

# Diffusing DeBias: Synthetic Bias Amplification for Model Debiasing

**Massimiliano Ciranni<sup>\*</sup> <sup>1</sup>, Vito Paolo Pastore<sup>\*</sup> <sup>1,2</sup>, Roberto Di Via<sup>\*</sup> <sup>1</sup>,  
Enzo Tartaglione <sup>3</sup>, Francesca Odone <sup>1</sup>, Vittorio Murino <sup>2,4</sup>**

<sup>1</sup>MaLGa-DIBRIS, University of Genoa, Italy

<sup>2</sup>AI For Good (AIGO), Istituto Italiano di Tecnologia (IIT), Genoa, Italy

<sup>3</sup>Télécom Paris, École Polytechnique, France

<sup>4</sup>Department of Computer Science, University of Verona, Italy

**<sup>\*</sup>Equal Contribution**

**NeurIPS 2025**

*The Thirty-Ninth Annual Conference on Neural Information Processing Systems*

*San Diego (CA), USA, 12/02/2025 – 12/07/2025*

# Bias definition in Image Classification

- Spurious correlations between class labels and samples;
- Shortcuts learned by models to minimize empirical risk;
- Present in most training samples (bias-aligned);
- Absent in a small percentage (bias-conflicting);
- A model learns these spurious correlations (instead of semantic attributes).

$$\mathcal{D}_{train} = \{x_i, y_i, b_i\}_{i=1}^N$$

$$y^{(j)}, j \in \{1, \dots, C\}$$

$$b^{(k)}, k \in \{1, \dots, B\}$$

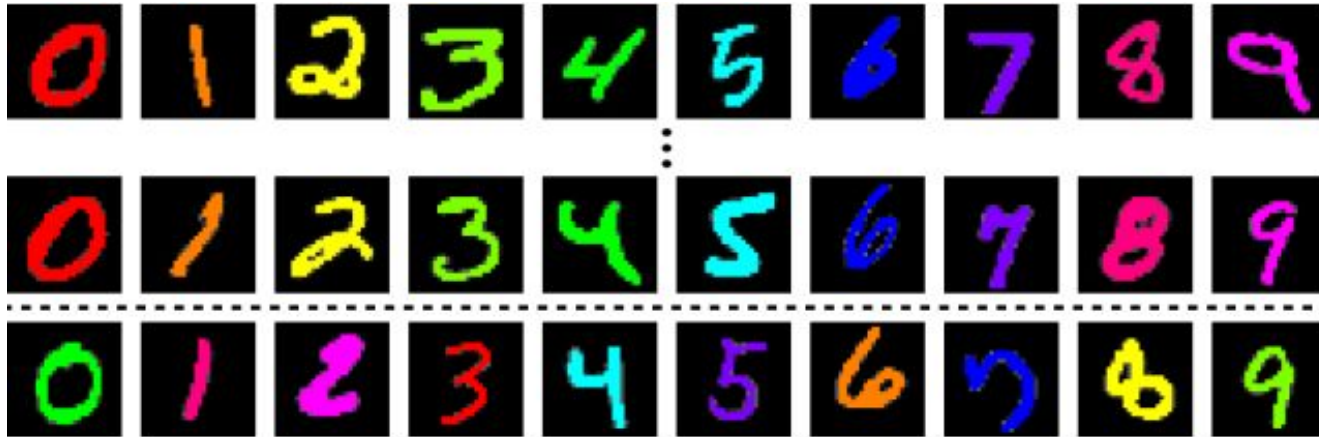
Problematic Bias when:

Most of the samples  $x_i$  belonging to class  $y^{(j)}$  share the same attribute  $b^{(k)}$ ,  
i.e.  $|\mathcal{D}_{bias-aligned}| \gg |\mathcal{D}_{bias-conflicting}|$ .

Training a model in this scenario often results in poor Generalization, i.e.

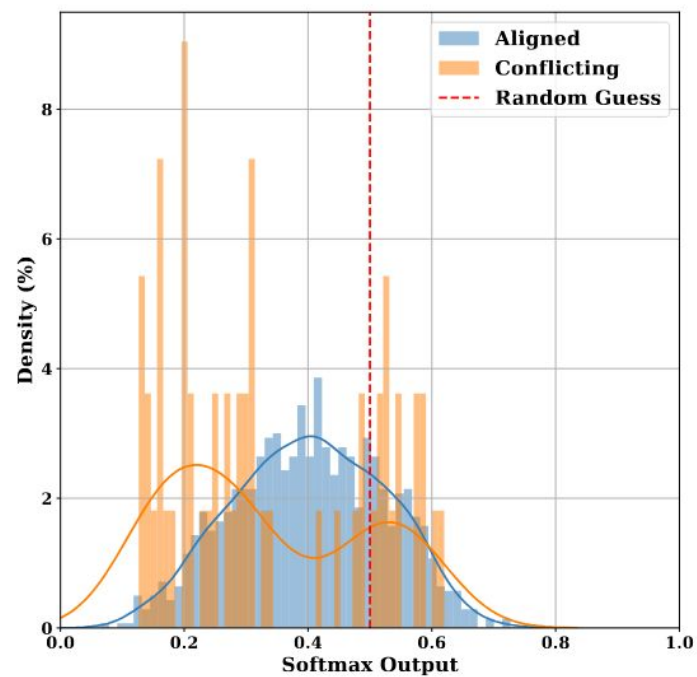
$$\text{Train Error}_{bias-conflicting} \ll \text{Test Error}_{bias-conflicting}$$

