

Decompose-and-Compose: A Compositional Approach to Mitigating Spurious Correlation

Fahimeh Hosseini Noohdani, Parsa Hosseini, Aryan Yazdanparast, Hamidreza Yaghoubi Araghi, Mahdiah Soleymani Baghshah

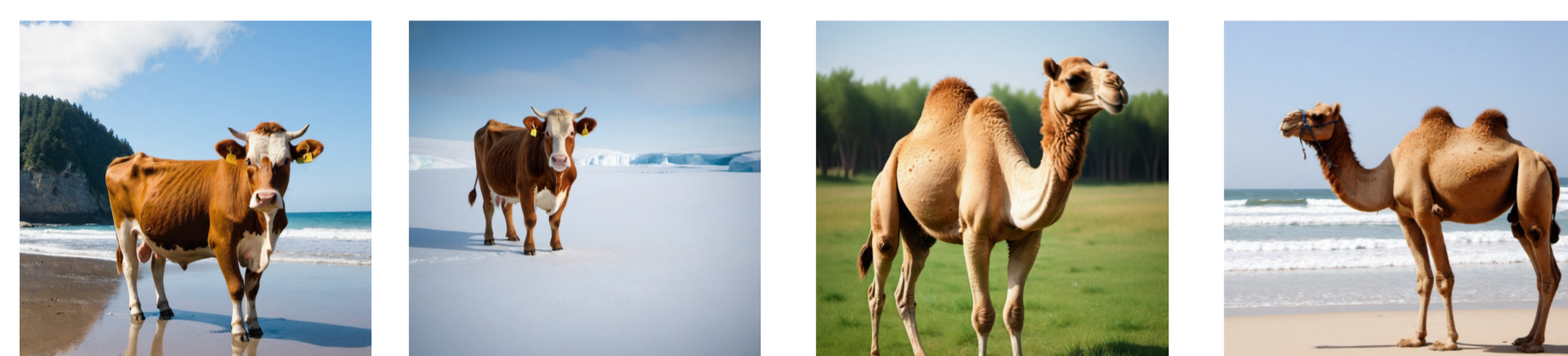
Sharif University of Technology

Spurious Correlation

When training a deep model on a classification task, there is a chance that the model relies on spurious correlation between parts of images and the label^[2].



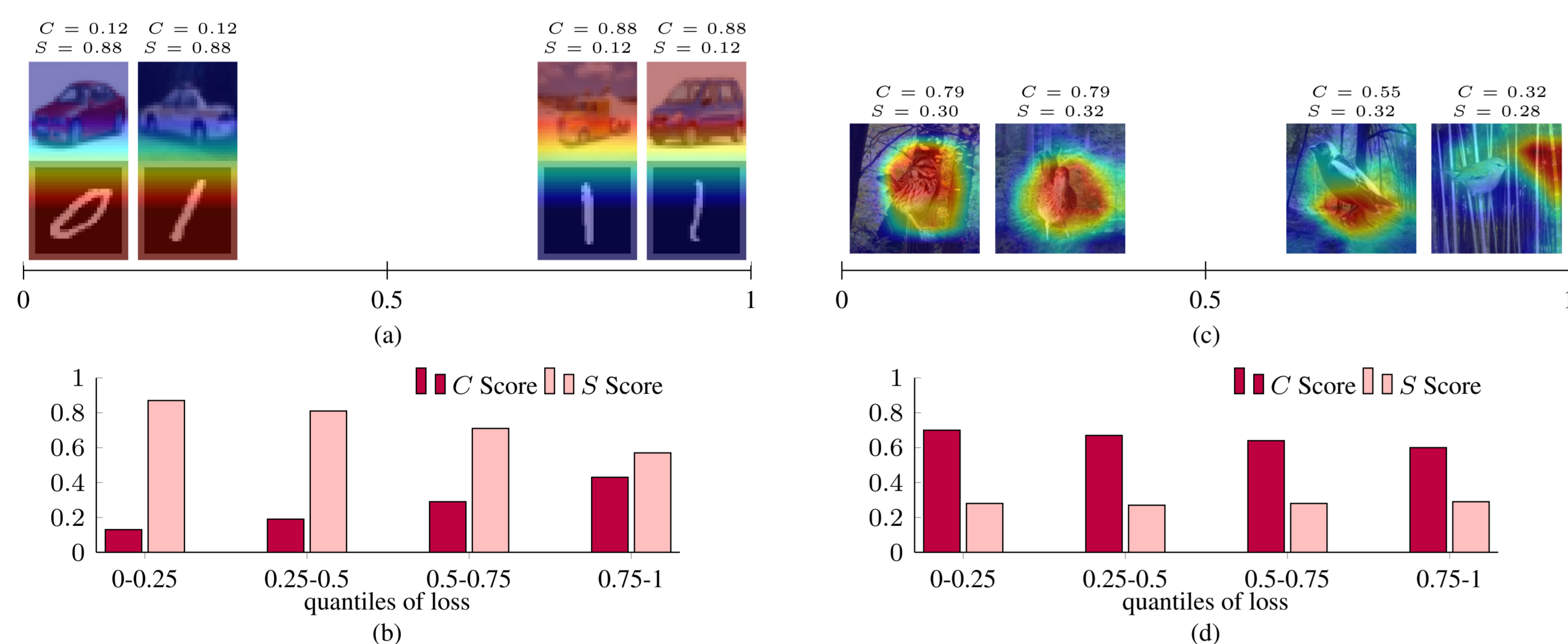
However, the test samples may not come from the train data distribution!



