

EFFICIENT BIAS MITIGATION WITHOUT PRIVILEGED INFORMATION



MATEO
ESPINOSA ZARLENGA



SWAMI
SANKARANARAYANAN



JERONE T. A.
ANDREWS



ZOHREH
SHAMS



MATEJA
JAMNIK



ALICE
XIANG

Sony AI



UNIVERSITY OF
CAMBRIDGE

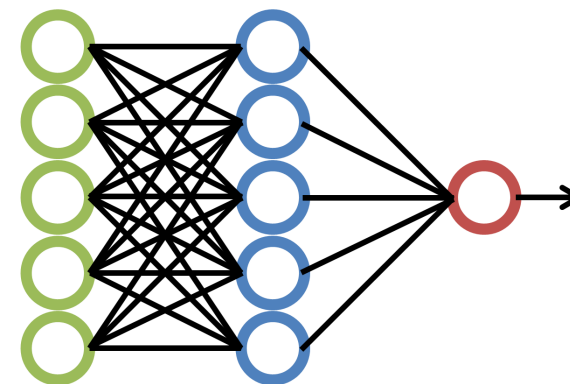
SUPERVISED LEARNING'S BREAD AND BUTTER

TODAY WE WILL CONSIDER THE TRADITIONAL SUPERVISED LEARNING SETUP:



TRAINING SAMPLES

$$\mathcal{D} = \{(x_i, y_i) \mid x_i \in \mathbb{R}^m, y_i \in \{1, 2, \dots, L\}\}_{i=1}^n$$



DEEP NEURAL NETWORK (DNN)

$$f_{\theta}: \mathbb{R}^m \rightarrow [0, 1]^L$$

SUPERVISED LEARNING'S BREAD AND BUTTER

IN PARTICULAR, WE WILL FOCUS ON INSTANCES WHEN A DNN f_θ IS TRAINED BY MINIMIZING THE EMPIRICAL MEAN LOSS $\ell(\cdot)$ OVER THE TRAINING SET:

$$J_{ERM}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i)$$

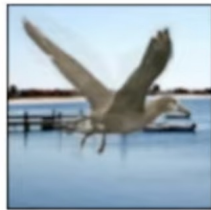
THIS IS WHAT IS REFERRED TO AS **EMPIRICAL RISK MINIMIZATION (ERM)**.

THE RISKS OF EMPIRICAL RISK MINIMIZATION

EVEN WHEN ERM LEADS TO HIGH PERFORMANCE ON AVERAGE, THIS CAN CHANGE WHEN WE LOOK AT SPECIFIC GROUPS:

Wildlife image classification (Wah et al., '11; Sagawa et al., '20)

Input: image of a bird



Label: bird type

water

land

THE RISKS OF EMPIRICAL RISK MINIMIZATION

EVEN WHEN ERM LEADS TO HIGH PERFORMANCE ON AVERAGE, THIS CAN CHANGE WHEN WE LOOK AT SPECIFIC GROUPS:

Wildlife image classification (Wah et al., '11; Sagawa et al., '20)

Input: image of a bird



Label: bird type

water

land

97.3% average test accuracy

THE RISKS OF EMPIRICAL RISK MINIMIZATION

EVEN WHEN ERM LEADS TO HIGH PERFORMANCE ON AVERAGE, THIS CAN CHANGE WHEN WE LOOK AT SPECIFIC GROUPS:

Wildlife image classification (Wah et al., '11; Sagawa et al., '20)

Input: image of a bird



Label: bird type

water

land

97.3% average test accuracy

72.6% on waterbirds on land backgrounds



THE RISKS OF EMPIRICAL RISK MINIMIZATION

EVEN WHEN ERM LEADS TO HIGH PERFORMANCE ON AVERAGE, THIS CAN CHANGE WHEN WE LOOK AT SPECIFIC GROUPS:

Wildlife image classification (Wah et al., '11; Sagawa et al., '20)

Input: image of a bird



Label: bird type

water

land

97.3% average test accuracy

72.6% on waterbirds on land backgrounds



WE ARE THEREFORE INTERESTED IN MAXIMIZING THE **WORST GROUP ACCURACY (WGA)**:

$$\text{WGA}(f_{\theta}, \mathcal{P}) := \min_{g \in \{1, 2, \dots, k\}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{P}(\mathbf{x}, y|g)} [\mathbb{1}(f_{\theta}(\mathbf{x}) = y)]$$

