

# medXGAN

Visual Explanations for Medical Classifiers through a Generative Latent Space  
CVPR 2022

Paweł Pawlik, Michał Siennicki

# Objective

We want to explain CNNs on covid dataset

Original  
Positive Image



Reconstruction



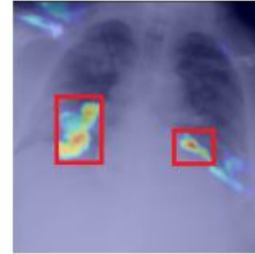
Negative  
Realization



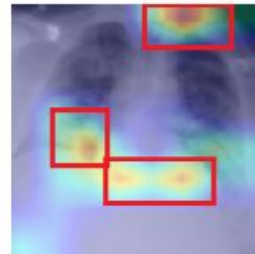
Difference  
Map



Colorized  
Map



Grad-CAM

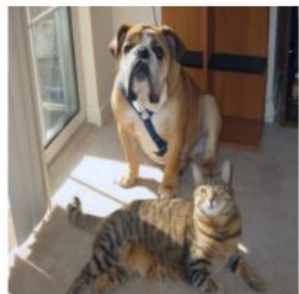


# Plan

- Background
  - Grad-CAM
  - GANs
  - GANs - reconstruction
- medXGAN
  - method overview
  - results

Background

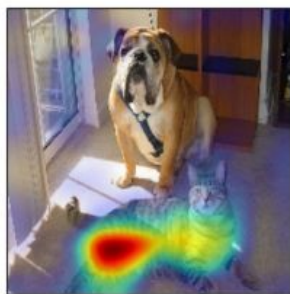
# Grad-CAM



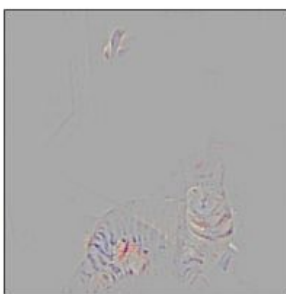
(a) Original Image



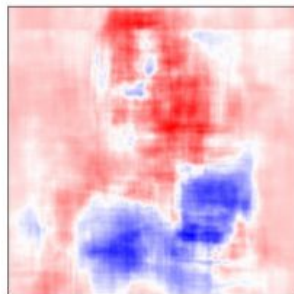
(b) Guided Backprop 'Cat'



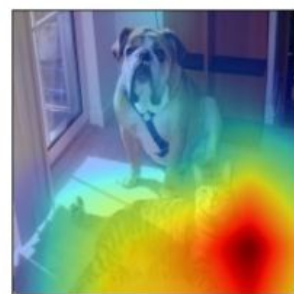
(c) Grad-CAM 'Cat'



(d) Guided Grad-CAM 'Cat'



(e) Occlusion map 'Cat'



(f) ResNet Grad-CAM 'Cat'



(g) Original Image



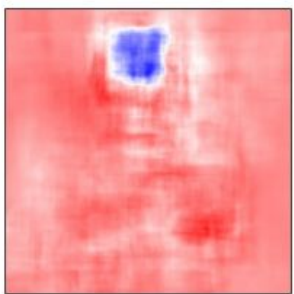
(h) Guided Backprop 'Dog'



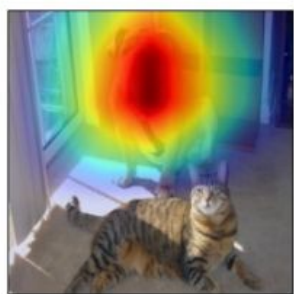
(i) Grad-CAM 'Dog'



(j) Guided Grad-CAM 'Dog'

