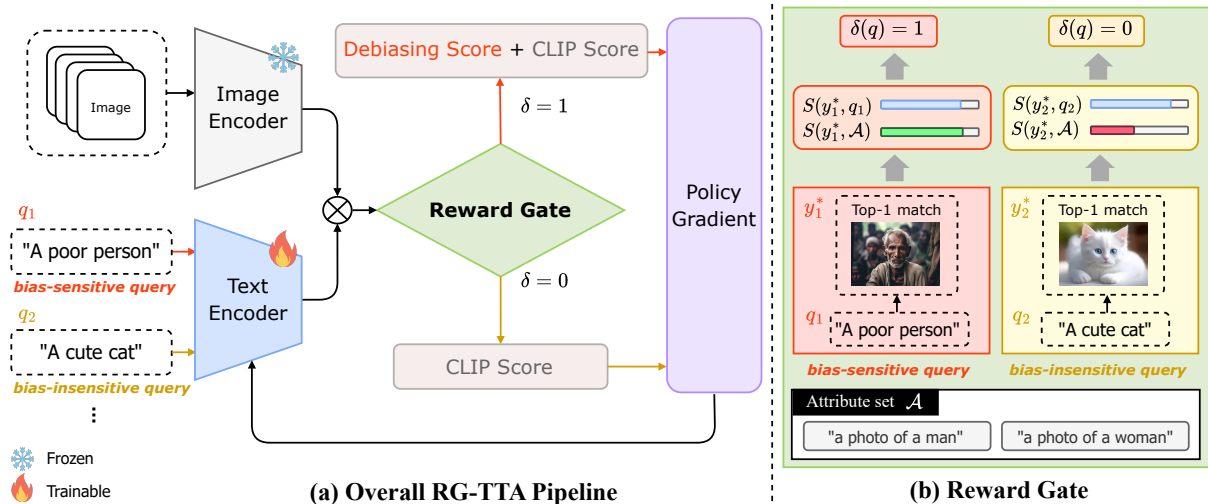


Figure 3: Overview of RG-TTA: RL-based episodic test-time adaptation with **reward gating** via the indicator  $\delta(q)$ . For each query, we update only the query-modality encoder (text encoder for text queries; image encoder for image queries) with a few policy-gradient steps on a truncated top- $K$  candidate set. The indicator  $\delta(q)$  controls the episode reward: when  $\delta(q) = 0$ , we optimize an alignment-only reward; when  $\delta(q) = 1$ , we add an attribute-balancing reward. The gate activates when the top-1 match  $y^*$  is sufficiently close to the attribute-alignment distribution, indicating elevated attribute entanglement.

137 methods remain *query-independent*, applying a single  
 138 globally-defined transformation to all queries.  
 139 Consequently, such fixed rules struggle to optimize  
 140 for both objectives simultaneously, forcing a compromise  
 141 between effective bias mitigation and the  
 142 preservation of general model utility.

143 **Test-Time Adaptation of CLIP via Reinforcement Learning.** Test-time Adaptation (TTA)  
 144 aims to improve model performance on unlabeled  
 145 test inputs by performing parameter updates at inference  
 146 time, enabling models to adapt to distribution shifts  
 147 without retraining (Sun et al., 2020; Wang et al., 2021a).  
 148 Early TTA methods relied on auxiliary self-supervised  
 149 objectives such as entropy minimization (Sun et al., 2020; Liu et al., 2021; Wang et al., 2021a).  
 150 However, these approaches often suffer from instability  
 151 from objective mismatch or prediction collapse (Park et al., 2025),  
 152 motivating recent interest in utilizing VLM’s internal  
 153 alignment signal as direct feedback. In particular,  
 154 treating CLIP similarity as explicit reward and  
 155 performing RL-based optimization (Zancato et al., 2023)  
 156 at inference time (Zhao et al., 2024) has been shown  
 157 to reduce the collapse behavior observed in entropy-based  
 158 updates and to improve zero-shot generalization. We  
 159 build upon the feedback-driven TTA paradigm and extend  
 160 it to address social bias in VLMs. Our key contributions  
 161 are (1) a bias-

sensitivity gate that activates debiasing per input,  
 and (2) a reward that combines cross-modal alignment  
 with a bias-subspace debiasing signal. With episodic  
 test-time updates and selective reward gating, RG-TTA  
 mitigates the fairness–utility trade-off, improving  
 fairness and general-purpose utility.

### 3 Method

We propose **Reward-Gated Test-Time Adaptation (RG-TTA)**, an RL-based framework for selective debiasing that adaptively updates model parameters during inference. RG-TTA is built on two key components. First, a **selective gating strategy** evaluates query sensitivity to trigger **fairness regularization** only when necessary. Second, an **adaptive reward function** balances CLIP alignment with an **attribute-balancing reward**, which encourages a uniform representation by favoring under-represented attributes. We update the model using a tractable approximation of the REINFORCE (Williams, 1992) algorithm. By focusing updates on the most relevant candidates through top- $K$  truncation, this approach ensures that the optimization stays centered on query-specific semantics and limits drift via episodic resets.

#### 3.1 Preliminaries

**CLIP.** A pretrained vision-language model (VLM) consists of an image encoder  $f(\cdot)$  and a text

192 encoder  $g(\cdot)$ , which project both modalities into  
 193 a shared embedding space (Radford et al., 2021).  
 194 Given a query  $q$  and a candidate  $y \in \mathcal{Y}$  selected  
 195 from the opposite modality, we define an alignment  
 196 score  $S(q, y)$  as the cosine similarity between their  
 197 embeddings. Specifically, in the text-to-image set-  
 198 ting we use  $S(q, y) = \cos(g(q), f(y))$ , whereas  
 199 in the image-to-text setting we use  $S(q, y) =$   
 200  $\cos(f(q), g(y))$ . We use the policy score  $S_\theta(q, y)$   
 201 for candidate ranking and parameterizing the policy.  
 202 We use the reference score  $S_{\text{ref}}(q, y)$  only to com-  
 203 pute the CLIP reward;  $S_{\text{ref}}$  is obtained from a fixed  
 204 CLIP ViT-L/14 model throughout all experiments.

