# Generative Modeling for LArTPC Images
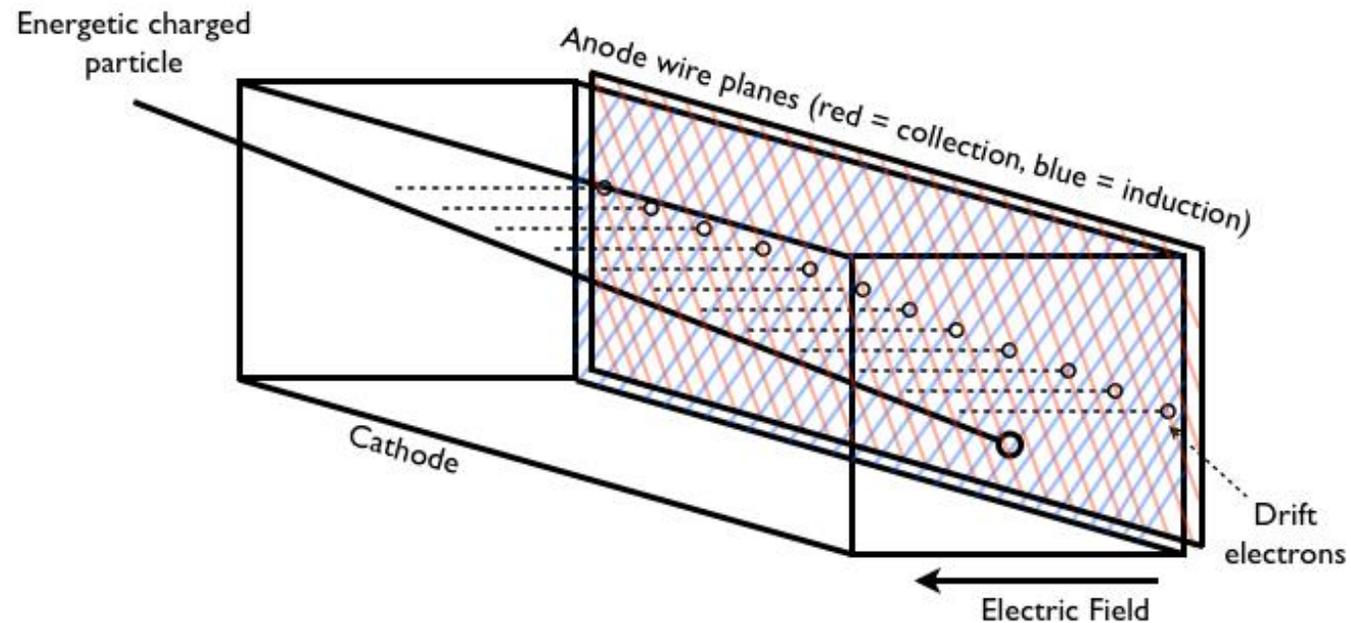