

# Background Invariance by Adversarial Learning

Ricardo Cruz<sup>1,2</sup>, Ricardo M. Prates<sup>3,4</sup>, Eduardo F. Simas Filho<sup>4</sup>,  
Joaquim F. Pinto Costa<sup>2</sup> and Jaime S. Cardoso<sup>1,2</sup>

<sup>1</sup> INESC TEC <sup>2</sup> University of Porto <sup>3</sup> Federal University of São Francisco Valley <sup>4</sup> Federal University of Bahia

Corresponding author: [rpcruz@inesctec.pt](mailto:rpcruz@inesctec.pt)

## Motivation

- CNNs fail to distinguish foreground (object) from background.
- Notice how a model with accuracy of **97.3%** can drop to as low as **11%** (random) just by changing the background.



Figure 1: Accuracy (%) of vanilla CNN trained for MNIST.

## Related Work

Little literature exists on making CNNs background invariant. One work proposes **an attention mechanism** to avoid artifacts, particularly irregular borders [1].

- Two classifiers: a *global* CNN  $G$ , and a *local* CNN  $L$ .
- $G$  is trained to classify the entire image  $x$ .
- A activation maps from a truncated  $G^T$  are then used to find a bounding box of the object and crop  $x'$ .
- $L$  is then trained using the smaller  $x'$

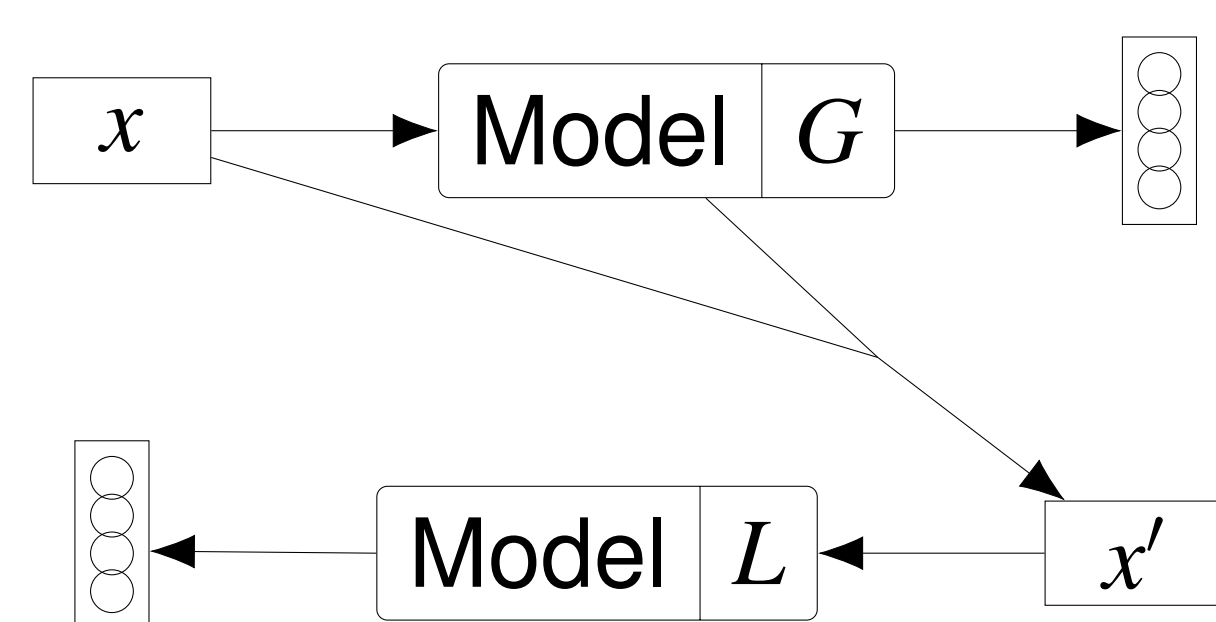


Figure 2: Attention mechanism diagram.

Two disadvantages are immediate:

- $L$  operates on a rectangular cropped version of the image and therefore is still influenced by artifacts that remain inside that rectangle.
- Model  $G$  is still influenced by artifacts because it did not have the benefit of being trained against the artifacts. While such artifacts are not presented in the training set, they could be generated in a controlled fashion, as our method now proposes:

