# Exploring data bias: from the known to the unknown

# Prelude

- Bias is a reflection of **real-world structure**, but they <u>affect the fairness</u> of the model and may also lead to a <u>decrease in the model's generalizability</u> in real-world applications.

- Format：
  - ✓ spurious correlation
  - ✓ distribution shift — Texture bias, background bias scene bias
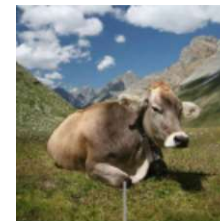  - ✓ Shortcut learning ....

- Way：
  - ✓ Inductive Bias [**No free lunch theorem**]

- Aspect：
  - Known bias
  - Unknown bias

Specific properties

Specific context

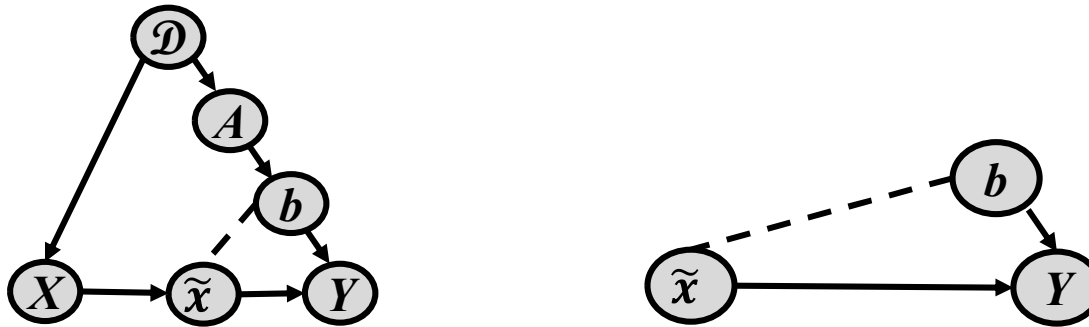Specific style

Image corruption

[1] Geirhos R et al. Shortcut learning in deep neural networks. (Nature Machine Intelligence, 2020)

$$P\left(y \mid do(x)\right) = \sum_{A} P\left(y \mid do(x), a\right) P(a \mid do(x)) = \sum_{A} P\left(y \mid x, a\right) P(a)$$

If $b$ is a complete bias from $X$, then the goal is to <u>make $b$ and $Y$ as independent as possible</u>

- Random experiment
- $P(Y|B) = P(Y)$
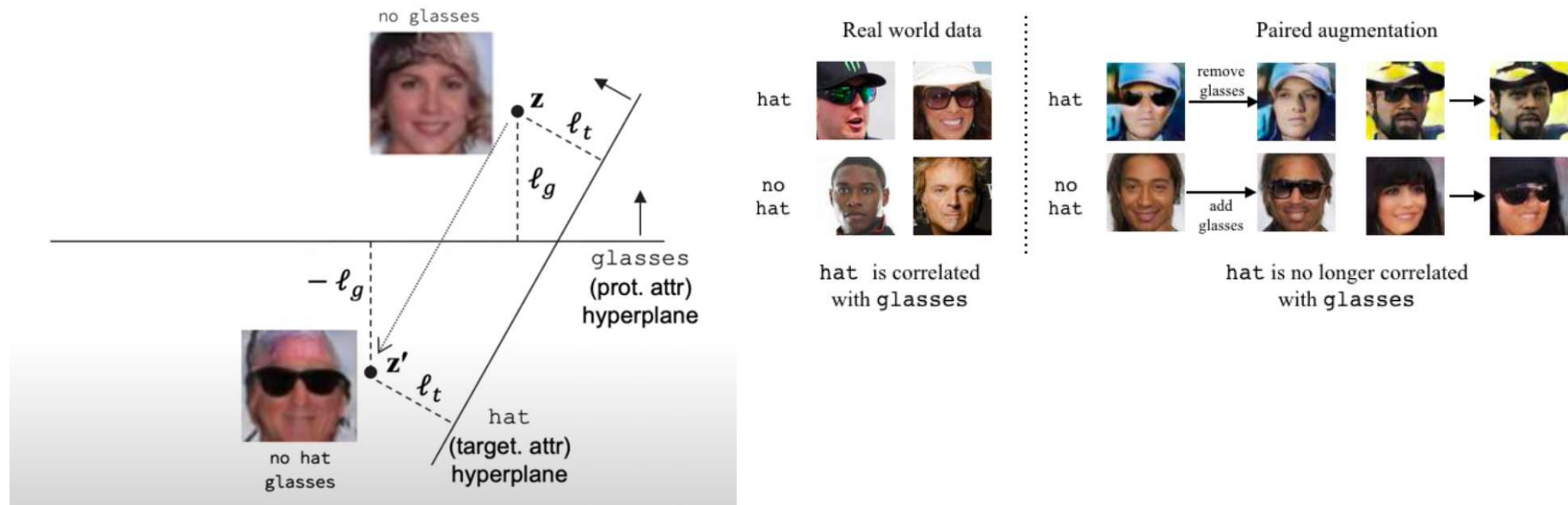- $\min I(\tilde{X}; B) \quad i.e. \min D_{KL}(p(\tilde{x}, b)||P(\tilde{x})P(b))$

# Fair data space

◦ Constructing a fair dataset by GAN-based data augmentation

- Assume: linearly differentiable in semantic properties

- Target: Learning interpretable operational directions, constructing complementary vectors

[2] Ramaswamy V V et al. Fair attribute classification through latent space de-biasing(CVPR2022).
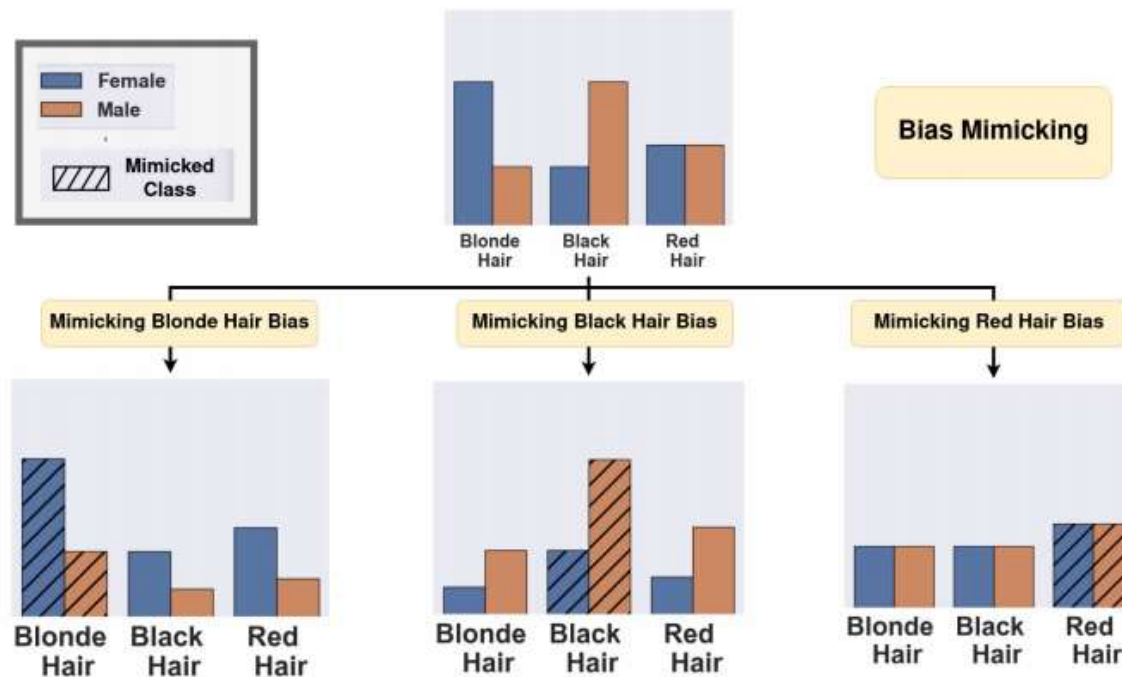
# Fair sampling distribution

