

Learning Debiased Representations via Conditional Attribute Interpolation

Yi-Kai Zhang, Qi-Wei Wang, De-Chuan Zhan, Han-Jia Ye[✉]



Biased Representation



(a) Orange lifeboat

97.8% lifeboat
0.5% beacon
0.4% container ship

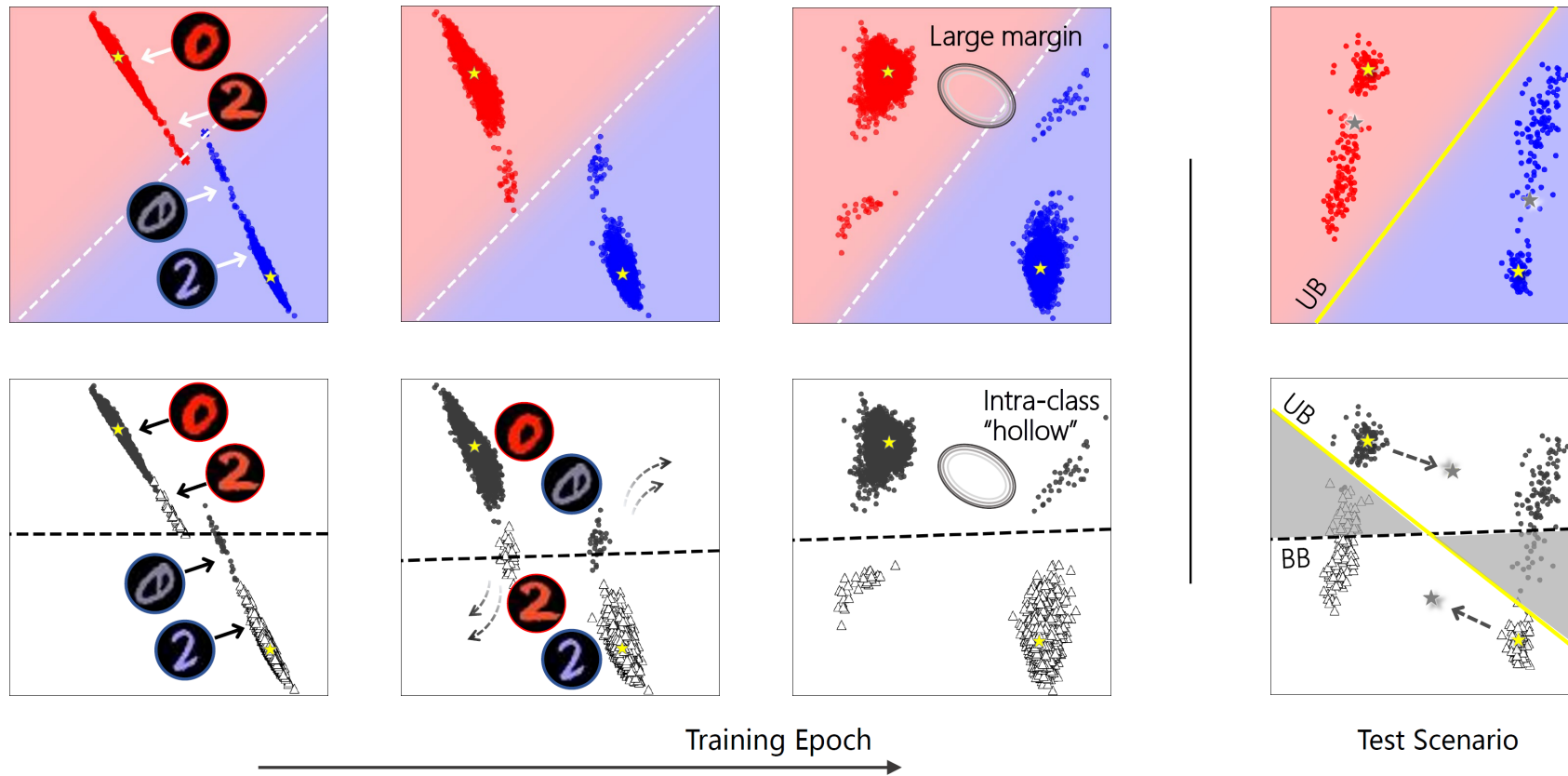


(b) Orange cyclists

57.0% lifeboat
13.8% bicycle-built-for-two
9.0% toyshop

- Classification of a standard **ImageNet pre-trained ResNet - 50** of
 - (a) an **orange lifeboat** for training (with **color** and **shape** attributes) ✓
 - (b) an **orange cyclist** for test (**aligned** with color, **conflicting** with shape) ✗

