



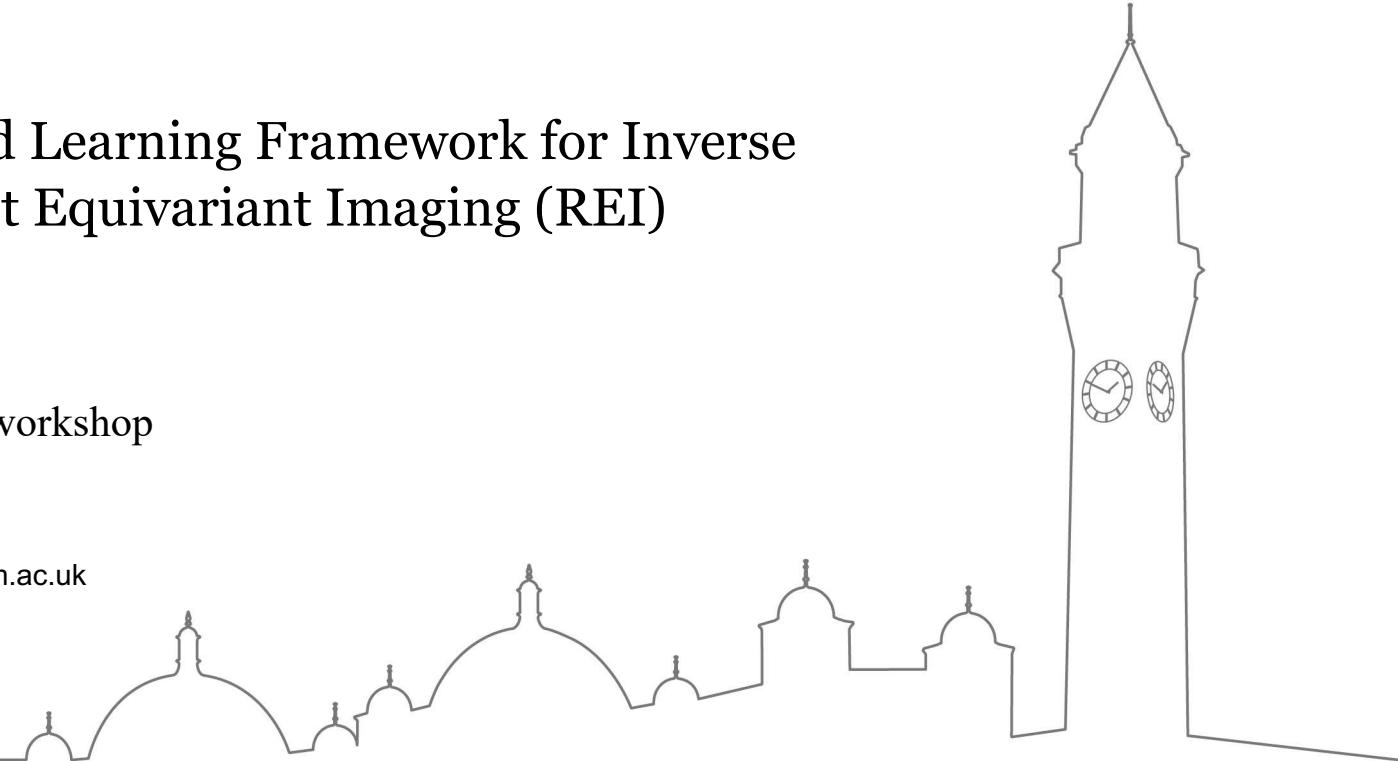
UNIVERSITY OF  
BIRMINGHAM

# An Unsupervised Learning Framework for Inverse Problem: Robust Equivariant Imaging (REI)

CVPR2022 digestion workshop

Name: Binglun Wang

Email: bxw135@student.bham.ac.uk



## Inverse Problem

Forward process: generate  $y$  from  $x$

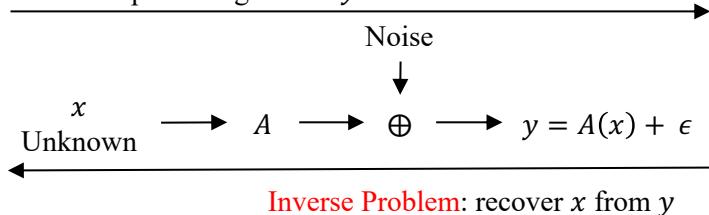


Fig 1. Definition of Inverse Problem [1, 2]

## Example of Inverse Problem

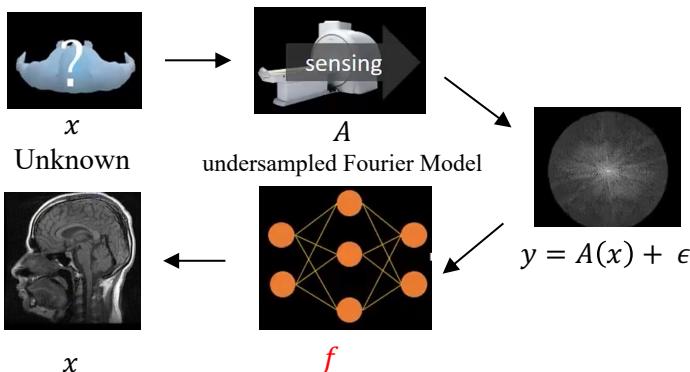
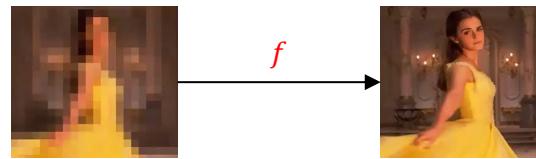