BIAS MITIGATION

THERE IS A PLETHORA OF BIAS MITIGATION METHODS

MaskTune: Mitigating Spurious Correlations by Forcing to Explore

Saeid Asgari Taghanaki*
Autodesk AI Lab

Aliasghar Khani*
Autodesk AI Lab

Fereshte Khani*
Stanford University

Ali Gholami*
Autodesk AI Lab

Linh Tran
Autodesk AI Lab

Ali Mahdavi-Amiri
Simon Fraser University

Ghassan Hamarneh
Simon Fraser University

Towards Last-layer Retraining for Group Robustness with Fewer Annotations

Tyler LaBonte¹ Vidya Muthukumar^{2,1} Abhishek Kumar³
¹H. Milton Stewart School of Industrial and Systems Engineering, Georgia Institute of Technology
²School of Electrical and Computer Engineering, Georgia Institute of Technology
³Google DeepMind
{tlabonte, vmuthukumar8}@gatech.edu abhishk@google.com

Multiaaccuracy: Black-Box Post-Processing for Fairness in Classification

Michael P. Kim*[†]
mpk@cs.stanford.edu
Stanford University

Amirata Ghorbani[†]
amiratag@stanford.edu
Stanford University

Learning from Failure: Training Debiased Classifier from Biased Classifier

Junhyun Nam¹ Hyuntak Cha² Sungsoo Ahn¹ Jaeho Lee¹ Jinwoo Shin^{1,2}
¹School of Electrical Engineering, KAIST
²Graduate School of AI, KAIST
{junhyun.nam, hyuntak.cha, sungsoo.ahn, jaeho-lee, jinwoos}@kaist.ac.kr

No Subclass Left Behind: Fine-Grained Robustness in Coarse-Grained Classification Problems

Nimit S. Sohoni, Jared A. Dunnmon, Geoffrey Angus, Albert Gu, Christopher Ré
Stanford University
nims@stanford.edu, jdunnmon@cs.stanford.edu, gdlangus@cs.stanford.edu, albertgu@stanford.edu, chrismre@cs.stanford.edu

DISTRIBUTIONALLY ROBUST NEURAL NETWORKS FOR GROUP SHIFTS: ON THE IMPORTANCE OF REGULARIZATION FOR WORST-CASE GENERALIZATION

Shiori Sagawa*
Stanford University
ssagawa@cs.stanford.edu

Tatsunori B. Hashimoto
Microsoft
tahashim@microsoft.com

Pang Wei Koh*
Stanford University
pangwei@cs.stanford.edu

Percy Liang
Stanford University
pliang@cs.stanford.edu

LAST LAYER RE-TRAINING IS SUFFICIENT FOR ROBUSTNESS TO SPURIOUS CORRELATIONS

James Zou
jamesz@stanford.edu
Stanford University

Polina Kirichenko*
New York University

Pavel Izmailov*
New York University

Andrew Gordon Wilson
New York University

Just Train Twice: Improving Group Robustness without Training Group Information

Evan Zheran Liu*¹ Behzad Haghgoo*¹ Annie S. Chen*¹ Aditi Raghunathan¹ Pang Wei Koh¹
Shiori Sagawa¹ Percy Liang¹ Chelsea Finn¹

BIAS MITIGATION

THERE IS A PLETHORA OF BIAS MITIGATION METHODS

THESE CAN BE:

1. **GROUP SUPERVISED**: WE ASSUME GROUP LABELS DURING TRAINING

BIAS MITIGATION

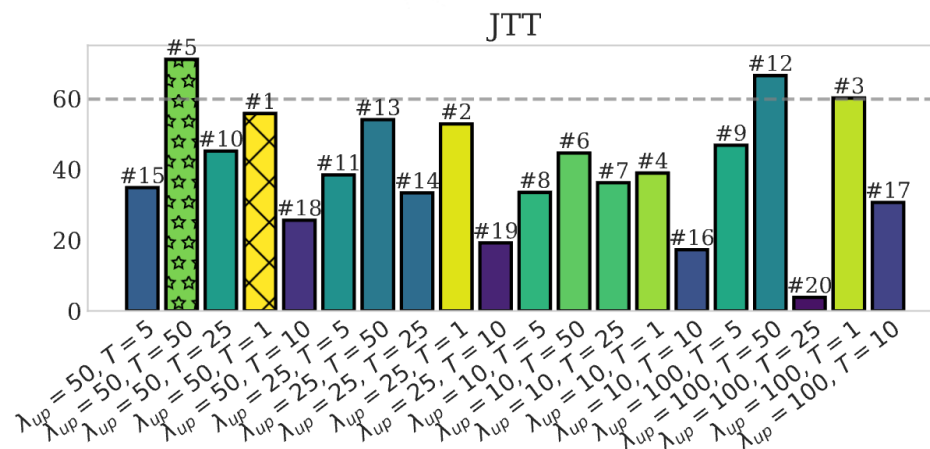
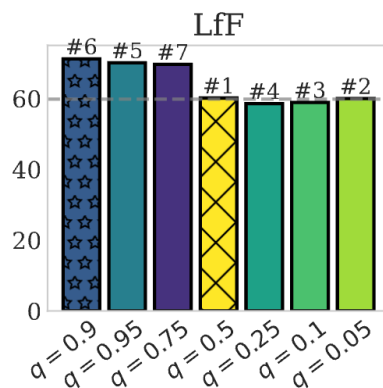
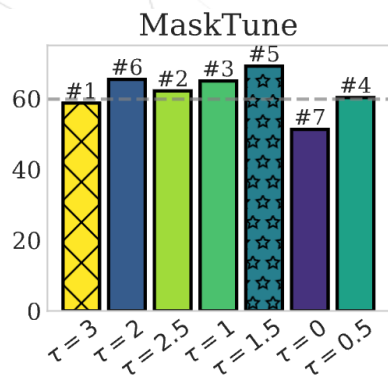
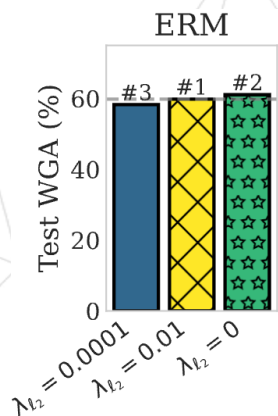
THERE IS A PLETHORA OF BIAS MITIGATION METHODS