Maan Qraitem[1], Kate Saenko[1,2], Bryan A. Plummer[1]
[1]Boston University   [2]MIT-IBM Watson AI Lab
{mqraitem,saenko,bplum}@bu.edu

- We want to guarantee $P_D(Y|B) = P_D(Y)$

- Bias Mimicking: Given class "c" :
  Ensure that is $P_D(B|Y=c)$
  "mimicked" in each other class

$$P_{d_c}(B = s|Y = c) = P_{d_c}(B = s|Y = c') \quad \forall s \in S.$$

$$P_D(B = s) = \sum_{c \in C} P_D(B = s|Y = c)P_D(Y = c)$$

$$P_D(B = s) = P_D(B = s|Y = c) \sum_{c' \in C} P_D(Y = c')$$

$$= P_D(B = s|Y = c)$$

[3] Qraitem M et al. Bias mimicking: A simple sampling approach for bias mitigation.(CVPR2023).

# Fair sampling distribution

Maan Qraitem[1], Kate Saenko[1,2], Bryan A. Plummer[1]
[1]Boston University   [2]MIT-IBM Watson AI Lab
{mqraitem,saenko,bplum}@bu.edu

How to Bias Mimick?

- constrain the solution space such that the solution retains **the most number** of samples.

- obtain the set of solutions using a linear program.

c: preserved class
c': mimicked class
s: bias
l: count

Set of biases

$$\max \quad \sum_s l_s^{c'}$$

$$\text{s.t.} \quad l_s^{c'} \leq |D_{c',s}| \qquad\qquad\qquad s \in S$$

$$\frac{l_s^{c'}}{\sum_s l_s^{c'}} = P_D(B = s | Y = c) \quad s \in S$$

[3] Qraitem M et al. Bias mimicking: A simple sampling approach for bias mitigation.(CVPR2023).

# Fair sampling distribution

Maan Qraitem[1], Kate Saenko[1,2], Bryan A. Plummer[1]
[1]Boston University  [2]MIT-IBM Watson AI Lab
{mqraitem,saenko,bplum}@bu.edu
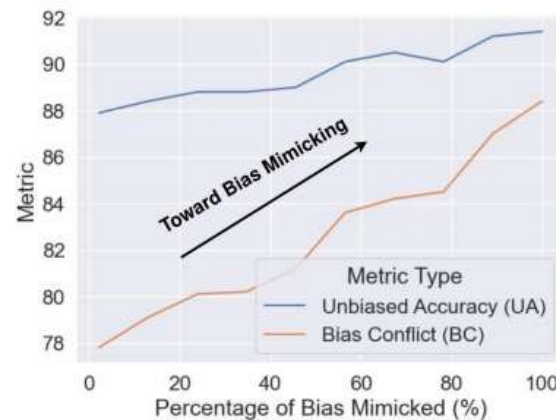
- How to train & inference?



✓ <u>Too many</u> additional parameters

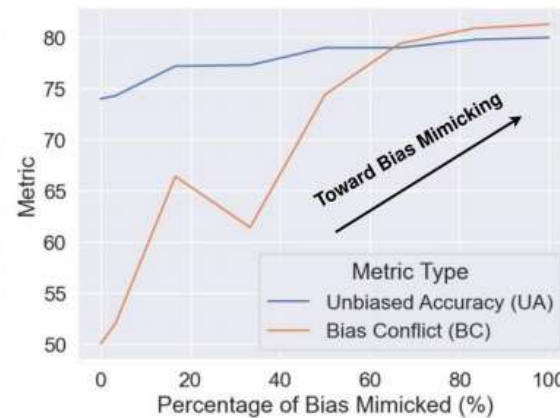✓ The scores may <u>not be calibrated</u> with respect to each other

[3] Qraitem M et al. Bias mimicking: A simple sampling approach for bias mitigation.(CVPR2023).

# Results

Maan Qraitem[1], Kate Saenko[1,2], Bryan A. Plummer[1]
[1]Boston University   [2]MIT-IBM Watson AI Lab
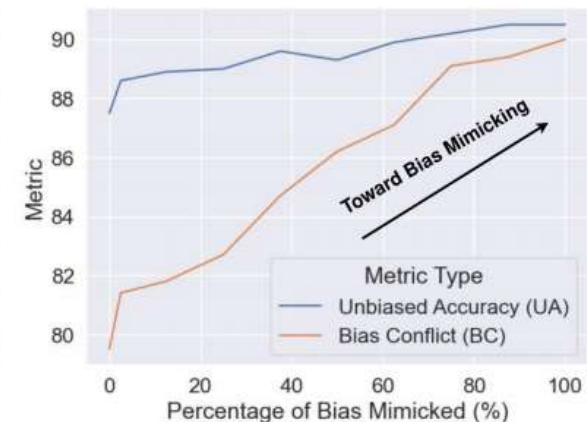{mqraitem,saenko,bplum}@bu.edu

- How sensitive is the model to the mimicking condition?