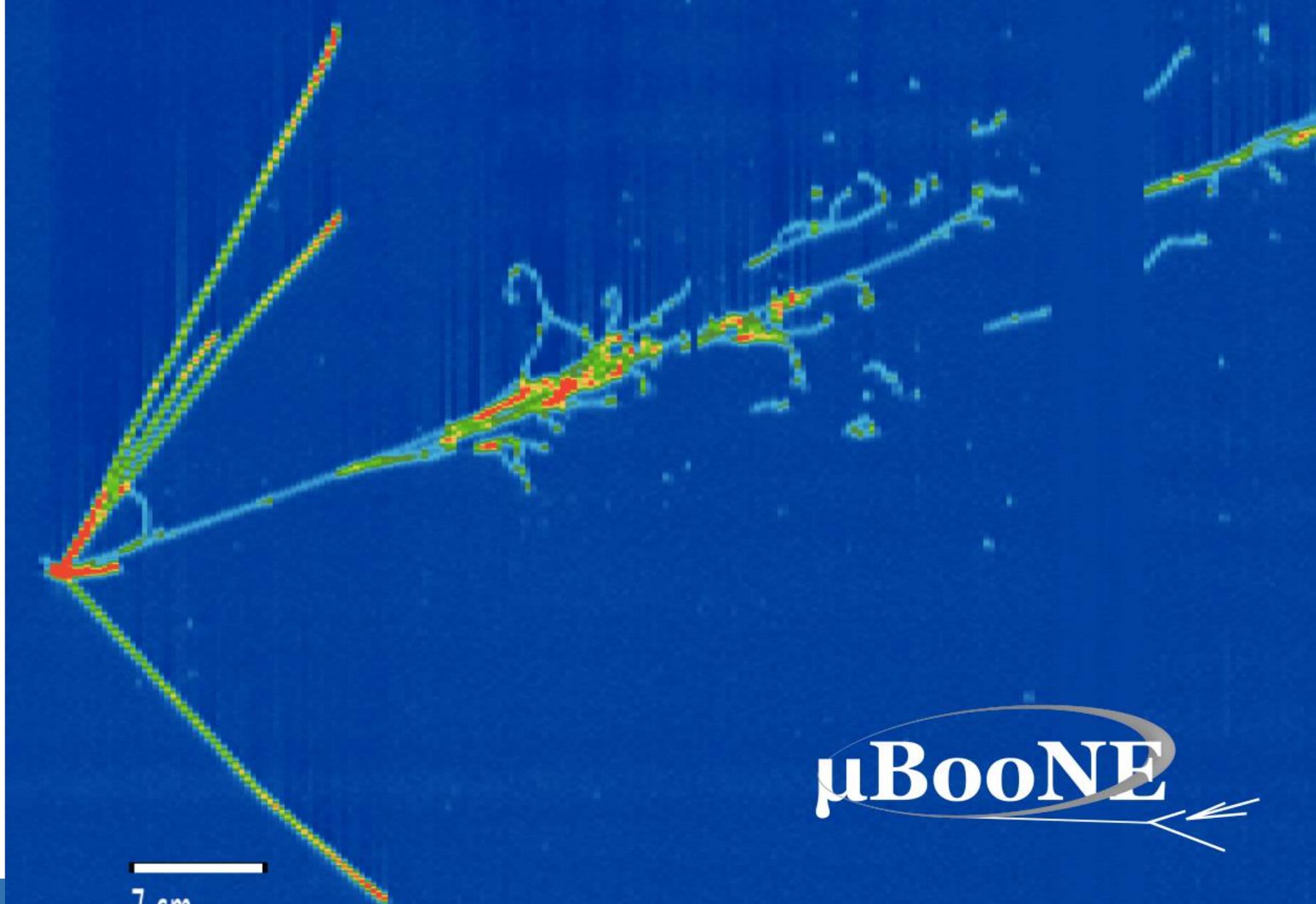
Zev Imani