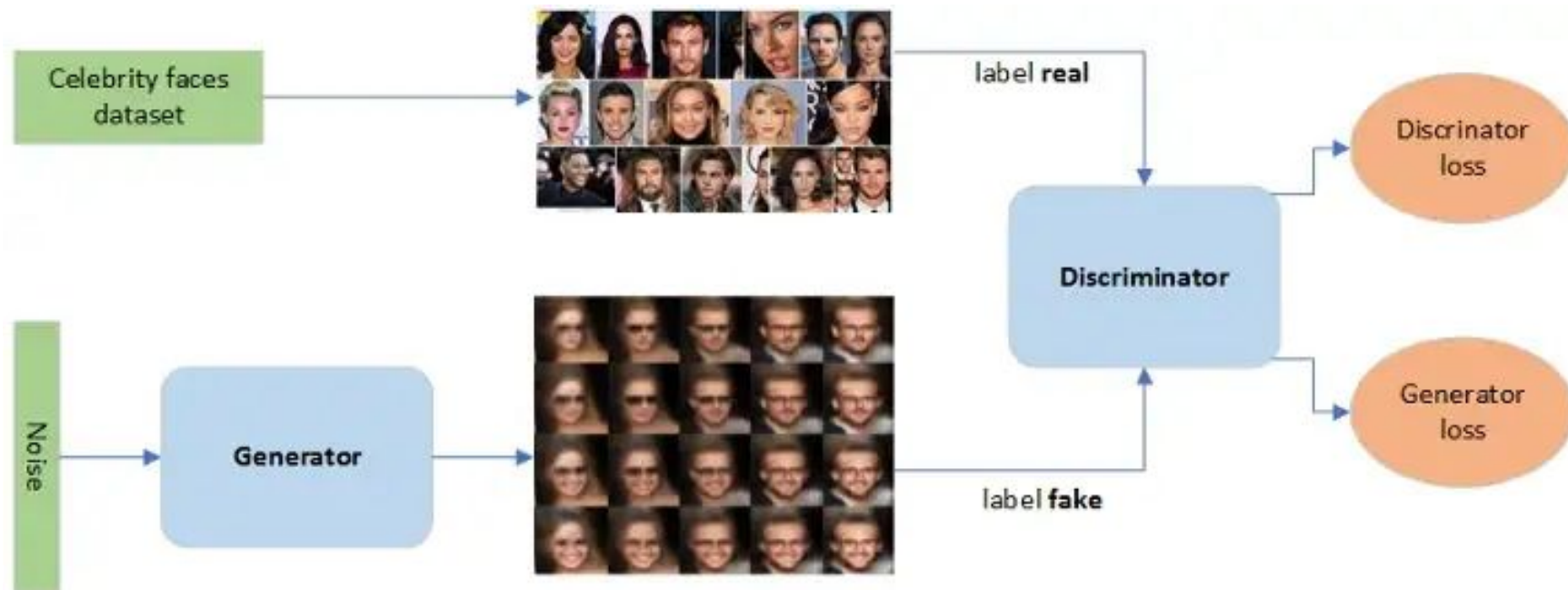


(k) Occlusion map 'Dog'



(l) ResNet Grad-CAM 'Dog'

# GANs



# GANs - reconstruction

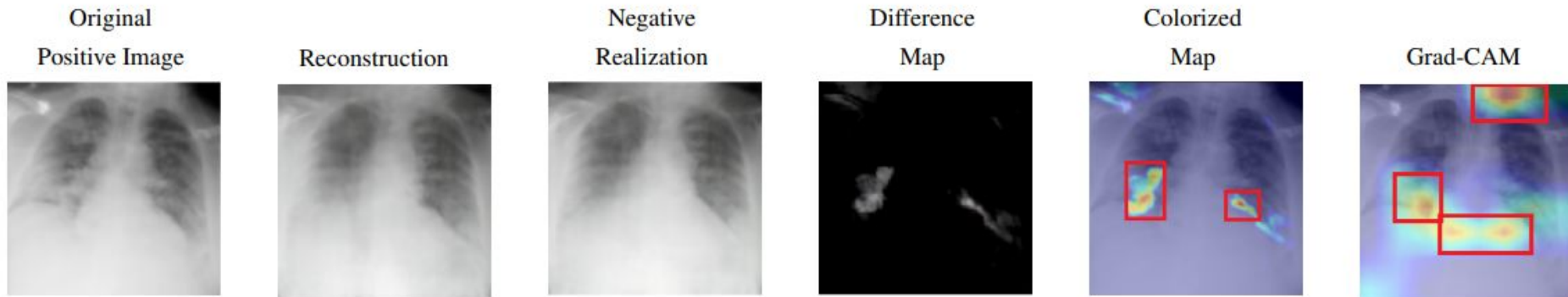
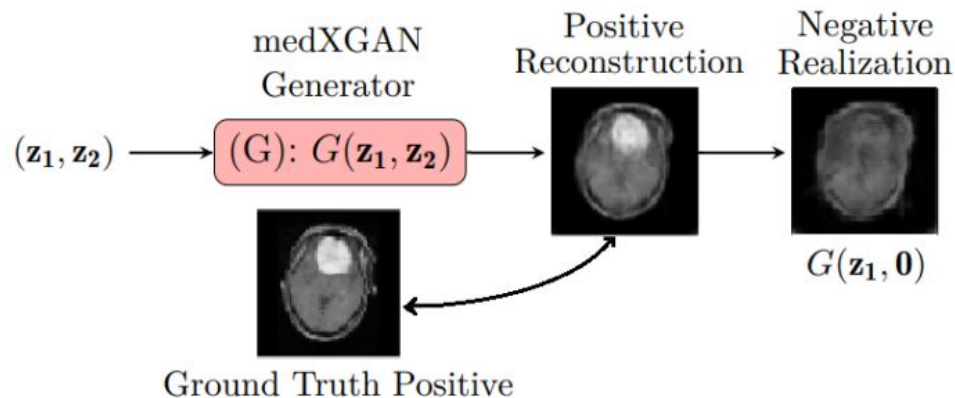


