



# **A Truly Unbiased Model**

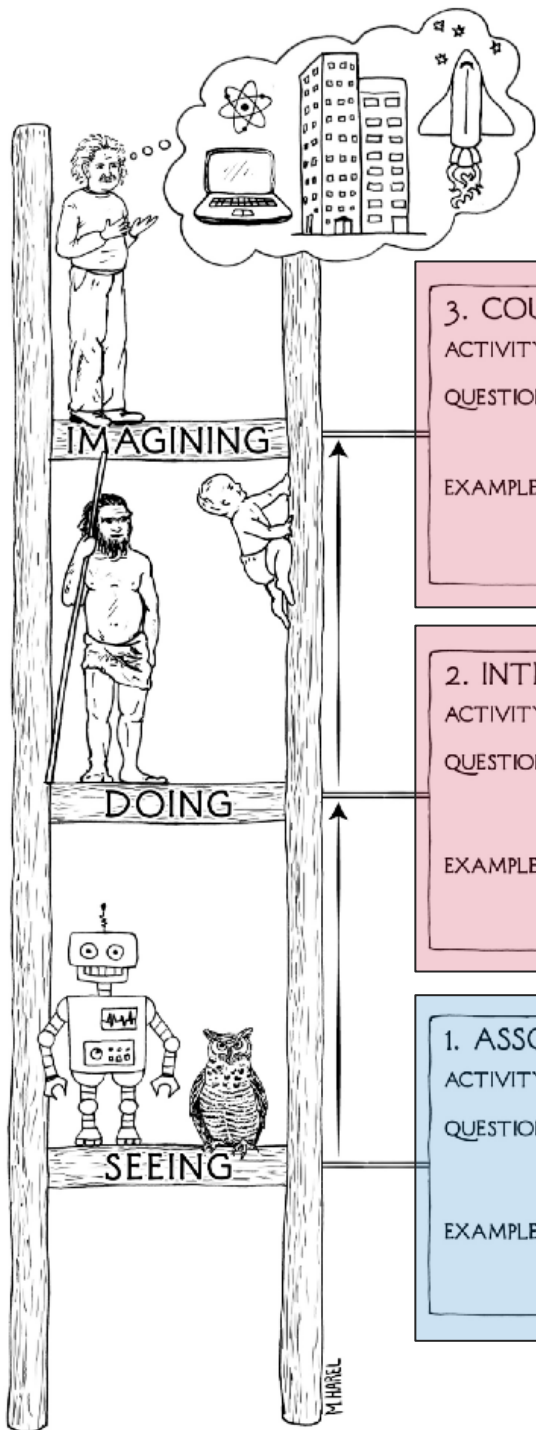
## **Recent Progress @ MReal**

Hanwang Zhang 张含望

<https://mreallab.github.io/>

hanwangzhang@ntu.edu.sg

# Causality vs. Correlation

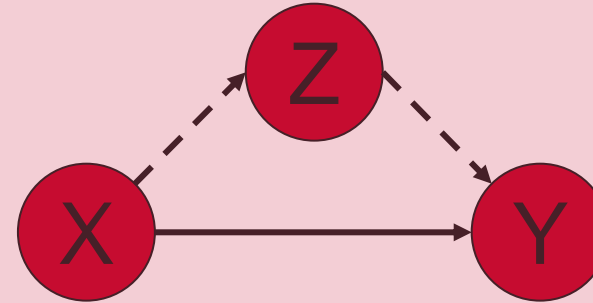


## 3. COUNTERFACTUALS

ACTIVITY: Imagining, Retrospection, Understanding

QUESTIONS: *What if I had done ...? Why?*  
(Was it X that caused Y? What if X had not occurred? What if I had acted differently?)

EXAMPLES: Was it the aspirin that stopped my headache?  
Would Kennedy be alive if Oswald had not killed him? What if I had not smoked for the last 2 years?

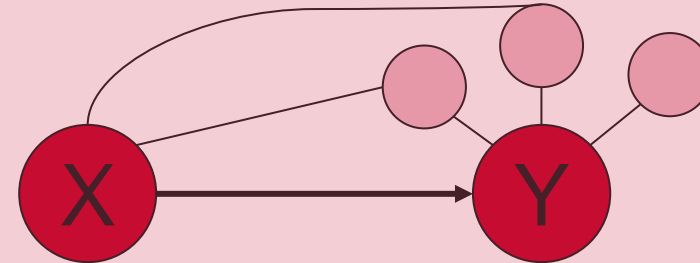


## 2. INTERVENTION

ACTIVITY: Doing, Intervening

QUESTIONS: *What if I do ...? How?*  
(What would Y be if I do X?  
How can I make Y happen?)

EXAMPLES: If I take aspirin, will my headache be cured?  
What if we ban cigarettes?



## 1. ASSOCIATION

ACTIVITY: Seeing, Observing

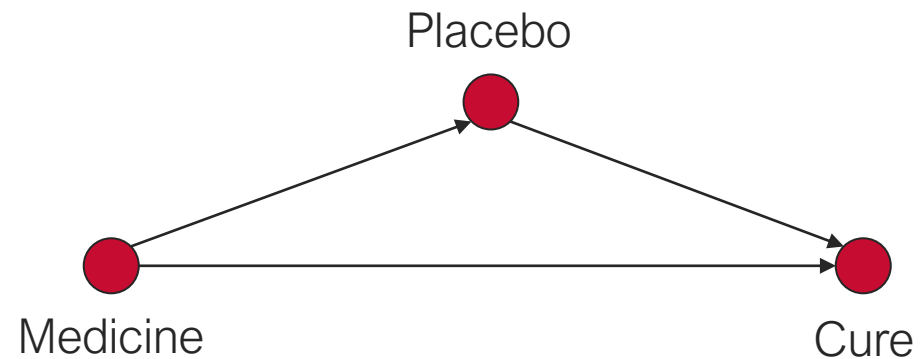
QUESTIONS: *What if I see ...?*  
(How are the variables related?  
How would seeing X change my belief in Y?)

EXAMPLES: What does a symptom tell me about a disease?  
What does a survey tell us about the election results?



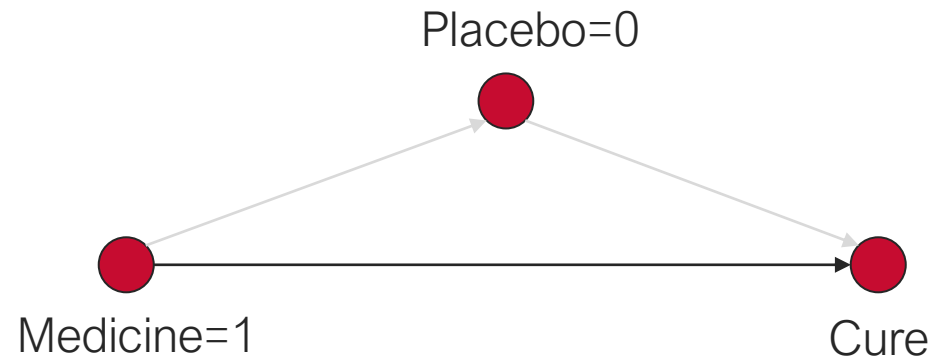
# Mediation Effect

- How to remove Placebo Effect



# Mediation Effect

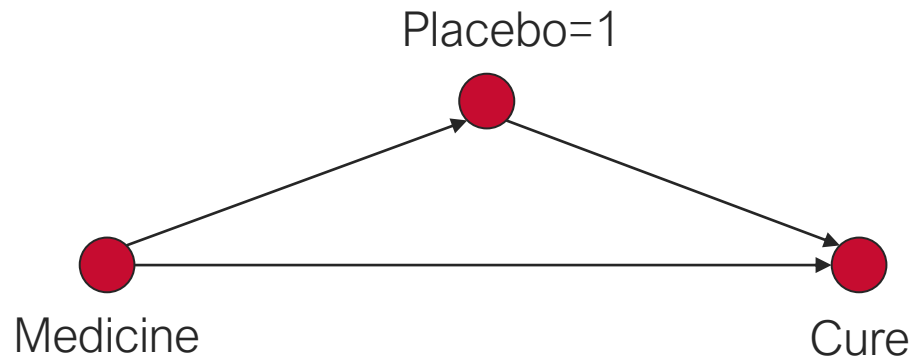
- How to remove Placebo Effect?
- Challenge: Med = 1 and Placebo = 1 always co-occur; or, illegal to realize the following graph



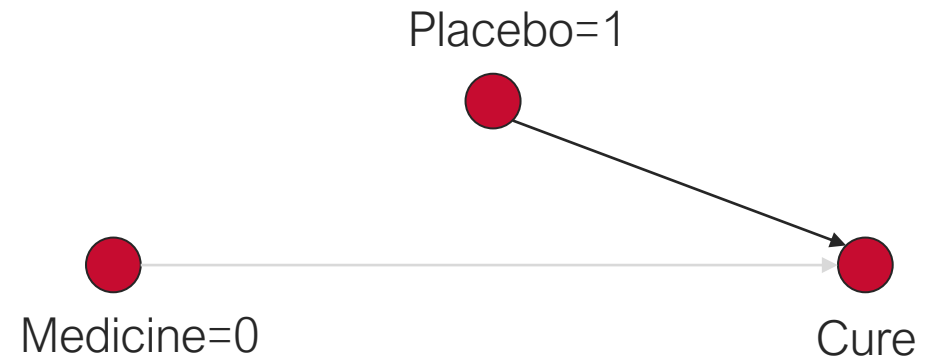
Ideal case

# Mediation Effect: TDE (the minus trick)

- How to remove Placebo Effect?
- Solution: counterfactual  $\rightarrow$  cheating  $\rightarrow$  Med = 0 but Placebo = 1



minus



# Do and CF in debiasing methods

- Assumption: train  $\neq$  test (OOD)
- Do: CSS, CVL, Re-weighting/Re-sample
- CF: RUBi, CF-VQA, LMH