  ○ 0%: distribution remains the same
  ○ 100%: Complete bias mimicking



a) CelebA/Blonde          b) Utk-Face/Age          a) Utk-Face/Race
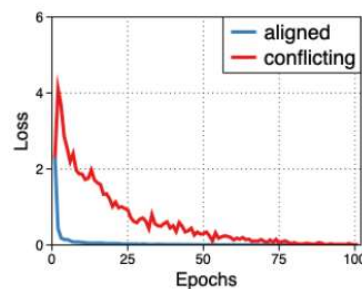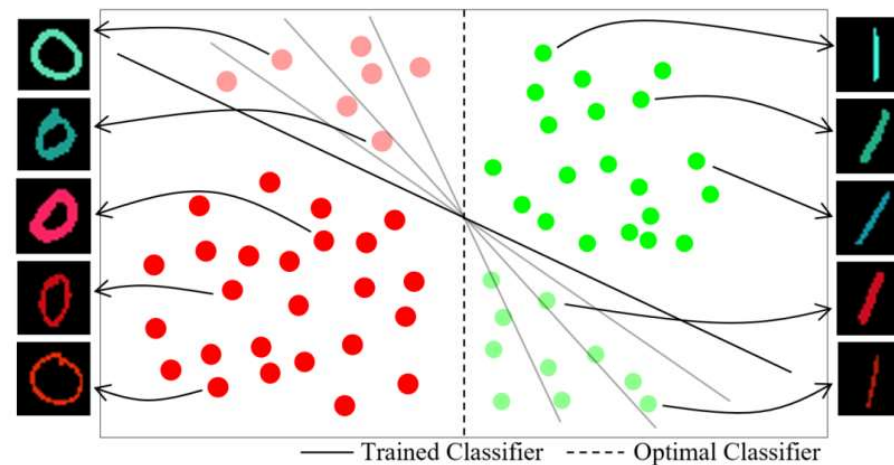
The performance of the model **is highly dependent** on the
degree of bias mimicry, particularly <u>the bias conflict groups</u>

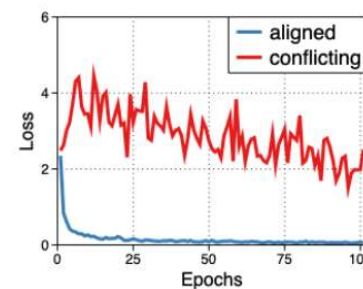[3] Qraitem M et al. Bias mimicking: A simple sampling approach for bias mitigation.(CVPR2023).

# Fair network

Byungju Kim[1]     Hyunwoo Kim[2]     Kyungsu Kim[3]     Sungjin Kim[3]     Junmo Kim[1]

School of Electrical Engineering, KAIST, South Korea[1]
Beijing Institute of Technology[2]
Samsung Research[3]

- How does the bias affect model training?

The **malignant bias** attributes **are easier to learn** than the original task[5]



Trained Classifier  ----- Optimal Classifier

(a) Colored MNIST, (Digit, Color)

(b) Corrupted CIFAR-10[1], (Object, Corruption)

[4] Nam J et al. Learning from failure: De-biasing classifier from biased classifier.(NeurIPS2020).
[5] Kim B et al. Learning not to learn: Training deep neural networks with biased data.(CVPR2020).
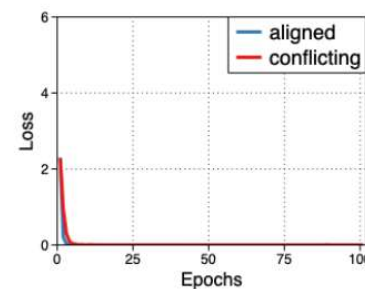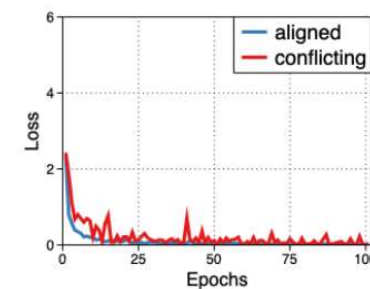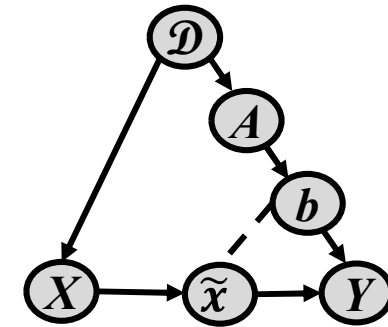
# Fair network

Byungju Kim[1]     Hyunwoo Kim[2]     Kyungsu Kim[3]     Sungjin Kim[3]     Junmo Kim[1]
School of Electrical Engineering, KAIST, South Korea[1]
Beijing Institute of Technology[2]
Samsung Research[3]

- How to unlearn the malignant bias？

  - **Minimize the mutual information** between feature embedding and target bias

    $$\min I(\tilde{X}; B) \quad i.e. \min D_{KL}(p(\tilde{x}, b) || P(\tilde{x})P(b))$$

  - malignant bias    $\mathcal{I}(b(X^{train}); Y) \gg \mathcal{I}(b(X^{test}); Y) \approx 0.$

  - Ultimate optimization goal

    $$\min_{\theta_f, \theta_g} \mathbb{E}_{\tilde{x} \sim P_X(\cdot)}[\mathcal{L}_c(y_{\tilde{x}}, g(f(\tilde{x})))] + \lambda \mathcal{I}(b(X); f(X))$$

- Issues: How to capture **the distribution of bias** in a network?

[5] Kim B et al. Learning not to learn: Training deep neural networks with biased data.(CVPR2020).

# Fair network

Byungju Kim[1]   Hyunwoo Kim[2]   Kyungsu Kim[3]   Sungjin Kim[3]   Junmo Kim[1]

School of Electrical Engineering, KAIST, South Korea[1]

Beijing Institute of Technology[2]

Samsung Research[3]

$$\min_{\theta_f,\theta_g} \mathbb{E}_{\tilde{x}\sim P_X(\cdot)}[\mathcal{L}_c(y_{\tilde{x}}, g(f(\tilde{x})))] + \lambda \mathcal{I}(b(X); f(X))$$

$$\mathcal{I}(b(X); f(X)) = H(b(X)) - H(b(X)|f(X)).$$

$Q$

$$\min_{\theta_f} \mathbb{E}_{\tilde{x}\sim P_X(\cdot)}[\mathbb{E}_{\tilde{b}\sim Q(\cdot|f(\tilde{x}))}[\log Q(\tilde{b}|f(\tilde{x}))]]$$

$$s.t. \quad Q(b(X)|f(X)) = P(b(X)|f(X)).$$

$$\mu D_{KL}(P(b(X)|f(X))||Q(b(X)|f(X)))$$

$$\mathcal{L}_\mathcal{B}(\theta_f,\theta_h) = \mathbb{E}_{\tilde{x}\sim P_X(\cdot)}[\mathcal{L}_c(b(\tilde{x}), h(f(\tilde{x})))]$$

$$\min_{\theta_f,\theta_g} \max_{\theta_h} \mathbb{E}_{\tilde{x}\sim P_X(\cdot)}[\mathcal{L}_c(y_{\tilde{x}}, g(f(\tilde{x})))$$
$$+ \lambda \mathbb{E}_{\tilde{b}\sim Q(\cdot|f(\tilde{x}))}[\log Q(\tilde{b}|f(\tilde{x}))]]$$
$$- \mu \mathcal{L}_\mathcal{B}(\theta_f,\theta_h)$$

[5] Kim B et al. Learning not to learn: Training deep neural networks with biased data.(CVPR2020).