Bias-aligned/biased

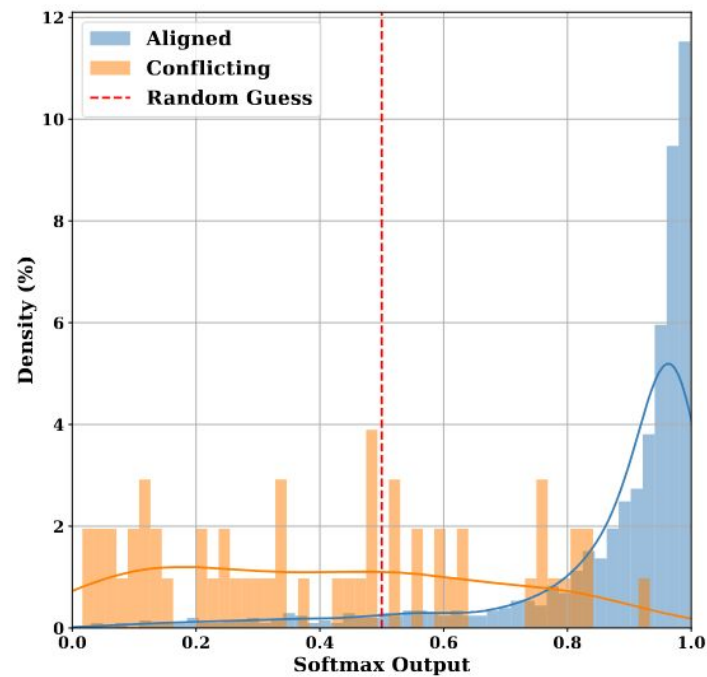


Bias-conflicting/unbiased

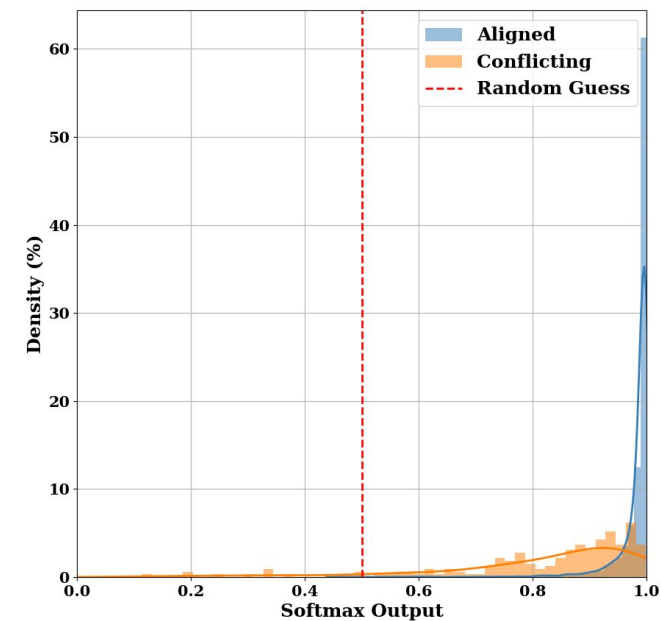




Epoch 0

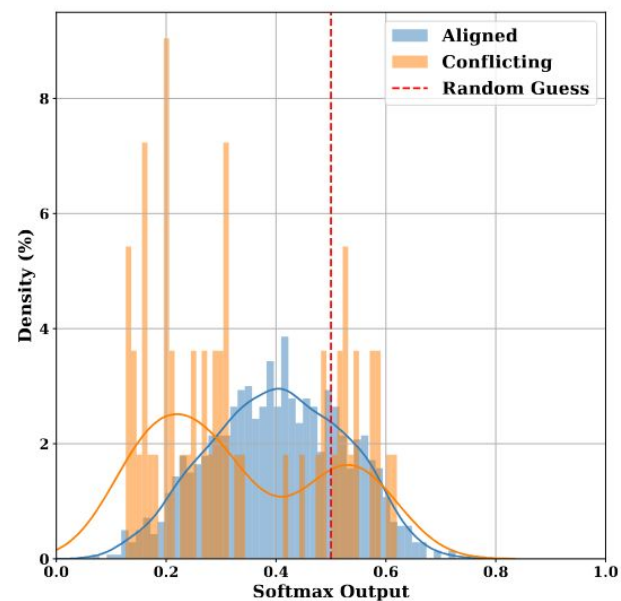


Epoch 6

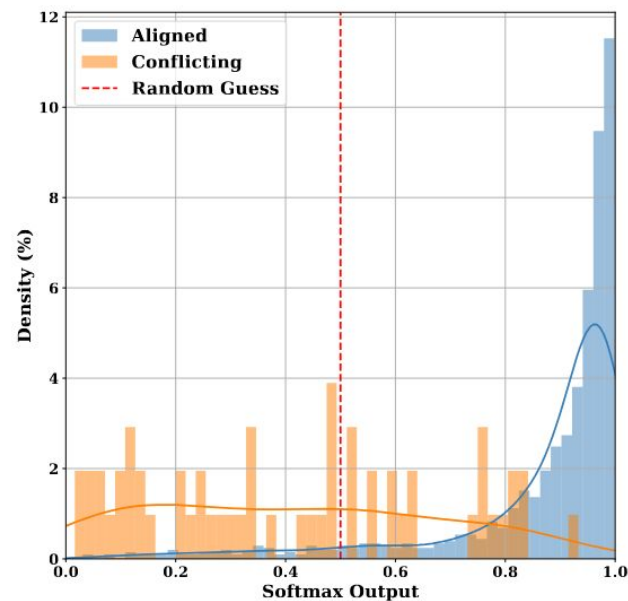
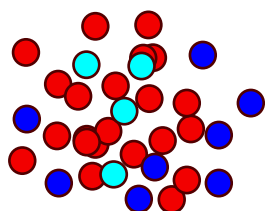


Epoch 10

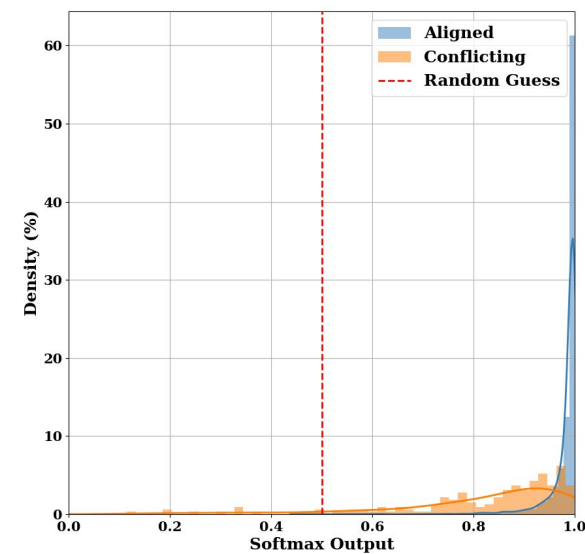
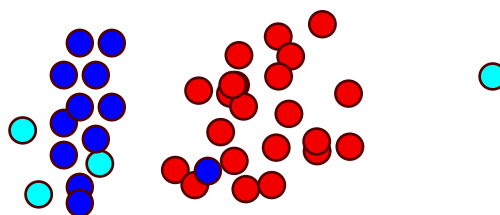
- *Aligned*
- *Conflicting*
- *Other class*