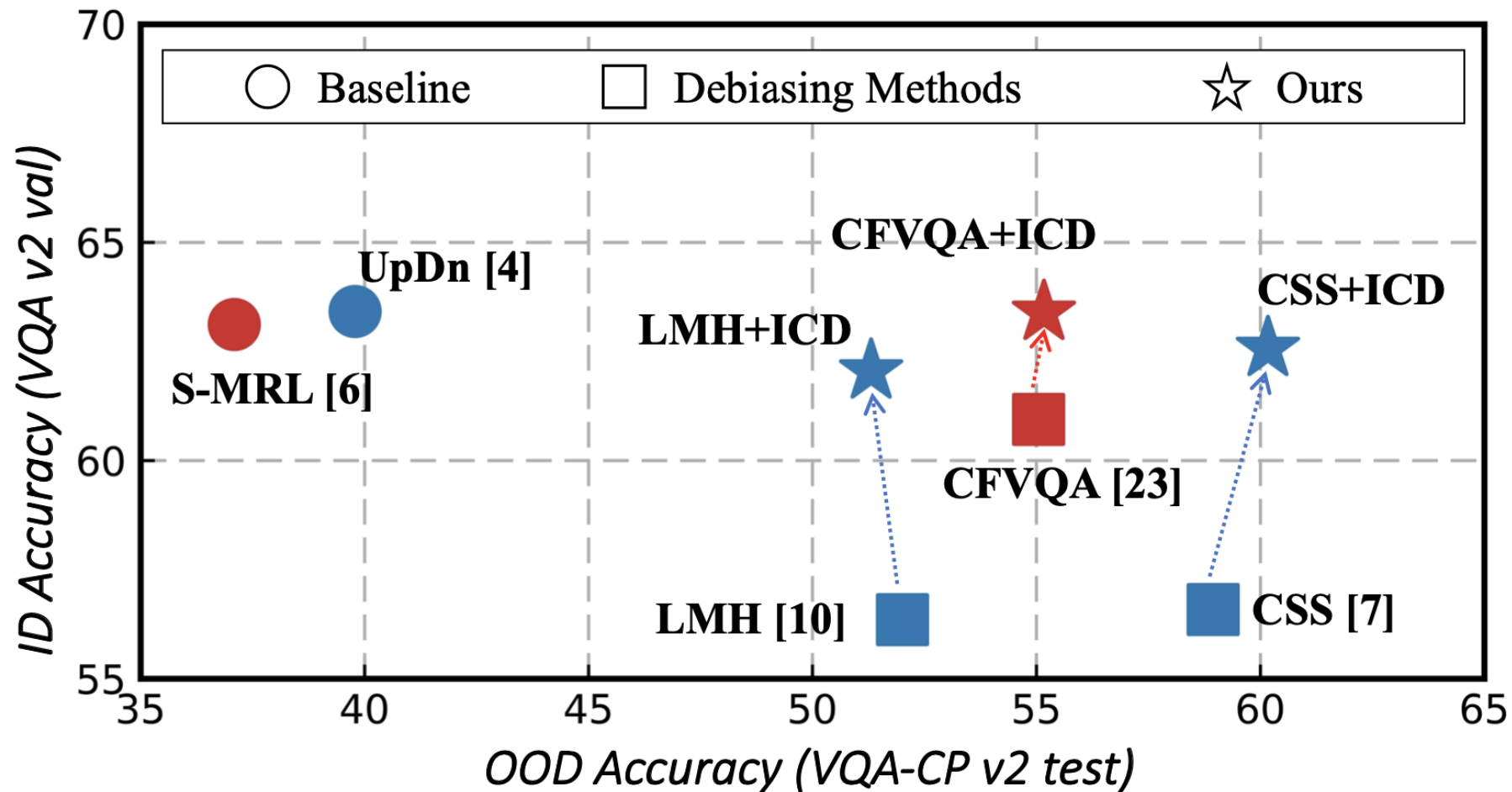
**CSS:** Chen et al. Counterfactual Samples Synthesizing for Robust Visual Question Answering. CVPR'20

**CVL:** Abbasnejad et al. Counterfactual Vision and Language Learning. CVPR'20

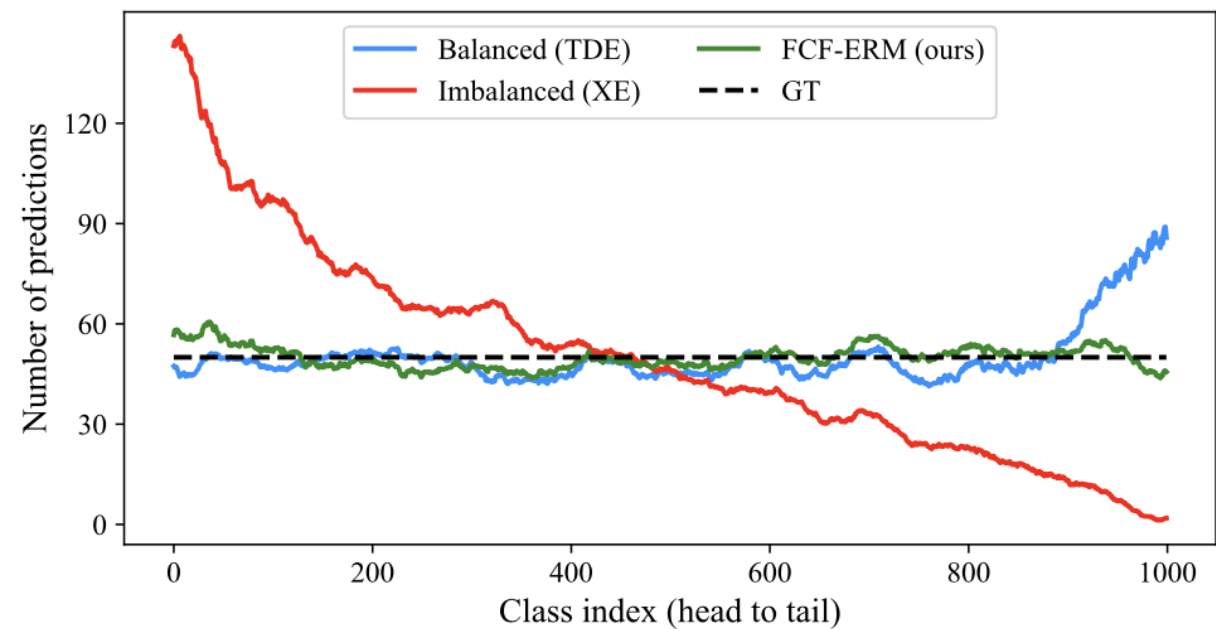
**RUBi:** Cadene et al. RUBi: Reducing Unimodal Biases in Visual Question Answering. NeurIPS'19