The data distribution is $p$. Optimize $\theta$ to find the model $M(\theta)$ that makes the evaluation function $l$ optimal

➤ **Ignoring uncertainty**

➤ **Stochastic optimization  (Known** about $p$ completely**)**

   Optimize $\theta$ to maximize the expected value of $l$ under $p$

➤ **Robust optimization   (Only basic information** about $p$**)**

   Find the optimal solution when $p$ is the worst case

➤ **Distributionally robust optimization (**Extraordinarily **known** distributional **characteristi**

   Finding the worst distribution function that satisfies the
   uncertain parameter features

ERM:   $\hat{\theta}_{\text{ERM}} := \arg\min_{\theta \in \Theta} \ \mathbb{E}_{(x,y)\sim\hat{P}}[\ell(\theta; (x, y))]$

DRO:   $\min_{\theta \in \Theta} \left\{ \mathcal{R}(\theta) := \sup_{Q \in \mathcal{Q}} \mathbb{E}_{(x,y)\sim Q}[\ell(\theta; (x, y))] \right\}$    e.g. $f$-divergence based

# What's more? ———Uncertainty Modeling and Optimization

ERM:
$$\hat{\theta}_{\text{ERM}} := \arg\min_{\theta \in \Theta} \mathbb{E}_{(x,y) \sim \hat{P}}[\ell(\theta; (x,y))].$$

DRO:
$$\min_{\theta \in \Theta}\left\{\mathcal{R}(\theta) := \sup_{Q \in \mathcal{Q}} \mathbb{E}_{(x,y) \sim Q}[\ell(\theta; (x,y))]\right\}$$

$\Delta^m$ is an $(m-1)$-dimensional probabilistic simplex

Group DRO:
$$\hat{\theta}_{\text{DRO}} := \arg\min_{\theta \in \Theta}\left\{\hat{\mathcal{R}}(\theta) := \max_{g \in \mathcal{G}} \mathbb{E}_{(x,y) \sim \hat{P}_g}[\ell(\theta; (x,y))]\right\}$$

$$\mathcal{Q} := \{\textstyle\sum_{g=1}^{m} q_g P_g : q \in \Delta_m\}$$

$$\min_{\theta \in \Theta} \sup_{q \in \Delta_m} \sum_{g=1}^{m} q_g \mathbb{E}_{(x,y) \sim P_g}[\ell(\theta; (x,y))]$$

---

**Algorithm 1:** Online optimization algorithm for group DRO

**Input:** Step sizes $\eta_q, \eta_\theta$; $P_g$ for each $g \in \mathcal{G}$
Initialize $\theta^{(0)}$ and $q^{(0)}$
**for** $t = 1, \ldots, T$ **do**
    $g \sim \text{Uniform}(1, \ldots, m)$      // Choose a group $g$ at random
    $x, y \sim P_g$      // Sample $x, y$ from group $g$
    $q' \leftarrow q^{(t-1)}; q'_g \leftarrow q'_g \exp(\eta_q \ell(\theta^{(t-1)}; (x,y)))$      // Update weights for group $g$
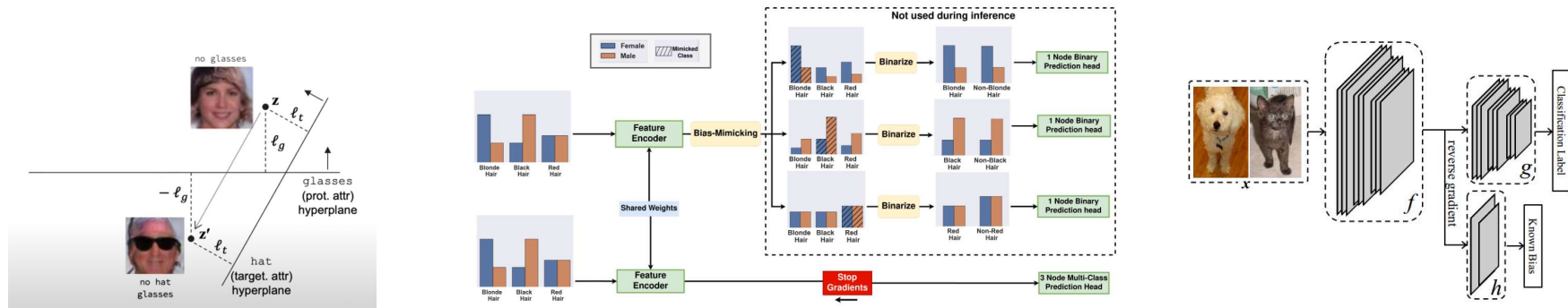    $q^{(t)} \leftarrow q' / \sum_{g'} q'_{g'}$      // Renormalize $q$
    $\theta^{(t)} \leftarrow \theta^{(t-1)} - \eta_\theta q_g^{(t)} \nabla \ell(\theta^{(t-1)}; (x,y))$      // Use $q$ to update $\theta$
**end**

[6] Sagawa S et al. Distributionally robust neural networks for group shifts: on the importance of regularization for worst-case generalization.(ICLR2020).

- Still need to <u>know a priori information</u> about the bias



- How to trick unknown bias?

  ○ Unsupervised learning

  ○ Self-supervised learning

  ○ MLLM-based

# Aspect 1 ——Clustering

Seonguk Seo[1]     Joon-Young Lee[3]     Bohyung Han[1,2]

[1]ECE & [1]ASRI & [1,2]IPAI, Seoul National University   [3]Adobe Research

{seonguk, bhhan}@snu.ac.kr    jolee@adobe.com

➢ **Observation**:

For a <u>particular attribute</u> (other than the target attribute), samples with **the same label** tend to have **similar representations** in the feature space of a fully trained model

➢ **Motivation** :

Define groups using **biased pseudo-attribute** information obtained through any **clustering algorithm** in the feature embedding space.

➢ **Optimization goals** :

$$\min_{\theta}\left\{\mathbb{R}_{\mathcal{K}}(\theta) := \mathbb{E}_{(\mathbf{x},y)\sim P}\left[\omega_{h((\mathbf{x},y);\tilde{\theta})}\ell((\mathbf{x},y);\theta)\right]\right\}$$

cluster membership

[7] Seo S, Lee J Y, Han B. Unsupervised learning of debiased representations with pseudo-attributes.(CVPR2022).