### 205 Test-time adaptation in vision–language tasks.

206 Test-time adaptation (TTA (Sun et al., 2020; Wang  
 207 et al., 2021a)) updates a trained model at inference  
 208 time using a few unlabeled steps. We follow an  
 209 *episodic* protocol: each query  $q$  is adapted inde-  
 210 pendently and the parameters are reset before the  
 211 next query, which mitigates negative transfer. For  
 212 vision–language retrieval, where  $q$  can be text or an  
 213 image, we adapt only the query-modality encoder  
 214 and keep the opposite-modality encoder fixed.

### 215 3.2 Selective Gating Strategy

216 Applying debiasing uniformly to all queries is  
 217 unnecessary and can induce parameter drift on  
 218 bias-insensitive queries. To mitigate the structural  
 219 fairness–utility trade-off, we introduce a gating  
 220 mechanism that measures the alignment discrep-  
 221 ancy between query semantics and demographic  
 222 attributes. For each query  $q$ , we select the top-1  
 223 candidate  $y^* = \arg \max_y S_\theta(q, y)$  as an **anchor**  
 224 and assess whether the query–candidate alignment  
 225 is disproportionately driven by demographic as-  
 226 sociations by comparing  $y^*$  to the mean attribute  
 227 similarity over a predefined attribute set  $\mathcal{A}$  of  
 228 protected-group exemplars. We instantiate  $\mathcal{A}$  in the  
 229 same modality as the query (e.g., attribute prompts  
 230 for text and exemplar images for vision), so that  
 231  $S_\theta(a, y^*)$  uses the same scoring function; concrete  
 232 instantiations are given in Sec. 4.4.

233 The gating logic is based on the following intu-  
 234 ition: a large discrepancy suggests that the match  
 235 is driven by general, attribute-independent seman-  
 236 tic alignment, in which case we deactivate fairness  
 237 regularization ( $\delta(q) = 0$ ) and focus on refining  
 238 standard cross-modal alignment to enhance utility.  
 239 Conversely, if  $y^*$  is close to the attribute alignment  
 240 distribution (within a threshold  $\epsilon$ ), the signal is  
 241 likely entangled with a specific attribute category,

242 triggering an attribute-balancing reward ( $\delta(q) = 1$ ).  
 243 Accordingly, the gate is defined as:

$$244 \delta(q) = \mathbb{I} \left[ S(q, y^*) - \frac{1}{|\mathcal{A}|} \sum_{a \in \mathcal{A}} S(a, y^*) < \epsilon \right]. \quad (1)$$

245 Here,  $\mathbb{I}[\cdot]$  denotes the indicator function and  $\epsilon$  is  
 246 a fixed threshold (0.02) kept constant throughout  
 247 the episode to prevent parameter drift.

### 248 3.3 Adaptive Reward Function

249 In this section, we define a candidate-wise reward  
 250  $r(q, y_k)$  for each candidate  $y_k \in \mathcal{Y}_K(q)$ , where  
 251  $\mathcal{Y}_K(q)$  denotes the candidate set used for adap-  
 252 tation. The construction of  $\mathcal{Y}_K(q)$  is deferred to  
 253 3.4. Our reward always includes a CLIP-based  
 254 alignment signal to preserve general cross-modal  
 255 matching ability, and adds an attribute-balancing  
 256 reward only when the gate is activated. This de-  
 257 sign maintains a balance between alignment and  
 258 debiasing within each episode.

259 **CLIP alignment reward.** We follow prior work  
 260 in defining a CLIP-based alignment reward (Zhao  
 261 et al., 2024). Concretely, for each candidate  $y_k$  we  
 262 compute a nonnegative alignment score as  $s_k =$   
 263  $\max(S_{\text{ref}}(q, y_k), 0)$ , and use the episode mean  $\bar{s} =$   
 264  $\frac{1}{K} \sum_{j=1}^K s_j$  as a baseline. The resulting baseline-  
 265 normalized alignment reward is

$$266 r_{\text{clip}}(q, y_k) = s_k - \bar{s}. \quad (2)$$

### 267 Bias subspace and attribute-balancing reward.

268 When  $\delta(q) = 1$ , we incorporate an attribute-  
 269 balancing reward computed in a predefined *bias*  
 270 *subspace* into the reward. Let  $\mu_c \in \mathbb{R}^D$  denote the  
 271 class-mean embedding for class  $c$  in the shared em-  
 272 bedding space, precomputed from a labeled source  
 273 dataset and kept fixed during test-time adaptation.  
 274 We further define a reference vector  $\bar{\mu}$ , which can  
 275 be instantiated as the average of  $\{\mu_c\}_{c=1}^C$ , and con-  
 276 struct the set of difference vectors  $\{\mu_c - \bar{\mu}\}_{c=1}^C$ . We  
 277 then perform PCA on these differences and extract  
 278 all nonzero principal components, yielding  $d$  com-  
 279 ponents in total. The resulting orthonormal basis  
 280  $U \in \mathbb{R}^{D \times d}$  defines the bias subspace.

281 For each selected candidate  $y_k$ , let  $z_k$  denote  
 282 the embedding produced by the frozen encoder of  
 283 the opposite modality. We then map both candi-  
 284 dates and class prototypes into the bias subspace  
 285 using the projection operator  $P = UU^\top$ , yielding  
 286  $\tilde{z}_k = P(z_k - \bar{\mu})$  and  $\tilde{\mu}_c = P(\mu_c - \bar{\mu})$ . Next, we

compute a soft assignment over classes by normalizing Gaussian-kernel similarities with a temperature parameter  $\gamma$ . Defining the squared distance between the projected candidate and class prototype as  $d_{kc} = \|\tilde{z}_k - \tilde{\mu}_c\|_2^2$ , the soft assignment is given by:

$$\alpha_c^{(k)} = \frac{\exp(-d_{kc}/\gamma)}{\sum_{c'=1}^C \exp(-d_{kc'}/\gamma)}. \quad (3)$$

To estimate the attribute distribution at the episode level, we compute the popularity for each class as  $p_c = \frac{1}{K} \sum_{j=1}^K \alpha_c^{(j)}$ . We then define the attribute-balancing reward for each candidate as

$$d(q, y_k) = \sum_{c=1}^C \alpha_c^{(k)} \left( p_c - \frac{1}{C} \right). \quad (4)$$

This score reflects how strongly candidate  $y_k$  is associated with attribute classes that are over- or under-represented in the current episode, encouraging a more balanced set of selected candidates.