Biased Training Dynamics



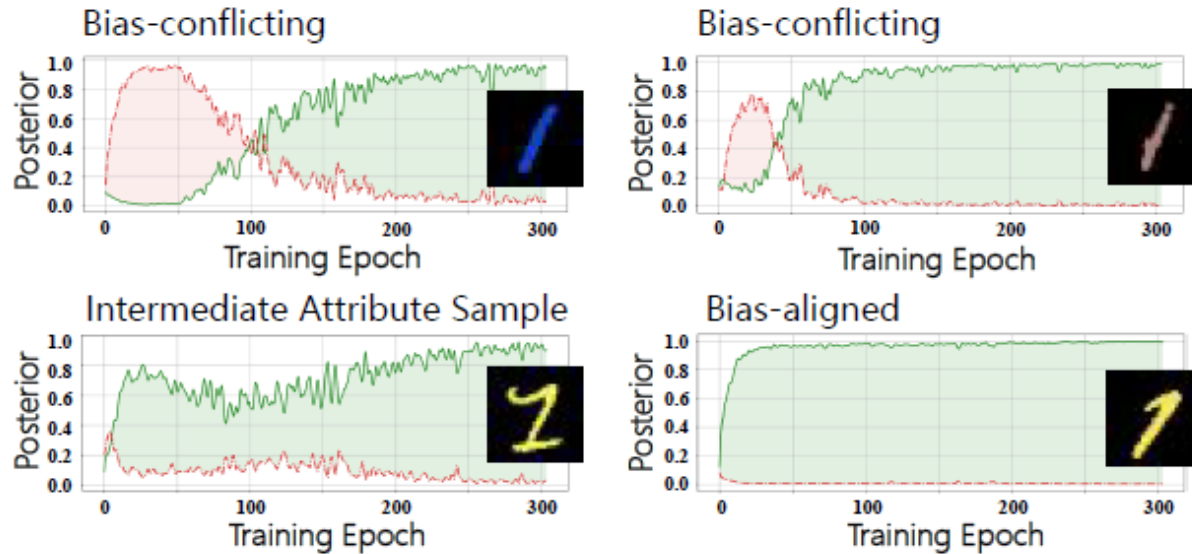
- Color** attribute classifier for ● in red and ● in blue, **Shape** attribute classifier for ○ ○ in ● and △ △ in △

(a) The target attribute **shape** is learned later in a “*lazy*” manner.

(b) The easier-to-learn **color** leaves a large margin and triggers the shape intra-class “*hollow*”.

(c) With the “hollow” **conditioned on** a particular shape, the true class representation is **deviated** toward color.