Fig 2. Example of MRI [3]

## Supervised Learning



$$\{y\} \arg \min_f \text{MSE}(x, f(y)) = \{x\}$$

Fig 3. Example of super-resolution [3]

## Unsupervised Learning

In some unsupervised methods [4, 5, 6], the setting is not often met in practice. However, Equivariant Imaging (EI):

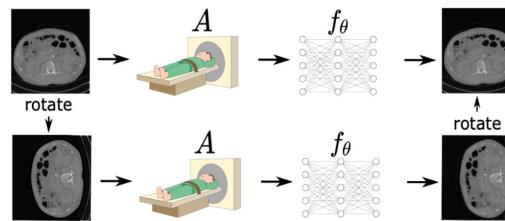


Fig 4. Idea of Equivariant Imaging [1, 2]

SOTA performance  
Based on  $\{y, x\}$  pairs

However, in lots of tasks.  
Obtaining targets  $\{x\}$  is  
**expensive or impossible**

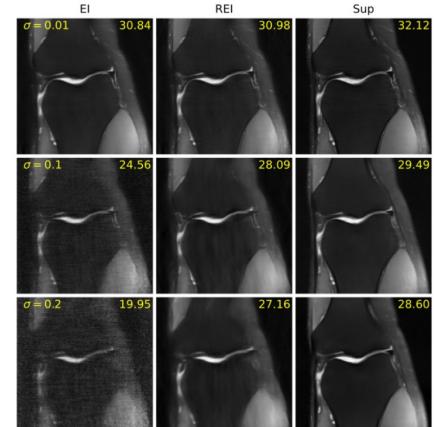


Fig 5. EI problem and REI performance [1]

[1] "Robust Equivariant Imaging: a fully unsupervised framework for learning to image from noisy and partial measurements.", [Chen et al., CVPR2022 Oral]

[2] "Equivariant imaging: Learning beyond the range space.", [Chen et al., ICCV2021]

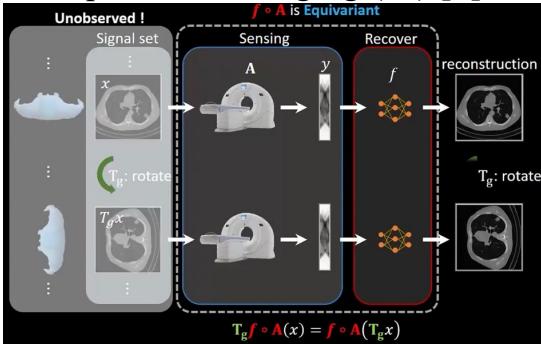
[3] "Robust Equivariant Imaging: A Fully Unsupervised Framework for Learning To Image From" [Wang, YouTube]

[4] "AmbientGAN: Generative models from lossy measurements." [Bora et al., ICLR2018]

[5] "Deep image reconstruction using unregistered measurements without groundtruth" [Gan et al., ISBI 2021]

[6] "Unsupervised adversarial image reconstruction." [Pajot et al., ICLR2021]

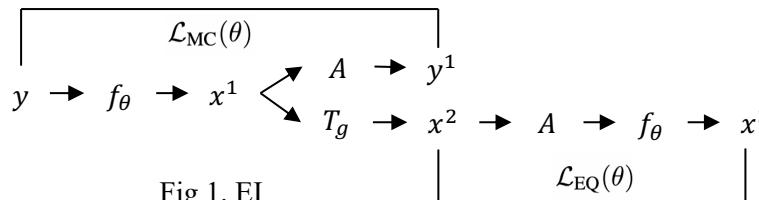
## Equivariant Imaging (EI) [1]



$$\mathcal{L}_{\text{EI}}(\theta) = \mathcal{L}_{\text{MC}}(\theta) + \alpha \mathcal{L}_{\text{EQ}}(\theta) \quad y \subseteq R^m, x \subseteq R^n$$

With:  $\mathcal{L}_{\text{MC}}(\theta) = \sum_{i=1}^N \frac{1}{m} \|y_i - A(f_\theta(y_i))\|^2,$

$$\mathcal{L}_{\text{EQ}}(\theta) = \sum_{i=1}^N \sum_{g=1}^{|\mathcal{G}|} \frac{1}{|\mathcal{G}|n} \|T_g f_\theta(y_i) - f_\theta(A(T_g f_\theta(y_i)))\|^2$$



- Self-supervised denoising via Stein's Unbiased Risk Estimator (SURE) [2]

If measurements contain Gaussian noise:

$y|u \sim \mathcal{N}(u, I\sigma^2)$   $h_\theta$  that maps  $y \mapsto u$

$$\mathcal{L}_{\text{SURE}}(\theta) = \sum_{i=1}^N \frac{1}{m} \|y_i - h_\theta(y_i)\|^2 - \sigma^2 + \frac{2\sigma^2}{m} \nabla \cdot h_\theta(y_i)$$

*Theorem 1:* SURE loss is an unbiased estimator of the supervised mean squared loss

$$\mathbb{E}_y \{\mathcal{L}_{\text{SURE}}(\theta)\} = \mathbb{E}_{y,u} \{\mathcal{L}_{\text{MSE}}(\theta)\}$$

*Theorem 2 [3]:* The divergence term can be approximated by sampling a single i.i.d.

$$\nabla \cdot h_\theta(y) \approx \frac{1}{\tau} b^\top (h_\theta(y + \tau b) - h_\theta(y))$$

Where  $\tau$  is small fixed number,  $b \sim \mathcal{N}(0, 1)$

The SURE framework [4, 5, 6] beyond simple additive Gaussian noise, including Poisson and mixed Poisson Gaussian noise, resulting in different loss expressions depending on the noise type.

