# Learning from Failure:

# Training Debiased Classficier from Biased Classifier

Jihyeon Lee, Oct 20th, 2020

# Motivation

## Learning debias without supervision

- Neural networks often learn to make predictions using the unintended decision rule

- Recent approaches focus on how to utilize various types of human supervision effectively

- Another line of research focuses on developing algorithms tailored to a domain-specific type of bias in the target dataset

- An approach to train a debiased classifier without relying on such expensive supervision is warranted



$$\min_{f} \left\{ \mathcal{L}(f) + \lambda \max_{g} \left( \mathrm{HSIC}_1(f,g) - \lambda_g \mathcal{L}(g) \right) \right\}$$

# Contribution

**Learning from Failure(LfF) – failure based debiasing scheme**

1. Propose debiasing network by simultaneously train two neural net, one to be biased and the other to be debiased

2. Classifier learns to fit samples aligned with the bias during the early stage of training and learns samples conflicting with

   the bias later

3. LfF Does not require expensive supervision on the bias or bias-tailored training technique

4. Show the effectiveness of LfF on various biased datasets
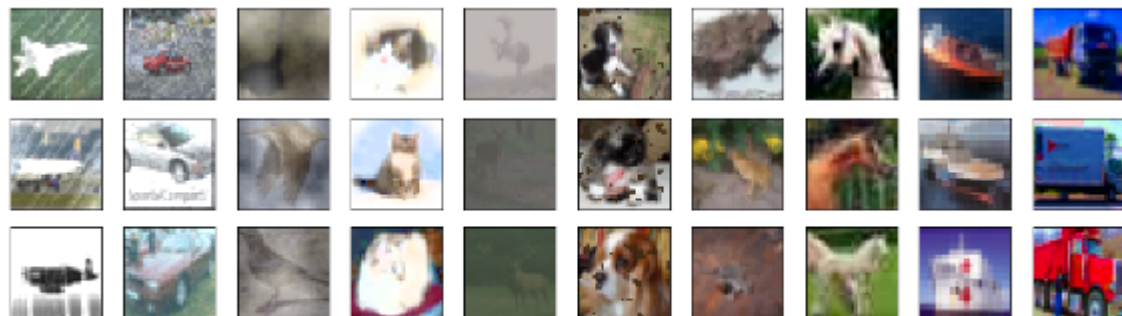
5. Constructed baised action recognition dataset (BAR)

# Biased Dataset

**Definition of bias**

Consider a dataset $\mathcal{D}$ where each input $x$ can be represented by a set of (possibly latent) *attributes* $\{a_1, \ldots, a_k\}$ for $a_i \in \mathcal{A}_i$ that describes the input. The goal is to train a predictor $f$ that belongs to a set of *intended decision rules* $\mathcal{F}_t$, consisting of decision rules that correctly predict the *target attribute* $y = a_t \in \mathcal{A}_t$. We say that a dataset $\mathcal{D}$ is *biased*, if (a) there exists another attribute $a_b \neq y$ that is highly correlated to the target attribute $y$ (i.e., $H(y|a_b) \approx 0$), and (b) one can settle an *unintended decision rule* $g_b \notin \mathcal{F}_t$ that correctly classifies $a_b$. We denote such an attribute $a_b$ by a *bias attribute*. In biased datasets with a bias attribute $a_b$, we say that a sample is *bias-aligned* whenever it can be
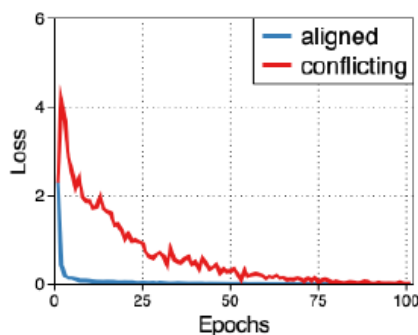


(a) Colored MNIST

(b) Corrupted CIFAR-10[1]
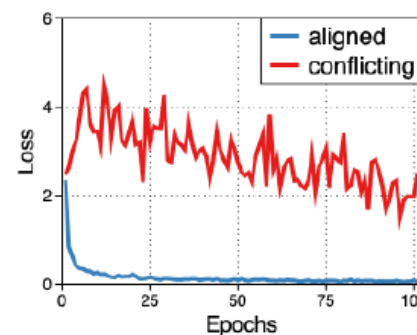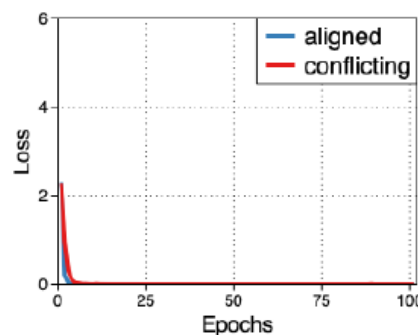
# Biased Dataset