NPML 2024

# Outline

1. Data Motivation

2. LArTPC Image Generation Attempts

3. ~~Diffusion Methodology~~

4. Quality Tests (Abridged)

5. Distance Metrics

6. Takeaways

# Liquid Argon Time Projection Chamber (LArTPC)

- Detector for HEP experiments

  - Ongoing neutrino research
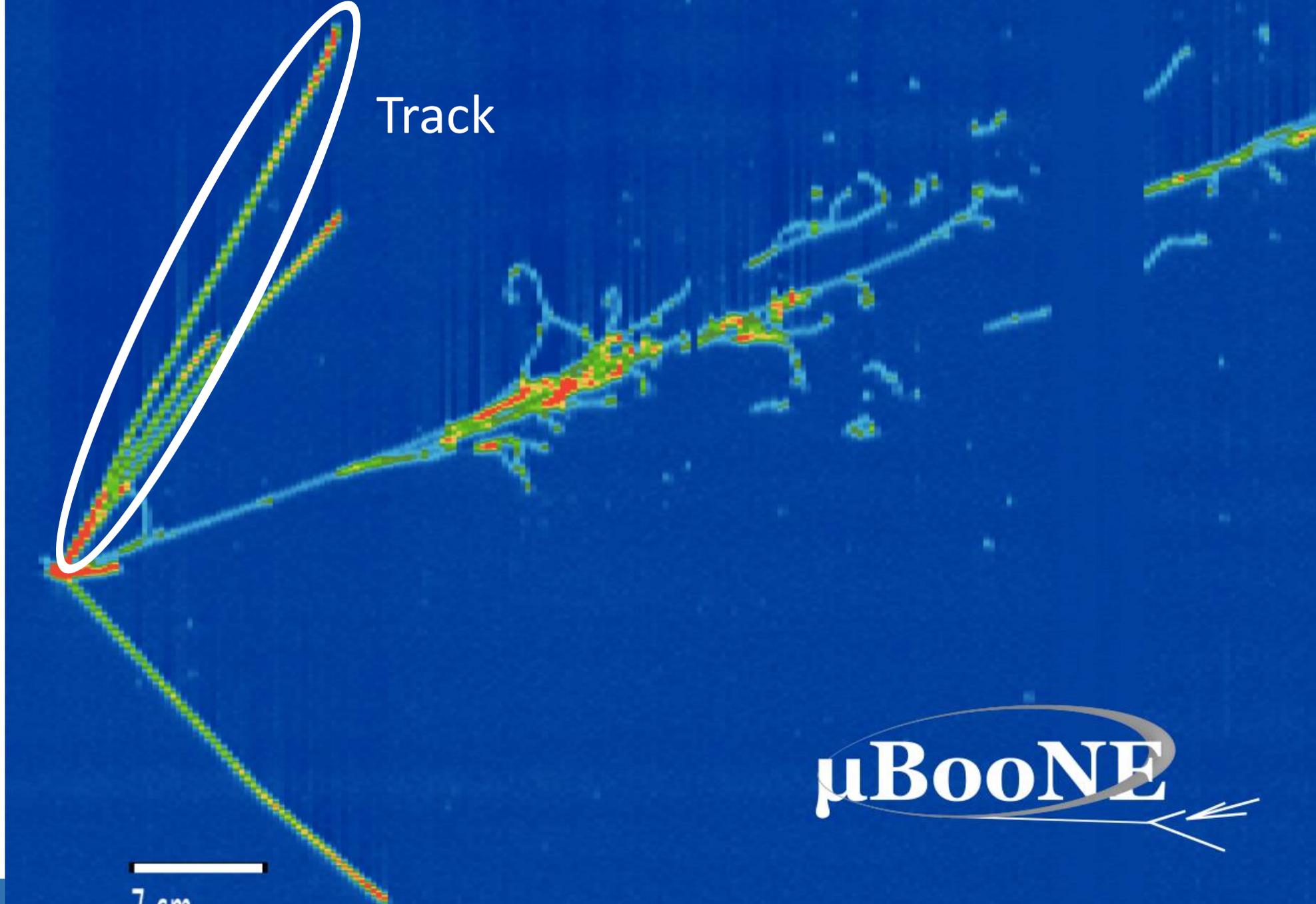
  - Particle interaction images