Find Intermediate Attribute Samples with χ -pattern



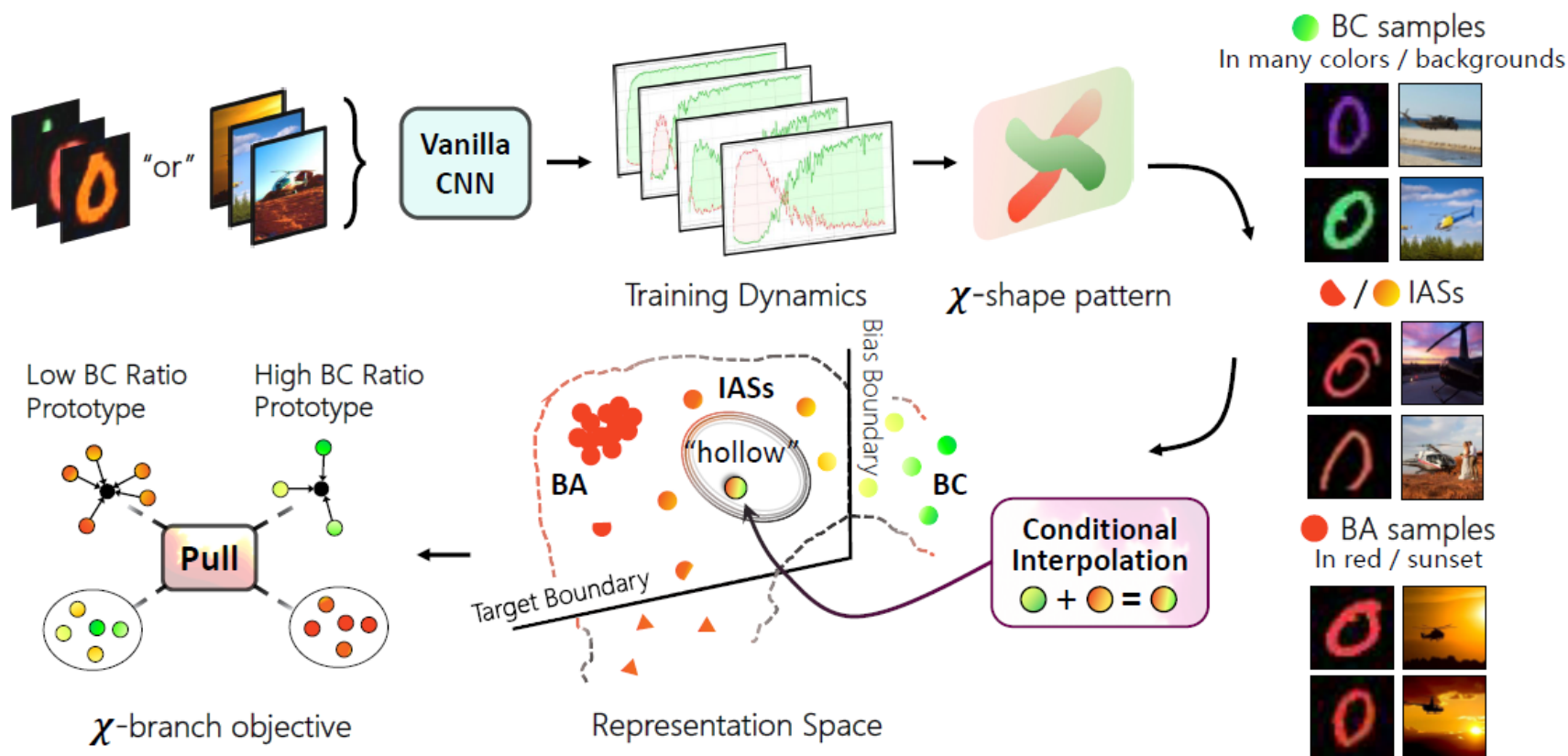
$$\mathcal{L}_{\text{CE}}(\mathbf{x}_i) = \begin{pmatrix} \mathcal{L}_{\text{CE}}^{gt}(\mathbf{x}_i) = \{-\log \text{Pr}^t(y_i | \mathbf{x}_i)\}_{t=1}^T \\ \mathcal{L}_{\text{CE}}^b(\mathbf{x}_i) = \{-\log \text{Pr}^t(b_i | \mathbf{x}_i)\}_{t=1}^T \end{pmatrix},$$

$$\chi_{\text{pattern}} = \begin{pmatrix} p^{gt} = \{e^{-A_1 t}\}_{t=1}^T \\ p^b = \{e^{A_2 t}\}_{t=1}^T \end{pmatrix},$$

$$\begin{aligned} \mathbf{s}(\mathbf{x}_i) &= \langle \mathcal{L}_{\text{CE}}(\mathbf{x}_i), \chi_{\text{shape}} \rangle \\ &= \langle \mathcal{L}_{\text{CE}}^{gt}(\mathbf{x}_i), p^{gt} \rangle + \langle \mathcal{L}_{\text{CE}}^b(\mathbf{x}_i), p^b \rangle \\ &= \sum_{t=1}^T (-e^{-A_1 t} \cdot \log \text{Pr}(h_{\theta}(\mathbf{x}_i) = y_i | \mathbf{x}_i) \\ &\quad - e^{A_2 t} \cdot \log \text{Pr}(h_{\theta}(\mathbf{x}_i) = b_i | \mathbf{x}_i)) . \end{aligned}$$

- The change of posterior over the **GT-class** as well as the **bias class**

Learning from χ -structured objective



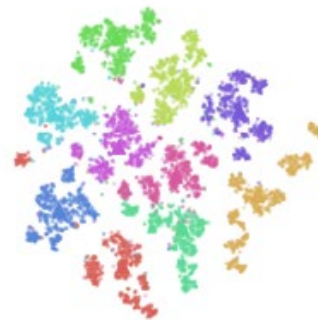
$$\mathbf{P}_{\gamma, c} = \frac{1}{K} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{B}_{\gamma}} f_{\phi}(\mathbf{x}_i) \cdot \mathbb{I}[y_i = c], \quad \Pr(y_i | \mathbf{x}_i) = \frac{\exp(-d(f_{\phi}(\mathbf{x}_i), \mathbf{P}_{\gamma, y_i}) / \tau)}{\sum_{c \in [C]} \exp(-d(f_{\phi}(\mathbf{x}_i), \mathbf{P}_{\gamma, c}) / \tau)}$$