[1] "Equivariant imaging: Learning beyond the range space.", [Chen et al., ICCV2021]

[2] "Estimation of the mean of a multivariate normal distribution." [Stein, The annals of Statistics 1981]

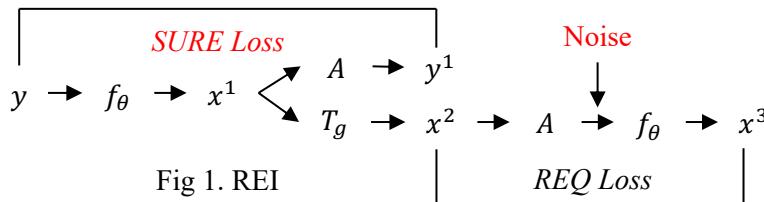
[3] "MonteCarlo SURE: A black-box optimization of regularization parameters for general denoising algorithms." [Ramani et al., ToIP2008]

[4] "Generalized SURE for exponential families: Applications to regularization." [Eldar et al., ToSP2008]

[5] "Image denoising in mixed Poisson-Gaussian noise." [Luisier et al., ToIP2011]

[6] "Least squares estimation without priors or supervision." [Raphan et al., Neural computation2011]

## ▪ REI



$$\mathcal{L}_{\text{REI}}(\theta) = \mathcal{L}_{\text{SURE}}(\theta) + \alpha \mathcal{L}_{\text{REQ}}(\theta)$$

REQ Loss:

$$\mathcal{L}_{\text{REQ}}(\theta) = \sum_{i=1}^N \sum_{g=1}^{|\mathcal{G}|} \frac{1}{|\mathcal{G}|n} \|T_g f_\theta(y_i) - f_\theta(\tilde{y}_i)\|^2$$

SURE Loss:

For Gaussian noise:

$$\begin{aligned} \mathcal{L}_{\text{SURE}}(\theta) &= \sum_{i=1}^N \frac{1}{m} \|y_i - A(f_\theta(y_i))\|_2^2 - \sigma^2 \\ &\quad + \frac{2\sigma^2}{m\tau} b_i^\top (A(f_\theta(y_i + \tau b_i)) - A(f_\theta(y_i))) \end{aligned}$$

where  $b_i \sim \mathcal{N}(0, I)$  and  $\tau$  is a small positive number.

For Poisson noise & Mixed Gaussian-Poisson (MGP): See Paper

## ▪ Experiments

Setup:  $f_\theta$  are residual U-Nets in all networks  
different loss in EI, REI, Sup.

In Accelerated MRI experiment:

1.  $A^\dagger$  is masked inverse Fourier transform.
2. Gaussian Noise based SURE loss
3.  $T_G$  are set of rotations.  $|\mathcal{G}|=360$  (degree)

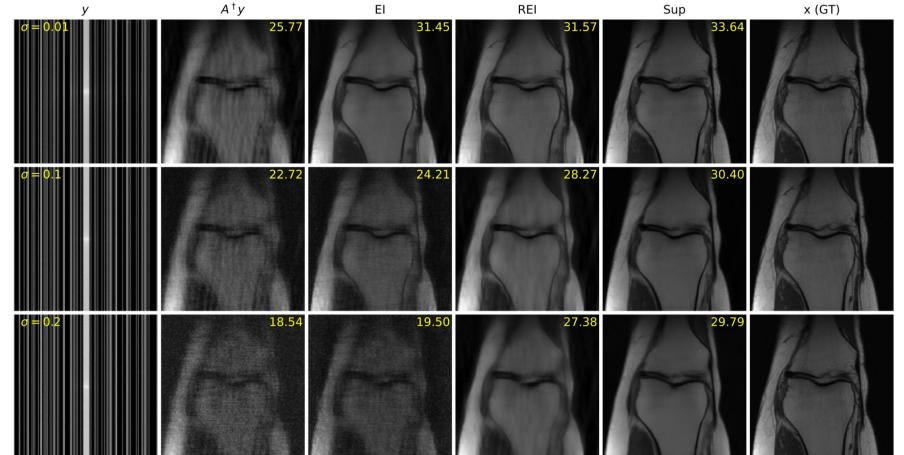
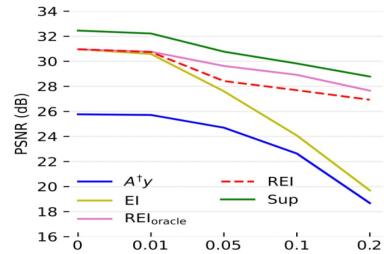


Fig 2. Reconstruction performance (PSNR) of different methods

Inpainting (Poisson noise) & Sparse-view CT (MGP) experiments: See Paper

## Method & Experiment



[1] “Robust Equivariant Imaging: a fully unsupervised framework for learning to image from noisy and partial measurements.”, [Chen et al., CVPR2022 Oral]  
[2] “Equivariant imaging: Learning beyond the range space.”, [Chen et al., ICCV2021]