**Final combined reward.** Finally, we define the overall reward for each selected candidate by combining the CLIP alignment reward with the query-conditioned attribute-balancing reward:

$$r(q, y_k) = r_{\text{clip}}(q, y_k) - \delta(q) \lambda d(q, y_k), \quad (5)$$

where  $\lambda$  is a hyperparameter that controls the strength of the attribute-balancing reward. When  $\delta(q) = 0$ , the reward reduces to the CLIP alignment term alone. When  $\delta(q) = 1$ , candidates associated with over-represented classes (i.e.,  $p_c > 1/C$ ) tend to have larger  $d(q, y_k)$  and thus incur a larger balancing penalty, while candidates linked to under-represented classes (i.e.,  $p_c < 1/C$ ) tend to have smaller or negative  $d(q, y_k)$  and are relatively favored. As a result, the episode-level attribute distribution is encouraged toward the uniform prior.

### 3.4 Optimization and Episodic Update

We optimize the proposed objective at test time via episodic policy-gradient updates (Wang et al., 2021a; Shu et al., 2022), treating each input query  $q$  as an episode. We first evaluate the gating indicator  $\delta(q)$  to determine whether to include the attribute-balancing term in the episode reward. Conditioned on this decision, we choose a candidate budget  $K$  (e.g.,  $K = 10$  when  $\delta(q) = 0$  and  $K = 1024$  when  $\delta(q) = 1$ ), and construct a truncated candidate set  $\mathcal{Y}_K(q)$  by selecting the top- $K$  candidates from the fixed pool  $\mathcal{Y}$  according to the alignment

score  $S_\theta(q, y)$ . We then compute the candidate-wise reward  $r(q, y_k)$  for each  $y_k \in \mathcal{Y}_K(q)$  and update the query-modality encoder parameters  $\theta$  with a small number of gradient steps. After the episode ends, we reset the parameters to their initial state before processing the next query, mitigating negative transfer across queries. We define the policy  $\pi_\theta(y | q)$  over the full candidate pool  $\mathcal{Y}$  using a softmax of the alignment score:

$$\pi_\theta(y | q) = \frac{\exp(S_\theta(q, y))}{\sum_{y' \in \mathcal{Y}} \exp(S_\theta(q, y'))}. \quad (6)$$

Although  $\pi_\theta$  is normalized over the full candidate pool  $\mathcal{Y}$  (i.e., the denominator is computed over  $\mathcal{Y}$ ), we approximate the policy-gradient objective by summing only over the truncated set  $\mathcal{Y}_K(q)$  for the current query. Concretely, for each episode we minimize the following REINFORCE-style objective:

$$\mathcal{L}(q) = -\frac{1}{K} \sum_{y_k \in \mathcal{Y}_K(q)} r(q, y_k) \log \pi_\theta(y_k | q). \quad (7)$$

This top- $K$  truncation based on  $S(q, y)$  yields a tractable approximation to the full policy-gradient objective while focusing updates on the most relevant candidates for the query.

## 4 Experiments

### 4.1 Datasets

To comprehensively evaluate both debiasing performance and generalization, we used a diverse set of benchmarks. For fairness evaluation, we consider in-domain settings on FairFace (val) (Kärkkäinen and Joo, 2021) and UTK-Face (Zhang et al., 2017), and an out-of-domain setting on FACET (Gustafson et al., 2023). To assess whether adaptation preserves general-purpose zero-shot utility, we additionally evaluate on ImageNet-1K (Deng et al., 2009) for image classification and Flickr1k (Plummer et al., 2015) for retrieval.

### 4.2 Metrics

**Fairness metrics.** Following prior work (Berg et al., 2022a; Seth et al., 2023a), we use retrieval-based metrics: MaxSkew@ $k$  and NDKL@ $k$ . These metrics quantify the disparity in the distribution of protected attributes (e.g., gender, age, and race) within the top- $k$  retrieved images for neutral queries. MaxSkew measures the maximum representation of a dominant group, while NDKL mea-

Table 1: Gender and age debiasing performance across different source datasets. We report results on in/out-of-domain fairness benchmarks and zero-shot utility tasks. ABLE is calculated based on in-domain metrics. Best results are shown in **bold**, and second-best results are underlined.