**LMH:** Clark et al. Don't Take the Easy Way Out: Ensemble based Methods for Avoiding Known Dataset Biases. EMNLP'19

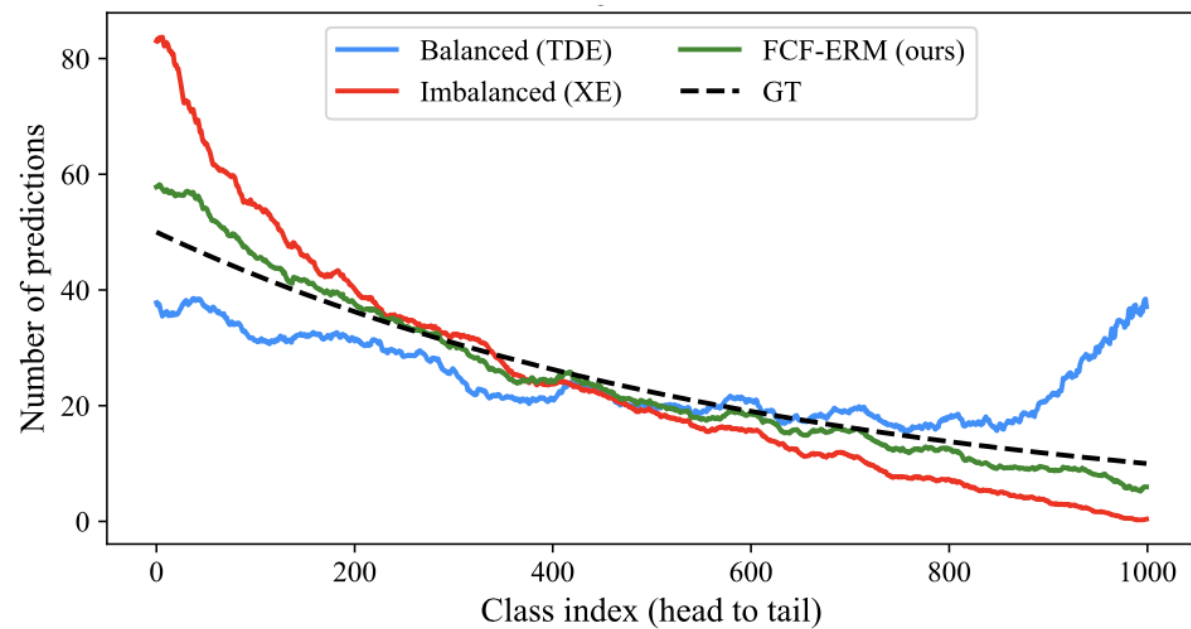
# VQA OOD



# Long-tail OOD



**(a) Balanced Test**



**(b) Imbalanced Test**



# What's new?

- A best of two worlds VQA model
- A best of two worlds long-tailed model

# Introspective Distillation for VQA: Key Idea

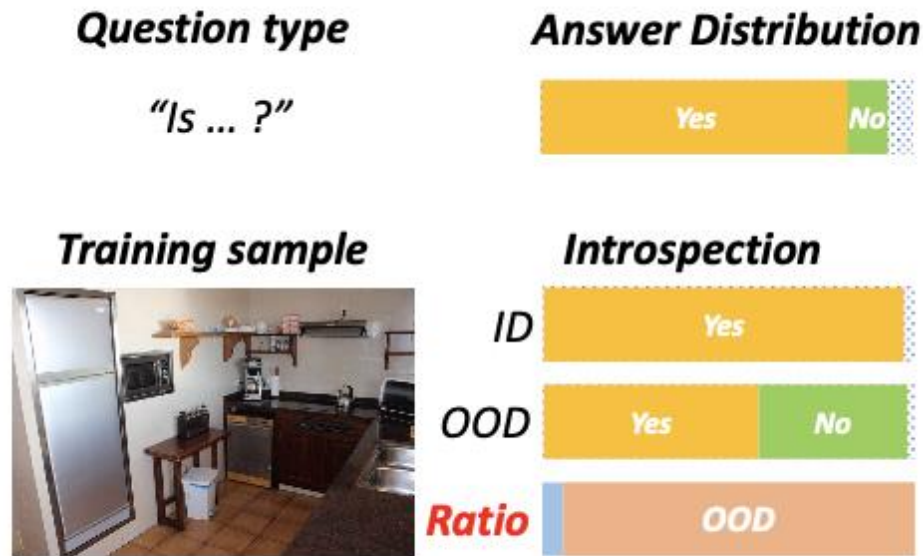
- **ID-Teacher**: Good @ Train = Test, Bad @ Train  $\neq$  Test
- **OOD-Teacher**: Good @ Train  $\neq$  Test, Bad @ Train = Test
- A **Student** learns the best of the two teachers
- By **ONLY** given the train, how does the student know to whom she should listen (**oracle**)?

# Introspective Distillation for VQA: Key Idea



# Introspection: Case 1

- if  $ID\text{-bias} > OOD\text{-bias}$ , then  $ID\text{-teacher} < OOD\text{-teacher}$

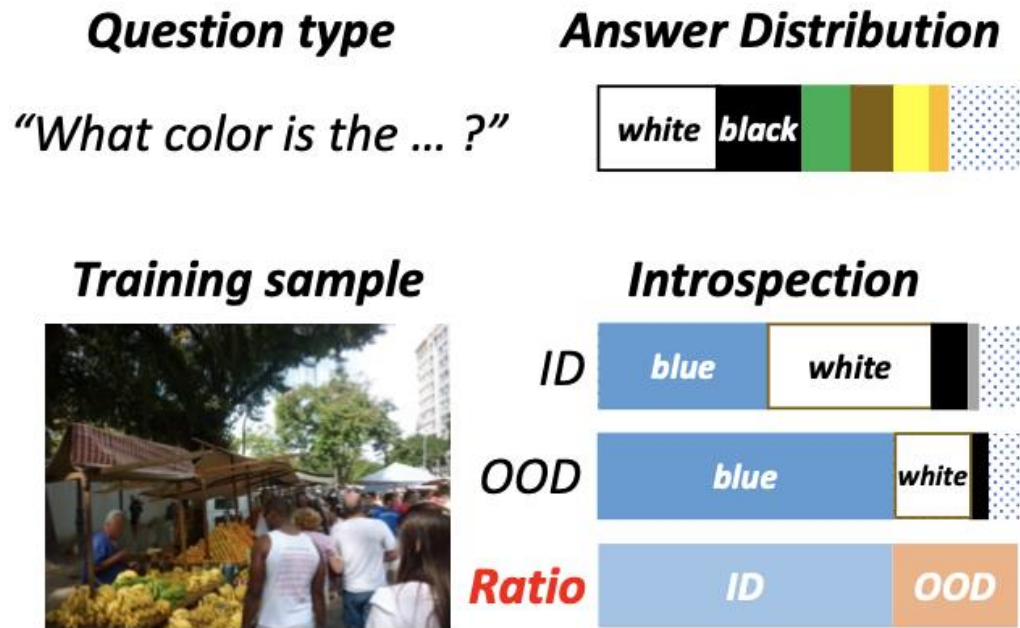


Q: Is that an electric oven? (GT: Yes.)

For each sample,  
If ID-Teacher is too good to be true  
OOD-Teacher not so good,  
 $W(OOD) \propto XE(OOD)/XE(ID)$

# Introspection: Case 2

- if  $ID\text{-bias} < OOD\text{-bias}$ , then  $ID\text{-teacher} > OOD\text{-teacher}$



For each sample,  
 If ID-Teacher is not so good,  
 OOD-Teacher is too good to be true,  
 $W(ID) \propto XE(ID)/XE(OOD)$

Q: What color is the older man's shirt? (GT: Blue.)

# Introspection: Case 3

- if  $ID\text{-bias} \approx OOD\text{-bias}$ , then  $ID\text{-teacher} \approx OOD\text{-teacher}$

## Question type

“How many ... ?”

## Answer Distribution



## Training sample



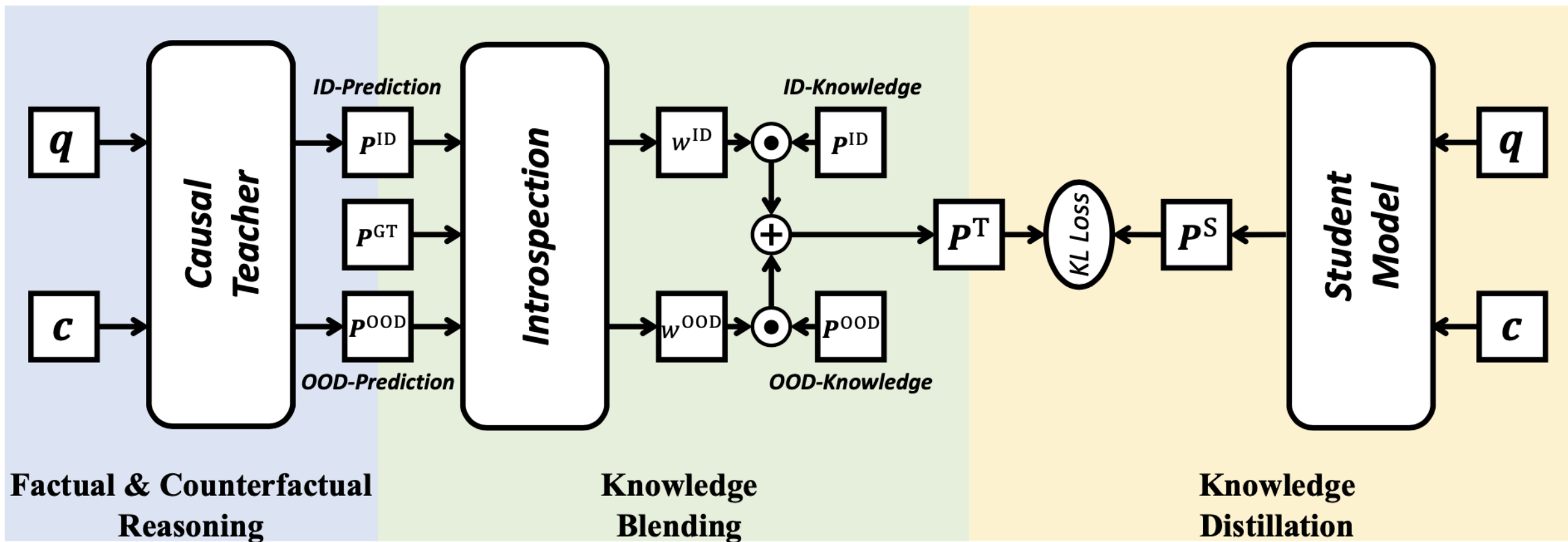
## Introspection



For each sample,  
 If ID/OOD-teachers are similar,  
 $W(ID) \approx W(OOD)$  as  
 $XE(ID) \approx XE(OOD)$

Q: How many skiers? (GT: 3.)