**Two types of bias: malignant and benign**

- Biased dataset does not necessarily lead the model to learn the unintended decision rule

- The bias negatively affects the model only when the bias attribute is "easier" to learn than the target attribute
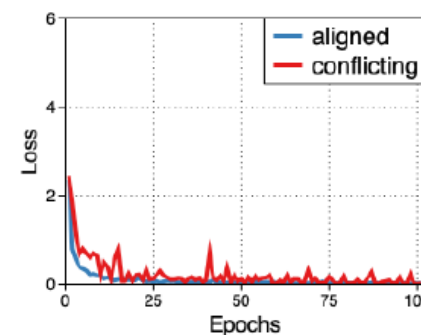
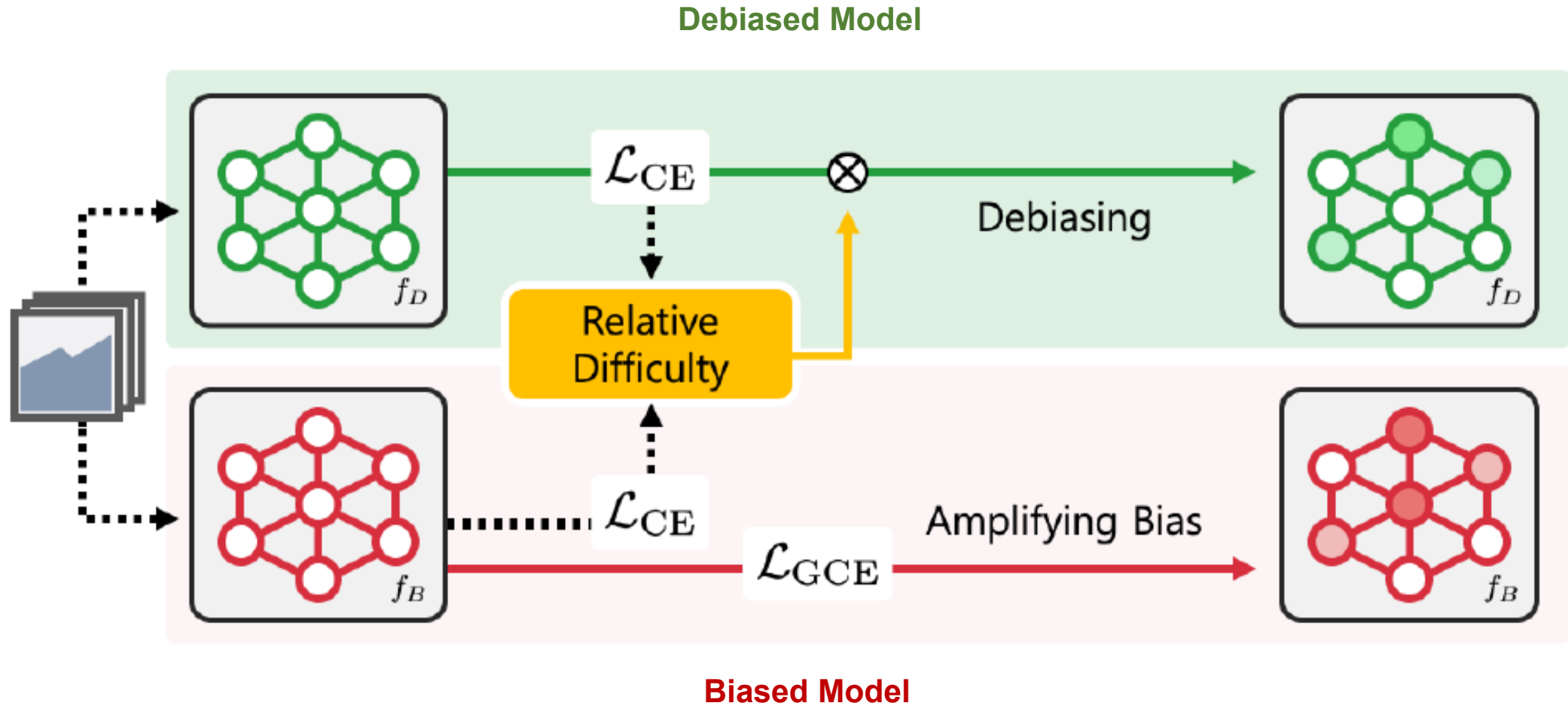| Dataset | Target | Bias | Accuracy | Accuracy* | Relative drop |
|---------|--------|------|----------|-----------|---------------|
| Colored MNIST | Color | Digit | $99.97_{\pm0.04}$ | $100.0_{\pm0.00}$ | -0.03% |
| | Digit | Color | $50.34_{\pm0.16}$ | $96.41_{\pm0.07}$ | -47.79% |
| Corrupted CIFAR-10[1] | Corruption | Object | $98.34_{\pm0.26}$ | $99.62_{\pm0.03}$ | -1.28% |
| | Object | Corruption | $22.72_{\pm0.87}$ | $80.00_{\pm0.01}$ | -71.60% |
| Corrupted CIFAR-10[2] | Corruption | Object | $98.64_{\pm0.20}$ | $99.80_{\pm0.01}$ | -1.16% |
| | Object | Corruption | $21.07_{\pm0.29}$ | $79.65_{\pm0.11}$ | -73.56% |