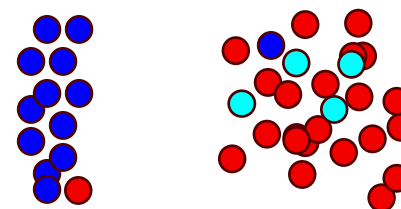
Epoch 0



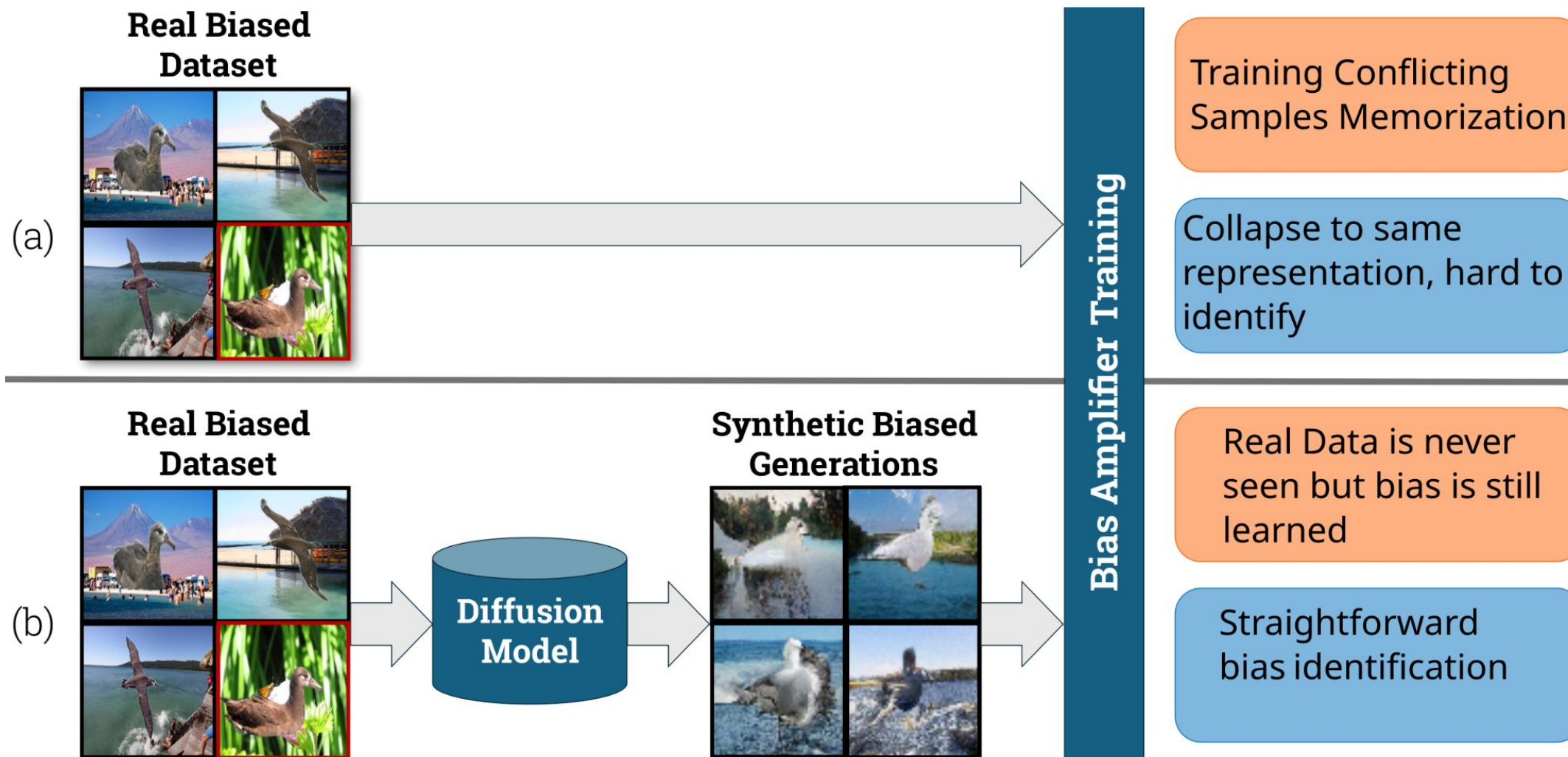
Epoch 6



Epoch 10







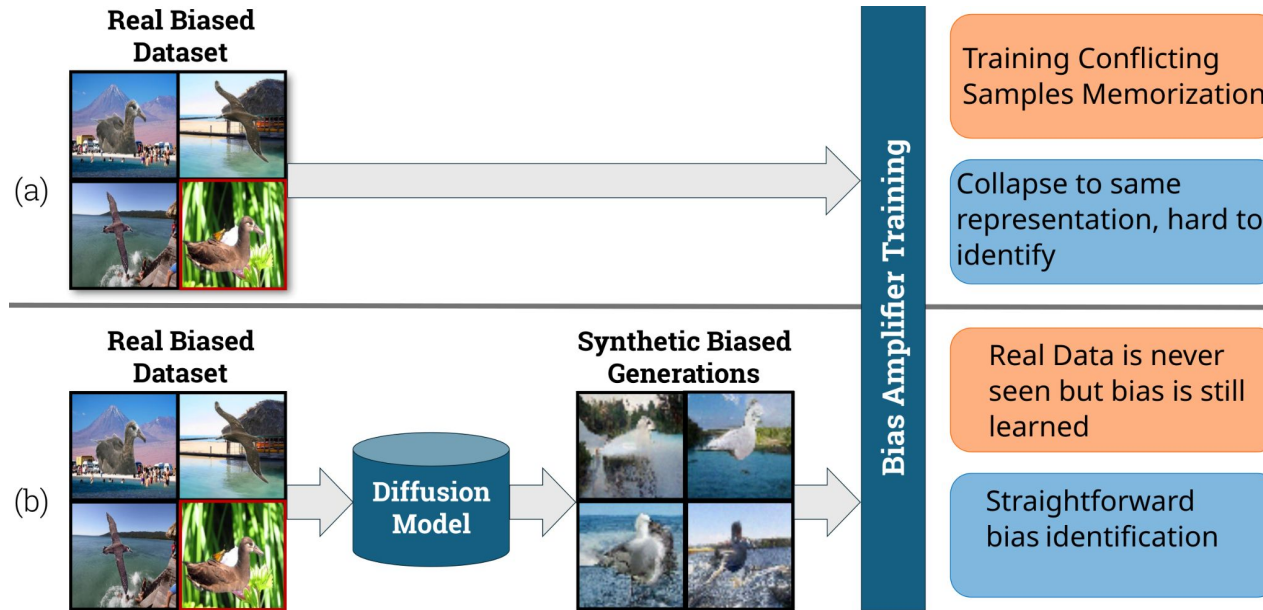
# Assume

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

$$\mathcal{D} = \mathcal{D}_{\text{unbiased}} \cup \mathcal{D}_{\text{biased}}$$

$$\mathcal{D}_{\text{unbiased}} \sim p_{\text{data}}$$

$$\mathcal{D}_{\text{biased}} \sim p_{\text{data}}(\mathbf{x}, y | b)$$



## Assume

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

$$\mathcal{D} = \mathcal{D}_{\text{unbiased}} \cup \mathcal{D}_{\text{biased}}$$

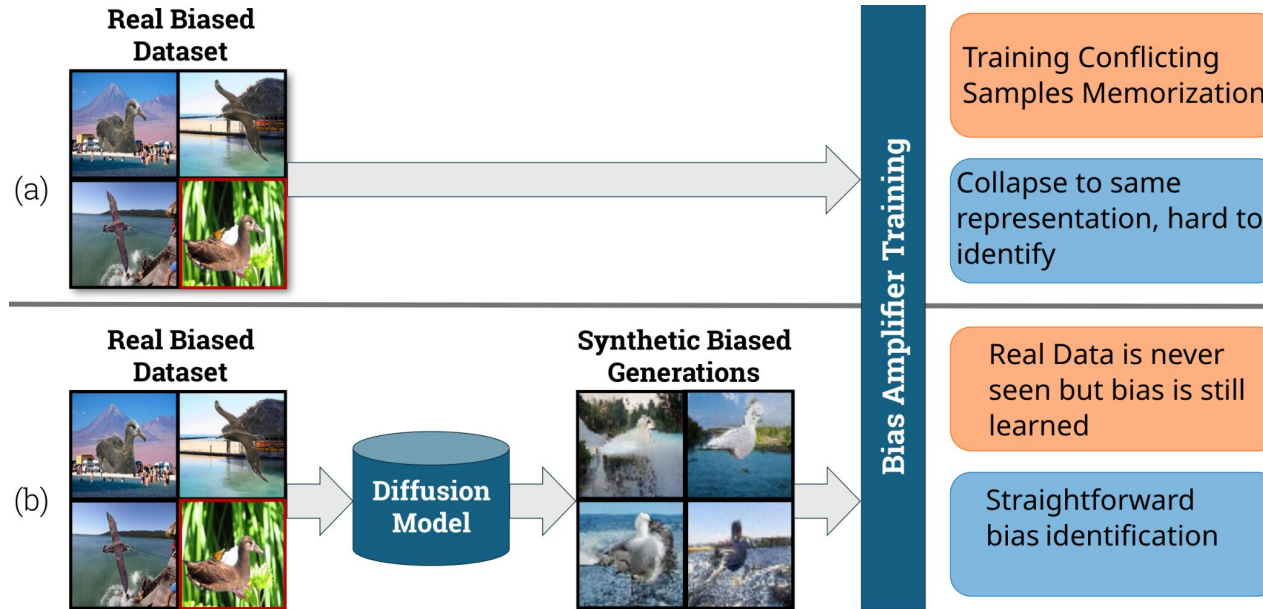
$$\mathcal{D}_{\text{unbiased}} \sim p_{\text{data}}$$

$$\mathcal{D}_{\text{biased}} \sim p_{\text{data}}(\mathbf{x}, y | b)$$