# Aspect 1——Clustering

Seonguk Seo[1]     Joon-Young Lee[3]     Bohyung Han[1,2]

[1]ECE & [1]ASRI & [1,2]IPAI, Seoul National University    [3]Adobe Research

{seonguk, bhhan}@snu.ac.kr    jolee@adobe.com

Group DRO:

$$\min_{\theta \in \Theta} \sup_{q \in \Delta_m} \sum_{g=1}^{m} q_g \mathbb{E}_{(x,y) \sim P_g} [\ell(\theta; (x,y))]$$

Revised version

$$\min_{\theta} \left\{ \mathbb{R}_{\mathcal{K}}(\theta) := \mathbb{E}_{(\mathbf{x},y) \sim P} \left[ \omega_{h((\mathbf{x},y);\tilde{\theta})} \ell((\mathbf{x},y);\theta) \right] \right\}$$

$$\omega_k = \frac{\mathbb{E}_{(\mathbf{x},y) \sim P_k}[\ell((\mathbf{x},y);\theta)]}{N_k}$$

$$= \frac{\mathbb{E}_{(\mathbf{x},y) \sim P}[\ell((\mathbf{x},y);\theta) \mid h((\mathbf{x},y);\tilde{\theta}) = k]}{\sum_i \mathbb{1}(h((\mathbf{x}_i,y_i);\tilde{\theta}) = k)}$$

**Step1：** **Clustering** the training samples in the feature embedding space of <u>the fully optimized base model</u>, assuming that each cluster corresponds to a bias pseudo-attribute

**Step2：** **Weighting** between groups considering the size and average difficulty of each cluster

**Algorithm 1: Debiasing with bias pseudo-attribute**

1 **Require:** step size $\eta_\theta$, momentum $m$, training steps $T$, batch size $B$, the number of clusters $K$

2 **Base model:**
3 **Initialize** $\tilde{\theta}$
4 **for** $t = 1, ..., T$ **do**
5      Sample $(\mathbf{x}_i, y_i) \sim P$ for $i = 1, ..., B$;
6      $\tilde{\theta} \leftarrow \tilde{\theta} - \eta_\theta \sum_{i=1}^{B} \nabla \ell((\mathbf{x}_i, y_i); \tilde{\theta})$;
7 **end**

8 **for** $k = 1, ..., K$ **do**
9      $P_k = \{(\mathbf{x}_n, y_n) \mid h((\mathbf{x}_n, y_n); \tilde{\theta}) = k \text{ for all } n\}$;
10      $N_k = |P_k|$;
11 **end**

12 **Target model:**
13 **Initialize** $\theta$ and $\omega_k$ for $k = 1, ..., K$
14 **for** $t = 1, ..., T$ **do**
15      $\omega_k \leftarrow (1-m)\omega_k + \frac{m}{N_k}\mathbb{E}_{(\mathbf{x},y) \sim P_k}[\ell((\mathbf{x},y); \theta)]$
                     for $k = 1, ..., K$;
17      Sample $(\mathbf{x}_i, y_i) \sim P$ for $i = 1, ..., B$;
18      $\alpha_i = \omega_{h((\mathbf{x}_i, y_i); \tilde{\theta})}$;
19      $\overline{\alpha}_i = \alpha_i / \sum_{i=1}^{B} \alpha_i$;
20      $\theta \leftarrow \theta - \eta_\theta \sum_{i=1}^{B} \overline{\alpha}_i \nabla \ell((\mathbf{x}_i, y_i); \theta)$;
21 **end**

[7] Seo S, Lee J Y, Han B. Unsupervised learning of debiased representations with pseudo-attributes.(CVPR2022).

# Aspect 2—Contrastive Approach

Correct-N-Contrast: A Contrastive Approach for
Improving Robustness to Spurious Correlations

Michael Zhang[†], Nimit S. Sohoni[†], Hongyang R. Zhang[‡], Chelsea Finn[†] & Christopher Ré[†]
[†]Stanford University, [‡]Northeastern University

{mzhang,nims,hongyang,cbfinn,chrismre}@cs.stanford.edu
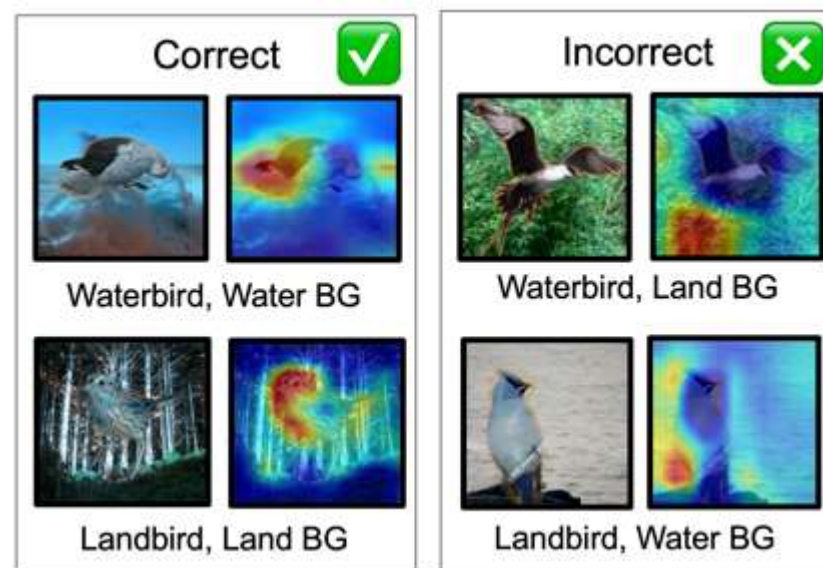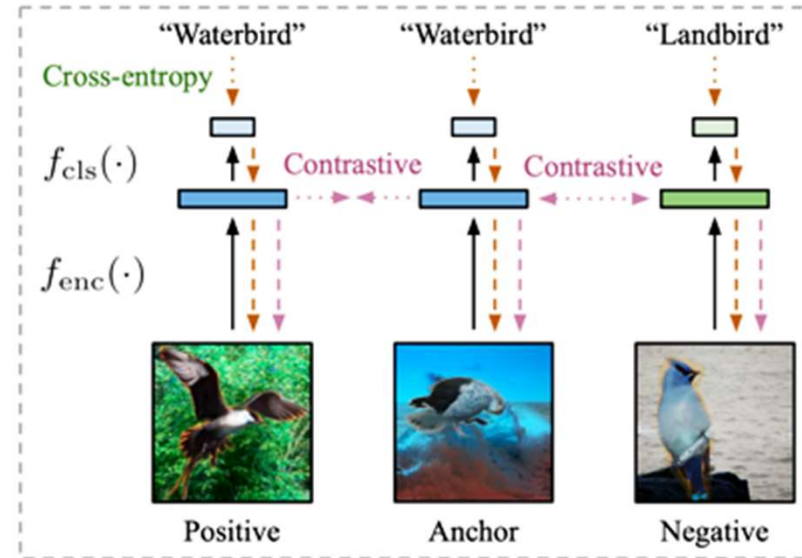
Corresponding author