(a) Colored MNIST, (Digit, Color)    (b) Corrupted CIFAR-10[1], (Object, Corruption)
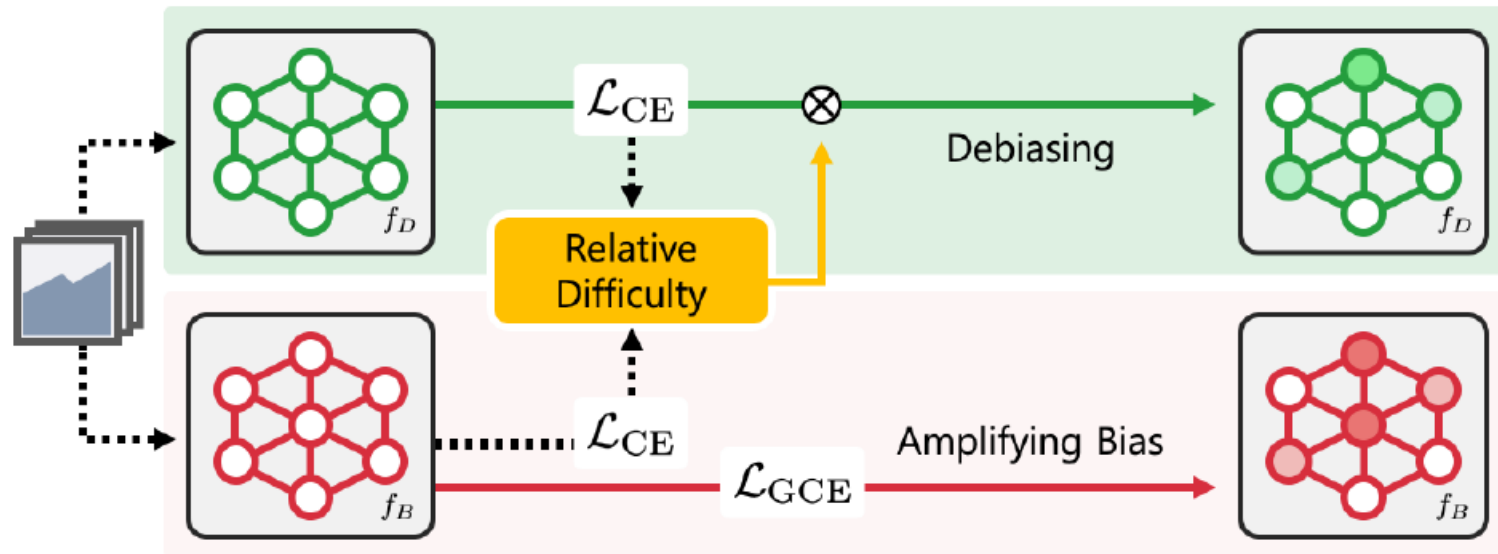
# Model

## Model Structure

# Model

## Training biased model – amplify the bias

- GCE loss up-weights the gradient of the CE loss for the samples with a high probability p of predicting the correct target

- emphasizing the 'easier' samples with the strong agreement between softmax output and the target

$$\text{GCE}(p(x;\theta),y) = \frac{1 - p_y(x;\theta)^q}{q} \qquad \lim_{q \to 0} \frac{1-p^q}{q} = -\log p$$

$$\frac{\partial \text{GCE}(p,y)}{\partial \theta} = p_y^q \frac{\partial \text{CE}(p,y)}{\partial \theta}$$