# The Introspective Pipeline



# How does Introspection look like?

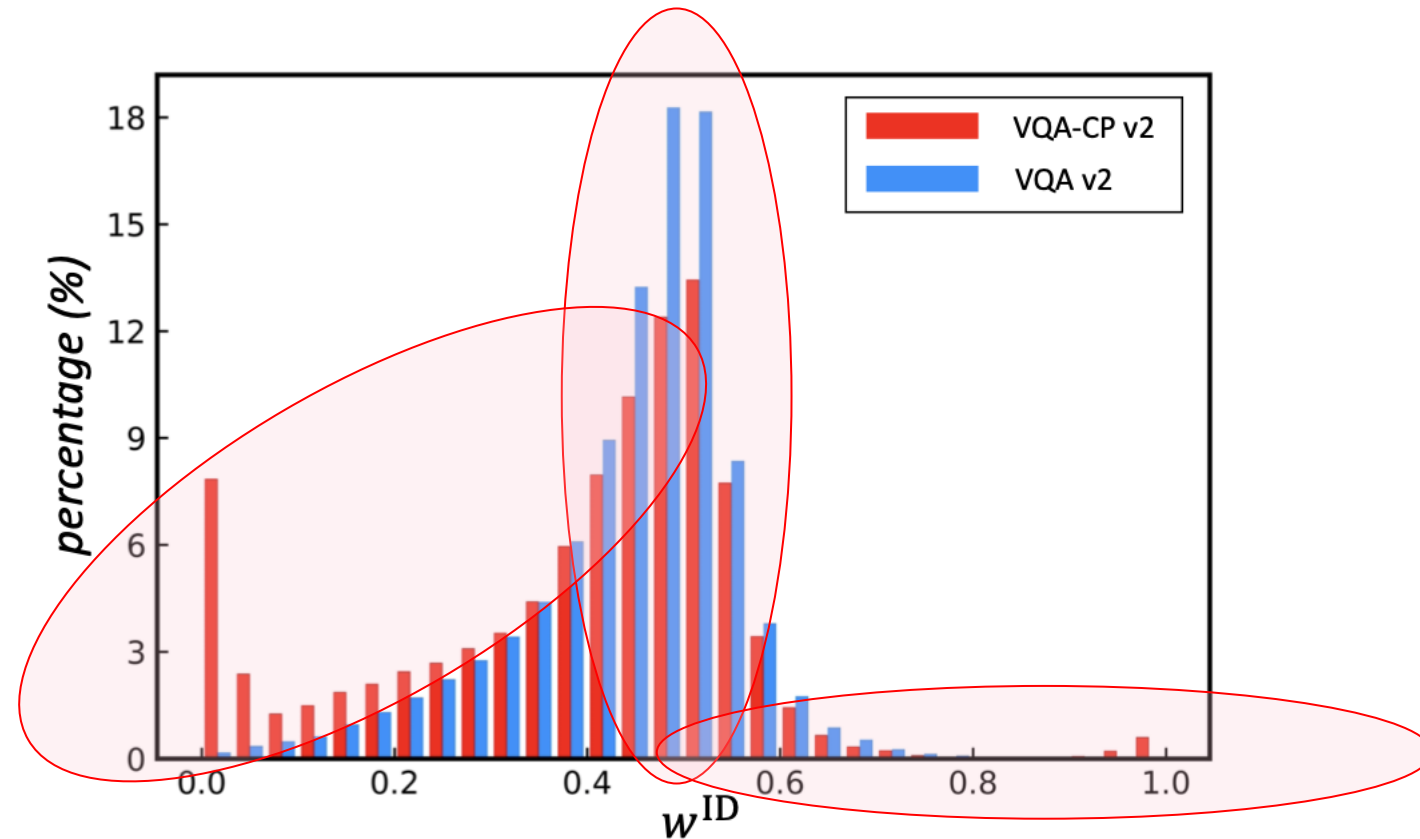


Figure 4: The distribution of  $w^{\text{ID}}$  on the VQA-CP v2 and VQA v2 training sets.



# How does Introspection look like? Both are mostly Case 3

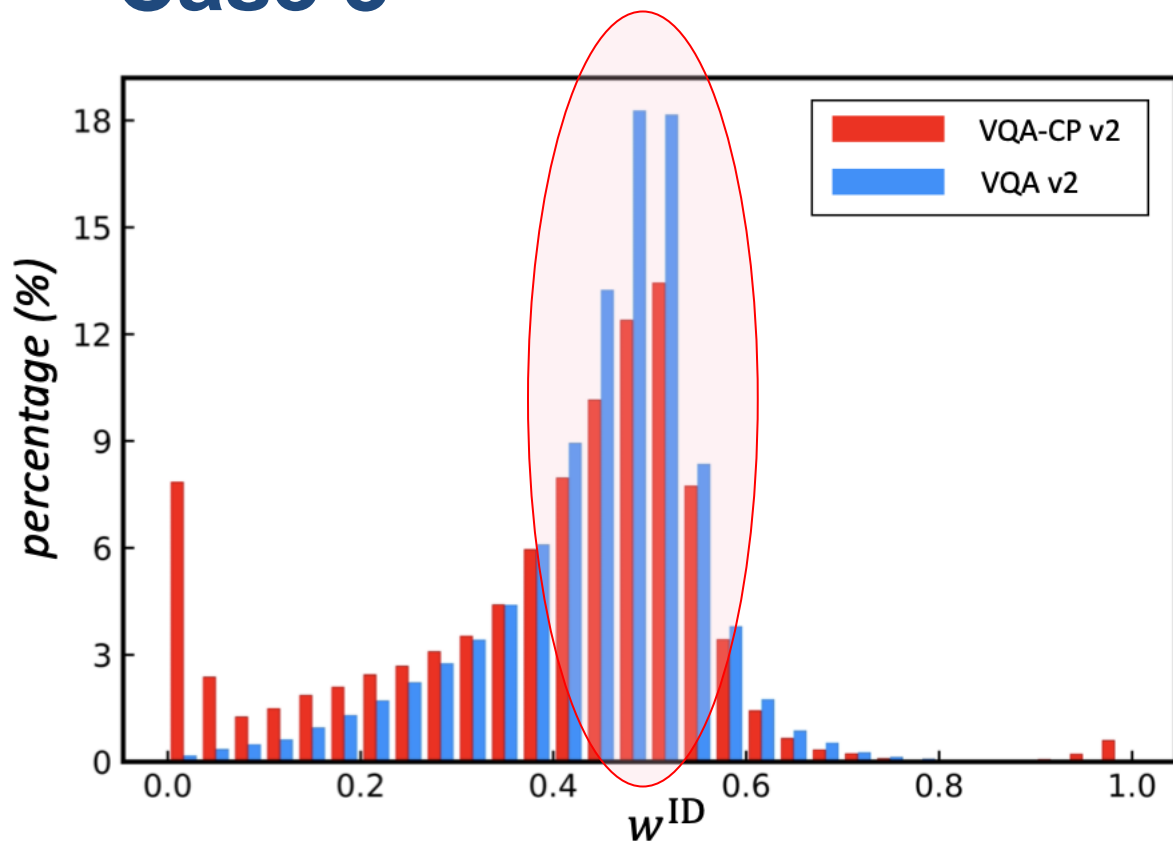
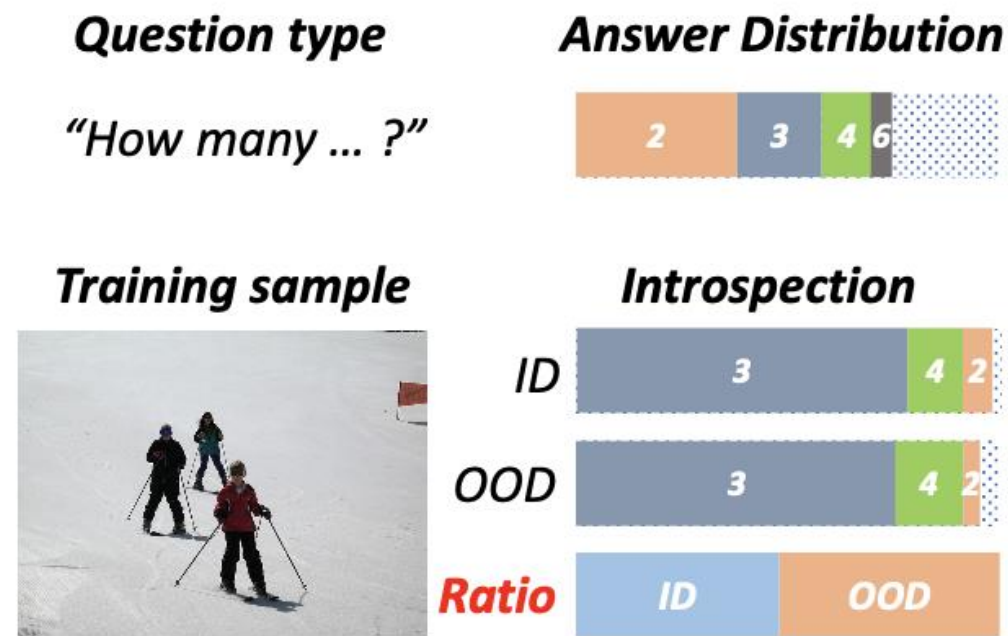


Figure 4: The distribution of  $w^{ID}$  on the VQA-CP v2 and VQA v2 training sets.

*ID-teacher  $\approx$  OOD-teacher*



Q: How many skiers? (GT: 3.)

# How does Introspection look like? VQA-CP has more Case 1 than VQA

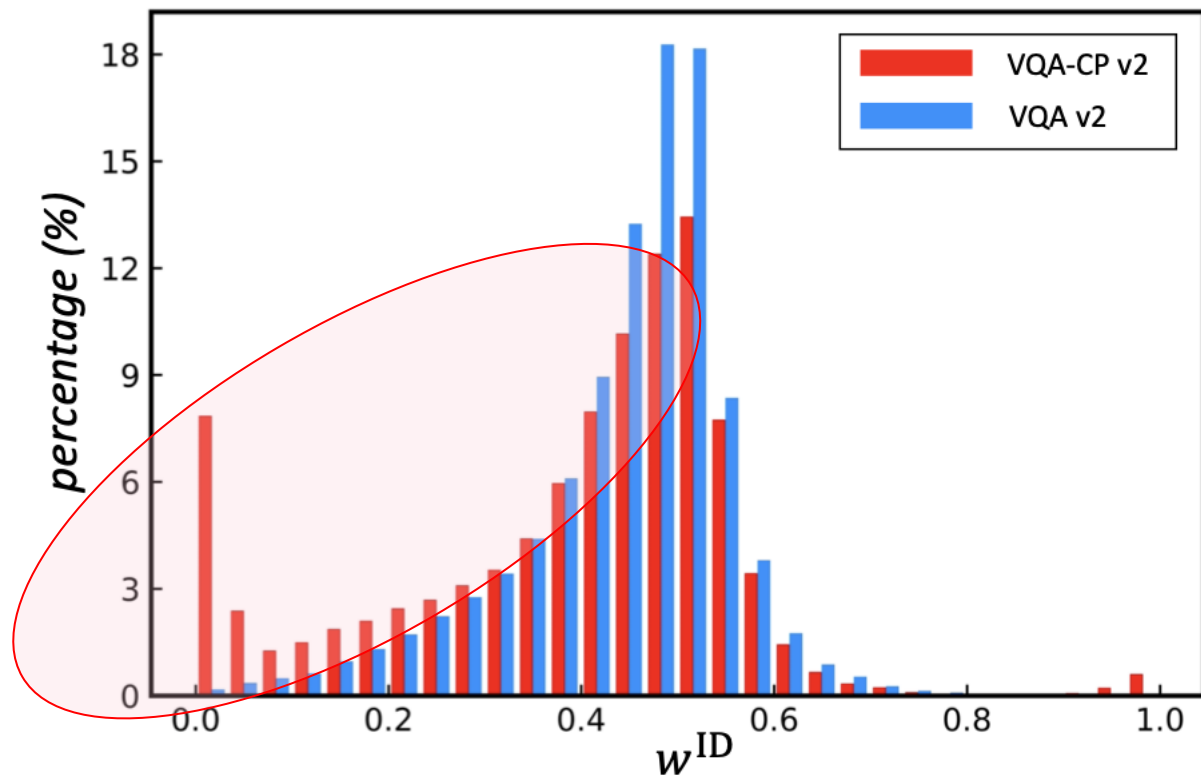


Figure 4: The distribution of  $w^{ID}$  on the VQA-CP v2 and VQA v2 training sets.