## Proposed Method

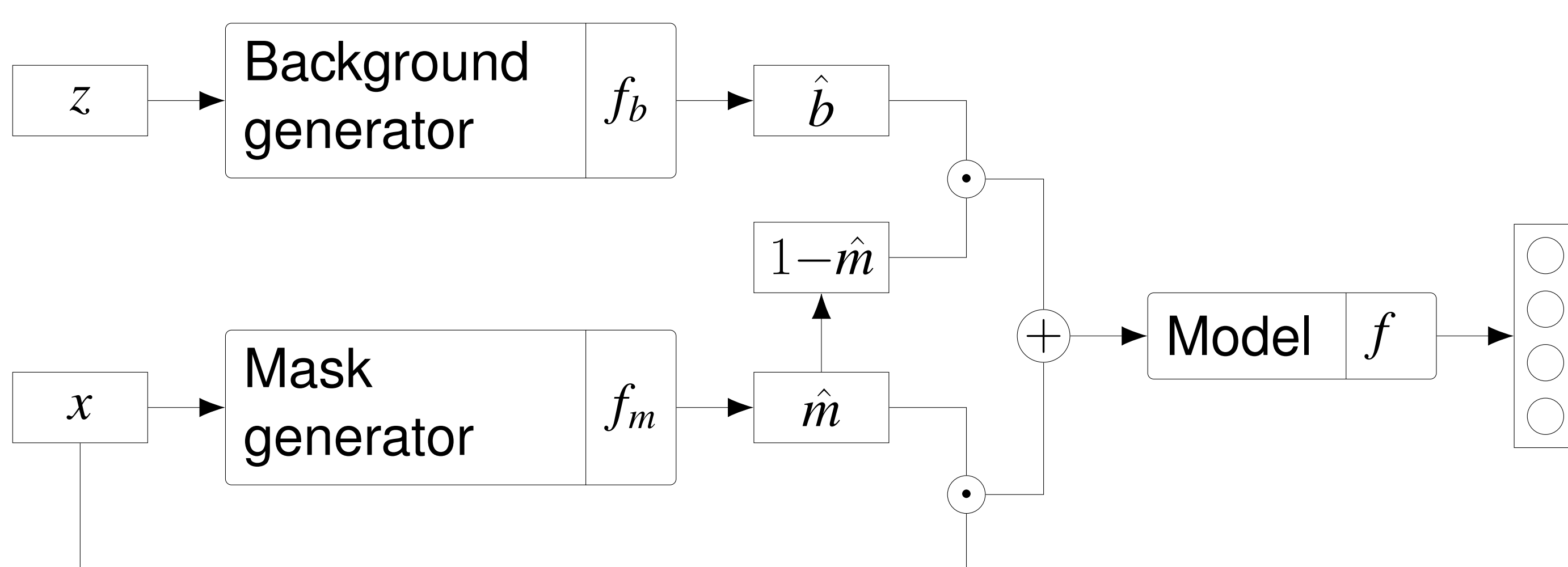


Figure 3: Proposed method.

- The goal is to (during training) be able to place the object in a multitude of contexts (backgrounds).
- Focus on “what” the object is rather than “where” the object is.
- The backgrounds are generated in an adversarial manner.
- However, the insertion of adversarial backgrounds in the image cannot be allowed to destroy the concept (class) one is trying to learn.

– A model  $f$  is optimized to minimize a loss  $\mathcal{L}(y, f(x))$ .

– A mask generator  $f_m$  is trained to produce a mask  $\hat{m} \in [0, 1]$ .

\* This U-Net is trained unsupervisedly by finding the best mask that minimizes the previous loss,  $\mathcal{L}(f(x \odot f_m(x)), y)$ .

– A background generator  $f_b$  transforms noise  $z$  into a background  $\hat{b}$  image.

\* The trick: it (adversarially) *maximize* the loss  $\mathcal{L}$ .

**Summary:** Model  $f$  tries to *minimize* a loss while  $f_b$  tries to find backgrounds that *maximize* it:

$$\min_{f, f_m} \max_{f_b} \sum_{i=1}^N \mathcal{L} \left( f \left( \hat{m}_i \odot x_i + \hat{b}_j \odot (1 - \hat{m}_i) \right), y_i \right).$$

This is inspired by literature in adversarial training and GANs.

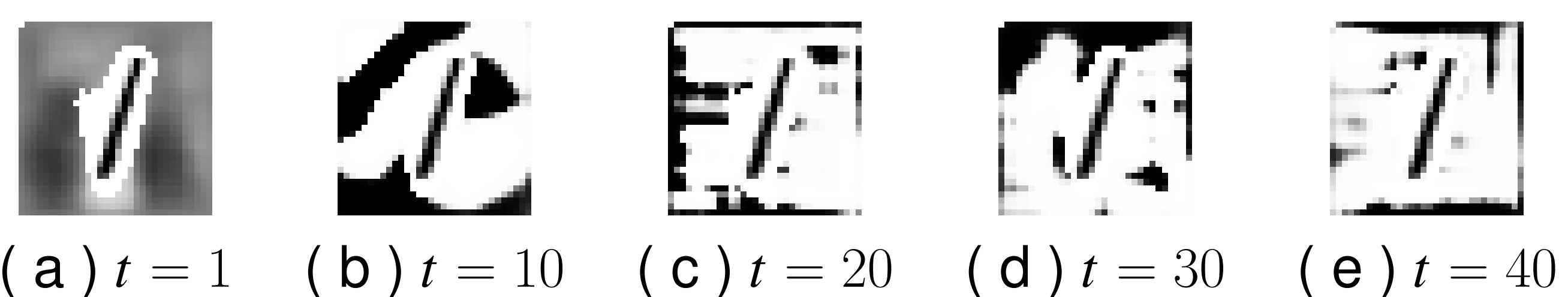


Figure 4: Background augmentation along the epochs.

## Experiments

- The attention mechanism only negligibly improves on the baseline classifier.
- The proposed method is resilient to a wide range of testing backgrounds.

	4	4	4	4	4	4	4
Vanilla CNN	97.3	38.0	24.3	61.4	32.9	19.7	11.2
Attention	93.4	28.1	26.8	57.3	40.1	29.3	25.1
Proposal	94.9	92.3	76.8	93.1	93.7	70.8	86.2

Table 1: Accuracy (%) using MNIST.

## Conclusion

- Sometimes it is easier to collect data inside a studio rather than in the real world – for example when training a drone. Unfortunately, a CNN does terrible when used in new backgrounds.
- An adversarially trained data augmentation method is proposed. The proposed method can be used for classification, regression, segmentation, reinforcement learning, etc.

## References

- [1] Qingji Guan, Yaping Huang, Zhun Zhong, Zhedong Zheng, Liang Zheng, and Yi Yang. Diagnose like a radiologist: Attention guided convolutional neural network for thorax disease classification. *arXiv preprint arXiv:1801.09927*, 2018.