

# Explaining an image classifier with a generative model conditioned by uncertainty

Adrien Le Coz<sup>1,2</sup>, Stéphane Herbin<sup>2</sup>, Faouzi Adjed<sup>1</sup>

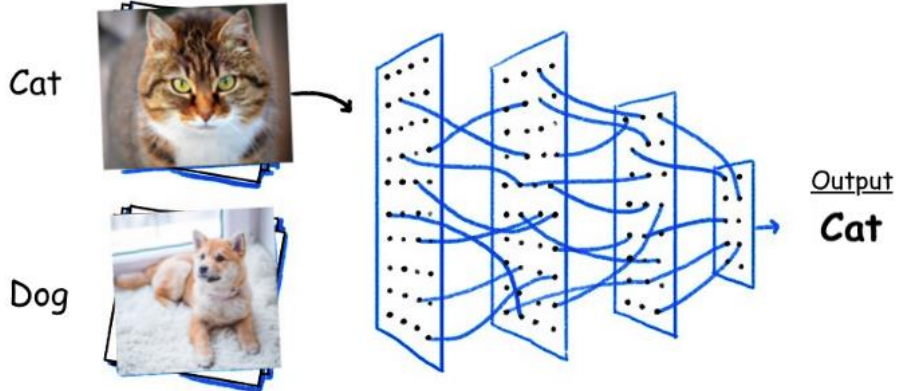
<sup>1</sup>IRT SystemX, Palaiseau, France

<sup>2</sup>DTIS, ONERA, Université Paris Saclay F-91123 Palaiseau - France

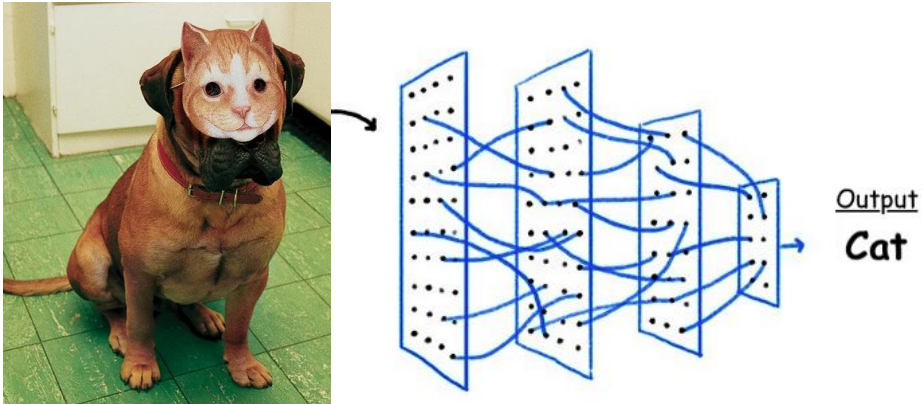


Context

Image classifiers work well most of the time...