➢ **Observation**:

The **worst population accuracy** of the neural network in the false correlation case
is closely related to its representation - i.e., the output of its last hidden layer - only
to **the extent that it relies on true labels** rather than false attributes

➢ **Motivation** :

By **improving alignment** while keeping the class mean error low, we can
help improve the worst group error for the class

➢ **New sampling scheme**

Samples <u>with the same category but
different spurious attributes </u>as **different
"views"** (anchors and **positive samples**)
of the same category, and samples
**negative samples** of data points with the
<u>same inferred spurious attributes </u>but
<u>different categories</u>



Correct ✓

Waterbird, Water BG

Landbird, Land BG

Incorrect ✗

Waterbird, Land BG

Landbird, Water BG

[8] Zhang M et al. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations.(CVPR2023).

# Aspect 2—Contrastive Approach

Correct-N-Contrast: A Contrastive Approach for Improving Robustness to Spurious Correlations

Michael Zhang[†], Nimit S. Sohoni[†], Hongyang R. Zhang[‡], Chelsea Finn[†] & Christopher Ré[†]
[†]Stanford University, [‡]Northeastern University

{mzhang,nims,hongyang,cbfinn,chrismre}@cs.stanford.edu
Corresponding author

**Stage 1**: ERM Training

o Train model with ERM to predict ground-truth labels from data

o Collect predictions of the trained ERM model on the *training data*

o Get contrastive batches from predictions:

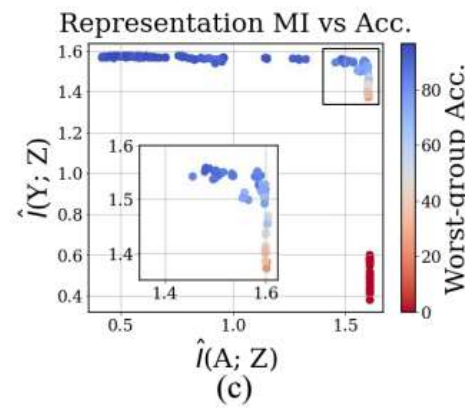anchor    positive    negative

$y =$ "Waterbird"   "Waterbird"   "Landbird"

$\hat{y} =$ "Waterbird" ✓   "Landbird" ✗   "Waterbird" ✗

Learn similar representations

**Stage 2**: Supervised Contrastive Learning

o Train a new model with supervised contrastive loss + classification loss using contrastive batches

"Waterbird"   "Waterbird"   "Landbird"

Cross-entropy

$f_{cls}(\cdot)$   Contrastive   Contrastive

$f_{enc}(\cdot)$

Positive    Anchor    Negative

$$\hat{\mathcal{L}}(f_\theta; x, y) = \lambda \hat{\mathcal{L}}_{con}^{sup}(f_{enc}; x, y) + (1 - \lambda)\hat{\mathcal{L}}_{cross}(f_\theta; x, y)$$

$$\hat{\mathcal{L}}_{con}^{sup}(f_{enc}) = \hat{\mathcal{L}}_{con}^{sup}\left(x_1, \{x_m^+\}_{m=1}^M, \{x_n^-\}_{n=1}^N; f_{enc}\right) + \hat{\mathcal{L}}_{con}^{sup}\left(x_1^+, \{x_i\}_{i=1}^M, \{x_n'^-\}_{n=1}^N; f_{enc}\right) \qquad -\frac{1}{M}\sum_{m=1}^M \log \frac{\exp(z_1^\top z_m^+/\tau)}{\sum_{m=1}^M \exp(z_1^\top z_m^+/\tau) + \sum_{n=1}^N \exp(z_1^\top z_n^+/\tau)}$$

[8] Zhang M et al. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations.(CVPR2023).
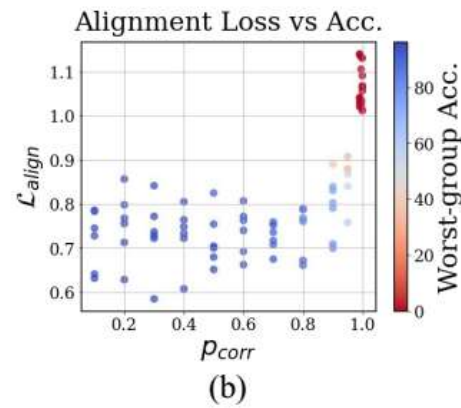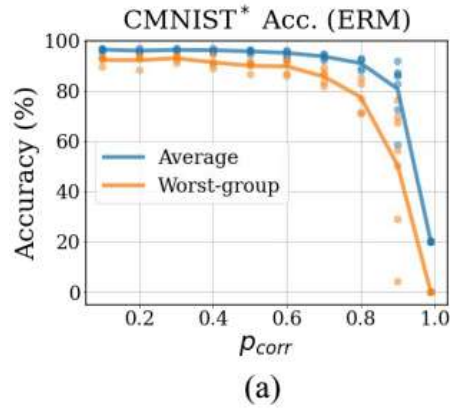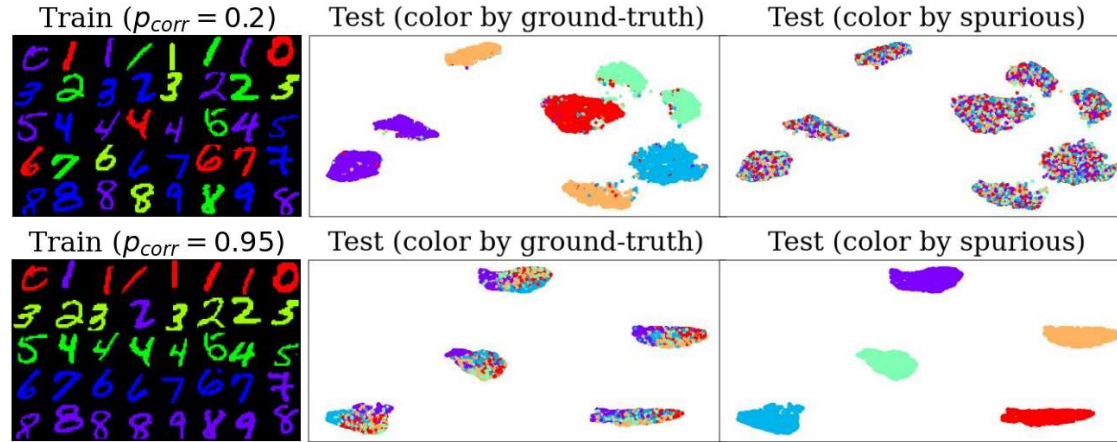
# Aspect 2—Theoretical proof