medXGAN



# Objective

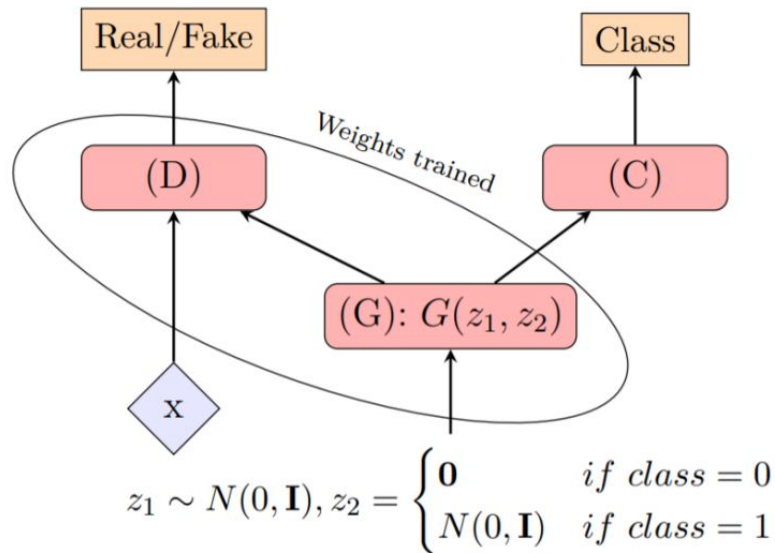
Find input image reconstruction and compare it with its negative realization



# The method

We want to learn the following disentangled representation:

- $z_1$ : lung features
- $z_2$ : covid pathologies
  - $z_2$  has high mutual information ( $I(z_2, \hat{y})$ ) with the CNN classifier



Reconstruction



Negative  
Realization



# The method

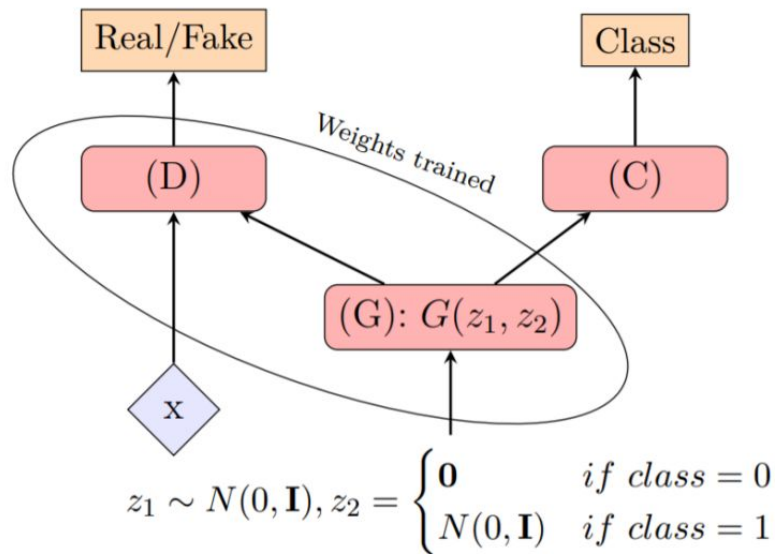
We want to learn the following disentangled representation:

- $z_1$ : lung features
- $z_2$ : covid pathologies
  - $z_2$  has high mutual information  $(I(z_2, \hat{y}))$  with the CNN classifier

$$\begin{aligned} \min_G \max_D \mathbb{E}_{x \sim p_x} [\log D(x)] \\ + \mathbb{E}_{z_1 \sim p_{z_1}, y \sim p_y} [\log(1 - D(G(z_1, y)))] \\ - \mathbb{E}_{z_1 \sim p_{z_1}, y \sim p_y} [\log(p_c(y|G(z_1, y)))] \end{aligned}$$

} normal GAN objective

} mutual information



# Sampling process

- Input image:
- Find reconstruction latent vector  $z_1:z_2$

$$\arg \min_{z_1, z_2} \text{MSE}(G(z_1, z_2), x) + \text{BCE}(C(G(z_1, z_2)), C(x))$$

- The reconstruction:  $G(z_1, z_2)$
- Negative realization:  $G(z_1, 0)$

Original  
Positive Image



Reconstruction

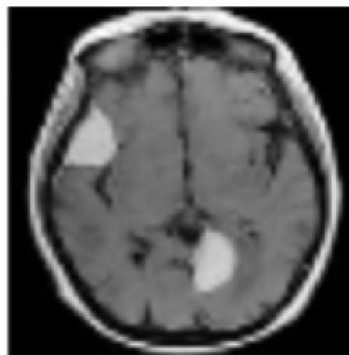


Negative  
Realization

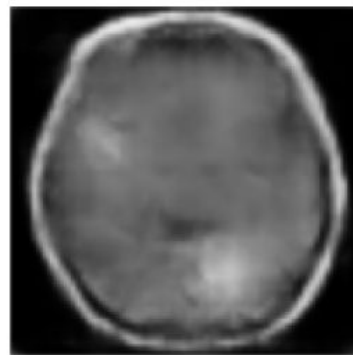


# medXGAN and Grad-CAM on Brain MRIs

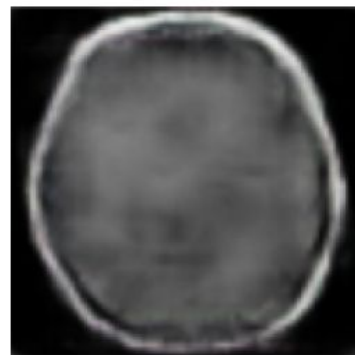
In this example, despite lacking a perfect reconstruction, the medXGAN method localizes two tumors, while Grad-CAM focuses on one tumor and an eye.