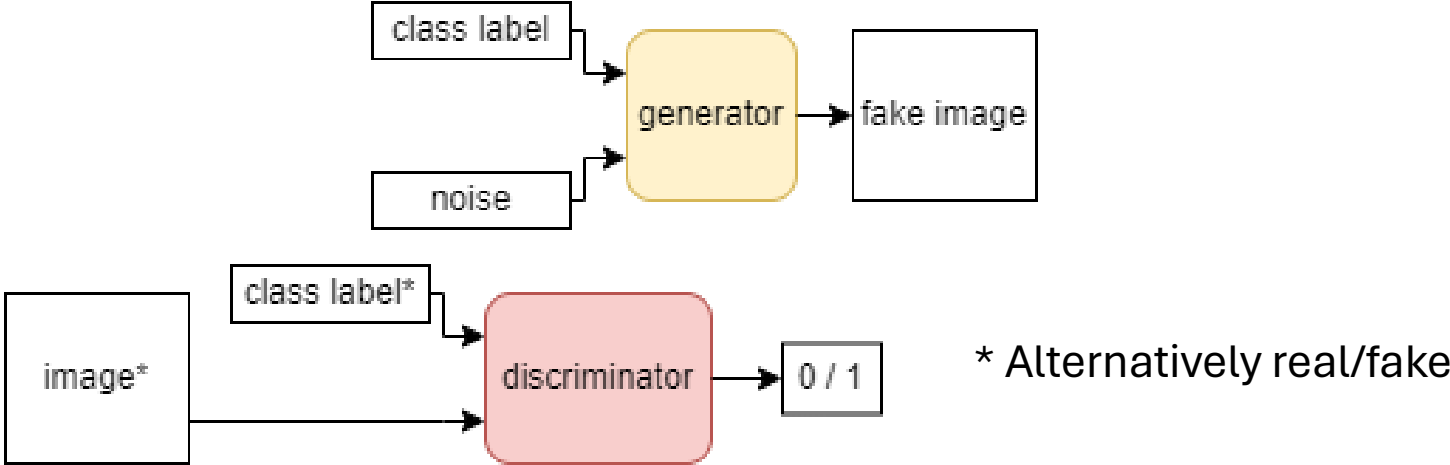


... but not always!

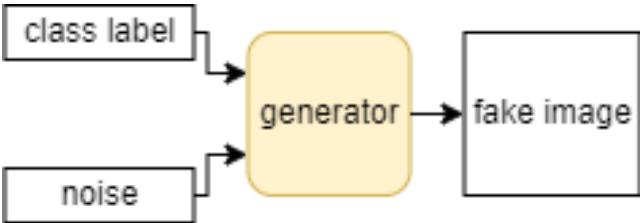


What makes them fail?

# Generative adversarial networks

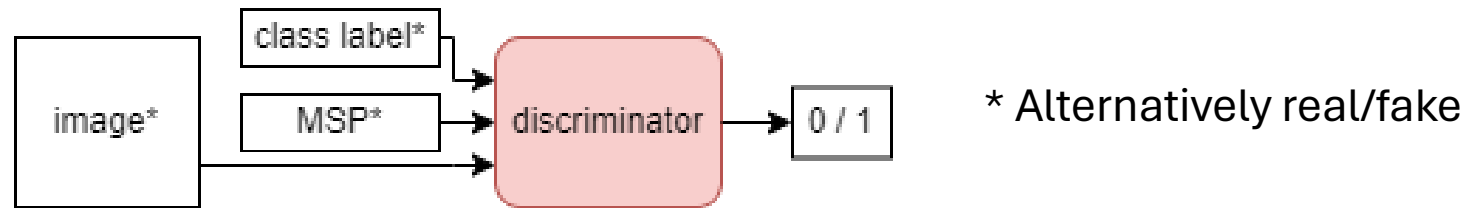
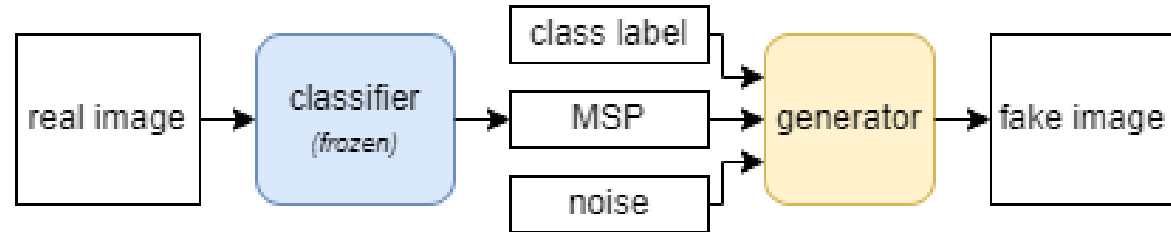


Can be used to generate data of a chosen class (vary label input)