## ▪ Problem

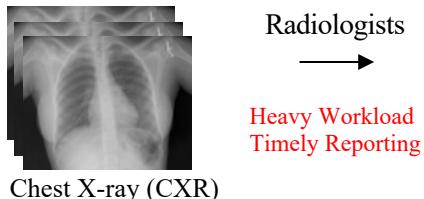


Fig 1. Radiologists' work [1]

## ▪ Anomaly Detection (AD)

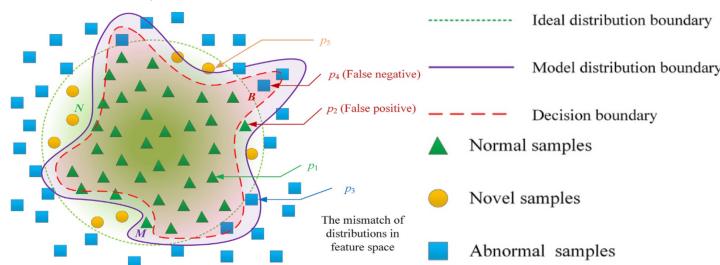


Fig 2. Definition of AD [2]

## ▪ Evaluation Metrics

Receiver Operating Characteristics (ROC) curve

$$TPR = \frac{\text{Number of correctly classified positive samples}}{\text{Number of positive samples}}$$

$$FPR = \frac{\text{Number of misclassified negative samples}}{\text{Number of negative samples}}$$

AUC-ROC,  
Sensitivity,  
Specificity  
etc. [4][5]

## ▪ Improve work efficiency

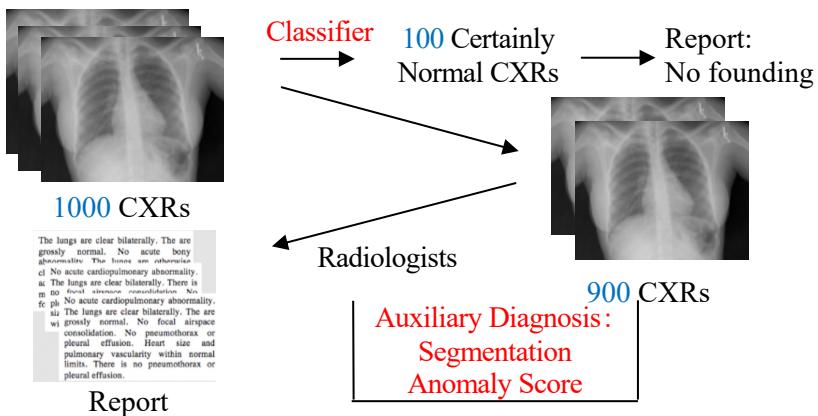


Fig 3. Example of Auxiliary Diagnosis [3]

## ▪ My Project:

Investigate and summarise several significant modern anomaly detection methods based on GANs in CXR

Create a GUI system to visually compare results from different networks and classify CXR in real time by trained models.

[1] “A Self-boosting Framework for Automated Radiographic Report Generation.”, [Wang et al., CVPR2021]

[2] “GAN-based anomaly detection: A review.”, [Xia et al., Neurocomputing2022]

[3] “Chestx-ray8: Hospital-scale chest x-ray database and...” [Wang et al., CVPR2017]

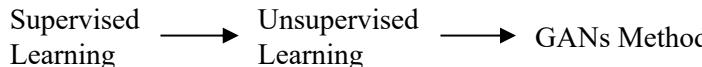
[4] “One-class classification: A survey.” [Perera et al., arXiv2021]

[5] “f-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks.” [Schlegl et al., MIA2019]

## ▪ My Recent work

## CXR AD based on GANs

Deep Learning      CXR Dataset [2]:  
 Good performance    Easily Accessible  
 Data hungry [1]      Already have open dataset



limited due to:  
 Abnormal samples      GAN can make abnormal inferences  
 Unpredictability [3]    using adversarial learning of the representation of samples [3]

## GAN [6]:

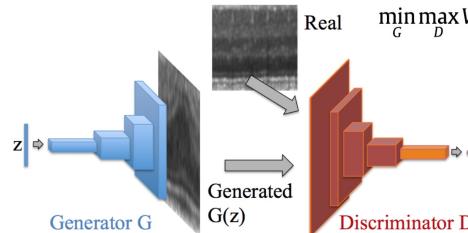


Fig 1. DCGAN Structure [4] [5]

$$\min_G \max_D V(D, G) = \mathbf{E}_{\mathbf{x} \sim P_{\text{data}}(\mathbf{x})}[\log D(\mathbf{x})] + \mathbf{E}_{\mathbf{z} \sim P_z(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$$

For G fixed, the optimal D:

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$$

$$C(G) = \max_D V(G, D) = -\log(4) + 2 \cdot JSD(p_{\text{data}} \| p_g)$$

$p_g = p_{\text{data}}$

## AD based on GAN

Theoretical Support [7][8]:

once the sample distribution and latent variable distribution were correlated through strong constraints, the changes of latent variables could be used to predict whether the sample was abnormal