Correct-N-Contrast: A Contrastive Approach for Improving Robustness to Spurious Correlations

Michael Zhang[†], Nimit S. Sohoni[†], Hongyang R. Zhang[‡], Chelsea Finn[†] & Christopher Ré[†]

[†]Stanford University, [‡]Northeastern University

{mzhang,nims,hongyang,cbfinn,chrismre}@cs.stanford.edu

Corresponding author

➢ **Motivation**:

**Reducing alignment** losses can **narrow the gap** between the worst and average group losses

➢ **Metrics design**:

○ Alignment loss

$$\hat{\mathcal{L}}_{\text{align}}(f_{\text{enc}}; g, g') = \frac{1}{|G|} \frac{1}{|G'|} \sum_{(x,y,a) \in G} \sum_{(x',y,a') \in G'} \|f_{\text{enc}}(x) - f_{\text{enc}}(x')\|_2$$

Degree to which samples with the same class but different spurious attributes map to neighborhood vectors

○ Mutual information

$$\hat{I}(Y; Z) = \frac{1}{|Z|} \sum_{z \in Z} \sum_{y \in Y} p(y \mid z) \log \frac{p(y \mid z)}{p(y)}$$

$$D_{KL}\{P(y|z)||P(y)\}$$

Approximate MI between model-learned representations and labels



$$f_\theta : \mathcal{X} \mapsto \mathcal{Y}$$

[8] Zhang M et al. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations.(CVPR2023).

# Aspect 2—Theoretical proof

Michael Zhang[†], Nimit S. Sohoni[†], Hongyang R. Zhang[‡], Chelsea Finn[†] & Christopher Ré[†]
[†]Stanford University, [‡]Northeastern University

{mzhang,nims,hongyang,cbfinn,chrismre}@cs.stanford.edu

Corresponding author

[8] Zhang M et al. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations.(CVPR2023).

# Aspect 2—Theoretical proof

Michael Zhang[†], Nimit S. Sohoni[†], Hongyang R. Zhang[‡], Chelsea Finn[†] & Christopher Ré[†]
[†]Stanford University, [‡]Northeastern University

{mzhang,nims,hongyang,cbfinn,chrismre}@cs.stanford.edu

Corresponding author

Let $\mathcal{L}_{\text{wg}}(f_\theta; y)$ be the worst-group loss among groups in $\mathcal{G}_y$:

$$\mathcal{L}_{\text{wg}}(f_\theta; y) := \max_{g \in \mathcal{G}_y} \mathbb{E}_{(x,\tilde{y},a) \sim P_g} [\ell(f_\theta(x), \tilde{y})].$$

Let $\mathcal{L}_{\text{avg}}(f_\theta; y)$ be the average loss among groups in $\mathcal{G}_y$:

$$\mathcal{L}_{\text{avg}}(f_\theta; y) := \mathbb{E}_{(x,\tilde{y},a) \sim P:\forall a \in \mathcal{A}} [\ell(f_\theta(x), \tilde{y})].$$

Additionally, let $\mathcal{L}_{\text{align}}(f_\theta; y)$ be the largest cross-group alignment loss among groups in $\mathcal{G}_y$:

$$\hat{\mathcal{L}}_{\text{align}}(f_\theta; y) := \max_{g \in \mathcal{G}_y, g' \in \mathcal{G}_y: g \neq g'} \hat{\mathcal{L}}_{\text{align}}(f_{\text{enc}}; g, g').$$

**Theorem 3.1.** *In the setting described above, suppose the weight matrix of the linear classification layer $W$ satisfies $\|W\|_2 \leq B$, for some $B > 0$. Suppose the loss function $\ell(x, y)$ is $C_1$-Lipschitz in $x$ and bounded from above by $C_2$, for some $C_1 >$ and $C_2 > 0$. Let $n_g$ be the size of any group $g \in \mathcal{G}$ in the training set. Then, for any $\delta > 0$, with probability $1 - \delta$, the following holds for any $y \in \mathcal{Y}$:*

$$\mathcal{L}_{wg}(f_\theta; y) - \mathcal{L}_{avg}(f_\theta; y)$$

$$\leq BC_1 \cdot \hat{\mathcal{L}}_{align}(f_\theta; y) + \max_{g \in \mathcal{G}_y} C_2 \sqrt{8 \log(|\mathcal{G}_y|/\delta)/n_g}.$$

The **upper bound** of the <u>loss gap between the worst and average groups</u> is **linearly and positively** correlated with the <u>largest cross-group alignment loss</u>

[8] Zhang M et al. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations.(CVPR2023).

Michael Zhang[†], Nimit S. Sohoni[†], Hongyang R. Zhang[‡], Chelsea Finn[†] & Christopher Ré[†]
[†]Stanford University, [‡]Northeastern University

{mzhang,nims,hongyang,cbfinn,chrismre}@cs.stanford.edu
Corresponding author

[8] Zhang M et al. Correct-n-contrast: A contrastive approach for improving robustness to spurious correlations.(CVPR2023).

# Aspect 3—Keyword Explanation

Younghyun Kim[*1]    Sangwoo Mo[*2]    Minkyu Kim[3]    Kyungmin Lee[1]    Jaeho Lee[4]    Jinwoo Shin

[1]KAIST    [2]University of Michigan    [3]KRAFTON    [4]POSTECH

younghyun.kim@kaist.ac.kr    swmo@umich.edu

➢ **Motivation**:   <u>Unknown visual bias can't be interpreted</u>

• Visualized spurious features that are **not human-readable**

• Thus, they are hard to be directly utilized for debiasing



(a) **class:** band aid, **spurious feature:** fingers, **-41.54%**  (b) **class:** space bar, **spurious feature:** keys, **-46.15%**  (c) **class:** plate, **spurious feature:** food, **-32.31%**  (d) **class:** butterfly, **spurious feature:** flowers, **-21.54%**  (e) **class:** potter's wheel, **spurious feature:** vase, **-21.54%**

[9] Kim Y et al. Discovering and Mitigating Visual Biases through Keyword Explanation.(CVPR2024).

# Aspect 3—Keyword Explanation

Younghyun Kim[*1]  Sangwoo Mo[*2]  Minkyu Kim[3]  Kyungmin Lee[1]  Jaeho Lee[4]  Jinwoo Shin
[1]KAIST  [2]University of Michigan  [3]KRAFTON  [4]POSTECH
younghyun.kim@kaist.ac.kr  swmo@umich.edu

➤ **Method**:  B2T: Bias-to-text



**Step1**: use ClipCap as our default **captioning model**



**Step2**: apply the YAKE algorithm to **extract keywords**

Text Preprocessing (Segmentation) --> Feature Extraction -->
Individual Word Weight Calculation --> Candidate Keyword Generation

[9] Kim Y et al. Discovering and Mitigating Visual Biases through Keyword Explanation.(CVPR2024).