7 cm
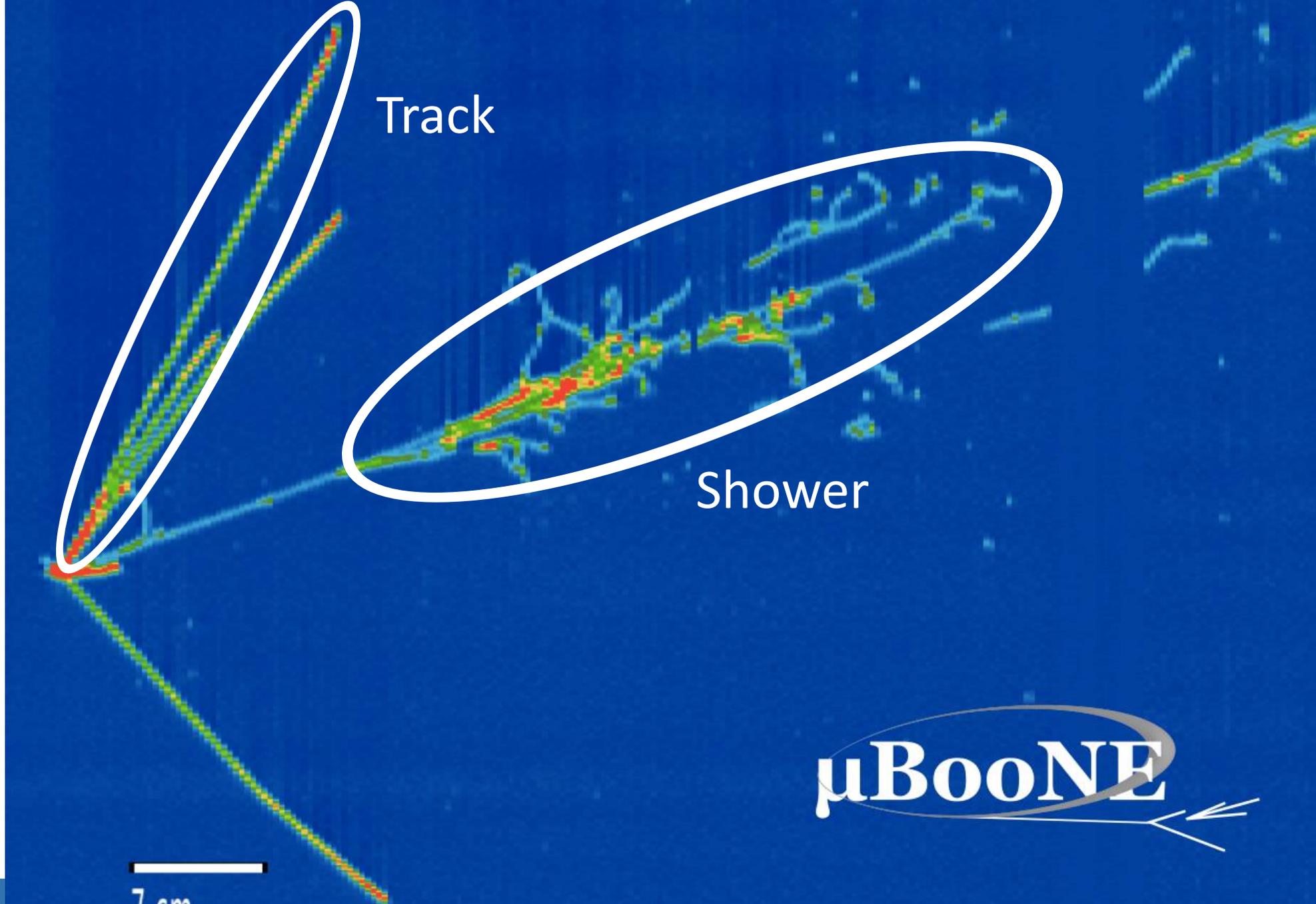
μBooNE

NuMI DATA: RUN 10811, EVENT 2549. APRIL 9, 2017.