# Model

## Training de-biased model – relative difficulty score

- Train simultaneously with the samples using the CE loss re-weighted by score

- Score indicates how much each sample is likely to be bias-conflicting

- For bias-alinged samples, biased model $f_B$ tends to have smaller loss compare to debiased model $f_D$ at the early stage, therefore having small weight for training debiased model

$$\mathcal{W}(x) = \frac{\text{CE}(f_B(x), y)}{\text{CE}(f_B(x), y) + \text{CE}(f_D(x), y)}$$

# Experiments

## Controlled experiments

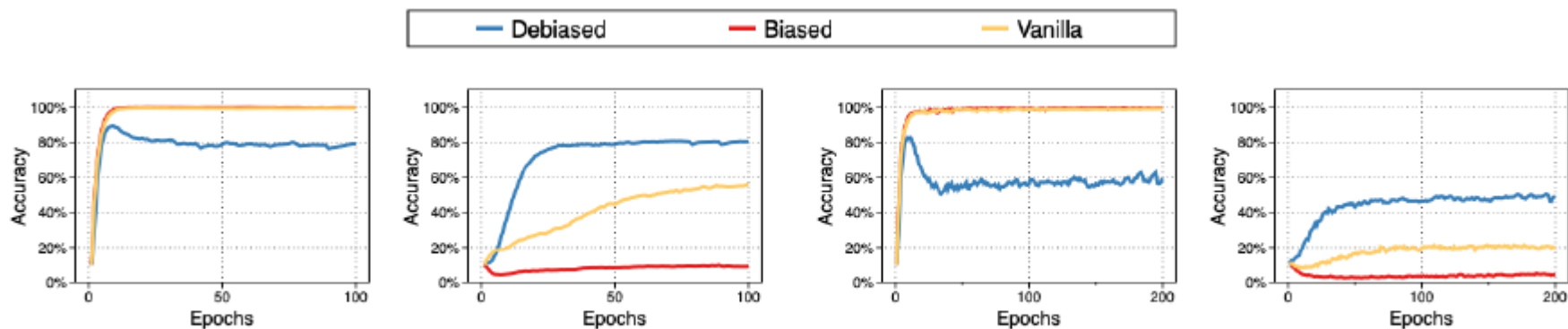| Dataset | Ratio (%) | Vanilla ○ | Ours ○ | HEX ◑ | REPAIR ● | Group DRO ● |
|---------|-----------|-----------|--------|-------|----------|-------------|
| Colored MNIST | 95.0 | $77.63_{\pm 0.44}$ | $\mathbf{85.39}_{\pm 0.94}$ | $70.44_{\pm 1.41}$ | $82.51_{\pm 0.59}$ | $84.50_{\pm 0.46}$ |
| | 98.0 | $62.29_{\pm 1.47}$ | $\mathbf{80.48}_{\pm 0.45}$ | $62.03_{\pm 0.24}$ | $72.86_{\pm 1.47}$ | $76.30_{\pm 1.53}$ |
| | 99.0 | $50.34_{\pm 0.16}$ | $\mathbf{74.01}_{\pm 2.21}$ | $51.99_{\pm 1.09}$ | $67.28_{\pm 1.69}$ | $71.33_{\pm 1.76}$ |
| | 99.5 | $35.34_{\pm 0.13}$ | $\mathbf{63.39}_{\pm 1.97}$ | $41.38_{\pm 1.31}$ | $56.40_{\pm 3.74}$ | $59.67_{\pm 2.73}$ |
| Corrupted CIFAR-10[1] | 95.0 | $45.24_{\pm 0.22}$ | $\mathbf{59.95}_{\pm 0.16}$ | $21.74_{\pm 0.27}$ | $48.74_{\pm 0.71}$ | $53.15_{\pm 0.53}$ |
| | 98.0 | $30.21_{\pm 0.82}$ | $\mathbf{49.43}_{\pm 0.78}$ | $17.81_{\pm 0.29}$ | $37.89_{\pm 0.22}$ | $40.19_{\pm 0.23}$ |
| | 99.0 | $22.72_{\pm 0.87}$ | $\mathbf{41.37}_{\pm 2.34}$ | $16.62_{\pm 0.80}$ | $32.42_{\pm 0.35}$ | $32.11_{\pm 0.83}$ |
| | 99.5 | $17.93_{\pm 0.66}$ | $\mathbf{31.66}_{\pm 1.18}$ | $15.39_{\pm 0.13}$ | $26.26_{\pm 1.06}$ | $29.26_{\pm 0.11}$ |
| Corrupted CIFAR-10[2] | 95.0 | $41.27_{\pm 0.98}$ | $\mathbf{58.57}_{\pm 1.18}$ | $19.25_{\pm 0.81}$ | $54.05_{\pm 1.01}$ | $57.92_{\pm 0.31}$ |
| | 98.0 | $28.29_{\pm 0.62}$ | $\mathbf{48.75}_{\pm 1.68}$ | $15.55_{\pm 0.84}$ | $44.22_{\pm 0.84}$ | $46.12_{\pm 1.11}$ |
| | 99.0 | $20.71_{\pm 0.29}$ | $\mathbf{41.29}_{\pm 2.08}$ | $14.42_{\pm 0.51}$ | $38.40_{\pm 0.26}$ | $39.57_{\pm 1.04}$ |
| | 99.5 | $17.37_{\pm 0.31}$ | $34.11_{\pm 2.39}$ | $13.63_{\pm 0.42}$ | $31.03_{\pm 0.42}$ | $\mathbf{34.25}_{\pm 0.74}$ |