| Backbone                                 | Biases         | Methods        | In-Domain       |                   | Out-of-Domain   |                   |                 |                   | IN1K                |              | Flickr             |              | ABLE (%) $\uparrow$ |
|--|----------------|----------------|-----------------|-------------------|-----------------|-------------------|-----------------|-------------------|---------------------|--------------|--------------------|--------------|---------------------|
|  |                |                | Source Dataset  |                   | Cross Dataset   |                   | FACET           |                   | Acc. (%) $\uparrow$ |              | R@5 (%) $\uparrow$ |              |                     |
|  |                |                | MS $\downarrow$ | NDKL $\downarrow$ | MS $\downarrow$ | NDKL $\downarrow$ | MS $\downarrow$ | NDKL $\downarrow$ | Top-1               | Top-5        | TR                 | IR           |                     |
| <b>Source: UTKFace   Cross: FairFace</b> |                |                |                 |                   |                 |                   |                 |                   |                     |              |                    |              |                     |
| ViT-B/16                                 | Gender         | Original CLIP  | 0.114           | 0.080             | 0.218           | 0.088             | 0.478           | 0.215             | 68.31               | 91.83        | 96.4               | 85.5         | 77.39               |
|  |                | CLIP-clip      | 0.070           | 0.055             | 0.133           | 0.038             | 0.459           | 0.190             | 67.81               | 91.42        | 95.4               | 83.0         | 78.52               |
|  |                | Biased-prompts | 0.179           | 0.062             | 0.161           | 0.048             | 0.460           | 0.215             | 65.07               | 89.38        | 94.3               | 86.1         | 73.18               |
|  |                | Joint V-L      | <b>0.048</b>    | 0.043             | 0.101           | <b>0.032</b>      | 0.456           | 0.181             | 67.99               | 91.64        | 95.8               | 84.6         | 79.36               |
|  |                | Ours           | <u>0.051</u>    | <b>0.029</b>      | <b>0.080</b>    | <u>0.035</u>      | <b>0.053</b>    | <b>0.064</b>      | <b>70.38</b>        | <b>93.00</b> | <b>97.2</b>        | <b>88.5</b>  | <b>80.87</b>        |
| Age                                      | Original CLIP  | 0.421          | 0.229           | 0.657             | 0.433           | 0.744             | 0.367           | 68.31             | 91.83               | 96.4         | 85.5               | 66.96        |                     |
|  | CLIP-clip      | 0.393          | <b>0.215</b>    | 0.643             | 0.430           | 0.745             | 0.364           | 67.93             | 91.58               | 96.1         | 84.5               | 67.70        |                     |
|  | Biased-prompts | 0.578          | 0.451           | 0.777             | 0.550           | <b>0.635</b>      | 0.355           | 66.43             | 90.28               | 94.1         | 85.2               | 60.83        |                     |
|  | Joint V-L      | 0.414          | 0.231           | <b>0.606</b>      | <b>0.410</b>    | 0.746             | 0.365           | 67.63             | 91.46               | 95.6         | 84.6               | 66.86        |                     |
|  | Ours           | <b>0.151</b>   | <u>0.226</u>    | <u>0.641</u>      | <u>0.420</u>    | <u>0.742</u>      | <b>0.330</b>    | <b>70.36</b>      | <b>92.99</b>        | <b>97.2</b>  | <b>88.3</b>        | <b>77.39</b> |                     |
| ViT-B/32                                 | Gender         | Original CLIP  | 0.066           | <b>0.032</b>      | 0.138           | 0.054             | 0.485           | 0.225             | 63.39               | 88.83        | 94.7               | 83.5         | 75.60               |
|  |                | CLIP-clip      | 0.098           | 0.045             | 0.253           | 0.105             | 0.500           | 0.240             | 62.21               | 88.23        | 93.0               | 81.0         | 73.79               |
|  |                | Biased-prompts | 0.089           | 0.036             | 0.094           | <b>0.027</b>      | 0.417           | 0.164             | 60.37               | 86.75        | 93.6               | 82.4         | 72.74               |
|  |                | Joint V-L      | <b>0.043</b>    | 0.033             | 0.108           | 0.039             | 0.469           | 0.212             | 62.46               | 88.23        | 94.7               | 82.9         | 75.60               |
|  |                | Ours           | <u>0.054</u>    | 0.038             | <b>0.088</b>    | <u>0.035</u>      | <b>0.050</b>    | <b>0.021</b>      | <b>69.76</b>        | <b>92.31</b> | <b>97.0</b>        | <b>86.6</b>  | <b>79.60</b>        |
| Age                                      | Original CLIP  | 0.412          | 0.253           | 0.617             | 0.416           | 0.752             | 0.388           | 63.39             | 88.83               | 94.7         | 83.5               | 64.77        |                     |
|  | CLIP-clip      | 0.415          | 0.264           | 0.659             | 0.435           | 0.754             | 0.397           | 62.70             | 88.31               | 94.2         | 83.2               | 64.34        |                     |
|  | Biased-prompts | 0.522          | 0.409           | 0.701             | 0.497           | <b>0.663</b>      | <b>0.366</b>    | 61.07             | 86.92               | 92.0         | 82.2               | 60.19        |                     |
|  | Joint V-L      | 0.407          | <b>0.252</b>    | 0.627             | 0.416           | 0.751             | 0.370           | 62.93             | 88.66               | 94.1         | 82.5               | 64.69        |                     |
|  | Ours           | <b>0.127</b>   | 0.283           | <b>0.385</b>      | <b>0.364</b>    | <u>0.741</u>      | <u>0.369</u>    | <b>69.75</b>      | <b>92.30</b>        | <b>96.9</b>  | <b>86.4</b>        | <b>77.85</b> |                     |
| <b>Source: FairFace   Cross: UTKFace</b> |                |                |                 |                   |                 |                   |                 |                   |                     |              |                    |              |                     |
| ViT-B/16                                 | Gender         | Original CLIP  | 0.218           | 0.088             | 0.114           | 0.080             | 0.478           | 0.215             | 68.31               | 91.83        | 96.4               | 85.5         | 73.87               |
|  |                | CLIP-clip      | 0.103           | 0.026             | 0.083           | 0.062             | 0.478           | 0.199             | 68.00               | 91.50        | 95.4               | 83.0         | 77.55               |
|  |                | Biased-prompts | 0.161           | 0.048             | 0.179           | 0.062             | 0.460           | 0.215             | 65.07               | 89.38        | 94.3               | 86.1         | 73.78               |
|  |                | Joint V-L      | <b>0.080</b>    | <b>0.025</b>      | 0.040           | 0.023             | 0.446           | 0.170             | 68.05               | 91.63        | 96.6               | 84.3         | 78.35               |
|  |                | Ours           | <u>0.082</u>    | 0.031             | <b>0.030</b>    | <b>0.022</b>      | <b>0.114</b>    | <b>0.040</b>      | <b>70.32</b>        | <b>92.98</b> | <b>97.2</b>        | <b>88.5</b>  | <b>79.75</b>        |
| Age                                      | Original CLIP  | 0.657          | 0.433           | 0.421             | 0.229           | 0.744             | 0.367           | 68.31             | 91.83               | 96.4         | 85.5               | 58.94        |                     |
|  | CLIP-clip      | 0.647          | 0.432           | 0.402             | 0.215           | 0.742             | 0.373           | 67.97             | 91.61               | 96.3         | 84.4               | 59.16        |                     |
|  | Biased-prompts | 0.777          | 0.550           | 0.578             | 0.451           | <b>0.635</b>      | 0.355           | 66.43             | 90.28               | 94.1         | 85.2               | 54.33        |                     |
|  | Joint V-L      | 0.608          | <b>0.294</b>    | 0.377             | <b>0.115</b>    | 0.738             | 0.341           | 68.34             | 91.74               | 96.0         | 84.0               | 60.61        |                     |
|  | Ours           | <b>0.526</b>   | <u>0.318</u>    | <b>0.245</b>      | 0.222           | 0.742             | <b>0.265</b>    | <b>70.36</b>      | <b>92.99</b>        | <b>97.2</b>  | <b>88.3</b>        | <b>64.24</b> |                     |
| ViT-B/32                                 | Gender         | Original CLIP  | 0.138           | 0.054             | 0.066           | 0.032             | 0.485           | 0.225             | 63.39               | 88.83        | 94.7               | 83.5         | 73.37               |
|  |                | CLIP-clip      | 0.107           | 0.030             | 0.061           | 0.023             | 0.492           | 0.215             | 59.62               | 86.29        | 90.9               | 76.2         | 71.68               |
|  |                | Biased-prompts | 0.094           | <b>0.027</b>      | 0.089           | 0.036             | 0.417           | 0.164             | 60.37               | 86.75        | 93.6               | 82.4         | 72.59               |
|  |                | Joint V-L      | 0.090           | 0.030             | <b>0.050</b>    | <b>0.021</b>      | 0.466           | 0.204             | 62.52               | 88.56        | 94.9               | 82.9         | 74.24               |
|  |                | Ours           | <b>0.078</b>    | 0.032             | <b>0.050</b>    | 0.029             | <b>0.137</b>    | <b>0.036</b>      | <b>69.76</b>        | <b>92.30</b> | <b>97.0</b>        | <b>86.6</b>  | <b>79.54</b>        |
| Age                                      | Original CLIP  | 0.617          | 0.416           | 0.412             | 0.253           | 0.752             | 0.388           | 63.39             | 88.83               | 94.7         | 83.5               | 58.29        |                     |
|  | CLIP-clip      | 0.635          | 0.425           | 0.400             | 0.252           | 0.749             | 0.387           | 62.40             | 88.30               | 94.5         | 82.5               | 57.32        |                     |
|  | Biased-prompts | 0.701          | 0.497           | 0.522             | 0.409           | <b>0.663</b>      | 0.366           | 61.07             | 86.92               | 92.0         | 82.2               | 54.76        |                     |
|  | Joint V-L      | 0.572          | 0.364           | 0.385             | <b>0.195</b>    | 0.750             | 0.381           | 63.13             | 88.71               | 94.1         | 82.8               | 59.60        |                     |
|  | Ours           | <b>0.523</b>   | <b>0.229</b>    | <b>0.245</b>      | <u>0.222</u>    | <u>0.743</u>      | <b>0.265</b>    | <b>69.75</b>      | <b>92.30</b>        | <b>96.9</b>  | <b>86.4</b>        | <b>64.09</b> |                     |