*ID-teacher < OOD-teacher*

Question type	Answer Distribution
"Is ... ?"	
Training sample	Introspection
	ID
	OOD
	Ratio
Q: <u>Is</u> that an electric oven? (GT: Yes.)	

# How does Introspection look like? VQA has more Case 2 than VQA-CP

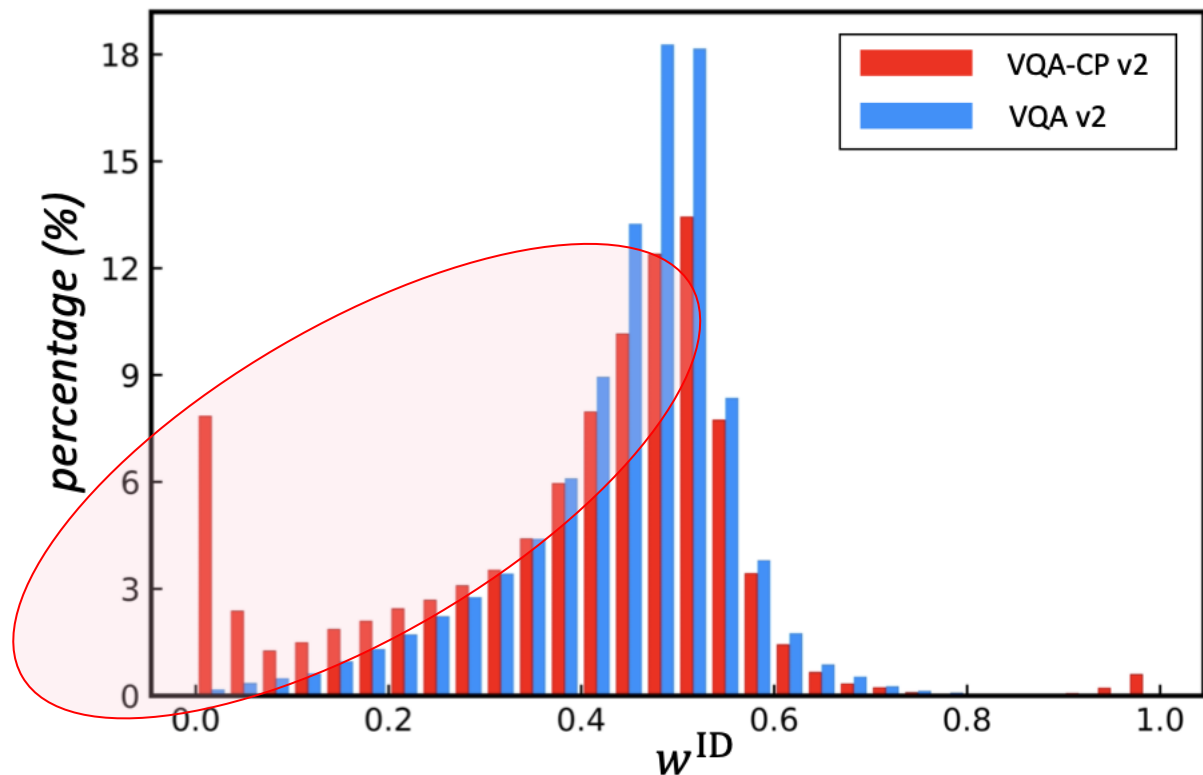
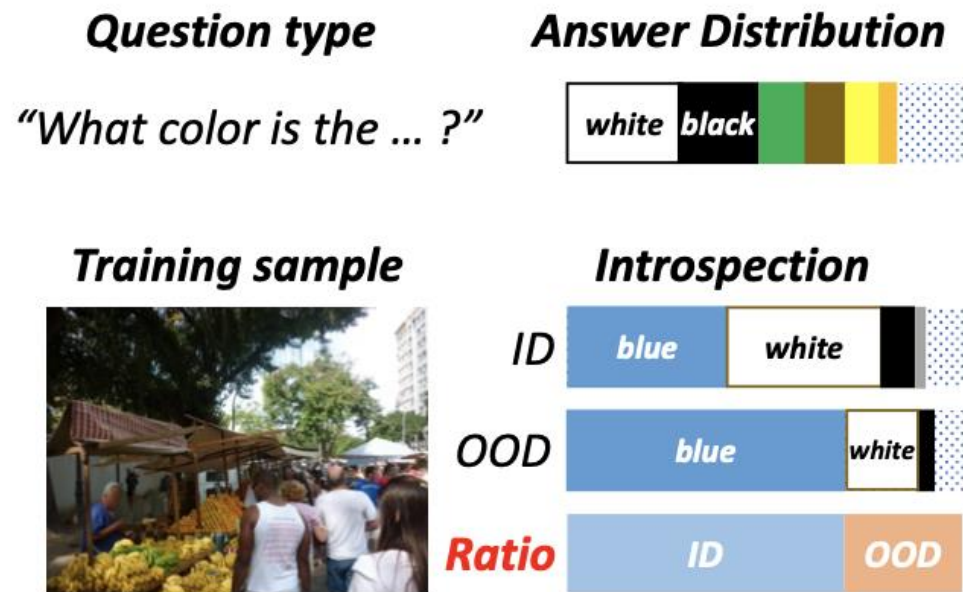


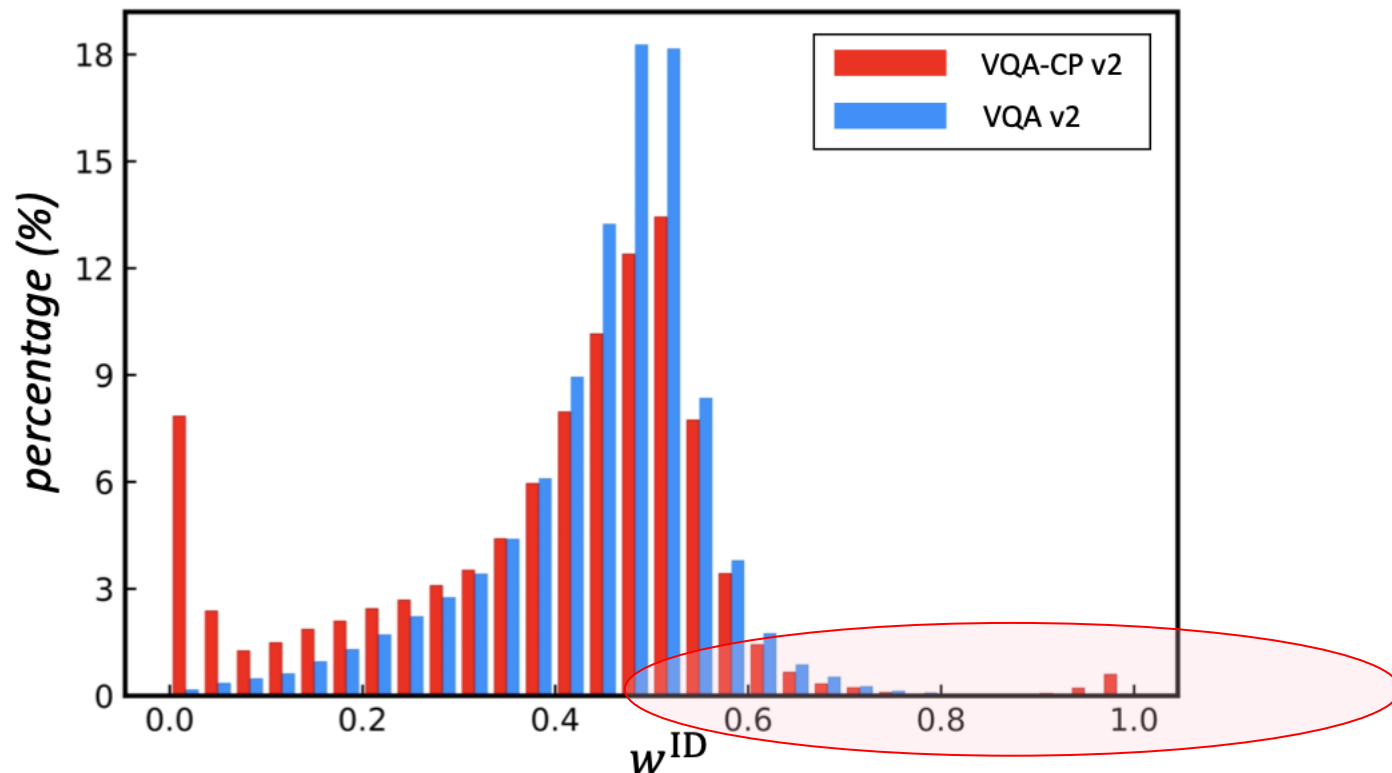
Figure 4: The distribution of  $w^{ID}$  on the VQA-CP v2 and VQA v2 training sets.

*ID-teacher > OOD-teacher*



Q: What color is the older man's shirt? (GT: Blue.)