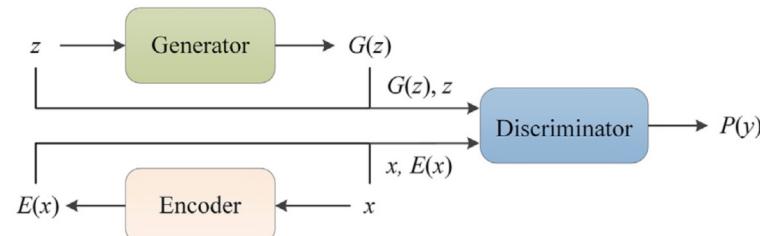


Fig 2. BiGAN and ALI Structure [7] [8]

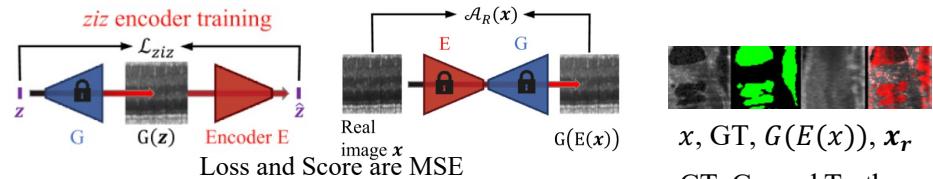


Fig 3. Example of f-AnoGAN [9]

[1] “The AI index 2021 annual report.”, [Zhang et al., Stanford University2021] [2] “Chestx-ray8: Hospital-scale chest x-ray database and...” [Wang et al., CVPR2017]

[3] “GAN-based anomaly detection: A review.”, [Xia et al., Neurocomputing2022]

[4] “Unsupervised representation learning with deep convolutional generative adversarial networks.” [Radford et al., arXiv2015]

[5] “Unsupervised anomaly detection with generative adversarial networks to guide marker discovery.” [Schlegl et al., IPMI2017]

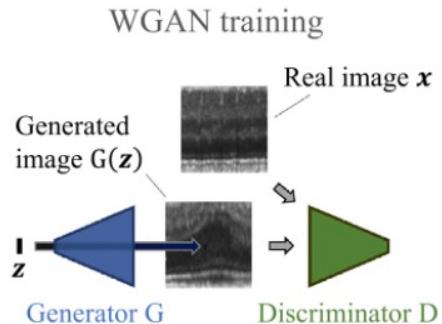
[6] “Generative adversarial nets.” [Goodfellow et al., NIPS2014]

[7] “Adversarial feature learning.” [Donahue et al., ICLR2017]

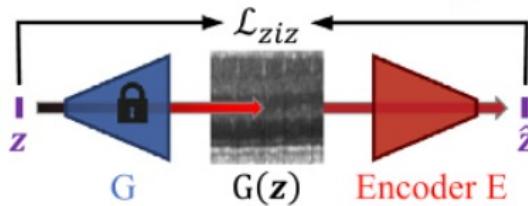
[8] “Adversarially learned inference” [Dumoulin et al., ICLR2017]

[9] “f-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks.” [Schlegl et al., MIA2019]

- Inverse Problem

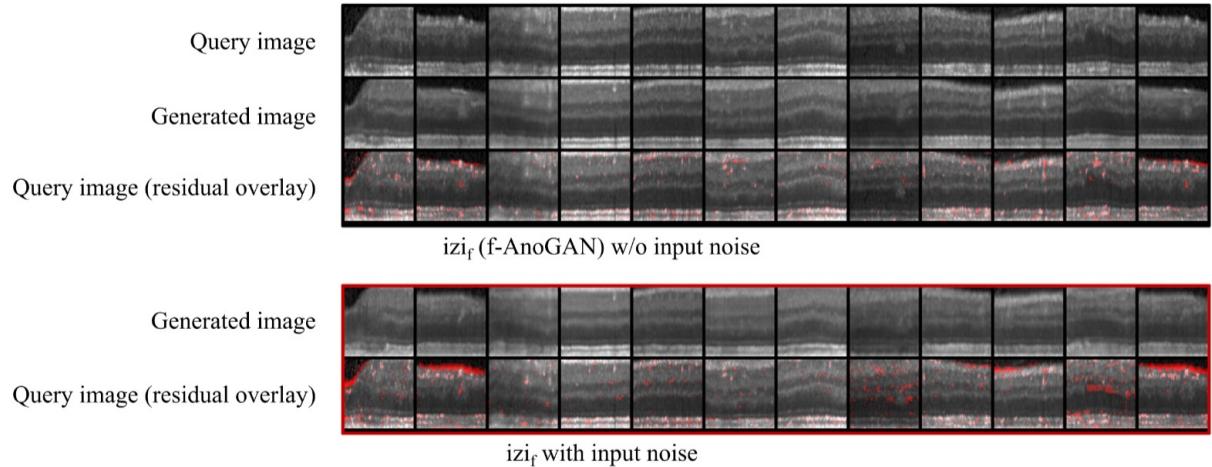


*ziz* encoder training



$x$  maps to  $E(x)$  is an Inverse Problem

- SURE Loss



The previous method is **sensitive to noise**.  
May be SURE can improve the robustness of methods

*Thank you!!!!*

- SURE Gaussian noise proof

$$\text{SURE}(h) = d\sigma^2 + \|g(x)\|^2 + 2\sigma^2 \sum_{i=1}^d \frac{\partial}{\partial x_i} g_i(x) = -d\sigma^2 + \|g(x)\|^2 + 2\sigma^2 \sum_{i=1}^d \frac{\partial}{\partial x_i} h_i(x),$$