THESE CAN BE:

1. **GROUP SUPERVISED**: WE ASSUME GROUP LABELS DURING TRAINING
2. **GROUP UNSUPERVISED**: WE DO NOT ASSUME GROUP LABELS DURING TRAINING*

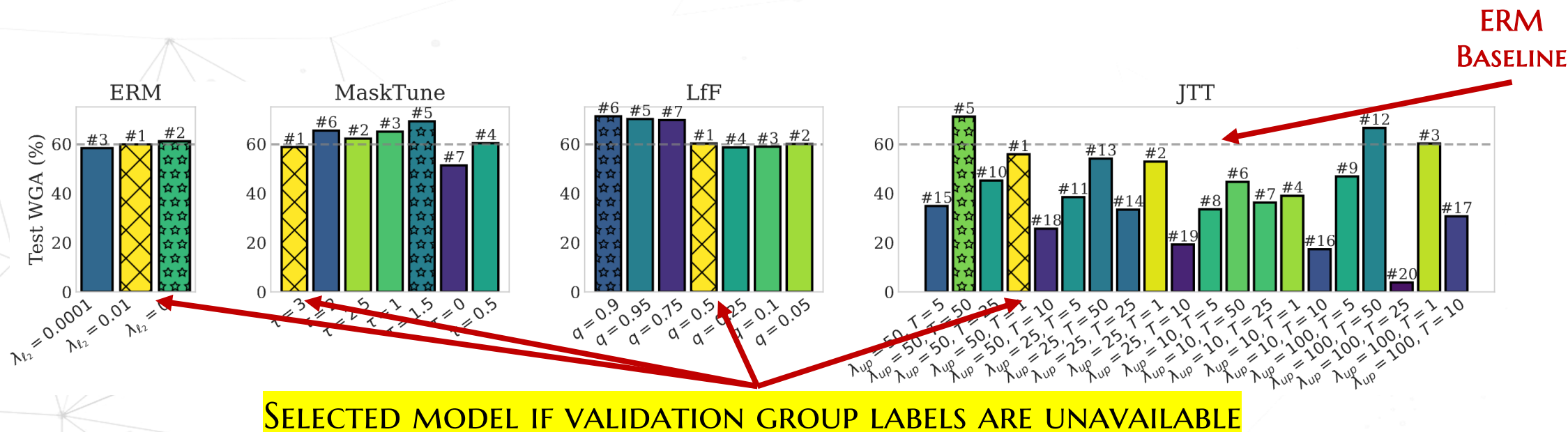
*THE REALITY OF UNSUPERVISED METHODS

IN PRACTICE, **UNSUPERVISED BIAS MITIGATION METHODS** NEED GROUP LABELS DURING MODEL SELECTION TO AVOID SELECTING A BIASED MODEL:



*THE REALITY OF UNSUPERVISED METHODS

IN PRACTICE, **UNSUPERVISED BIAS MITIGATION METHODS** NEED GROUP LABELS DURING MODEL SELECTION TO AVOID SELECTING A BIASED MODEL:



THE SELECTED HYPERPARAMETERS LEAD TO A MODEL **NO BETTER THAN AN ERM MODEL!**

OUR WORK



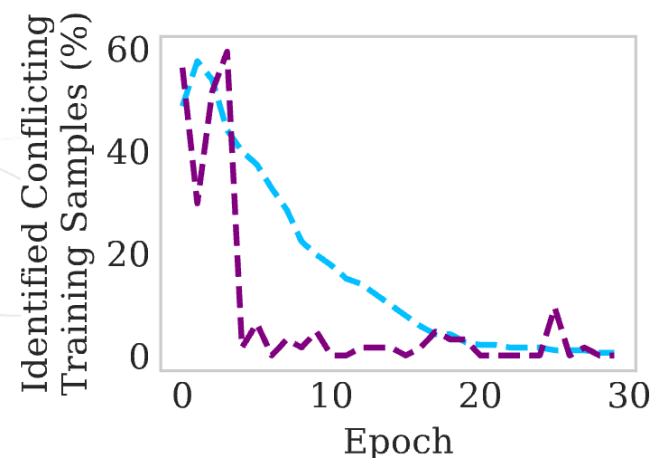
HOW CAN WE DESIGN A BIAS MITIGATION METHOD THAT DOES NOT REQUIRE GROUP LABELS FOR EITHER TRAINING OR MODEL SELECTION?

CIRCUMVENTING PRIVILEGED INFORMATION: INSIGHTS



CIRCUMVENTING PRIVILEGED INFORMATION: INSIGHTS

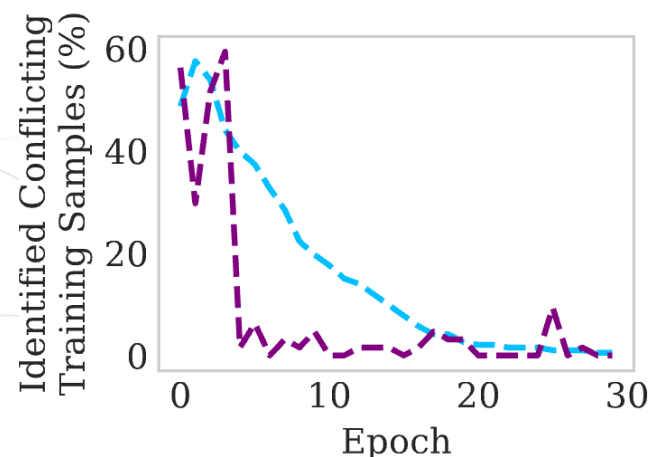
INSIGHT #1 (NAM ET AL. AND LIU ET AL.): SAMPLES WITH SPURIOUS CORRELATIONS ARE LEARNT BEFORE SAMPLES WITHOUT THE SPURIOUS CORRELATION



--- Bias-conflicting (Waterbirds) — Bias-aligned (Waterbirds) - - - Bias-conflicting (CelebA) — Bias-aligned (CelebA)

CIRCUMVENTING PRIVILEGED INFORMATION: INSIGHTS

INSIGHT #1 (NAM ET AL. AND LIU ET AL.): SAMPLES WITH SPURIOUS CORRELATIONS ARE LEARNT BEFORE SAMPLES WITHOUT THE SPURIOUS CORRELATION



--- Bias-conflicting (Waterbirds) — Bias-aligned (Waterbirds) - - - Bias-conflicting (CelebA) — Bias-aligned (CelebA)

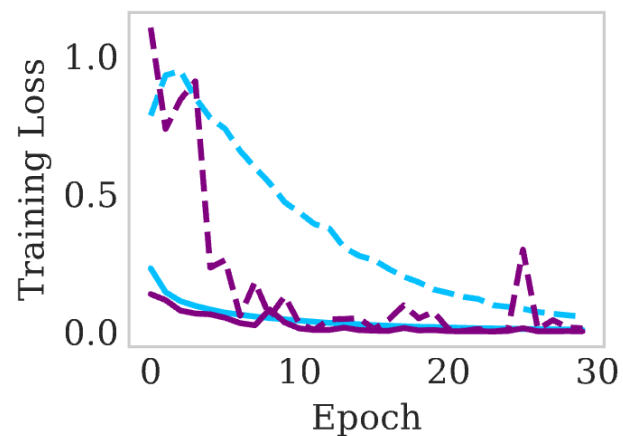
PREVIOUS WORKS [1, 2] EXPLOIT THIS BY LOOKING AT **A PRE-DETERMINED "TIME SLICE"**

[1] NAM ET AL. "LEARNING FROM FAILURE: DE-BIASING CLASSIFIER FROM BIASED CLASSIFIER." *NEURIPS 2020*.

[2] LIU ET AL. "JUST TRAIN TWICE: IMPROVING GROUP ROBUSTNESS WITHOUT TRAINING GROUP INFORMATION." *ICML, 2021*.