When such correlation is absent in the test data, the models accuracy drops. The Goal is to train a classifier that is robust to spurious correlation.

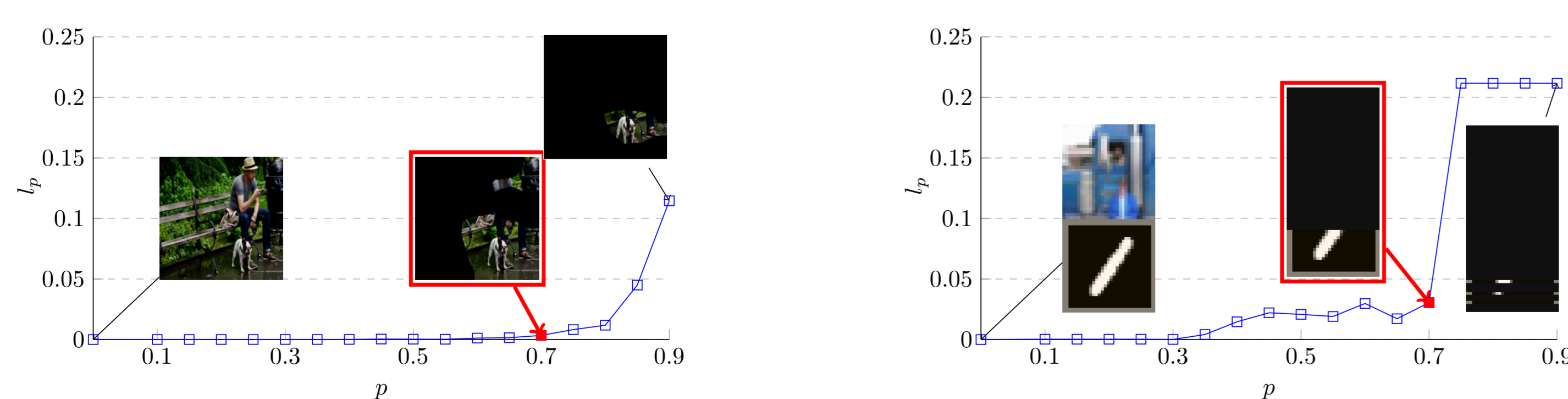
The Focal Regions of Models

Based on the easiness of inferring the label from the spurious parts, the models may attend more either to the spurious parts or the parts that are the real causes of the label.