We wish to show that

$$\mathbb{E}_\mu \|h(x) - \mu\|^2 = \mathbb{E}_\mu \{\text{SURE}(h)\}.$$

We start by expanding the MSE as

$$\begin{aligned}\mathbb{E}_\mu \|h(x) - \mu\|^2 &= \mathbb{E}_\mu \|g(x) + x - \mu\|^2 \\ &= \mathbb{E}_\mu \|g(x)\|^2 + \mathbb{E}_\mu \|x - \mu\|^2 + 2\mathbb{E}_\mu g(x)^T(x - \mu) \\ &= \mathbb{E}_\mu \|g(x)\|^2 + d\sigma^2 + 2\mathbb{E}_\mu g(x)^T(x - \mu).\end{aligned}$$

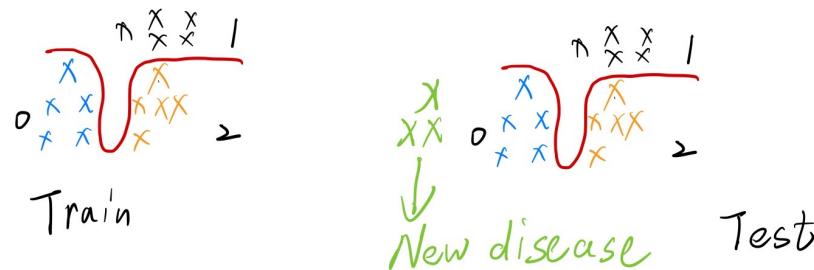
Now we use [integration by parts](#) to rewrite the last term:

$$\begin{aligned}\mathbb{E}_\mu g(x)^T(x - \mu) &= \int_{\mathbb{R}^d} \frac{1}{\sqrt{2\pi\sigma^{2d}}} \exp\left(-\frac{\|x - \mu\|^2}{2\sigma^2}\right) \sum_{i=1}^d g_i(x)(x_i - \mu_i) d^d x \\ &= \sigma^2 \sum_{i=1}^d \int_{\mathbb{R}^d} \frac{1}{\sqrt{2\pi\sigma^{2d}}} \exp\left(-\frac{\|x - \mu\|^2}{2\sigma^2}\right) \frac{dg_i}{dx_i} d^d x \\ &= \sigma^2 \sum_{i=1}^d \mathbb{E}_\mu \frac{dg_i}{dx_i}.\end{aligned}$$

Substituting this into the expression for the MSE, we arrive at

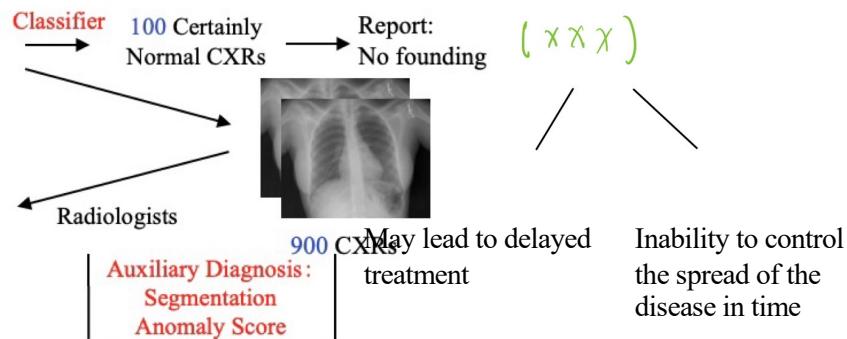
$$\mathbb{E}_\mu \|h(x) - \mu\|^2 = \mathbb{E}_\mu \left( d\sigma^2 + \|g(x)\|^2 + 2\sigma^2 \sum_{i=1}^d \frac{dg_i}{dx_i} \right).$$

## ▪ Supervised Learning Problem



The lungs are clear bilaterally. The airways are grossly normal. No acute heart gaseous. The bones are otherwise normal. No findings of pulmonary embolism. The lungs are clear bilaterally. There is no evidence of cardiomegaly. No findings of pulmonary embolism. The lungs are clear bilaterally. The airways are grossly normal. No focal findings of pneumonia or other acute respiratory distress syndrome. There are no pleural effusions. Heart size and position are within normal limits. There are no pneumothorax or pleural effusion.

Report



## Theory & Technology

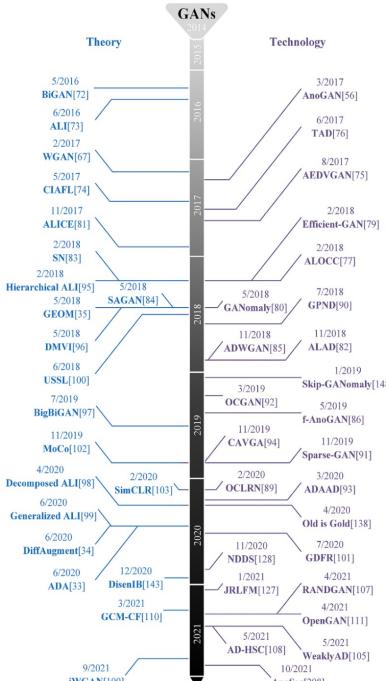


Fig 1. Time series diagram of the theory and application of GAN-based anomaly detection [2]