CIRCUMVENTING PRIVILEGED INFORMATION: INSIGHTS

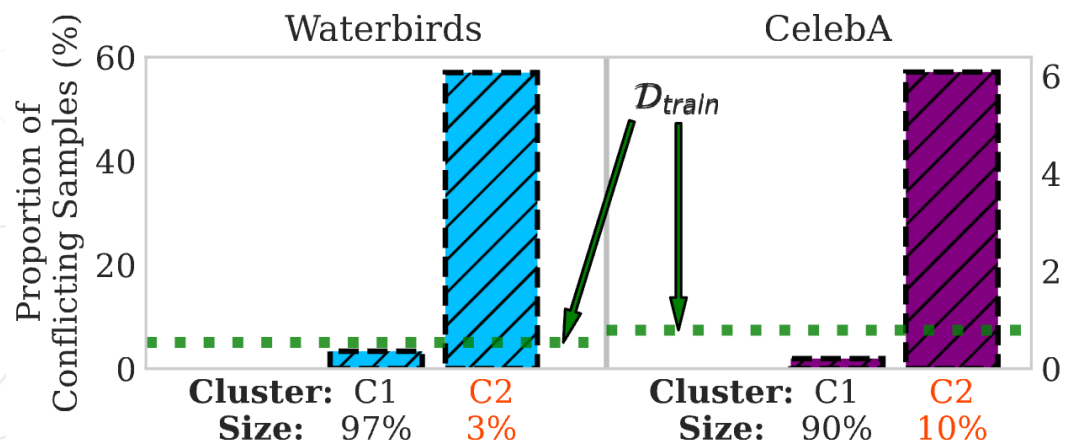
INSIGHT #2: TRAINING LOSS HISTORIES ARE VERY INFORMATIVE SIGNALS



--- Bias-conflicting (Waterbirds) — Bias-aligned (Waterbirds) - - - Bias-conflicting (CelebA) — Bias-aligned (CelebA)

CIRCUMVENTING PRIVILEGED INFORMATION: INSIGHTS

INSIGHT #2: TRAINING LOSS HISTORIES ARE VERY INFORMATIVE SIGNALS



--- Bias-conflicting (Waterbirds) — Bias-aligned (Waterbirds) --- Bias-conflicting (CelebA) — Bias-aligned (CelebA)

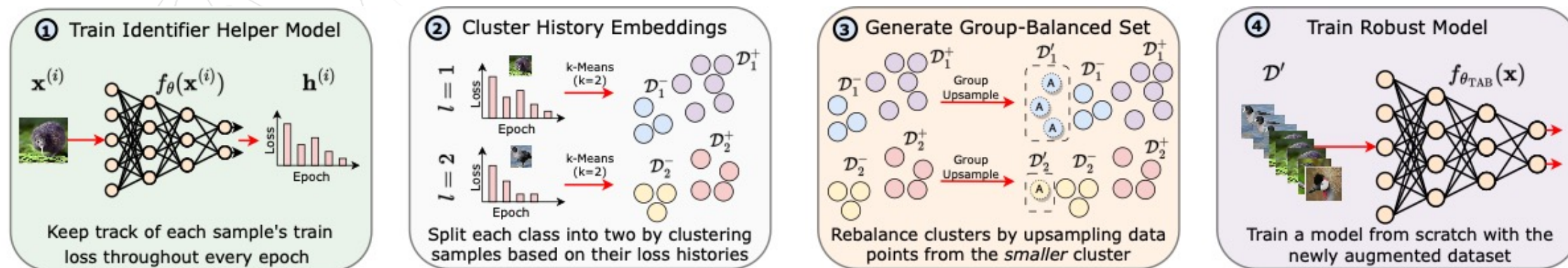
CLUSTERING SAMPLES BASED ON THEIR TRAINING HISTORIES PRODUCES **A DATA SUBSET WITH A HIGHER PROPORTION OF SAMPLES WITHOUT THE SPURIOUS CORRELATION!**