# How does Introspection look like? Both ID-Teachers are weaker (more biased than OOD-Teachers)



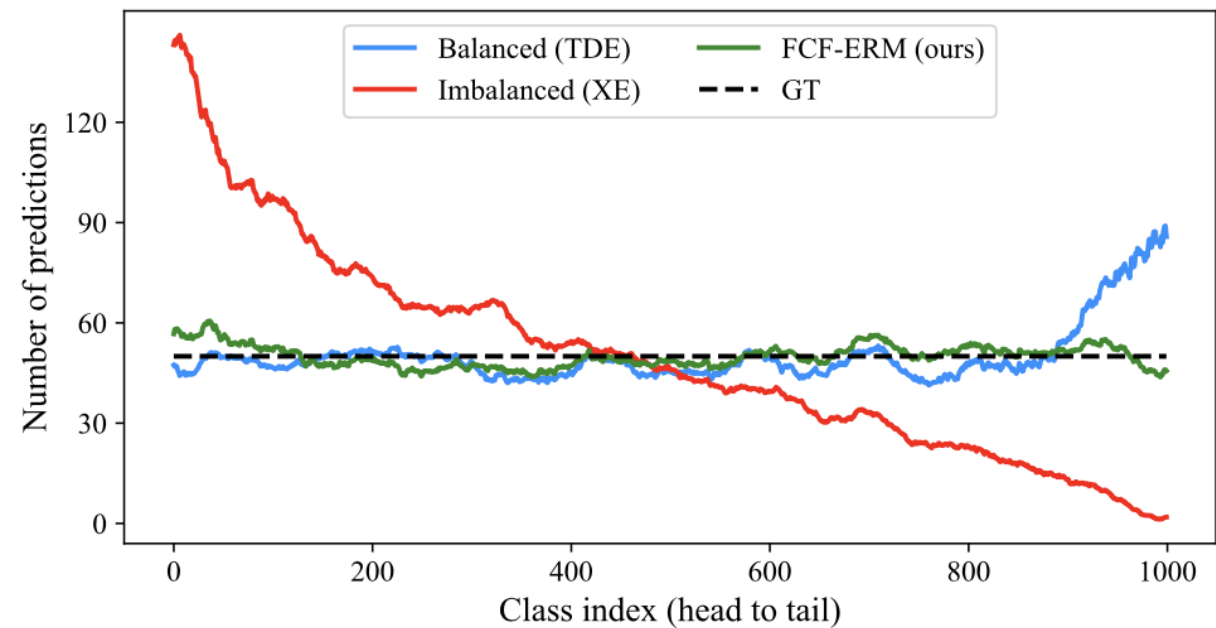
Homework

Figure 4: The distribution of  $w^{\text{ID}}$  on the VQA-CP v2 and VQA v2 training sets.

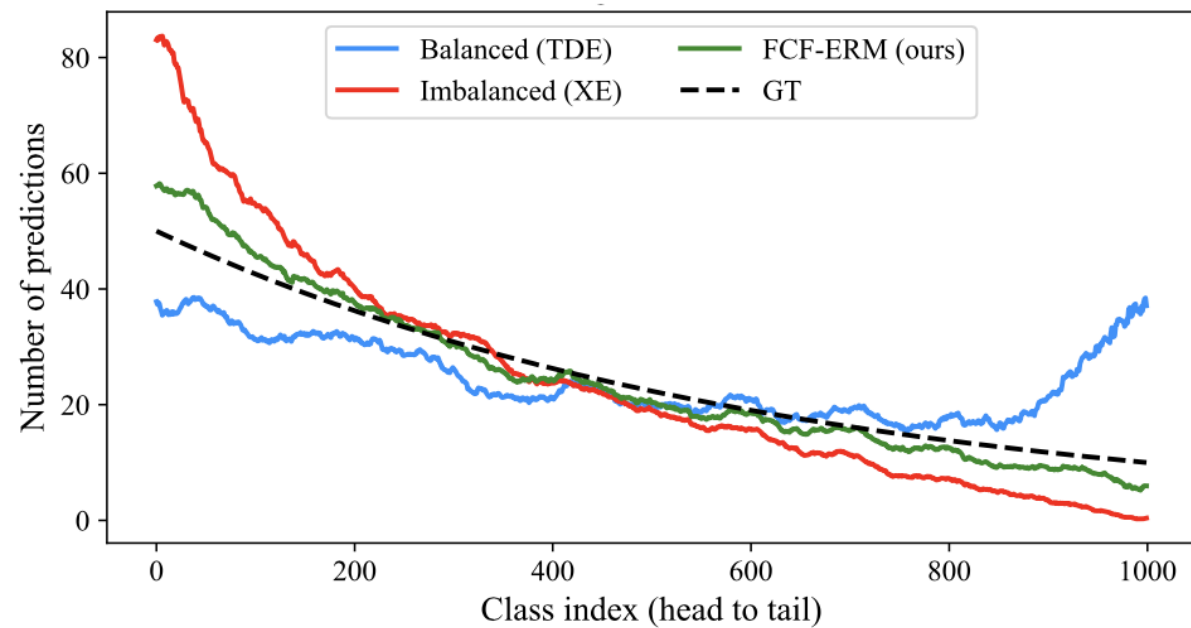
# The best of the two worlds

Methods	VQA-CP v2 test (OOD)				VQA v2 val (ID)				HM
	All	Y/N	Num.	Other	All	Y/N	Num.	Other	
UpDn [4]	39.79	43.23	12.28	45.54	63.42	81.19	42.43	55.47	48.90
LMH [10]	<b>52.01</b>	<b>72.58</b>	<b>31.12</b>	46.97	56.35	65.06	37.63	54.69	54.09
+ IntroD	51.31 <sup>-0.70</sup>	71.39	27.13	<b>47.41</b>	<b>62.05</b> <sup>+5.70</sup>	<b>77.65</b>	<b>40.25</b>	<b>55.97</b>	<b>56.17</b> <sup>+2.08</sup>
CSS [7]	58.95	84.37	<b>49.42</b>	48.21	56.98	65.90	38.19	55.18	57.95
+ IntroD	<b>60.17</b> <sup>+1.22</sup>	<b>89.17</b>	46.91	<b>48.62</b>	<b>62.57</b> <sup>+5.59</sup>	<b>78.57</b>	<b>41.42</b>	<b>56.00</b>	<b>61.35</b> <sup>+3.40</sup>
S-MRL [6]	37.09	41.39	12.46	41.60	63.12	81.83	45.95	53.43	46.72
RUBi [6]	47.60	70.48	<b>20.33</b>	43.09	61.16	81.97	44.86	49.65	53.53
+ IntroD	<b>48.54</b> <sup>+0.96</sup>	<b>73.94</b>	19.43	<b>43.21</b>	<b>61.86</b> <sup>+0.70</sup>	<b>82.40</b>	<b>45.40</b>	<b>50.58</b>	<b>54.40</b> <sup>+0.87</sup>
RUBi-CF [23]	54.90	90.26	<b>34.33</b>	42.01	60.53	81.39	42.87	49.34	57.58
+ IntroD	<b>54.92</b> <sup>+0.02</sup>	<b>90.84</b>	25.17	<b>44.26</b>	<b>63.15</b> <sup>+2.62</sup>	<b>82.44</b>	<b>45.12</b>	<b>53.25</b>	<b>58.75</b> <sup>+1.17</sup>
CF-VQA [23]	55.05	90.61	<b>21.50</b>	45.61	60.94	81.13	43.86	50.11	57.85
+ IntroD	<b>55.17</b> <sup>+0.12</sup>	<b>90.79</b>	17.92	<b>46.73</b>	<b>63.40</b> <sup>+2.46</sup>	<b>82.48</b>	<b>46.60</b>	<b>54.05</b>	<b>58.99</b> <sup>+1.14</sup>