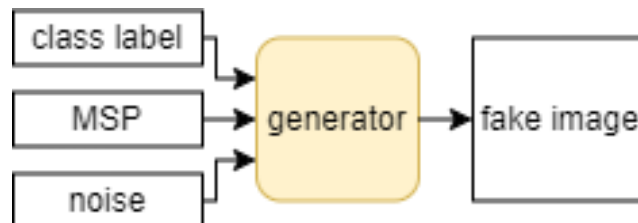


# Conditioning with uncertainty

Here, Maximum Softmax Probability (MSP) = classifier confidence

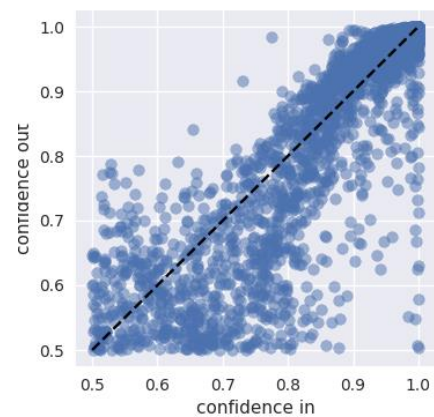
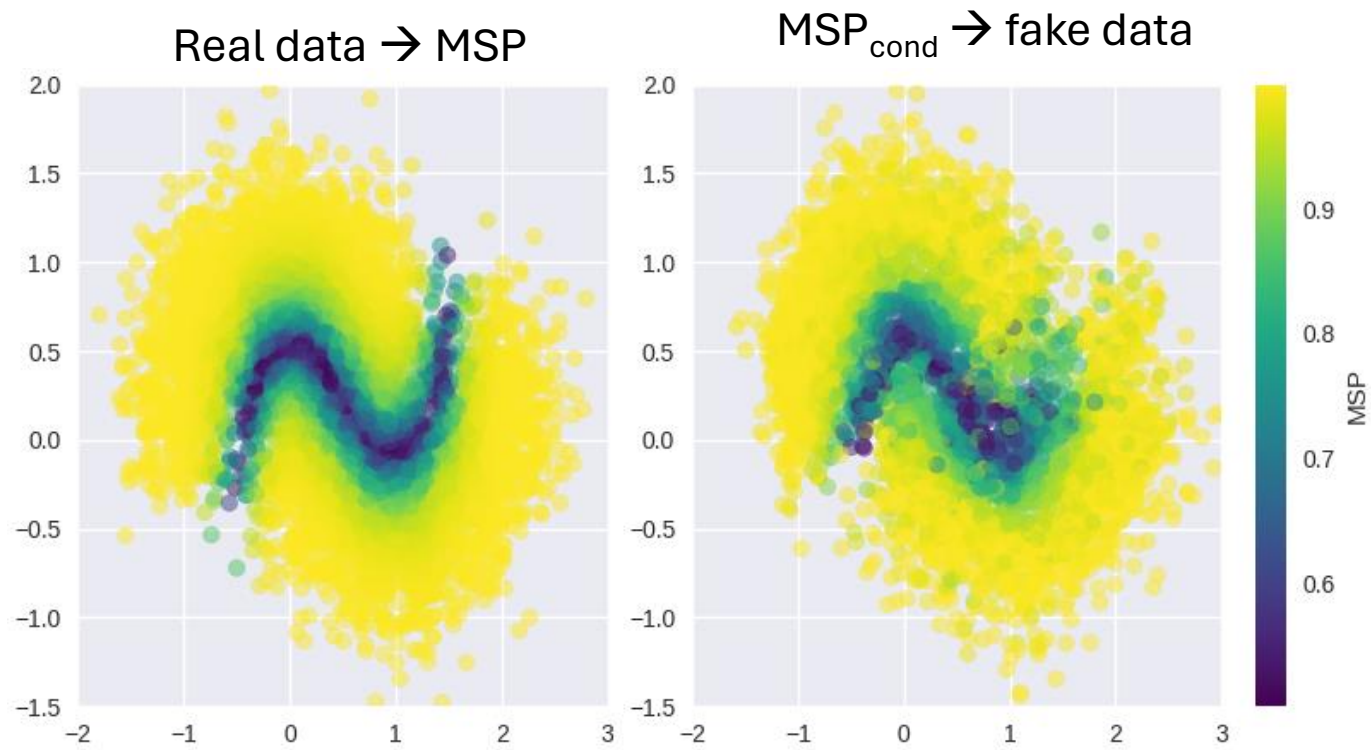


Can be used to generate uncertain images (vary MSP input)