**Our hypothesis is that a Conditional Diffusion Model can naturally approximate the biased distribution of the training set.**

If  $|\mathcal{D}_{\text{biased}}| \gg |\mathcal{D}_{\text{unbiased}}|$

$$\tilde{g}_{\phi}(\mathbf{x} | y) \approx p(\mathbf{x} | y)$$



## Assume

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

$$\mathcal{D} = \mathcal{D}_{\text{unbiased}} \cup \mathcal{D}_{\text{biased}}$$

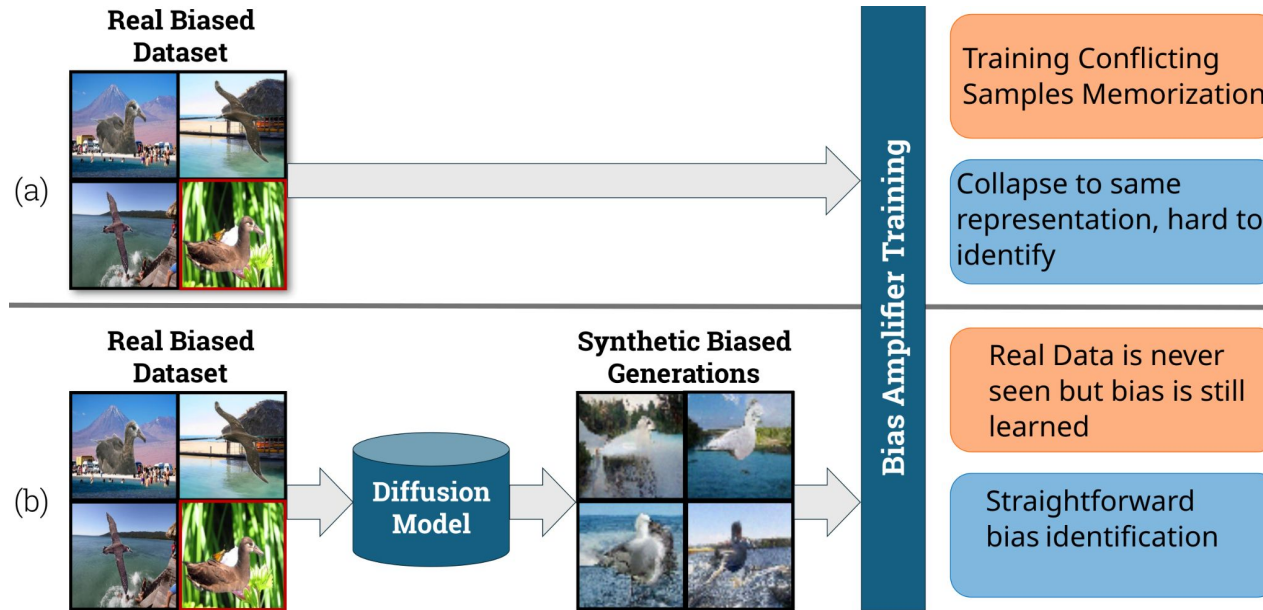
$$\mathcal{D}_{\text{unbiased}} \sim p_{\text{data}}$$

$$\mathcal{D}_{\text{biased}} \sim p_{\text{data}}(\mathbf{x}, y | b)$$

**Our hypothesis is that a Conditional Diffusion Model can naturally approximate the biased distribution of the training set.**