374 sures the divergence from uniform distribution. For  
375 both metrics, lower values indicate a fairer model.

376 **Utility (V-L alignment) metrics.** To ensure that  
377 our selective debiasing does not compromise the  
378 intrinsic V-L alignment of the pre-trained model,  
379 we evaluated zero-shot performance on standard  
380 benchmarks. We report Top-1 and Top-5 accu-  
381 racy on ImageNet-1K(Deng et al., 2009) for  
382 classification, and Recall@5 for both Image-to-  
383 Text (TR) and Text-to-Image (IR) retrieval on  
384 Flickr1k(Plummer et al., 2015).

385 **Alignment and Bias Level Evaluation (ABLE).**  
386 Single metrics capture either fairness or utility, but  
387 not their trade-off. We use ABLE proposed by

Zhang et al.(Zhang et al., 2025) for a holistic as-  
388 sessment. ABLE is defined as the harmonic mean  
389 of the zero-shot accuracy and the fairness score:  
390

$$391 \text{ABLE} = \frac{2}{\frac{1}{acc} + \frac{1}{\exp(-\text{MaxSkew}@k)}} \quad (8)$$

392 where  $acc$  denotes the ImageNet Top-1 accuracy.  
393 Higher ABLE indicates a better balance between  
394 mitigating social bias and retaining zero-shot accu-  
395 racy; we report  $\text{ABLE} \times 100$  (%) in tables.

### 396 4.3 Baselines

397 We compare RG-TTA against the **Original**  
398 **CLIP** and three representative debiasing base-  
399 lines using ViT-B/16 and ViT-B/32 backbones.

Table 2: Race debiasing performance using UTKFace as the source.

| Backbone | Methods        | UTKFace         |                   | IN1K                |              | Flickr             |             | ABLE (%) $\uparrow$ |
|----------|----------------|-----------------|-------------------|---------------------|--------------|--------------------|-------------|---------------------|
|          |                | MS $\downarrow$ | NDKL $\downarrow$ | Acc. (%) $\uparrow$ |              | R@5 (%) $\uparrow$ |             |                     |
|          |                |                 |                   | Top-1               | Top-5        | TR                 | IR          |                     |
| ViT-B/16 | Original CLIP  | 0.575           | 0.137             | 68.31               | 91.83        | 96.4               | 85.5        | 63.31               |
|          | CLIP-clip      | 0.613           | 0.157             | 67.74               | 91.48        | 95.8               | 85.1        | 60.20               |
|          | Biased-prompts | 0.604           | 0.208             | 67.00               | 90.72        | 94.1               | 85.8        | 60.21               |
|          | Joint V-L      | 0.378           | 0.069             | 68.07               | 91.64        | 96.5               | 83.8        | 68.30               |
|          | Ours           | <b>0.150</b>    | <b>0.022</b>      | <b>70.35</b>        | <b>93.00</b> | <b>97.2</b>        | <b>88.3</b> | <b>77.42</b>        |
| ViT-B/32 | Original CLIP  | 0.698           | 0.213             | 63.39               | 88.83        | 94.7               | 83.5        | 55.75               |
|          | CLIP-clip      | 0.840           | 0.426             | 62.90               | 88.32        | 93.6               | 81.74       | 51.20               |
|          | Biased-prompts | <b>0.317</b>    | <b>0.140</b>      | 61.80               | 87.46        | 91.9               | 83.34       | 60.21               |
|          | Joint V-L      | 0.638           | 0.230             | 63.01               | 88.56        | 94.6               | 82.4        | 57.48               |
|          | Ours           | 0.351           | 0.281             | <b>69.74</b>        | <b>92.34</b> | <b>97.0</b>        | <b>86.5</b> | <b>70.07</b>        |

CLIP-clip (Wang et al., 2021b) removes bias-correlated embedding dimensions identified via mutual information with attribute labels. Biased-prompts (Chuang et al., 2023) neutralizes bias directions by projecting embeddings using prompt-derived attribute subspaces without retraining. Joint V-L (Zhang et al., 2025) jointly debiases image and text representations to mitigate over-debiasing effects. All methods are evaluated under identical settings for Gender, Age, and Race.