# Current LT is just a “bias flip” game



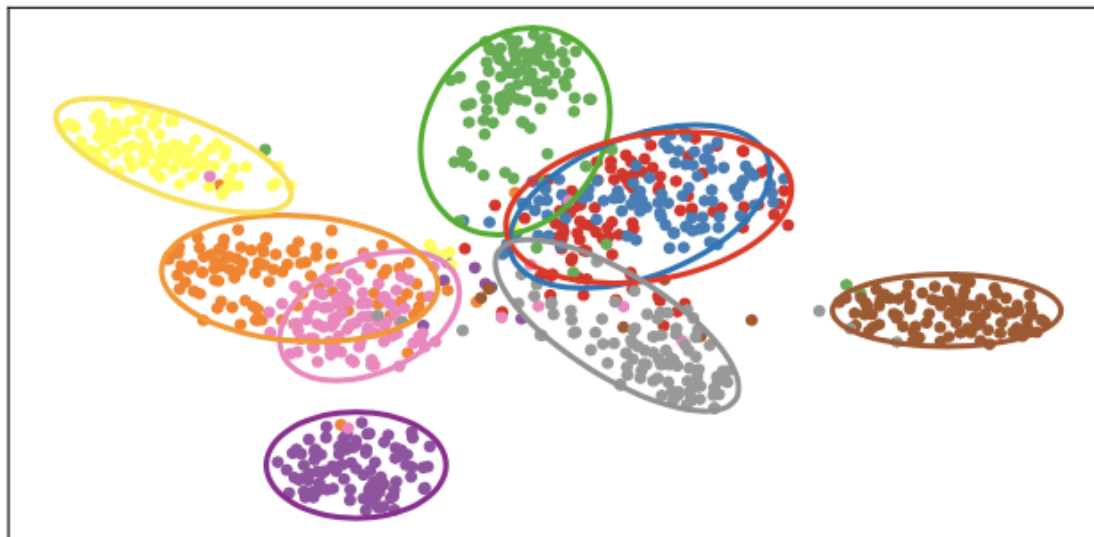
(a) **Balanced Test**



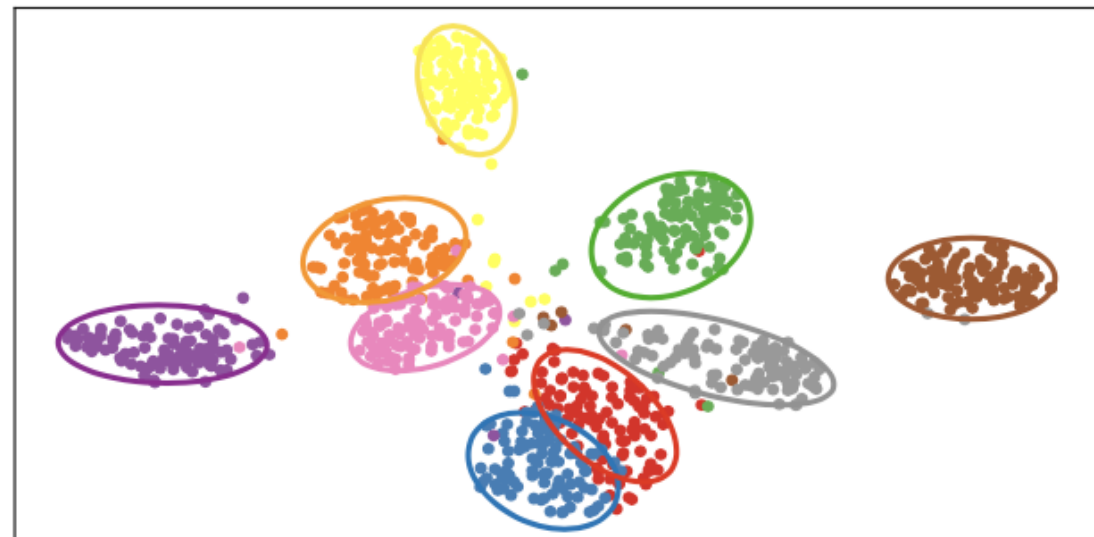
(b) **Imbalanced Test**



# So, it does not truly improve the feature



(c) t-SNE of balanced model (TDE)



(d) t-SNE of FCF-ERM (ours)

# Factual and Counterfactual ERMs Blend: 3 Steps

## Step 1

- Learn a conventional classifier on the imbalanced training data as the ***factual*** model
- Learn a balanced classifier as the ***counterfactual*** model



# Factual and Counterfactual ERMs Blend: 3 Steps

## Step 2: ER Weights

(Factual ER weight)

$$w^f = \frac{(XE^f)^\gamma}{(XE^f)^\gamma + (XE^{cf})^\gamma},$$

(Counterfactual ER weight)

$$w^{cf} = 1 - w^f = \frac{(XE^{cf})^\gamma}{(XE^f)^\gamma + (XE^{cf})^\gamma}.$$

# Factual and Counterfactual ERMs Blend: 3 Steps

## Step 3: Blended ERM

(Factual ER)

$$\mathcal{R}^f(f) = -w^f \sum_i y_i \log f_i(x),$$

where  $y_i$  and  $f_i$  are the ground-truth and the predicted label for  $i$ -th class, respectively.

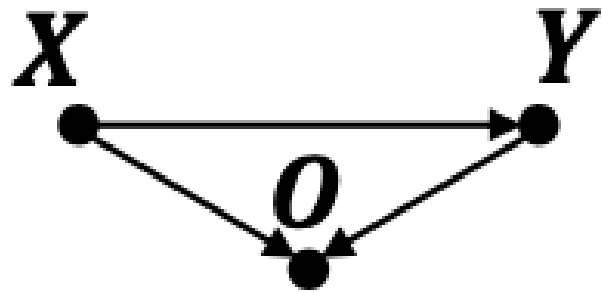
(Counterfactual ER)

$$\mathcal{R}^{\text{cf}}(f) = -w^{\text{cf}} \sum_i \hat{y}_i \log f_i(x),$$

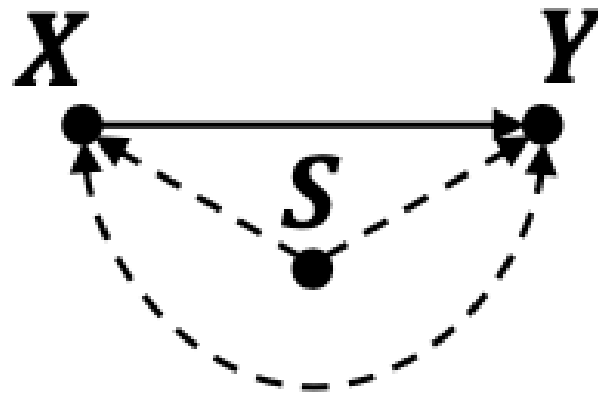
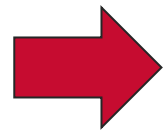
where  $\hat{y}_i = p^{\text{cf}}(y_i|x)$  denotes the balanced prediction for  $i$ -th class.  
The overall empirical risk minimization:

$$\mathcal{R}(f) = \mathcal{R}^f(f) + \mathcal{R}^{\text{cf}}(f).$$

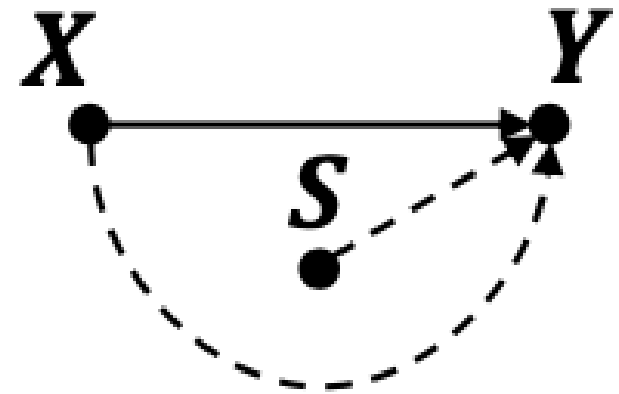
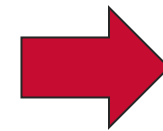
# Why? Selection Bias Removal



(a)



(b)

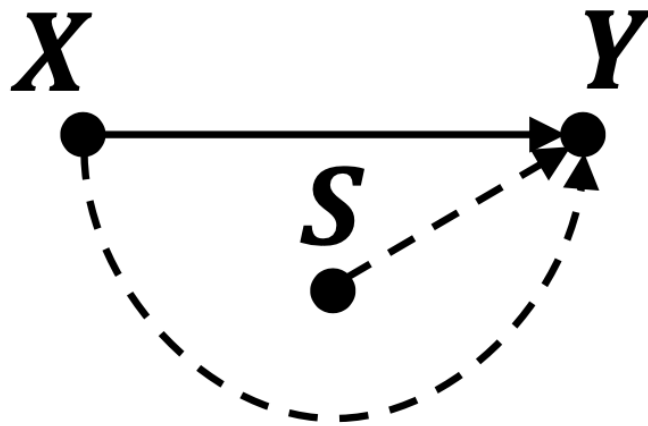


(c)

Reichenbach Principle [raikin-ba:k]

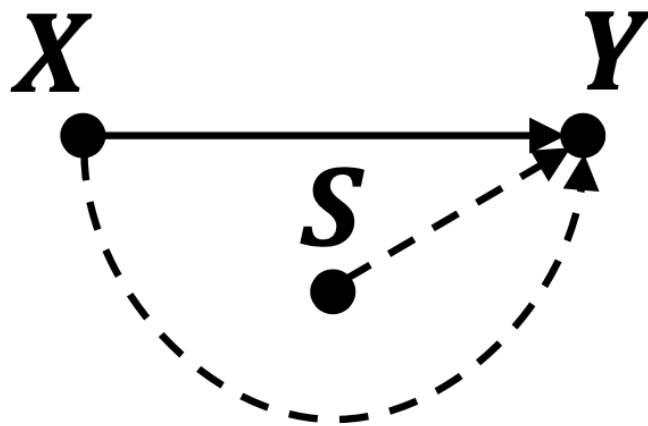
Do-operator

# ERM on the Do-modified graph



$$\mathcal{R}(f) = \mathbb{E}_{x \sim P(X), y \sim P(Y|do(X=x))} \mathcal{L}(y, f(x)) = \sum_x \sum_y \mathcal{L}(y, f(x)) P(y|do(x)) P(x)$$

# Backdoor Adjustment: from “interventional” distribution to “observational” distribution



$$P(y|do(x)) = \sum_{S=s \in \{0,1\}} P(y|x, S = s)P(S = s) = \frac{P(x, y, S = 1)}{P(x|S = 1)} + \frac{P(x, y, S = 0)}{P(x|S = 0)}.$$

## More math

$$\mathcal{R}(f) = \mathbb{E}_{x \sim P(X), y \sim P(Y|do(X=x))} \mathcal{L}(y, f(x)) = \sum_x \sum_y \mathcal{L}(y, f(x)) P(y|do(x)) P(x)$$

$$P(y|do(x)) = \sum_{S=s \in \{0,1\}} P(y|x, S=s) P(S=s) = \frac{P(x, y, S=1)}{P(x|S=1)} + \frac{P(x, y, S=0)}{P(x|S=0)}.$$

# Overall ERM

$(x, y, 1)$  means factual sample, drawn from training data

$(x, y, 0)$  means cf sample, drawn from balanced model

$$\begin{aligned}\mathcal{R}(f) &= \sum_{(x,y)} \sum_{s \in \{0,1\}} \mathcal{L}(y_s, f(x)) \frac{P(x)}{P(x|S=s)} P(x, y, S=s) \\ &= \frac{1}{N} \sum_{(x,y)} \left[ \underbrace{\mathcal{L}(y_{s=1}, f(x)) \frac{P(x)}{P(x|S=1)}}_{\text{factual ER}} + \underbrace{\mathcal{L}(y_{s=0}, f(x)) \frac{P(x)}{P(x|S=0)}}_{\text{counterfactual ER}} \right]\end{aligned}$$

