

➤ What is Model Reprogramming and Why?

- ❖ **Input-agnostic** perturbation to manipulate model behavior.
- ❖ Sharing similar ideas with text and visual prompting^[1, 2].
- ❖ Model weights are sometimes not alterable or inaccessible.

- large computational
- storage costs
- low data efficiency
- model privacy issues

➤ Motivations

Can an unfair model be reprogrammed to fair one?
 If so, why and how would it work?

➤ Contributions

- ❖ FairReprogram: a novel generic fairness-enhancing paradigm in a min-max fashion through model reprogramming for both NLP and CV tasks.
- ❖ Theoretically and empirically show **why and how** fairness can be promoted using an input-agnostic fairness trigger.
- ❖ Compared to retraining-based methods^[3], FairReprogram promotes model fairness with far less trade-off in accuracy across various NLP and CV datasets with in-the-wild biases.

➤ Fairness Metrics

- ❖ Equalized Odds (EO): $\hat{Y} \perp Z|Y, \sum_{z \in Z} (|FPR - FPR_z| + |FNR - FNR_z|)/2$,
- ❖ Demographic Parity (DP): $\hat{Y} \perp Z, \sum_{z \in Z} |p(\hat{Y} = 1) - p(\hat{Y} = 1|Z = z)|$,

➤ Related Work

[1] Li X L, Liang P. Prefix-tuning: Optimizing continuous prompts for generation.
 [2] Bahng H, et al. Visual Prompting: Modifying Pixel Space to Adapt Pre-trained Models.
 [3] Zhang B H, Lemoine B, Mitchell M. Mitigating unwanted biases with adversarial learning

➤ Fairness Reprogramming

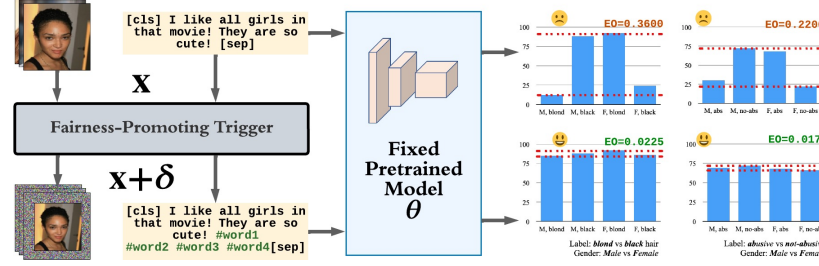


Figure 1. An example of fairness reprogramming in CV and NLP tasks. The input-agnostic trigger can promote fairness without altering the pretrained model.

❖ Optimization objective:

- Optimizable variables: fairness trigger: δ discriminator: ϕ
- Overall optimization objective:

$$\min_{\delta, \theta} \mathcal{L}_{util}(\mathcal{D}_{tune}, f^* \circ m) + \lambda \mathcal{L}_{fair}(\mathcal{D}_{tune}, f^* \circ m),$$

- Task utility loss:

$$\mathcal{L}_{util}(\mathcal{D}_{tune}, f^* \circ m) = \mathbb{E}_{\mathbf{X}, Y \sim \mathcal{D}_{tune}} [\text{CE}(Y, f^*(m(\mathbf{X})))]$$

- Discriminator-based fairness promoting loss:

$$\mathcal{L}_{fair}(\mathcal{D}_{tune}, f^* \circ m) = \max_{\phi} \mathbb{E}_{\mathbf{X}, Y, Z \sim \mathcal{D}_{tune}} [-\text{CE}(Z, d(f^*(m(\mathbf{X})), Y; \phi))]$$

❖ Why does fairness trigger work?

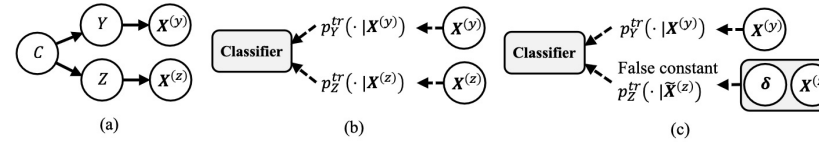


Figure 2. Illustration of why fairness trigger works. (a) Data generation. (b) Information flow from data to the classifier through the sufficient statistics. (c) Fairness trigger **strongly indicative of a demographic group** can confuse the classifier with a false demographic posterior, and thus preventing the classifier from using the correct demographic information.