## CXR AD based on GAN

Model	Dataset	AUC
AnoGAN	COVIDx (balanced test set)	0.54
AnoGAN	Segmented COVIDx (balanced test set)	0.71
RANDGAN	Segmented COVIDx (balanced test set)	<b>0.77</b>
RANDGAN	Segmented COVIDx (imbalanced test set)	0.76

Fig 2. RANDGAN (proposed in 2021) [1]  
Only compared AnoGAN (proposed in 2017)

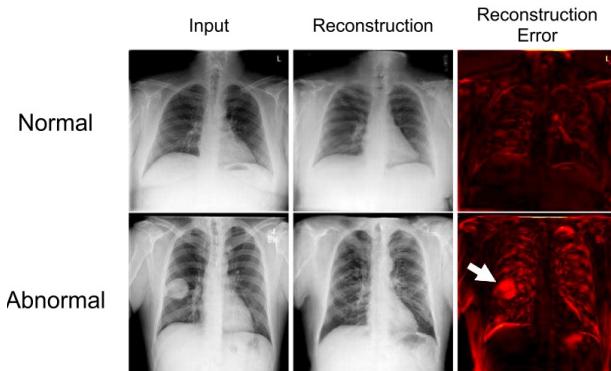


Fig 3.  $\alpha$ GAN (proposed in 2021) [3]  
No compare & similar with f-AnoGAN (proposed in 2019)

## Supplement

### Research Gap

Not much research & Lack of comparisons in modern anomaly detection methods based on GANs.

#### Reason:

1. Cross field. At present, many methods are still based on traditional CNN.
2. Many experts work on general anomaly detection methods.
3. Current anomaly detection methods lack explainable research.

But these are no reasons not to research, and my project can make a small contribution to this research gap

SOTA performance is not high

Adjusting parameters, modifying models, designing for losses

New Technologies, like: Classifier free guidance

### Significance

1. Research significance:  
Research Gap

2. Application significance:  
GUI

- [1] "RANDGAN: randomized generative adversarial network for detection of COVID-19 in chest X-ray." [Motamed et al., Scientific Reports 2021]
- [2] "GAN-based anomaly detection: A review.", [Xia et al., Neurocomputing 2022]
- [3] "Unsupervised deep anomaly detection in chest radiographs." [Nakao et al., Journal of Digital Imaging 2021]

- Challenge

- Instability of Training

(Proof details in supplement slides)

[1][2][3][4]

If D is optimal, loss is:

$$C(G) = \max_D V(G, D)$$

$$= -\log(4) + 2 \cdot JSD(p_{\text{data}} \| p_g)$$

When  $p_g$  and  $p_{\text{data}}$  do not overlap:

$$JSD = 0 \text{ (No Gradient)}$$

No overlapping:

- Image are 2D manifold of 3D
- $p_g \sim G(z)$  It is a mapping from low dimension to high dimension
- Each training is a sampling of the data

WGAN-GP: There are still problems

$$\min_G \max_D V(D, G) = \mathbf{E}_{z \sim P_z(\mathbf{z})} [D(G(\mathbf{z}))] - \mathbf{E}_{x \sim P_{\text{data}}(\mathbf{x})} [D(\mathbf{x})] + \lambda_{gp} \mathbf{E}_{\hat{\mathbf{x}} \sim P(\hat{\mathbf{x}})} \left[ (\|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}})\|_2 - 1)^2 \right]$$

- Challenge (Cont.)

- Dataset

CXR data are high-resolution images, but training and detection on high-resolution images can take up a lot of computational resources

- Evaluation Metrics

ROC-curve metrics can only measure the degree of anomaly in a way that is uninterpretable. [4]

Previous works had [5][6][7][8]:

Turing test, Residual image, Distribution

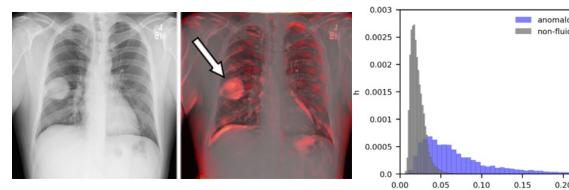


Fig 1. Left: Residual image, Right: Distribution

[1] “Wasserstein generative adversarial networks.”, [Arjovsky et al., ICLR2017]

[2] “Towards principled methods for training generative adversarial networks.”, [Arjovsky et al., arXiv2017]

[3] “Improved training of wasserstein gans.”, [Gulrajani et al., NIPS2017]

[4] “GAN-based anomaly detection: A review.”, [Xia et al., Neurocomputing2022]

[5] “f-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks.” [Schlegl et al., MIA2019]

[6] “Unsupervised representation learning with deep convolutional generative adversarial networks.” [Radford et al., arXiv2015]

[7] “RANDGAN: randomized generative adversarial network for detection of COVID-19 in chest X-ray.” [Motamed et al., Scientific Reports 2021]

[8] “Unsupervised deep anomaly detection in chest radiographs.” [Nakao et al., Journal of Digital Imaging 2021]

## GAN Global Minimum

GAN Global minimum

$$V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{x \sim p_g} [\log (1 - D(x))]$$