XE loss

# Overall ERM: it explains all

$$\begin{aligned}
 \mathcal{R}(f) &= \sum_{(x,y)} \sum_{s \in \{0,1\}} \mathcal{L}(y_s, f(x)) \frac{P(x)}{P(x|S=s)} P(x, y, S=s) \\
 &= \frac{1}{N} \sum_{(x,y)} \left[ \underbrace{\mathcal{L}(y_{s=1}, f(x))}_{\text{factual ER}} \underbrace{\frac{P(x)}{P(x|S=1)}}_{\text{factual ER}} + \underbrace{\mathcal{L}(y_{s=0}, f(x))}_{\text{counterfactual ER}} \underbrace{\frac{P(x)}{P(x|S=0)}}_{\text{counterfactual ER}} \right] \\
 w^f &= \frac{P(x)}{P(x|S=1)} \propto \frac{(XE^f)^\gamma}{(XE^{cf})^\gamma}, \quad w^{cf} = \frac{P(x)}{P(x|S=0)} \propto \frac{(XE^{cf})^\gamma}{(XE^f)^\gamma}.
 \end{aligned}$$



## The best of the two worlds: balanced test

Methods	Acc	Recall			Precision			F1		
		Many	Med	Few	Many	Med	Few	Many	Med	Few
XE	49.0	68.6	42.9	15.0	46.9	<b>59.1</b>	<b>60.7</b>	55.7	49.7	24.1
$\tau$ -Norm [17]	49.6	61.8	46.2	27.4	52.2	48.5	43.7	56.6	47.3	33.7
LWS [17]	49.9	60.2	47.2	30.3	53.0	49.1	41.3	56.4	48.1	35.0
LADE [13]	51.7	62.6	49.0	30.4	55.3	50.5	41.2	58.7	49.7	34.9
DiVE [11]	53.1	64.1	50.4	<b>31.5</b>	-	-	-	-	-	-
DisAlign [40]	53.4	61.3	<b>52.2</b>	31.4	-	-	-	-	-	-
PC [13]	48.9	60.4	46.7	23.8	56.3	49.7	32.0	58.3	48.2	27.3
TDE [16]	51.8	62.7	49.0	31.4	<b>57.3</b>	52.3	39.5	59.9	50.6	35.0
<b>FCF-ERM<sub>PC</sub></b>	53.2	67.6	49.8	24.0	53.1	55.0	52.4	59.3	51.9	33.0
<b>FCF-ERM<sub>TDE</sub></b>	<b>54.1</b>	<b>68.6</b>	50.0	27.5	53.5	57.3	52.0	<b>60.1</b>	<b>53.4</b>	<b>36.0</b>

## The best of the two worlds: imbalanced test

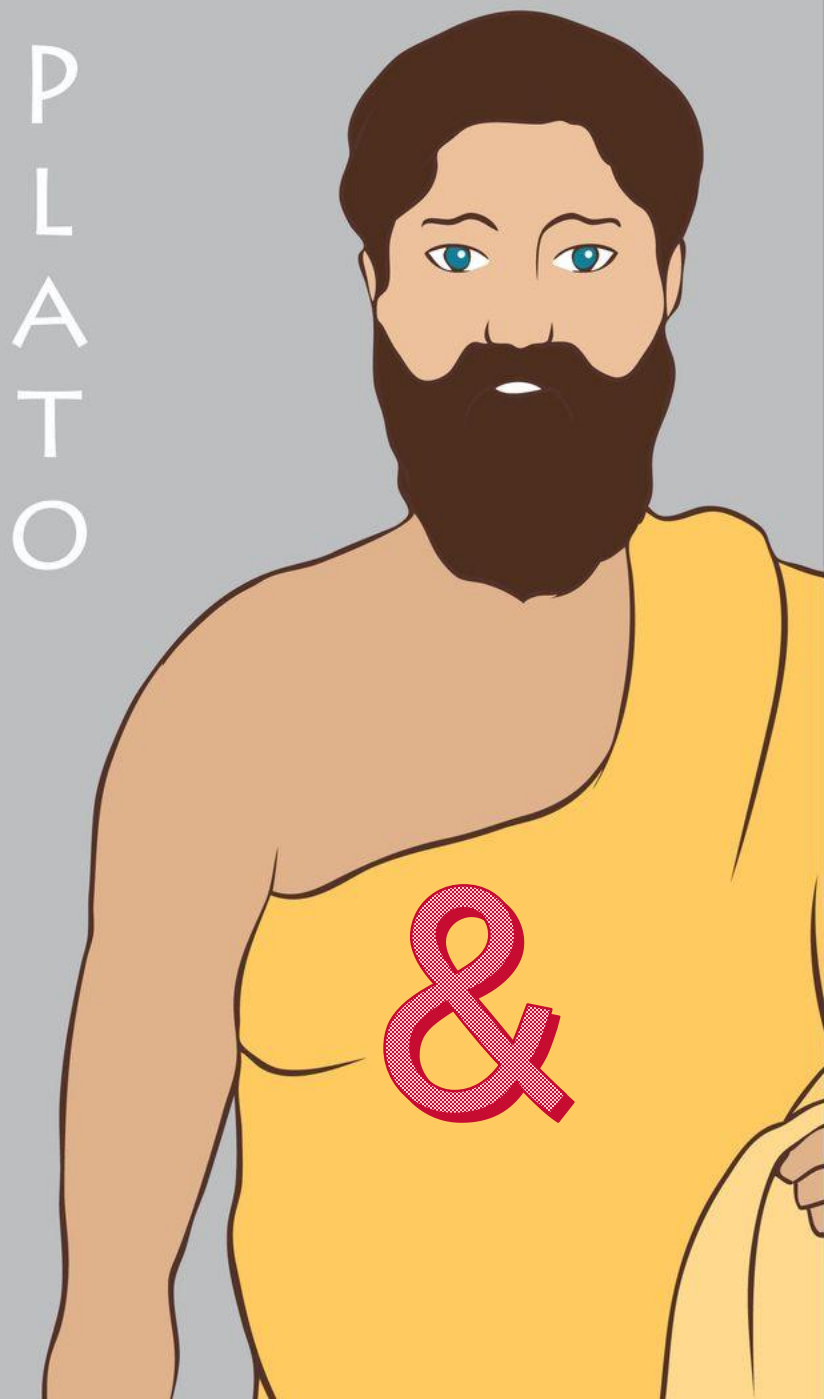
Imbalanced ratio	50	25	10	5
$\tau$ -Norm [17]	59.6	58.2	56.2	54.6
LWS [17]	60.6	59.2	57.0	55.0
PC [13]	58.2	56.8	54.5	52.7
LADE [13]	61.8	60.6	58.6	56.8
TDE [16]	63.0	61.6	59.5	57.6
<b>XE</b>	<b>67.7</b>	65.2	61.4	58.0
<b>FCF-ERM<sub>PC</sub></b>	66.8	65.3	62.5	60.1
<b>FCF-ERM<sub>TDE</sub></b>	<b>67.7</b>	<b>66.0</b>	<b>63.5</b>	<b>60.9</b>

# The best of two worlds: improved feature (LT data trained backbone. Normal classification on balanced data

Backbone	Acc	Recall			Precision			F1		
		Many	Med	Few	Many	Med	Few	Many	Med	Few
CIFAR100										
XE (PC [13])	52.6	60.3	51.9	44.4	59.6	51.1	44.4	60.0	51.5	44.4
TDE [16]	52.6	60.4	51.7	44.4	59.5	51.0	44.5	60.0	51.4	44.5
LADE [13]	53.9	58.7	53.8	47.8	60.2	54.5	47.1	59.4	54.1	47.4
<b>FCF-ERM<sub>TDE</sub></b>	55.1	62.8	54.5	46.7	61.7	53.9	48.1	62.3	54.2	47.4
<b>FCF-ERM<sub>PC</sub></b>	<b>55.3</b>	<b>60.9</b>	<b>56.0</b>	<b>48.0</b>	<b>63.7</b>	<b>54.3</b>	<b>48.3</b>	<b>62.3</b>	<b>55.1</b>	<b>48.1</b>
Places365										
XE (PC [13])	43.8	43.8	44.0	43.5	39.9	43.5	49.3	41.7	43.7	46.2
TDE [16]	43.8	43.8	43.9	43.6	39.7	43.6	48.7	41.6	43.8	46.0
LADE [13]	44.3	42.9	45.9	43.1	43.4	45.1	45.7	43.1	45.5	44.4
<b>FCF-ERM<sub>TDE</sub></b>	44.6	44.1	45.3	44.0	40.4	44.9	49.5	42.1	45.1	46.6
<b>FCF-ERM<sub>PC</sub></b>	<b>46.6</b>	<b>45.1</b>	<b>48.2</b>	<b>46.0</b>	<b>44.2</b>	<b>49.0</b>	<b>53.3</b>	<b>44.6</b>	<b>48.6</b>	<b>49.4</b>
ImageNet										
XE (PC [13])	56.5	64.5	53.8	43.2	59.8	55.1	50.6	62.1	54.4	46.6
TDE [16]	56.5	64.4	53.8	43.7	60.2	55.2	49.8	62.2	54.5	46.6
LADE [13]	57.9	62.6	55.7	52.2	62.4	56.5	52.9	62.5	56.1	52.5
<b>FCF-ERM<sub>TDE</sub></b>	58.9	66.5	56.4	46.2	62.1	57.8	<b>63.2</b>	64.2	57.1	49.4
<b>FCF-ERM<sub>PC</sub></b>	<b>60.2</b>	<b>64.8</b>	<b>58.2</b>	<b>53.8</b>	<b>64.9</b>	<b>58.3</b>	53.9	<b>64.8</b>	<b>58.2</b>	<b>53.8</b>



S  
O  
C  
R  
A  
T  
E  
S



P  
L  
A  
T  
O



A  
R  
I  
S  
T  
O  
T  
L  
E