# Aspect 3—Keyword Explanation

Younghyun Kim[*1]   Sangwoo Mo[*2]   Minkyu Kim[3]   Kyungmin Lee[1]   Jaeho Lee[4]   Jinwoo Shin
[1]KAIST   [2]University of Michigan   [3]KRAFTON   [4]POSTECH
younghyun.kim@kaist.ac.kr  swmo@umich.edu

➢ **Method**:   B2T: Bias-to-text



**Step3**: **verify** that keywords represent bias by CLIP score

To **measures the similarity** between the <u>keywords</u> and the <u>incorrectly predicted images</u>

$$s_{\text{CLIP}}(a; \mathcal{D}) := \text{sim}(a, \mathcal{D}_{\text{wrong}}) - \text{sim}(a, \mathcal{D}_{\text{correct}}). \qquad \text{sim}(a, \mathcal{D}) := \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} f_{\text{image}}(x) f_{\text{text}}(a).$$



**Effect of the CLIP score (waterbird class)**

[9] Kim Y et al. Discovering and Mitigating Visual Biases through Keyword Explanation.(CVPR2024).

# Aspect 3—Keyword Explanation

Younghyun Kim[*1]  Sangwoo Mo[*2]  Minkyu Kim[3]  Kyungmin Lee[1]  Jaeho Lee[4]  Jinwoo Shin

[1]KAIST  [2]University of Michigan  [3]KRAFTON  [4]POSTECH

younghyun.kim@kaist.ac.kr  swmo@umich.edu

We first extract B2T keywords, then use them to various applications:
- ➢ **Debiased training**          ➢ Model comparison
- ➢ CLIP prompting               ➢ Label diagnosis



☐ B2T discovers spurious correlations and distributions shifts

- e.g.)
"man" for CelebA blond、
"forest" and "ocean" for Waterbirds、
"illustration" and "drawing" for IN-R、
"snow" and "window" for IN-C

[9] Kim Y et al. Discovering and Mitigating Visual Biases through Keyword Explanation.(CVPR2024).

25

# Aspect 3—Keyword Explanation

Younghyun Kim[*1]  Sangwoo Mo[*2]  Minkyu Kim[3]  Kyungmin Lee[1]  Jaeho Lee[4]  Jinwoo Shin
[1]KAIST  [2]University of Michigan  [3]KRAFTON  [4]POSTECH
younghyun.kim@kaist.ac.kr  swmo@umich.edu

We first extract B2T keywords, then use them to various applications:

- ➢ Debiased training
- ➢ **CLIP prompting**
- ➢ Model comparison
- ➢ Label diagnosis

Table 8. Prompt designs for debiaisng zero-shot classifiers.

| Dataset | Dataset-wise Template | Class Name |
|---|---|---|
| CelebA | • [class name]<br>• [class name] man<br>• [class name] player<br>• [class name] person<br>• [class name] artist<br>• [class name] comedy<br>• [class name] film<br>• [class name] actor<br>• [class name] face | 1. Blond<br>• blond hair<br>• celebrity of blond hair<br><br>2. Non blond<br>• non blond hair<br>• celebrity of non blond hair |
| Waterbirds | • [class name]<br>• [class name] on the forest<br>• [class name] with woods<br>• [class name] on a tree<br>• [class name] on a branch<br>• [class name] in the forest<br>• [class name] on the tree<br>• [class name] on the ocean<br>• [class name] on a beach<br>• [class name] on the lake<br>• [class name] with a surfer<br>• [class name] on the water<br>• [class name] on a boat<br>• [class name] on the dock<br>• [class name] on the rocks<br>• [class name] in the sunset<br>• [class name] with a kite<br>• [class name] on the sky<br>• [class name] is on flight<br>• [class name] is on flies | 1. Landbird<br>• landbird<br><br>2. Waterbird<br>• waterbird |

☐ **Modify** the cue by adding a keyword, e.g., "[class]'s photo" in [group], where the keyword represents the name of the group

- **Obtaining** the average prompts embedding for a class in all groups
- **Comparing** broader class embeddings for image classification

[9] Kim Y et al. Discovering and Mitigating Visual Biases through Keyword Explanation.(CVPR2024).

# Aspect 3—Keyword Explanation

Younghyun Kim[*1]   Sangwoo Mo[*2]   Minkyu Kim[3]   Kyungmin Lee[1]   Jaeho Lee[4]   Jinwoo Shin
[1]KAIST   [2]University of Michigan   [3]KRAFTON   [4]POSTECH
younghyun.kim@kaist.ac.kr   swmo@umich.edu

We first extract B2T keywords, then use them to various applications:
- ➤ Debiased training
- ➤ CLIP prompting
- ➤ **Model comparison**
- ➤ Label diagnosis

| Keyword | Work | | Supermarket | |
|---|---|---|---|---|
| Samples |  | | | |
| ViT-B | O | O | O | O |
| RN50 | O | X | O | X |
| Actual (RN50) | dumbbell | dumbbell | shopping basket | shopping basket |
| Pred (RN50) | dumbbell | horizontal bar | shopping basket | grocery store |
| Caption | a set of dumbbells with weights. | person **works** out in the gym. | a basket full of food. | woman shopping in a **supermarket**. |

☐ Bias keywords can be used to <u>analyze and compare different classifiers</u> based on their keywords

e.g.) architecture: ResNet vs. ViT

[9] Kim Y et al. Discovering and Mitigating Visual Biases through Keyword Explanation.(CVPR2024).

# Aspect 3—Keyword Explanation

Younghyun Kim[*1]   Sangwoo Mo[*2]   Minkyu Kim[3]   Kyungmin Lee[1]   Jaeho Lee[4]   Jinwoo Shin

[1]KAIST   [2]University of Michigan   [3]KRAFTON   [4]POSTECH

younghyun.kim@kaist.ac.kr   swmo@umich.edu

We first extract B2T keywords, then use them to various applications:

- ➤ Debiased training
- ➤ CLIP prompting
- ➤ Model comparison
- ➤ **Label diagnosis**



☐ B2T can diagnose common labeling errors, such as mislabeling and label ambiguities

[9] Kim Y et al. Discovering and Mitigating Visual Biases through Keyword Explanation.(CVPR2024).

南京航空航天大学
Nanjing University of Aeronautics and Astronautics

模式分析与机器智能
工业和信息化部重点实验室
MIIT Key Laboratory of
Pattern Analysis & Machine Intelligence

# THANKS