If  $|\mathcal{D}_{\text{biased}}| \gg |\mathcal{D}_{\text{unbiased}}|$

$$\tilde{g}_{\phi}(\mathbf{x} | y) \approx p(\mathbf{x} | y)$$





## Assume

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$$

$$\mathcal{D} = \mathcal{D}_{\text{unbiased}} \cup \mathcal{D}_{\text{biased}}$$

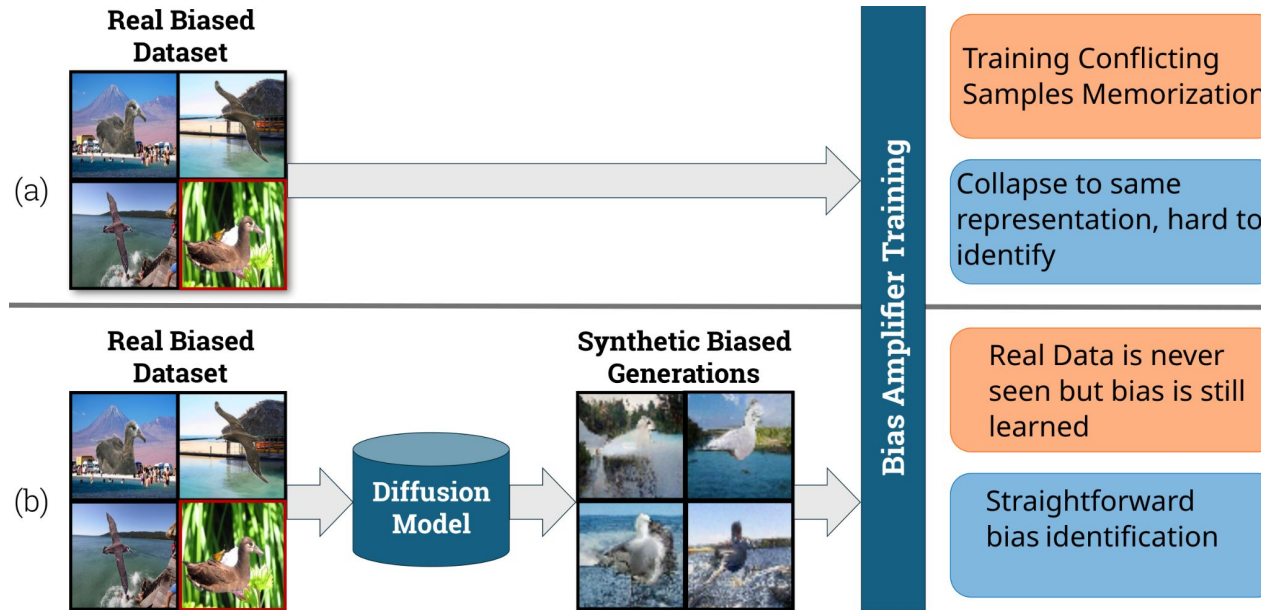
$$\mathcal{D}_{\text{unbiased}} \sim p_{\text{data}}$$