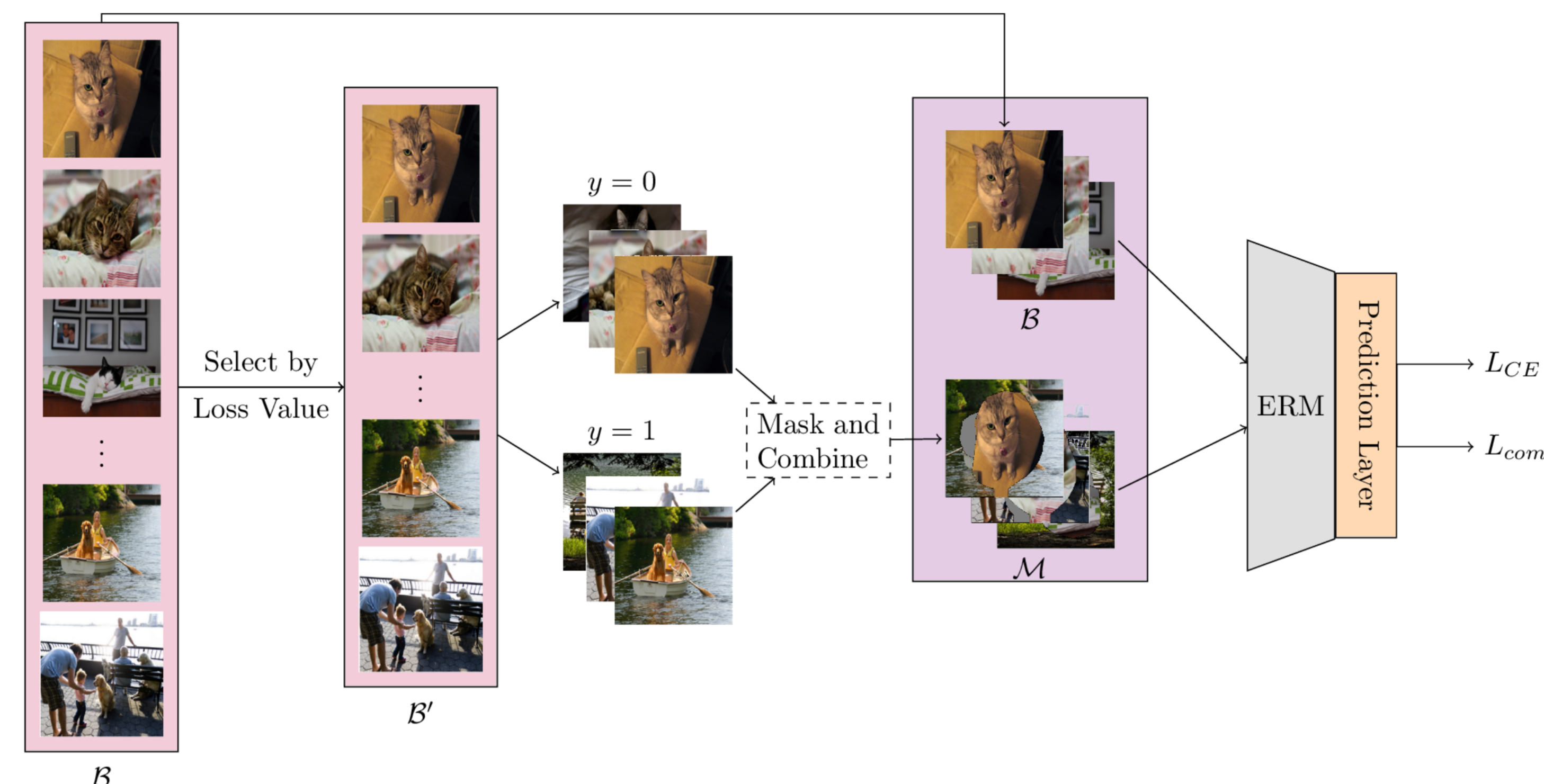
The amount of a model's attention on the C (core) and S (spurious) parts varies significantly among two different datasets exhibiting spurious correlation. This is more evident in samples on which the model has a low loss.

Adaptive Masking



When the most predictive pixels are masked, the loss increases rapidly. Therefore, the optimal masking portion is just before this point!