(a) Bias-{aligned, conflicting} Colored MNIST   (b) Bias-{aligned, conflicting} Corrupted CIFAR-10[1]

# Experiments

**Real-world experiments**

| Target attribute | Unbiased | | | Bias-conflicting | | |
|---|---|---|---|---|---|---|
| | Vanilla | Ours | Group DRO | Vanilla | Ours | Group DRO |
| HairColor | $70.25_{\pm0.35}$ | $84.24_{\pm0.37}$ | $85.43_{\pm0.53}$ | $52.52_{\pm0.19}$ | $81.24_{\pm1.38}$ | $83.40_{\pm0.67}$ |
| HeavyMakeup | $62.00_{\pm0.02}$ | $66.20_{\pm1.21}$ | $64.88_{\pm0.42}$ | $33.75_{\pm0.28}$ | $45.48_{\pm4.33}$ | $50.24_{\pm0.68}$ |

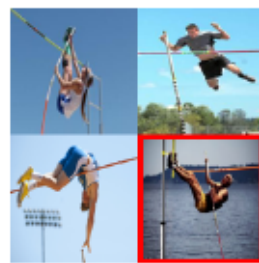| Action | Climbing | Diving | Fishing | Racing | Throwing | Vaulting | Average |
|---|---|---|---|---|---|---|---|
| Vanilla | $59.05_{\pm17.48}$ | $16.56_{\pm1.58}$ | $62.69_{\pm3.64}$ | $77.27_{\pm2.62}$ | $28.62_{\pm2.95}$ | $66.92_{\pm7.25}$ | $51.85_{\pm5.92}$ |
| Ours | $\mathbf{79.36_{\pm4.79}}$ | $\mathbf{34.59_{\pm2.26}}$ | $\mathbf{75.39_{\pm3.63}}$ | $\mathbf{83.08_{\pm1.90}}$ | $\mathbf{33.72_{\pm0.68}}$ | $\mathbf{71.75_{\pm3.32}}$ | $\mathbf{62.98_{\pm2.76}}$ |



(a) Climbing  (b) Diving  (c) Fishing  (d) Racing  (e) Throwing  (f) Vaulting

Figure 5: Illustration of BAR images of six typical action-place pairs settled as (Climbing, RockWall), (Diving, Underwater), (Fishing, WaterSurface), (Racing, APavedTrack), (Throwing, PlayingField), and (Vaulting, Sky). The images with red border lines belong to BAR evaluation set, and others belong to BAR training set.
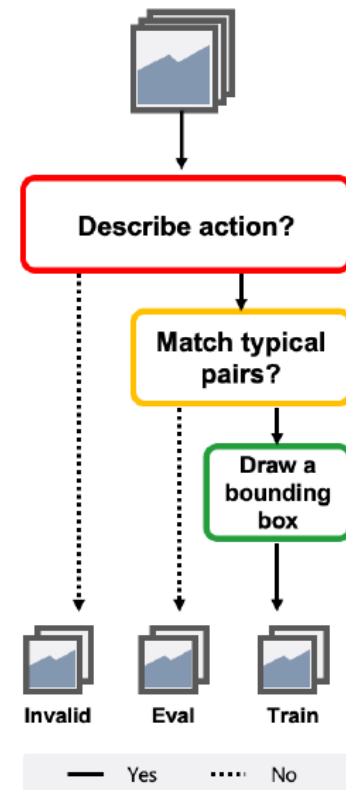


Figure 6: Illustration of BAR reasoning process.

# Conclusion

- We consider general properties of bias from observations on training dynamics of bias-aligned and bias-conflicting samples, in a more straightforward approach

- We propose a simple yet widely applicable debiasing scheme free from the choice for form and amount of supervision on the bias

- Besides, our scheme introduces only a single additional hyperparameter(GCE), where one can enjoy simplicity in use

- Using this approach can increase awareness of underexplored biases