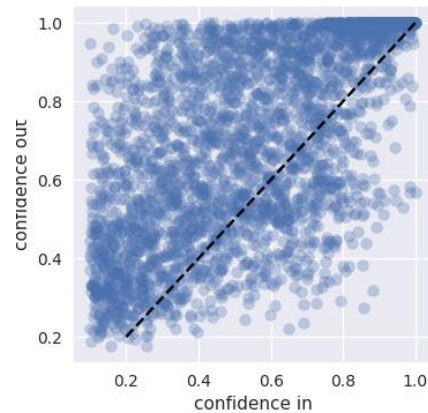
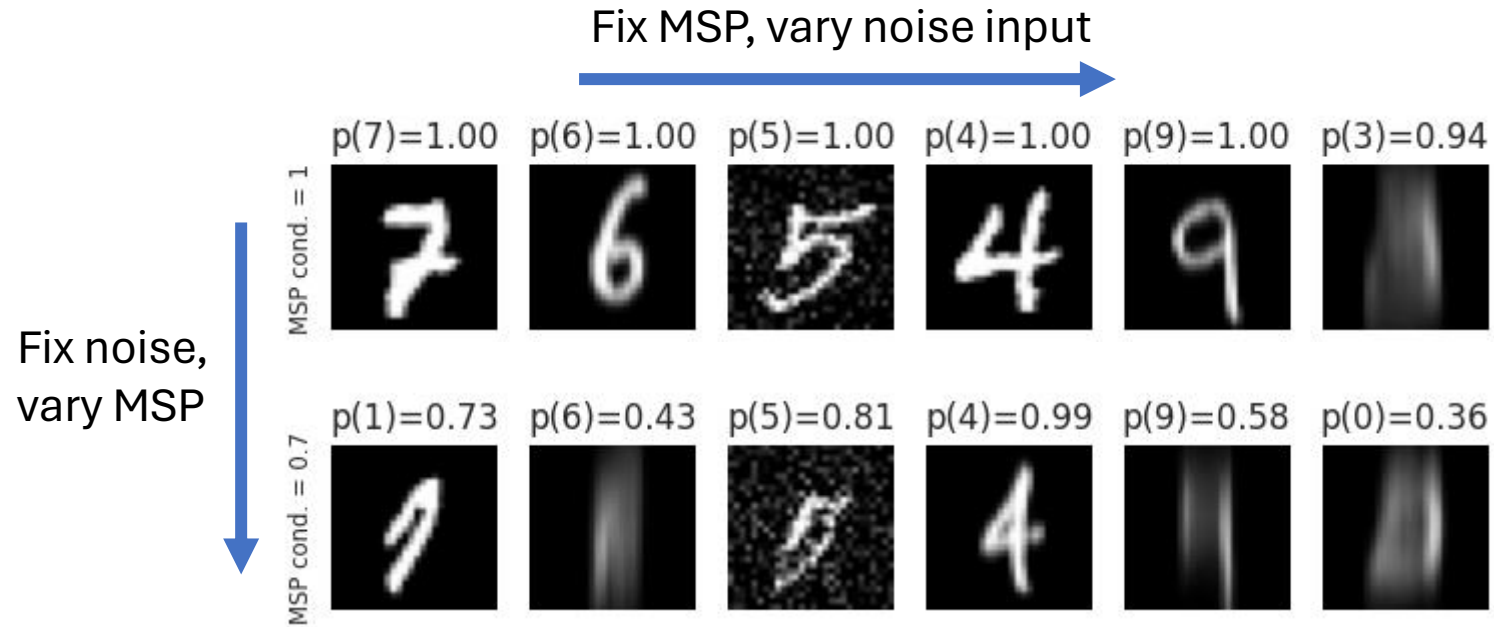


# Results

# Toy data



# Corrupted MNIST





Conclusion

- GAN conditioned by uncertainty can:
  - Corrupt images (fix noise, vary confidence)
  - Generate uncertain images (fix confidence, vary noise)
- Identification of uncertainty sources is manual
- Still limited results for images

Sources of images:

<https://towardsdatascience.com/10-papers-you-should-read-to-understand-image-classification-in-the-deep-learning-era-4b9d792f45a7>

[https://www.reddit.com/r/photoshopbattles/comments/3hpcvq/psbattle\\_dog\\_with\\_cat\\_mask/](https://www.reddit.com/r/photoshopbattles/comments/3hpcvq/psbattle_dog_with_cat_mask/)