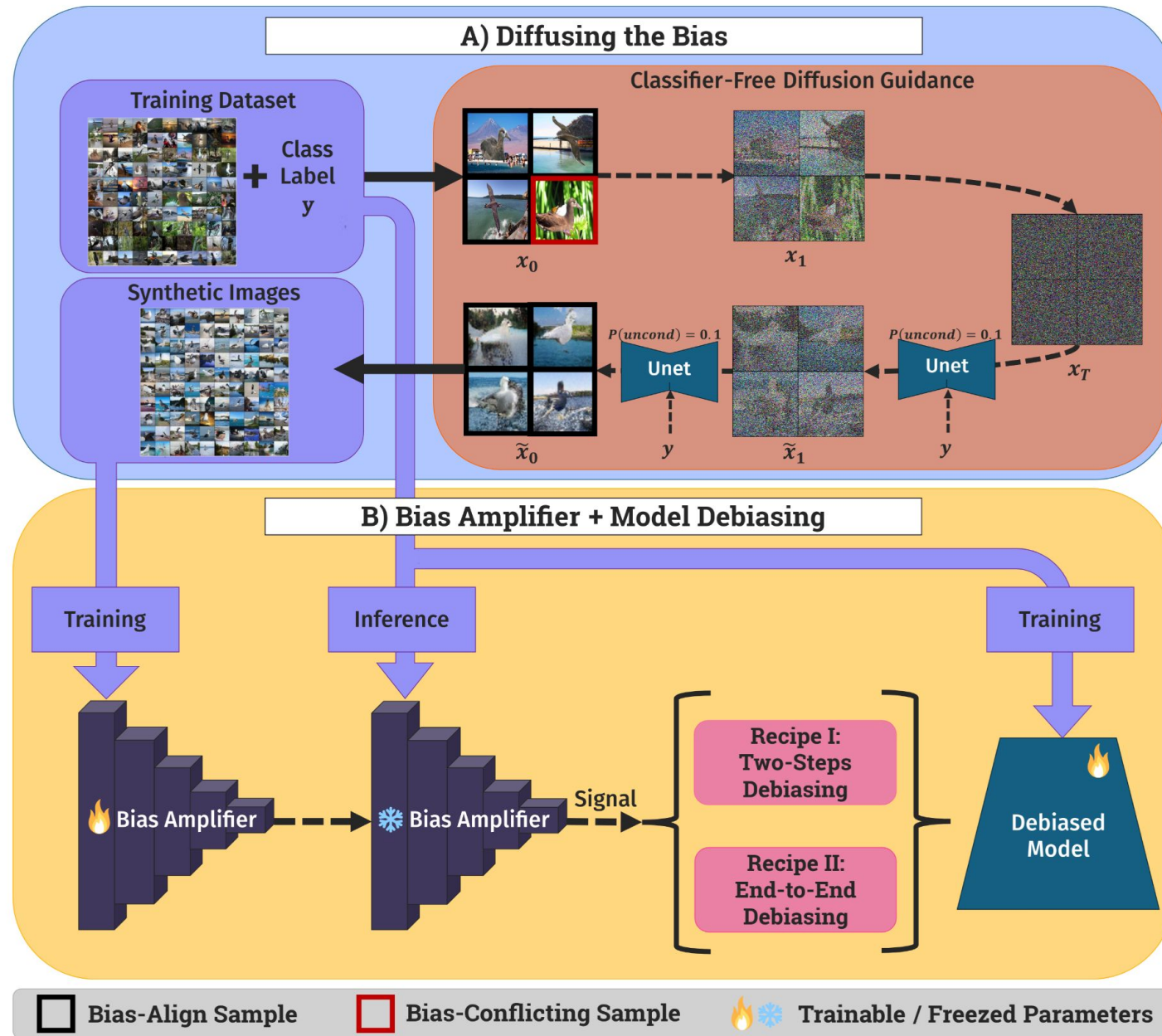
$$\mathcal{D}_{\text{biased}} \sim p_{\text{data}}(\mathbf{x}, y | b)$$

**Our hypothesis is that a Conditional Diffusion Model can naturally approximate the biased distribution of the training set.