For fixed  $G$ ,

$$\max_D V(D, G)$$

$$\text{When } P_d = \frac{P_d + P_g}{2}$$

$$= \int (P_{\text{data}} \cdot \log D + P_g \cdot \log (1 - D)) dx$$

$$\frac{\partial}{\partial D} = \frac{P_{\text{data}}}{D} + \frac{-P_g}{1-D}$$

$$D_g^* = \frac{P_{\text{data}}}{P_{\text{data}} + P_g}$$

Let  $D_g^*$

$$\min_G = \mathbb{E}_{x \sim p_{\text{data}}} [\log \frac{P_g}{P_d + P_g}] + \mathbb{E}_{x \sim p_g} [\log (\frac{P_d}{P_d + P_g})]$$

$$\log \frac{P_{\text{data}}}{(P_{\text{data}} + P_g)/2} \cdot \frac{1}{2} \left( \log \frac{P_d}{P_d + P_g} \right) \text{ flog}$$

$$\begin{aligned} & \text{KL}(P_d \parallel \frac{P_d + P_g}{2}) \log_2 \frac{P_d}{\frac{P_d + P_g}{2}} \\ & \quad + \mathbb{E}_{x \sim p_g} [\log (\frac{P_g}{P_d + P_g})] \end{aligned}$$

## GAN Training Problem

$$\text{In } D^*, \min_a = 2JS[P_d || Pg] - 2\log_2 P_d || Pg$$

$$KL(P_1 || P_2) = \mathbb{E}_{P_1} \log \frac{P_1}{P_2} \quad P_1 \in [0, 1], P_2 \in [0, 1]$$

$$JS(P_1 || P_2) = \frac{1}{2} KL(P_1 || \frac{P_1 + P_2}{2}) + \frac{1}{2} KL(P_2 || \frac{P_1 + P_2}{2})$$

$$= \frac{1}{2} \left[ \int p_1 \log \frac{2p_1}{P_1 + P_2} dx + \int p_2 \log \frac{2p_2}{P_1 + P_2} dx \right]$$

$$= \log_2 + \frac{1}{2} \int p_1 \log \frac{P_1}{P_1 + P_2} dx + \frac{1}{2} \int p_2 \log \frac{P_2}{P_1 + P_2} dx$$

When  $P_1$   $P_2$ ,  $J_S = \log_2$   
are not overlapping

# WGAN

Wasserstein distance      Earth - Mover distance

$$W(P_d, P_g) = \inf_{\pi \sim \Pi(P_d, P_g)} \mathbb{E}_{(x, y) \sim \pi} [ \|x - y\| ]$$



$$= \frac{1}{K} \sup_{\|f\|_L \leq K} \mathbb{E}_x \sim p_d [f(x)] - \mathbb{E}_x \sim p_g [f(x)]$$

(Lipschitz

$$|f(x_1) - f(x_2)| \leq K |x_1 - x_2|$$

Any  $x_1, x_2$ ,

## WGAN

$$K \cdot W(\Pr, \Pr_g) \approx \max_{\mathbb{E}_{x \sim p_d} [f_w(x)] - \mathbb{E}_{x \sim p_g} [f_w(x)]}$$

$$\|w\|_2 \leq k$$

Expr  $\mathbb{E}_{x \sim p_d} [f_w(x)] \rightarrow D(x)$   
 $\mathbb{E}_{x \sim p_g} [f_w(x)] \rightarrow D(x)$

Trick: Clipping  
 $[-c, c] \cap [-1, 1]$

$$\min_{G} \max_{D} \mathbb{E}_{x \sim p_d} D(x) - \mathbb{E}_{x \sim p_g} D(G(z))$$

$$G \text{ Loss} = -\mathbb{E}_{x \sim p_g} D(G(z))$$

$$D \text{ Loss} = \mathbb{E}_{x \sim p_g} D(G(z)) - \mathbb{E}_{x \sim p_d} D(x)$$

**Algorithm 1** WGAN, our proposed algorithm. All experiments in the paper used the default values  $\alpha = 0.00005$ ,  $c = 0.01$ ,  $m = 64$ ,  $n_{\text{critic}} = 5$ .

**Require:** :  $\alpha$ , the learning rate.  $c$ , the clipping parameter.  $m$ , the batch size.  
 $n_{\text{critic}}$ , the number of iterations of the critic per generator iteration.

**Require:** :  $w_0$ , initial critic parameters.  $\theta_0$ , initial generator's parameters.

```

1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$ 
12: end while
```

*Adam*

① remove sigmoid

② No log

③ clip

④ momentum  
Not good

## WGAN - GP

$$V(G, D) = \max_{D \text{ Lipschitz}} \left\{ \mathbb{E}_{x \sim p_{\text{data}}} [D_x] - \mathbb{E}_{x \sim p_G} [D_x] \right\}$$

Diagram illustrating the WGAN-GP loss function. It shows two overlapping distributions,  $p_D$  (blue) and  $p_G$  (red), with a boundary between them. A point  $x$  is sampled from the combined distribution. A blue circle labeled "penalty" indicates the distance from the boundary to  $x$ , with the condition  $\|\nabla_x D(x)\| \leq 1$ .

$$V(G, D) \approx \max_D \left\{ \mathbb{E}_{x \sim p_{\text{data}}} [D_x] - \mathbb{E}_{x \sim p_G} [D_x] \right\}$$

Derivation of the WGAN-GP loss:

$$\hat{x}_i = x_{d_i} \cdot \beta + x_{g_i} \cdot (1-\beta)$$

$$\int_X \max(0, \|\nabla_x D_x\| - 1) dx$$

$$- \mathbb{E}_{x \sim \text{penalty}} (\nabla_x D_x - 1)$$