Decompose and Compose

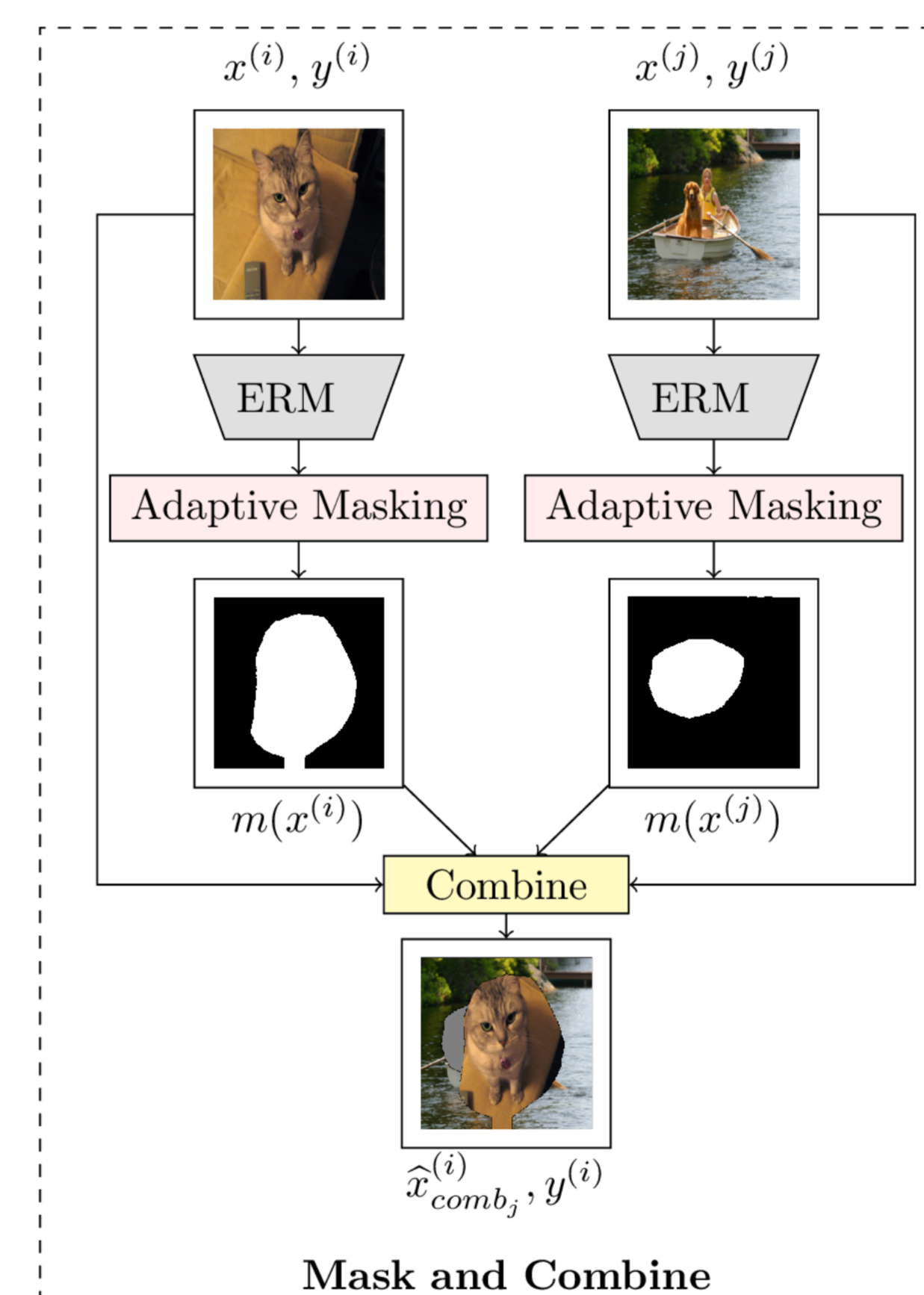


DaC

- Choose one of these two assumptions: The mask, or its inverse capture the core parts more.
 - The correct choice is determined by the worst validation group accuracy.
- Select low-loss images.
- Mask and Combine: Merge the core part of one image with the spurious part of another.
- Retrain Prediction Layer with the augmented data.

Justification

- By combining low-loss (majority) samples, we create samples representing minority groups.
- The augmented data is more group-balanced.
- retraining the last layer on the balanced augmented data makes the model more robust to spurious correlation.