TARGETED AUGMENTATIONS FOR BIAS MITIGATION



TARGETED AUGMENTATIONS FOR BIAS MITIGATION

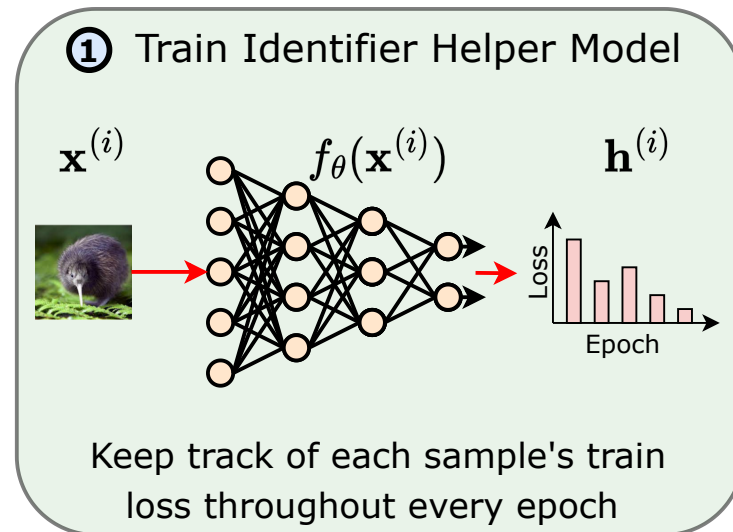
WE PROPOSE **TARGETED AUGMENTATIONS FOR BIAS MITIGATION (TAB)**, A NEW **HYPERPARAMETER-FREE** GROUP-UNSUPERVISED BIAS MITIGATION PIPELINE:



OUR APPROACH EXPLOITS THE **TRAINING HISTORY** OF AN IDENTIFICATION MODEL TO GENERATE **A GROUP-BALANCED DATASET** FROM WHICH A ROBUST MODEL CAN BE TRAINED

TARGETED AUGMENTATIONS FOR BIAS MITIGATION

TAB FIRST TRAINS AN ERM MODEL WHILE **KEEPING TRACK OF THE LOSS** ACROSS ALL TRAINING SAMPLES AND EPOCHS:

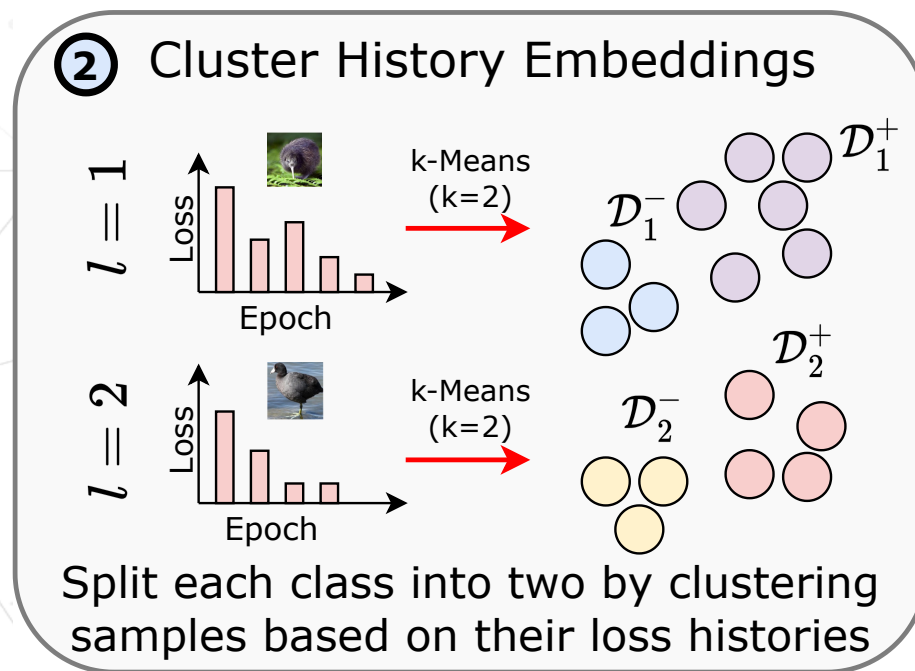


$$\mathbf{h}_i := [\ell(f_{\theta^{(1)}}(\mathbf{x}_i), y_i), \ell(f_{\theta^{(2)}}(\mathbf{x}_i), y_i), \dots, \ell(f_{\theta^{(T)}}(\mathbf{x}_i), y_i)]^T$$

THE LOSS DURING TRAINING PROVIDES VALUABLE INSIGHTS INTO WHICH CONCEPTS ARE PRESENT OR MISSING IN A SAMPLE.