#### 4.4 Implementation Details

We use the official pretrained CLIP checkpoints (Radford et al., 2021) with ViT-B/16, ViT-B/32, and ViT-L/14 backbones. Consistent with the episodic TTA protocol, we update only the encoder corresponding to the input query modality (e.g., the text encoder for text-to-image retrieval) while keeping the target modality encoder frozen. Optimization is performed using AdamW (Loshchilov and Hutter, 2019). Crucially, to balance efficiency and performance, we dynamically adjust the computational budget based on the gating decision  $\delta(q)$ : we perform a lightweight update with  $T=3$  steps and  $K=10$  candidates for bias-insensitive queries ( $\delta(q)=0$ ), while expanding to  $T=10$  steps and  $K=1024$  candidates for bias-sensitive ones ( $\delta(q)=1$ ). Detailed hyperparameters, including learning rates, gating thresholds, and prompt templates, are provided in Appendix A.

#### 4.5 Results

Table 1 summarizes bias-mitigation results across different source datasets (UTKFace and FairFace) for Gender and Age. Overall, our method improves fairness in both in-domain and out-of-domain settings while enhancing ImageNet zero-shot performance. For Gender, we observe a clear in-domain improvement on UTKFace (MaxSkew@k:

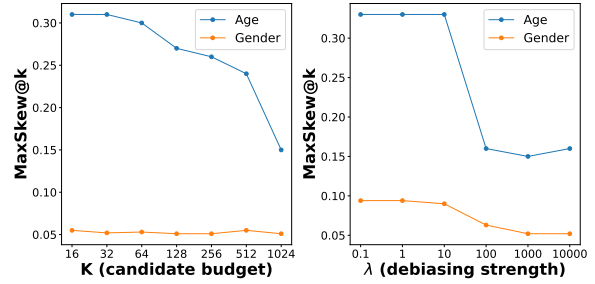


Figure 4: Ablation of  $K$  and  $\lambda$  showing their effect on MaxSkew@k.

Table 3: Ablation of reward variants for selective debiasing.

| Reward       | MS $\downarrow$ | IN1K Top-1 $\uparrow$ | ABLE (%) $\uparrow$ |
|--------------|-----------------|-----------------------|---------------------|
| CLIP         | 0.604           | 70.41                 | 61.54               |
| CLIP+Debias  | 0.149           | 53.98                 | 58.23               |
| RG-TTA(Ours) | 0.151           | 70.36                 | 77.39               |

0.114 $\rightarrow$ 0.051), with particularly large mitigation under distribution shift on FACET (MaxSkew@k: 0.478 $\rightarrow$ 0.053). Meanwhile, ImageNet Top-1 accuracy increases from 68.31 to 70.38, and ABL also rises from 77.39 to 80.87, indicating that fairness gains do not come at the expense of utility (see Table 1 for full results).

This trend is consistent across attributes and source configurations: for Age and Race, we observe substantial in-domain fairness gains while maintaining (or slightly improving) zero-shot utility (Tables 1, 2). Switching the source dataset to FairFace shows the same behavior across in-domain and out-of-domain settings, aligning with our design—gating with episodic resets—that focuses updates on bias-sensitive queries while preventing drift on others.

## 5 Discussion

### 5.1 Ablation Study

We ablate key design choices to examine how selective reward-gating mitigates the fairness-utility trade-off. Figure 4 ablates the hyperparameters used when debiasing is active, varying the candidate budget  $K$  and the debiasing strength  $\lambda$  in the  $\delta(q) = 1$  regime. We observe that increasing  $K$  generally improves MaxSkew@1000 (notably for Age), suggesting that a larger truncated set yields a more stable estimate of the episode-level attribute distribution and thus a more reliable balancing signal. Varying  $\lambda$  shows a similar trend: weak values yield limited mitigation, while larger values

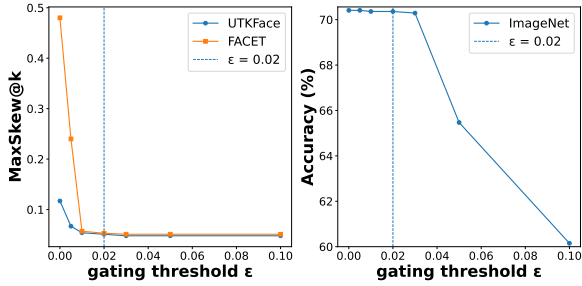


Figure 5: Sensitivity to the gating threshold  $\epsilon$ . Left: MaxSkew@ $k$  (gender) on UTKFace/FACET. Right: ImageNet-1K Top-1 accuracy.

468 deliver consistent gains with diminishing returns,  
 469 suggesting performance is not brittle once  $\lambda$  passes  
 470 a minimal effective threshold.

471 We next examine when to apply the debiasing  
 472 signal versus how strongly to apply it. Table 3 com-  
 473 pares (i) an alignment-only reward, (ii) adding the  
 474 attribute-balancing term uniformly to all queries,  
 475 and (iii) activating the balancing term only when  
 476 the gating indicator  $\delta(q)$  triggers. While uniform  
 477 debiasing can substantially reduce MaxSkew, it  
 478 risks unnecessary adaptation and notably degrades  
 479 zero-shot utility; in contrast, selective activation  
 480 preserves utility while retaining the fairness gains,  
 481 yielding a markedly better overall balance.

482 Finally, we characterize the gate through thresh-  
 483 old sensitivity and activation patterns under a  
 484 dataset-level proxy setting. Figure 5 shows that  
 485 across a broad range of the threshold  $\epsilon$ , fairness  
 486 improves rapidly and remains largely stable, and  
 487 utility is stable around the default  $\epsilon=0.02$ , with no-  
 488 ticeable degradation only when  $\epsilon$  is overly large and  
 489 activates debiasing too aggressively. We also evalu-  
 490 ate the gate’s activation behavior (Table 4) under a  
 491 proxy labeling scheme: treating UTKFace/FACET  
 492 queries as debiasing-needed and ImageNet-1K  
 493 queries as debiasing-not-needed, the gate activates  
 494 for almost all UTKFace/FACET queries while  
 495 remaining inactive for nearly all ImageNet-1K  
 496 queries, yielding low false negative/positive rates  
 497 in this proxy setting.