Benchmark Results of χ^2 -model

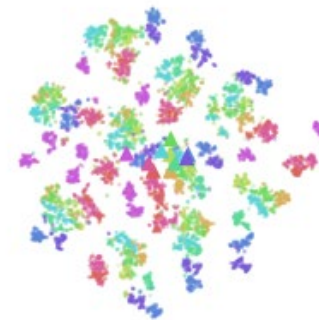
- State-of-the-art on Colored MNIST, Corrupted CIFAR-10, Biased CelebA and Biased NICO.

Dataset	Ratio (%)	Colored MNIST				Corrupted CIFAR-10				Data Source	Biased CelebA			NICO All
		99.9	99.5	99.0	95.0	99.9	99.5	99.0	95.0		BA	BC	All	
Vanilla	○	28.58	59.29	74.42	87.13	26.91	30.16	37.71	41.60	LfF [36]	73.69	70.41	72.05	34.44
+ p [44]	○	31.01	64.82	76.84	87.86	26.55	29.48	38.07	42.30	DFA [27]	<u>94.01</u>	58.98	<u>76.50</u>	33.10
RUBi [7]	●	27.82	70.80	86.58	96.77	33.70	34.70	34.59	47.23	χ^2 -model	97.66	<u>60.79</u>	79.23	36.99
ReBias [4]	●	27.71	72.89	85.95	96.87	33.65	34.40	35.82	47.45					
End [46]	●	28.19	<u>81.81</u>	88.10	96.99	31.30	33.83	34.02	38.77					
DI [50]	●	33.18	80.63	86.28	98.36	32.09	33.37	37.65	<u>51.27</u>					
LfF [36]	○	30.24	68.90	76.69	96.81	29.89	33.68	35.28	45.38					
LFA [27]	○	22.31	64.13	81.83	95.45	32.49	35.74	39.63	47.25					
χ -pattern + p	●	<u>60.33</u>	64.15	93.53	<u>98.30</u>	<u>35.33</u>	39.31	41.32	53.37					
χ^2 -model (Ours)	○	66.91	88.73	<u>92.15</u>	97.87	35.67	<u>37.61</u>	<u>40.74</u>	49.04					

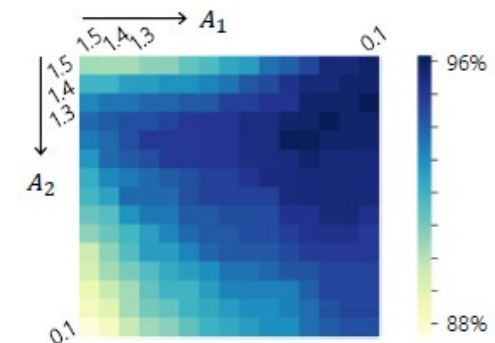
Measure	Acc. \uparrow	98%- σ \downarrow	AP \uparrow
Entropy [20]	78.33	632	83.52
Confidence [28]	80.33	590	85.61
Loss [36]	<u>94.39</u>	<u>418</u>	<u>98.22</u>
Pleiss <i>et al.</i> [39]	82.67	686	89.24
Zhao <i>et al.</i> [55]	90.33	451	96.04
χ -pattern (Ours)	95.84	372	98.44



(a) t-SNE on a_t

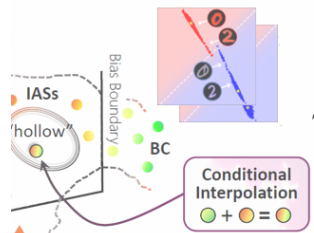


(b) t-SNE on a_b



(c) ablation study of A_1, A_2

Feel free to discuss with us~

- 

χ^2 -model: a practical method to learn debiased representation
- Code is available at github.com/ZhangYikai/chi-square

“Thanks!”



Paper



Code



Yi-Kai Zhang's Homepage
lamda.nju.edu.cn/zhangyk



Qi-Wei Wang's Homepage
lamda.nju.edu.cn/wangqiwei