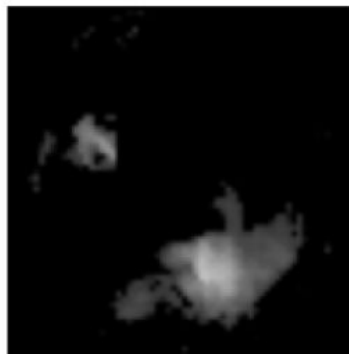
Original Positive



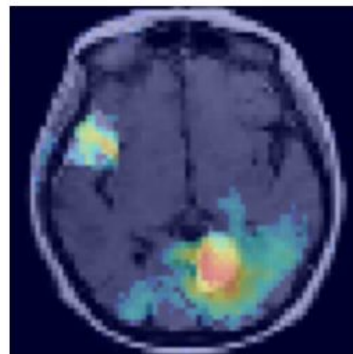
Reconstruction



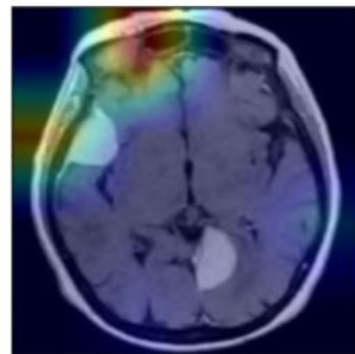
Negative Realization



Difference Map

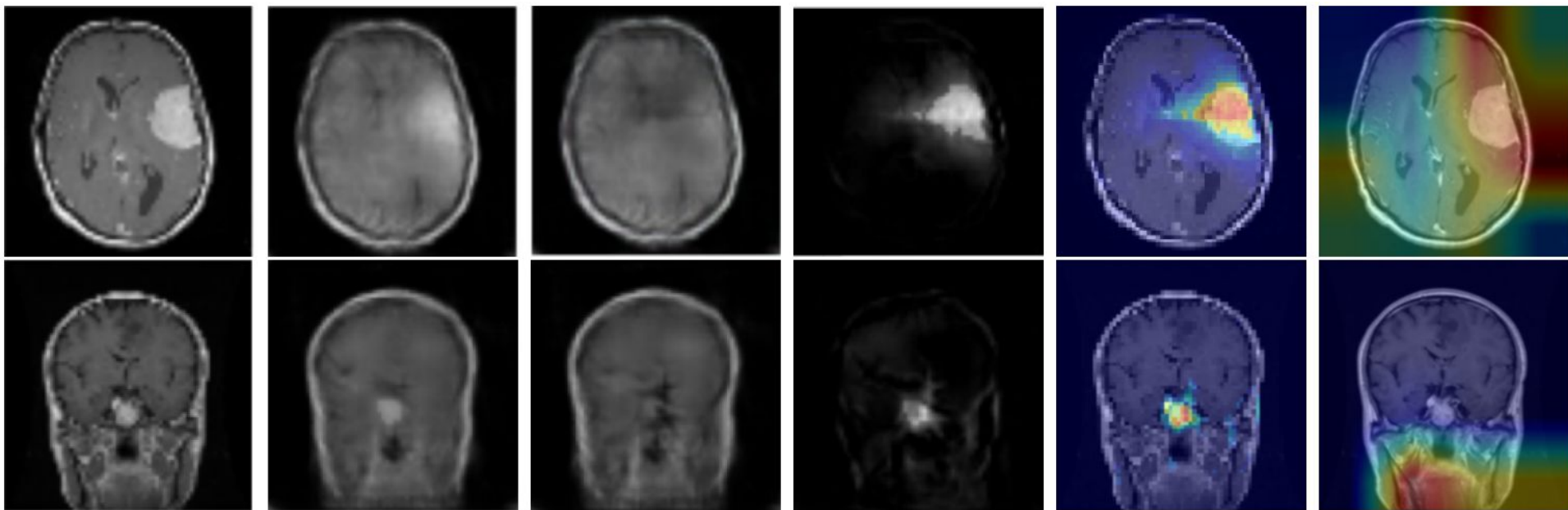


Colorized Map



Grad-CAM

# medXGAN and Grad-CAM on Brain MRIs



(a) Original

(b) Reconstruction

(c) Negative  
Realization

(d) Difference  
Map

(e) Colorized Map

(f) Grad-CAM

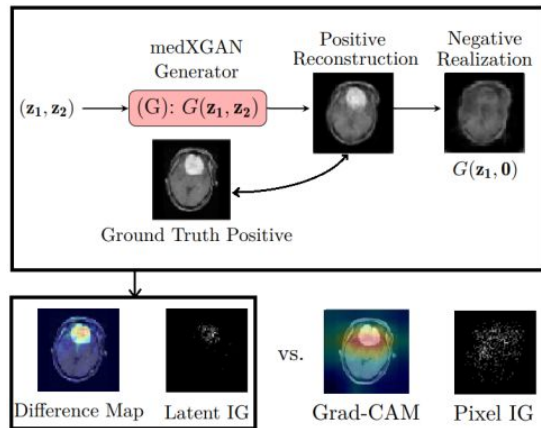
# Summary

- medXGAN can be used to replace

- Grad-Cam
- IG (integrated gradients)

- limitations

- the space along which we are explaining the model must be continuous and independent
- can only explain positive samples
- GANs require significant amount of data to train
- GANs have unstable training



Thank you!