**

If  $|\mathcal{D}_{\text{biased}}| \gg |\mathcal{D}_{\text{unbiased}}|$

$$\tilde{g}_{\phi}(\mathbf{x} | y) \approx p(\mathbf{x} | b)$$

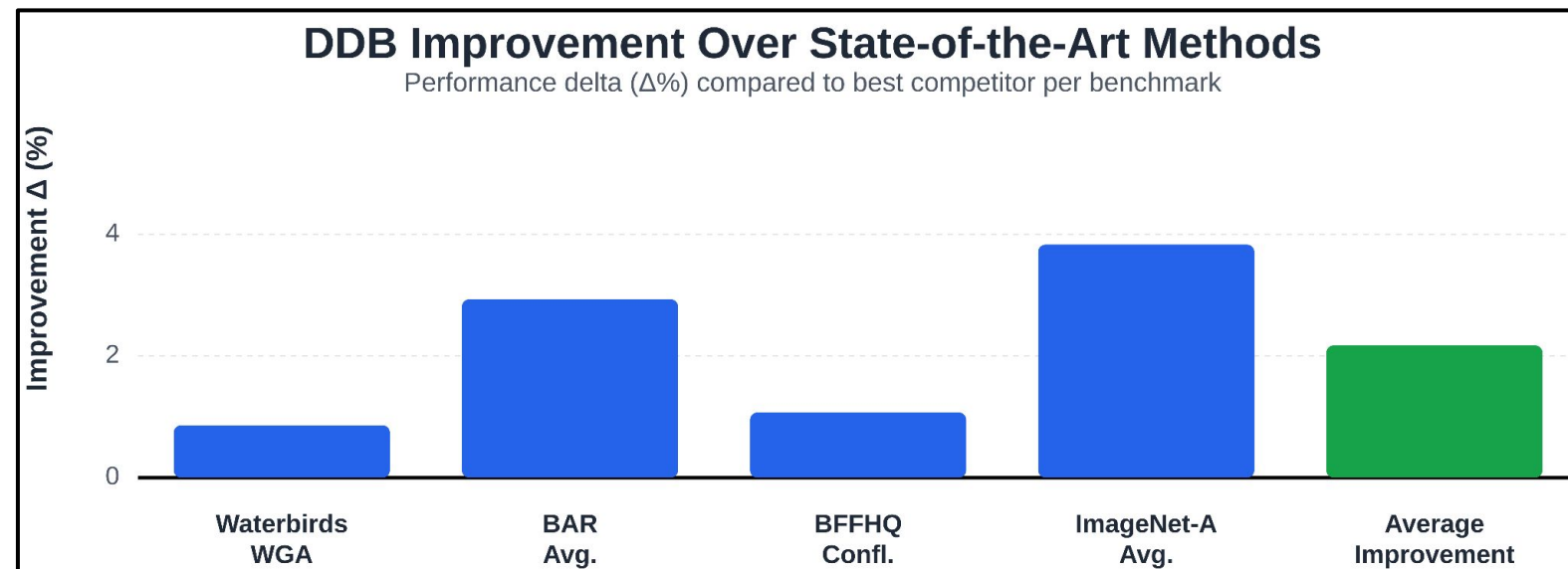
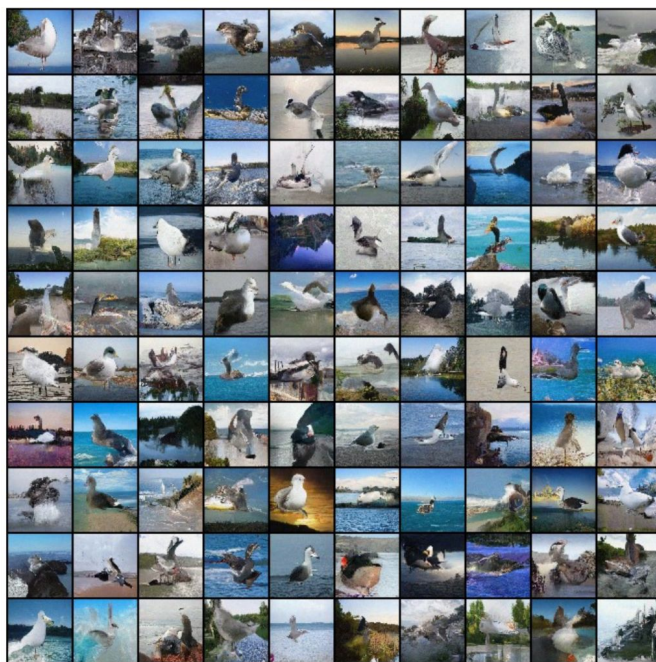