➤ Experiment Results Highlights

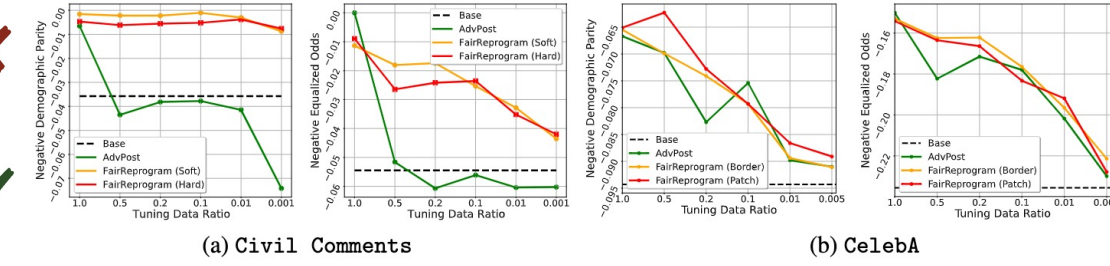


Figure 4. (Main Results) Results on (a) Civil Comments and (b) CelebA. We report the negative DP (left) and the negative EO (right) scores. We vary the trade-off parameter λ to record the performance. The closer a dot to the upper-right corner, the better the model is.

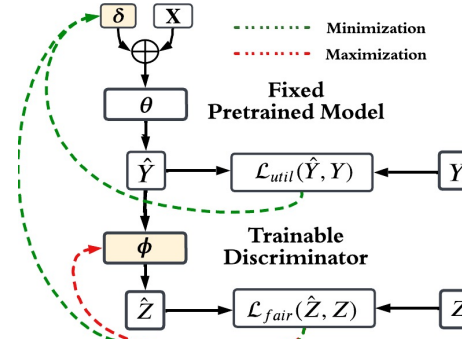


Figure 3. (Algorithm Illustration) An illustration of the FairReprogram algorithm pipelines formulated in a min-max fashion.

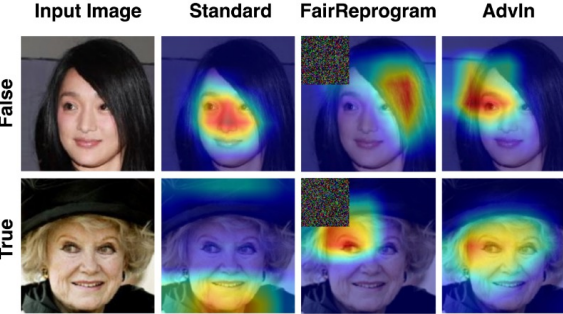


Figure 5. (Input Saliency Analysis) Gradient-based saliency map visualized with GradCAM. Highlighted zones (in red) are with major influence on predicted labels.

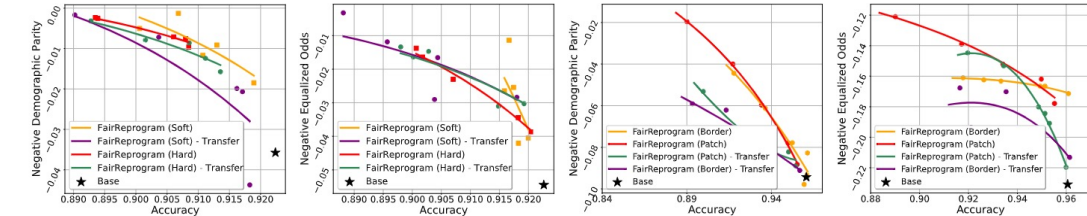


Figure 6. (Transfer Setting) We report negative DP (left) and negative EO (right) scores. The triggers are firstly trained in a BASE model. Then we evaluate the triggers based on another unseen BASE model.