TARGETED AUGMENTATIONS FOR BIAS MITIGATION

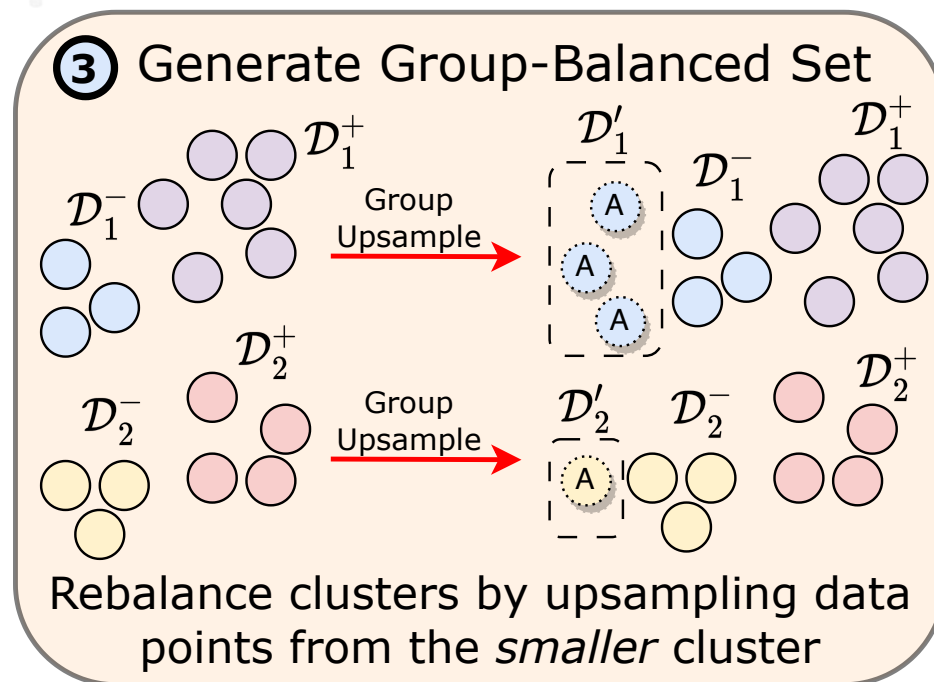
WE THEN IDENTIFY ERROR GROUPS BY **CLUSTERING THE LOSS HISTORY EMBEDDING** SPACE FOR EACH CLASS LABEL:



$g_i := \text{ClusterLabel}(\mathbf{h}_i, 2\text{-Means}(H_l))$ FOR ALL $\mathbf{h}_i \in H_l$
WHERE $H_l := \{h_j \mid y_j = l\}$ IS THE SET OF HISTORY EMBEDDINGS FOR SAMPLES IN CLASS l

TARGETED AUGMENTATIONS FOR BIAS MITIGATION

NEXT, WE GENERATE A **GROUP-BALANCED TRAINING SET** BY **UPSAMPLING** EACH MINORITY CLUSTER TO MATCH THE SIZE OF THE MAJORITY CLUSTER.

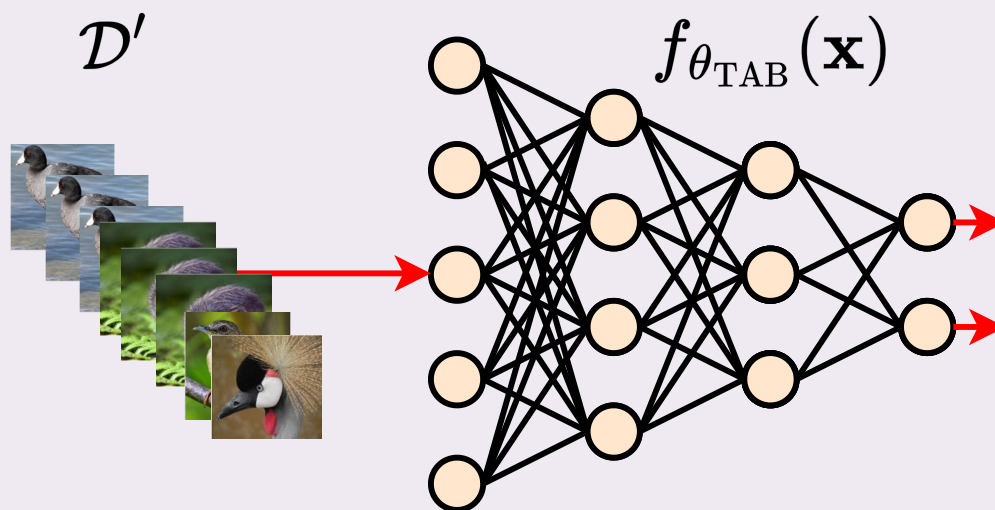


WE DO SO BY **RANDOMLY UPSAMPLING ELEMENTS FROM THE MINORITY CLUSTER** .

TARGETED AUGMENTATIONS FOR BIAS MITIGATION

FINALLY, WE **TRAIN A ROBUST MODEL** USING ERM ON THIS GROUP-BALANCED DATASET:

④ Train Robust Model



Train a model from scratch with the newly augmented dataset

TARGETED AUGMENTATIONS FOR BIAS MITIGATION

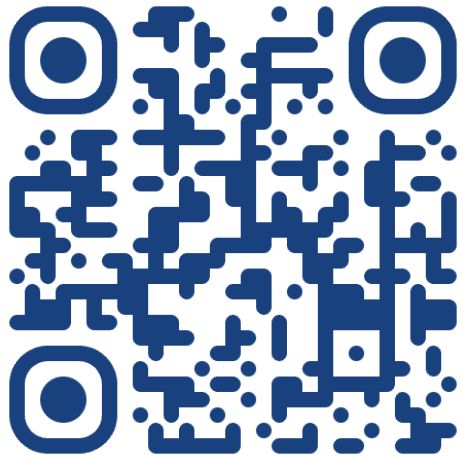
SO HOW DOES TAB PERFORM IN PRACTICE?