Results

Method	Group Info	Waterbirds		CelebA		Metashift		Dominoes	
		Worst	Average	Worst	Average	Worst	Average	Worst	Average
DFR ^[3]	X/✓✓	92.3±0.2	93.3±0.5	88.3±1.1	91.3±0.3	72.8±0.6	77.5±0.6	90±0.4	92.3±0.2
Group DRO ^[5]	✓/✓	91.4±1.1	93.5±0.3	88.9±2.3	92.9±0.2	66.0±3.8	73.6±2.1	-	-
LISA ^[6]	✓/✓	89.2±0.6	91.8±0.3	89.3±1.1	92.4±0.4	59.8±2.3	70.0±0.7	-	-
MaskTune ^[1]	X/X	86.4±1.9	93.0±0.7	79.4	89.5	66.3±6.3	73.1±2.2	65.8±4.7	85.6±0.7
CnC ^[7]	X/✓	88.5±0.3	90.9±0.1	88.8±0.9	89.9±0.5	-	-	-	-
JTT ^[4]	X/✓	86.7	93.3	81.1	88.0	64.6±2.3	74.4±0.6	-	-
Base (ERM)	X/X	70.8±0.5	91.6±0.1	41.7	96.0	61.3±3.4	73.9±1.5	72.8±1.6	88.5±0.3
DaC-C	X/✓	92.6±0.2	94.9±0.2	76.11±0	91.35±0.2	76.0±0.8	80.0±1.4	89.0±0.7	92.2±0.2
DaC	X/✓	92.3±0.4	95.3±0.4	81.9±0.7	91.4±1.1	78.3±1.6	79.3±0.1	89.2±0.1	92.2±0.3

Algorithm

Algorithm: Decompose-and-Compose (DaC)

Input: Model $f_{\theta}(\cdot) = w \circ g_{\theta}(\cdot)$; Dataset \mathcal{D}_{tr} ; Loss function $l(\cdot, \cdot)$; Hyperparameters α, q , causalflag

```

1 for epoch=1, 2, ... K do
2   for batch  $\mathcal{B}$  in  $\mathcal{D}_{tr}$  do
3      $b \leftarrow \text{mean}(\mathcal{B})$ 
4      $\mathcal{B}' \leftarrow q$  portion of samples in  $\mathcal{B}$ 
       with the lowest loss
5      $\mathcal{M} \leftarrow \{\}$ 
6     for each image  $(x, y) \in \mathcal{B}'$  do
7       Pick  $(x', y') \in \mathcal{B}'$  s.t.  $y \neq y'$ 
8        $m \leftarrow \text{AdaptiveMasking}(f, x, y, l)$ 
9        $m' \leftarrow \text{AdaptiveMasking}(f, x', y', l)$ 
10      if causalflag=False then
11         $m \leftarrow 1 - m$ 
12         $m' \leftarrow 1 - m'$ 
13      end
14       $\hat{x}_{comb} = m \odot x$ 
15       $+ (1 - m) \odot (1 - m')x' + (1 - m) \odot m'b$ 
16       $\mathcal{M} \leftarrow \mathcal{M} \cup \{(\hat{x}_{comb}, y)\}$ 
17    end
18     $L_{CE} \leftarrow \frac{1}{|\mathcal{B}'|} \sum_{(x,y) \in \mathcal{B}'} l(f_{\theta}(x), y)$ 
19     $L_{comb} \leftarrow \frac{1}{|\mathcal{M}|} \sum_{(x,y) \in \mathcal{M}} l(f_{\theta}(x), y)$ 
20     $L_{total} \leftarrow L_{CE} + \alpha L_{comb}$ 
21     $w \leftarrow \text{UpdateWeights}(L_{total})$ 
22  end
23 end

```

References

- Saeid Asgari, Aliasghar Khani, Fereshte Khani, Ali Gholami, Linh Tran, Ali Mahdavi-Amiri, and Ghassan Hamarneh. Masktune: Mitigating spurious correlations by forcing to explore. In *Advances in Neural Information Processing Systems*, 2022.
- Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. Springer-Verlag, 2018.
- Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient for robustness to spurious correlations. In *International Conference on Learning Representations*, 2023.
- Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training group information. In *International Conference on Machine Learning*, 2021.
- Shiori Sagawa*, Pang Wei Koh*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2020.
- Huaxiu Yao, Yu Wang, Sai Li, Linjun Zhang, Weixin Liang, James Zou, and Chelsea Finn. Improving out-of-distribution robustness via selective augmentation. In *International Conference on Machine Learning*, 2022.
- Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Re. Correct-n-contrast: a contrastive approach for improving robustness to spurious correlations. In *International Conference on Machine Learning*, 2022.