First we leverage a Conditional Diffusion Model (CDPM), to create a **substitute** of the original training set with a synthetic and bias-amplified set of images.

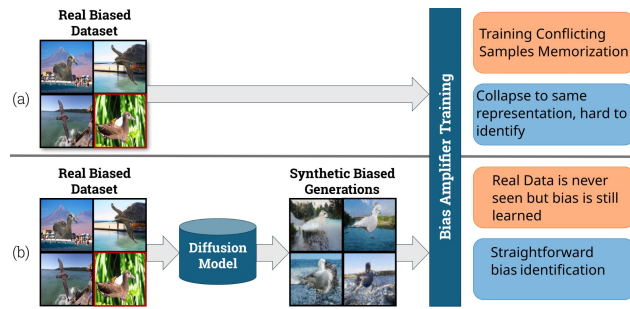
With this new dataset, we train a **Bias-Amplifier**, which can act as a plug-in for existing debiasing methods.





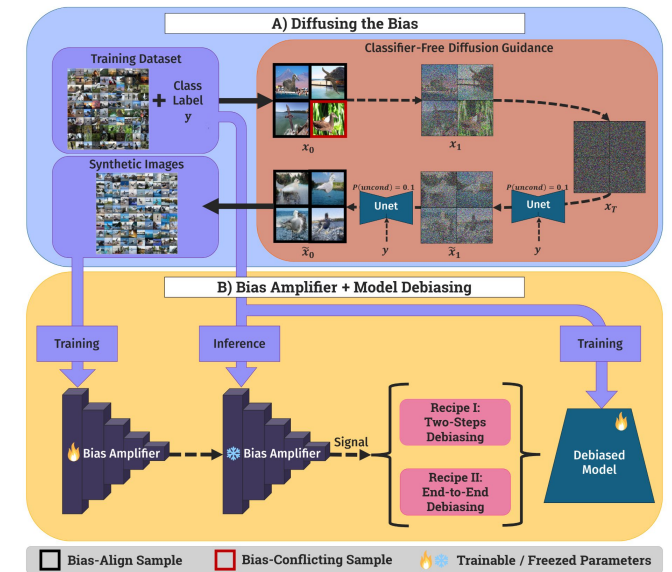
Our Bias Amplifier, regardless of the employed *recipe*, shows improvements over SOTA in all the explored benchmarks, surpassing other plug-in approaches employing our same debiasing methods.

# Conclusions and takeaways



In Diffusing DeBias, we show:

- **Diffusion Models** capture **biases in the training distribution**.
- bias-conflicting memorization issue solved by construction using aligned synthetically generated images to train a bias amplifier.
- Our Bias Amplifier improves **Bias Identification capabilities**
- **Acting as a plug-in**, provides improvements over SOTA





UniGe



ISTITUTO ITALIANO  
DI TECNOLOGIA  
AI FOR GOOD



# Diffusing DeBias: Synthetic Bias Amplification for Model Debiasing

**Massimiliano Ciranni\*<sup>1</sup>, Vito Paolo Pastore\*<sup>1,2</sup>, Roberto Di Via\*<sup>1</sup>,  
Enzo Tartaglione<sup>3</sup>, Francesca Odone<sup>1</sup>, Vittorio Murino<sup>2,4</sup>**

<sup>1</sup>MaLGa-DIBRIS, University of Genoa, Italy

<sup>2</sup>AI For Good (AIGO), Istituto Italiano di Tecnologia (IIT), Genoa, Italy

<sup>3</sup>Télécom Paris, École Polytechnique, France

<sup>4</sup>Department of Computer Science, University of Verona, Italy

**\*Equal Contribution**

**GitHub Project Page**



**ArXiv Preprint**