KEY RESULTS TL;DR

	Method - {Hypers}	Even-Odd ($p = 99\%$)	cMNIST ($p = 98\%$)	Waterbirds	CelebA	BAR	CUB
WGA (%)	G-DRO - $\{\eta, \lambda_{\ell_2}\}$	57.66 ± 6.76	59.29 ± 3.27	68.54 ± 1.75	85.74 ± 0.69	N/A	N/A
	ERM - $\{\eta, \lambda_{\ell_2}\}$	55.98 ± 13.85	46.97 ± 8.71	44.86 ± 1.11	34.81 ± 0.26	29.56 ± 1.78	16.67 ± 0.00
	LfF - $\{q\}$	2.97 ± 3.36	48.45 ± 5.83	51.14 ± 1.08	40.00 ± 0.00	29.56 ± 2.35	14.44 ± 3.14
	JTT - $\{T, \lambda_{\text{up}}\}$	79.32 ± 1.76	57.21 ± 3.59	44.50 ± 0.45	37.78 ± 2.83	30.98 ± 2.00	12.22 ± 1.57
	MaskTune - $\{\tau\}$	72.82 ± 3.08	13.94 ± 7.37	35.67 ± 1.75	37.04 ± 1.14	17.61 ± 1.54	10.00 ± 7.20
	TAB (ours) - \emptyset	81.85 ± 2.39	63.26 ± 2.50	55.92 ± 1.80	40.00 ± 1.20	38.94 ± 1.03	18.89 ± 1.57
Mean Acc. (%)	G-DRO - $\{\eta, \lambda_{\ell_2}\}$	58.97 ± 6.79	94.83 ± 0.55	97.19 ± 0.28	92.67 ± 0.14	N/A	N/A
	ERM - $\{\eta, \lambda_{\ell_2}\}$	85.52 ± 12.09	91.22 ± 0.26	97.68 ± 0.06	95.45 ± 0.04	56.93 ± 1.13	74.81 ± 0.29
	LfF - $\{q\}$	60.29 ± 13.53	90.48 ± 1.17	97.46 ± 0.12	95.22 ± 0.02	55.96 ± 1.25	74.00 ± 0.67
	JTT - $\{T, \lambda_{\text{up}}\}$	93.12 ± 4.74	92.13 ± 1.13	97.71 ± 0.11	94.77 ± 0.05	58.00 ± 2.34	69.92 ± 0.10
	MaskTune - $\{\tau\}$	92.60 ± 5.02	83.25 ± 3.26	98.15 ± 0.04	95.32 ± 0.07	50.66 ± 1.38	70.07 ± 0.97
	TAB (ours) - \emptyset	94.98 ± 3.37	93.28 ± 1.09	97.52 ± 0.09	94.67 ± 0.05	61.11 ± 0.94	72.98 ± 0.34

TAB ACHIEVES **BETTER WORST-GROUP ACCURACIES** THAN COMPETING APPROACHES WHILE MAINTAINING **A COMPETITIVE MEAN ACCURACY COMPARED TO ERM MODELS**

PAPER, POSTER, AND CONTACT INFORMATION



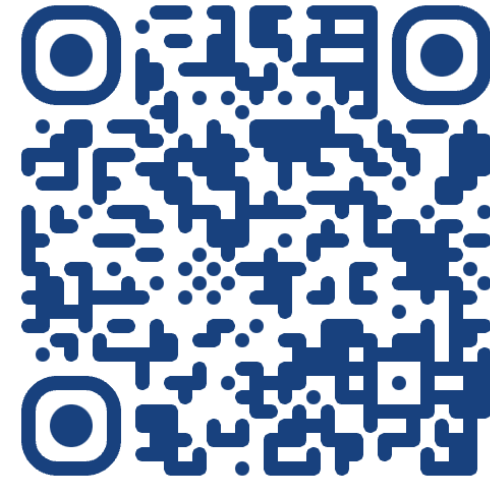
PAPER

POSTER INFORMATION

TODAY, TUESDAY OCT 1ST

4:30 P.M. — 6:30 P.M.

POSTER #27



PROJECT PAGE

IF YOU WANT TO DISCUSS FURTHER, CONTACT ME AT ME466@CAM.AC.UK