## 498 5.2 Out-of-Domain Analysis on FACET

499 On FACET (out-of-domain), our method transfers  
 500 well for Gender, maintaining lower MaxSkew@ $k$   
 501 than prior methods. We hypothesize that Gender  
 502 is associated with multiple robust visual cues (e.g.,  
 503 face, hairstyle, clothing), which preserve group se-  
 504 paration even after bias-subspace projection and lead  
 505 to more stable debiasing signals. In contrast, Age

Table 4: FP/FN are computed under a proxy labeling scheme where UTKFace and FACET are treated as positives (debiasing-needed) and ImageNet-1K as negatives (debiasing-not-needed).

| Dataset     | $\delta = 1$ (%) | $\delta = 0$ (%) | FP/FN    |
|-------------|------------------|------------------|----------|
| UTKFace     | 99.9             | 0.1              | FN = 0.1 |
| FACET       | 99.8             | 0.2              | FN = 0.2 |
| ImageNet-1K | 0.1              | 99.9             | FP = 0.1 |



Figure 6: **Gating failure cases on FACET.** Bias-sensitive queries are incorrectly gated off ( $\delta(q) = 0$ ), preventing fairness regularization from being activated.

506 relies on fine-grained facial cues that are sensitive  
 507 to distance, resolution, and occlusion; in FACET’s  
 508 unconstrained images, these cues are weakened,  
 509 reducing the reliability of the debiasing signal and  
 510 ultimately limiting gains.

## 511 5.3 Gating Failure Case Analysis

512 Figure 6 shows *false-negative* cases on FACET,  
 513 where bias-sensitive queries are gated off ( $\delta(q) = 0$ )  
 514 and thus do not trigger fairness regularization.  
 515 These failures often arise when salient object and  
 516 scene semantics in the top-1 anchor  $y^*$  overwhelm  
 517 demographic-correlated cues, causing the semantic-  
 518 attribute discrepancy to exceed the threshold.

## 519 6 Conclusion

520 We introduced Reward-Gated Test-Time Adapta-  
 521 tion (RG-TTA), a selective test-time debiasing  
 522 framework for CLIP-style vision–language models.  
 523 RG-TTA uses an input-dependent reward gate to ac-  
 524 tivate an attribute-balancing term only when neces-  
 525 sary, together with episodic updates and parameter  
 526 resets to limit unintended drift. Experiments across  
 527 in-domain and out-of-domain fairness benchmarks  
 528 as well as standard zero-shot utility tasks show that  
 529 RG-TTA consistently reduces demographic skew  
 530 while maintaining competitive utility. Overall, our  
 531 results highlight selective, input-conditioned adap-  
 532 tation as a practical design principle for mitigating  
 533 bias without broadly disrupting model behavior.

## 534 Limitations

535 Our approach relies on test-time adaptation (TTA)  
536 with per-query parameter updates, which can in-  
537 crease computation and latency, especially when  
538 using many update steps or large candidate sets.  
539 Because offline-trained debiasing methods amortize  
540 cost during training whereas TTA incurs cost  
541 online, direct runtime comparisons across these  
542 paradigms are not always apples-to-apples and can  
543 vary with the deployment scenario. Our exper-  
544 iments focus on a single protected attribute; ex-  
545 tending the framework to multiple attributes re-  
546 quires more complex reward/constraint design and  
547 may introduce conflicting objectives. The method  
548 also assumes access to an external reward signal  
549 from a stronger model, which raises availability  
550 and cost considerations and may transfer the re-  
551 ward model’s own biases into the adaptation sig-  
552 nal. Moreover, our fairness objective implicitly  
553 targets proximity to a chosen reference distribution  
554 (e.g., uniform), and both evaluation and reward de-  
555 pend on the accuracy and domain robustness of  
556 attribute estimators; TTA behavior can further be  
557 sensitive to hyperparameters and may be unstable  
558 for some queries. Future work will develop more  
559 efficient update schemes to reduce online overhead.  
560 We will also explore scalable multi-attribute ob-  
561 jectives/constraints and robustness techniques to  
562 mitigate sensitivity to reward sources and attribute  
563 estimators.

## 564 Ethics Statement

565 This work proposes a selective test-time adapta-  
566 tion (TTA) approach to mitigate distributional bi-  
567 ases over protected attributes (e.g., gender, age,  
568 and race) that can arise for person-centric queries.  
569 Our experiments use publicly available fairness  
570 evaluation datasets (e.g., FairFace, UTKFace, and  
571 FACET) together with standard utility benchmarks,  
572 and quantify bias using distribution-based fair-  
573 ness metrics. Because face images and protected-  
574 attribute annotations can be sensitive, our study  
575 does not aim to identify individuals and assumes  
576 use strictly in accordance with the datasets’ licenses  
577 and usage conditions. Our method further relies  
578 on an external reward signal from a stronger model  
579 (e.g., a fixed CLIP ViT-L/14 reference), which in-  
580 troduces practical considerations about the avail-  
581 ability and cost of such signals and raises the pos-  
582 sibility that biases present in the reward model  
583 could be propagated through the adaptation pro-

cess. In addition, our fairness objective implicit- 584  
ly assumes a chosen target distribution (e.g., a 585  
uniform prior), which may not be appropriate for 586  
all tasks or domains. We therefore recommend 587  
that any real-world deployment be accompanied 588  
by careful auditing of the reward source and at- 589  
tribute estimators for bias and error, and that use in 590  
high-stakes decision-making contexts be avoided or 591  
subjected to additional, domain-specific validation 592  
and oversight. 593

## References 594

- 595 Salma Abdel Magid, Jui-Hsien Wang, Kushal Kafle, and  
596 Hanspeter Pfister. 2024. *They’re all doctors: Synthe-  
597 sizing diverse counterfactuals to mitigate associative  
598 bias*. *arXiv preprint arXiv:2406.11331*.
- 599 Ibrahim Alabdulmohsin, Xiao Wang, Andreas Steiner,  
600 Priya Goyal, Alexander D’Amour, and Xiaohua Zhai.  
601 2024. Clip the bias: How useful is balancing data  
602 in multimodal learning? In *Proceedings of the In-  
603 ternational Conference on Learning Representations  
604 (ICLR)*.
- 605 Hugo Berg, Siobhan Hall, Yash Bhalgat, Hannah Kirk,  
606 Aleksandar Shtedritski, and Max Bain. 2022a. A  
607 prompt array keeps the bias away: Debiasing vision-  
608 language models with adversarial learning. In *Pro-  
609 ceedings of the 60th Annual Meeting of the Associ-  
610 ation for Computational Linguistics (ACL)*, pages  
611 806–822. Association for Computational Linguistics.
- 612 Hugo Berg, Siobhan Mackenzie Hall, Yash Bhalgat,  
613 Hannah Rose Kirk, Aleksandar Shtedritski, and Max  
614 Bain. 2022b. A prompt array keeps the bias away:  
615 Debiasing vision-language models with adversarial  
616 learning. In *Proceedings of the 2nd Conference of  
617 the Asia-Pacific Chapter of the Association for Com-  
618 putational Linguistics*, pages 806–822. Association  
619 for Computational Linguistics.
- 620 Abeba Birhane, Vinay Uday Prabhu, and Emmanuel  
621 Kahembwe. 2021. *Multimodal datasets: Misog-  
622 yny, pornography, and malignant stereotypes*. *arXiv  
623 preprint arXiv:2110.01963*.
- 624 Ching-Yao Chuang, Varun Jampani, Yuanzhen Li, Anto-  
625 nio Torralba, and Stefanie Jegelka. 2023. Debiasing  
626 vision-language models via biased prompts. In *Pro-  
627 ceedings of the IEEE/CVF Conference on Computer  
628 Vision and Pattern Recognition (CVPR)*, pages 4007–  
629 4016.
- 630 Sepehr Dehdashtian, Lan Wang, and Vishnu Naresh  
631 Boddeti. 2024. Fairerclip: Debiasing clip’s zero-  
632 shot predictions using functions in rkhss. In *Inter-  
633 national Conference on Learning Representations*.  
634 ArXiv:2403.15593.
- 635 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai  
636 Li, and Li Fei-Fei. 2009. Imagenet: A large-scale





Table 5: Race debiasing performance using FairFace as the source.

| Backbone | Methods        | Fairface        |                   | IN1K                |              | Flickr             |             | ABLE (%) $\uparrow$ |
|----------|----------------|-----------------|-------------------|---------------------|--------------|--------------------|-------------|---------------------|
|          |                | MS $\downarrow$ | NDKL $\downarrow$ | Acc. (%) $\uparrow$ |              | R@5 (%) $\uparrow$ |             |                     |
|          |                |                 |                   | Top-1               | Top-5        | TR                 | IR          |                     |
| ViT-B/16 | Original CLIP  | 0.528           | 0.182             | 68.31               | 91.83        | 96.4               | 85.5        | 63.31               |
|          | CLIP-clip      | 0.544           | 0.161             | 67.97               | 91.62        | 95.4               | 85.3        | 62.62               |
|          | Biased-prompts | 0.518           | 0.219             | 67.00               | 90.72        | 94.1               | 85.8        | 63.07               |
|          | Joint V-L      | <b>0.353</b>    | <b>0.125</b>      | 68.07               | 91.64        | 96.5               | 83.8        | 69.14               |
|          | Ours           | 0.372           | 0.167             | <b>70.32</b>        | <b>92.96</b> | <b>97.2</b>        | <b>88.4</b> | <b>69.62</b>        |
| ViT-B/32 | Original CLIP  | 0.568           | 0.165             | 63.39               | 88.83        | 94.7               | 83.5        | 59.84               |
|          | CLIP-clip      | 0.713           | 0.227             | 62.51               | 88.33        | 92.6               | 81.2        | 54.95               |
|          | Biased-prompts | 0.595           | 0.282             | 61.80               | 87.46        | 91.9               | 83.3        | 58.29               |
|          | Joint V-L      | 0.503           | 0.149             | 63.07               | 88.61        | 94.2               | 83.0        | 61.74               |
|          | Ours           | <b>0.389</b>    | <b>0.148</b>      | <b>69.73</b>        | <b>92.26</b> | <b>97.0</b>        | <b>86.5</b> | <b>68.74</b>        |

- **Bias-Sensitive** ( $\delta(q) = 1$ ): We increase the budget to  $T = 10$  steps and  $K = 1024$  candidates. In this case, candidates are ranked by their CLIP alignment scores to filter the most relevant samples for the update.

## A.5 Bias Subspace Construction Details

The bias subspace is constructed offline using the training split of the dataset:

- **Text Queries:** We use attribute-specific prompts formatted as “*a photo of a {attribute class} person*”. One prompt is generated per attribute class.
- **Image Queries:** We construct the reference set  $A$  by sampling  $M = 5$  images per attribute class uniformly from the **UTKFace** dataset. To ensure reproducibility and consistency, this reference set is fixed once and reused across all test episodes. The total size of the reference set is  $|A| = C \times 5$ , where  $C$  is the number of attribute classes.

## B Additional Results

### B.1 Additional Race Debiasing Results using FairFace

In the main text (Table 2), we utilized UTKFace as the source dataset for constructing the bias subspace to mitigate Race bias. To verify the robustness of our framework across different source domains, we conducted an additional experiment using **FairFace** as the source dataset.

The results are presented in Table 5. Consistent with the findings in the main text, our RG-TTA framework demonstrates a superior capability to balance fairness and utility. Regarding fairness, our method effectively mitigates racial



Figure 7: **Gating failure cases on ImageNet.** In ImageNet, queries are expected to be bias-insensitive and thus gated off ( $\delta(q) = 0$ ). Shown are false-positive cases where the gate is incorrectly activated ( $\delta(q) = 1$ ), causing debiasing and alignment rewards to be jointly applied.

bias compared to the Original CLIP; while **Joint V-L** shows competitive scores on ViT-B/16, our method achieves the best performance on ViT-B/32 (MaxSkew: 0.389). Crucially, in terms of utility preservation, our method consistently outperforms all baselines in zero-shot tasks (ImageNet and Flickr) across both backbones, confirming that our selective routing mechanism successfully prevents the over-debiasing observed in static approaches. Consequently, our method achieves the highest ABL scores for both backbones (69.62% and 68.74%), proving that it maintains the optimal trade-off between fairness and utility regardless of the source dataset used.

### B.2 False-positive gating failures on ImageNet

Figure 7 illustrates *false-positive* cases on ImageNet-1K, where bias-insensitive object queries are incorrectly gated on ( $\delta(q) = 1$ ), triggering unnecessary fairness regularization. These failures typically arise when the queried object strongly co-occurs with humans or human-like features in the top-1 retrieved anchor  $y^*$ . For instance, as shown in Figure 7, human figurines on a cake (left) or bystanders in the background (right) provide strong demographic signals that reduce the semantic-attribute discrepancy below the threshold  $\epsilon$ . This misleads the gate into treating the object query as bias-sensitive. Although our hybrid reward design minimizes semantic drift even when the gate is mistakenly active, these cases represent a computational inefficiency.