Track

7 cm

μBooNE

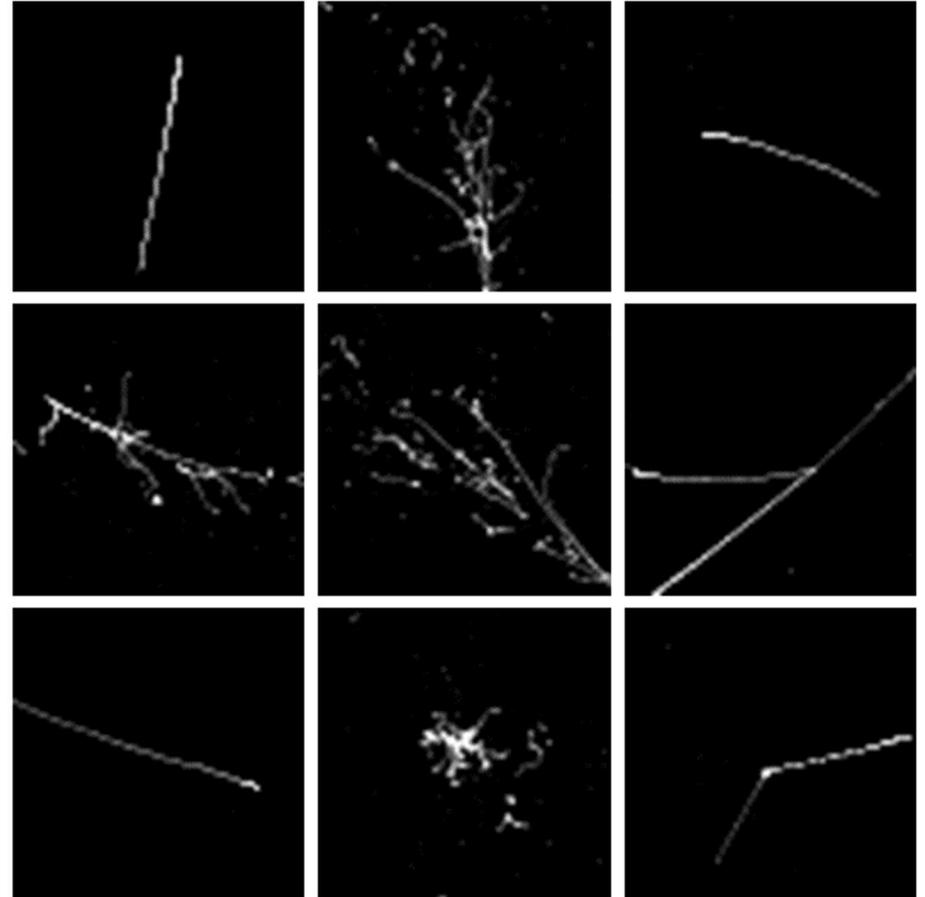NuMI DATA: RUN 10811, EVENT 2549. APRIL 9, 2017.

Track

Shower

μBooNE

NuMI DATA: RUN 10811, EVENT 2549. APRIL 9, 2017.

7 cm

# LArTPC Images

- Cropped image from detector
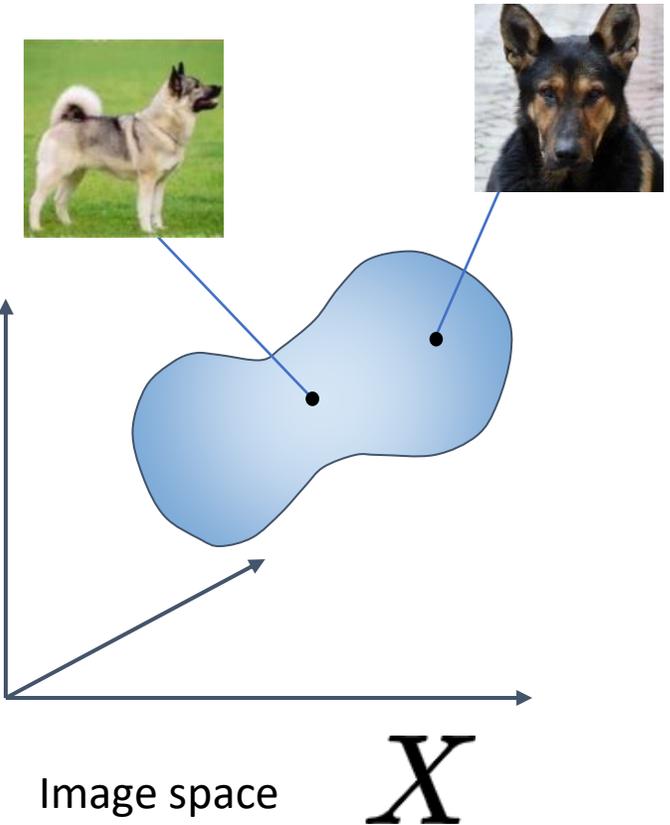
- Globally sparse, but locally dense

# Why Generative Modeling

- Observing rare neutrino events requires analyzing large datasets

- Potential to be faster than traditional simulation methods

- New tool for reconstruction and analyses

- Another way of understanding our data
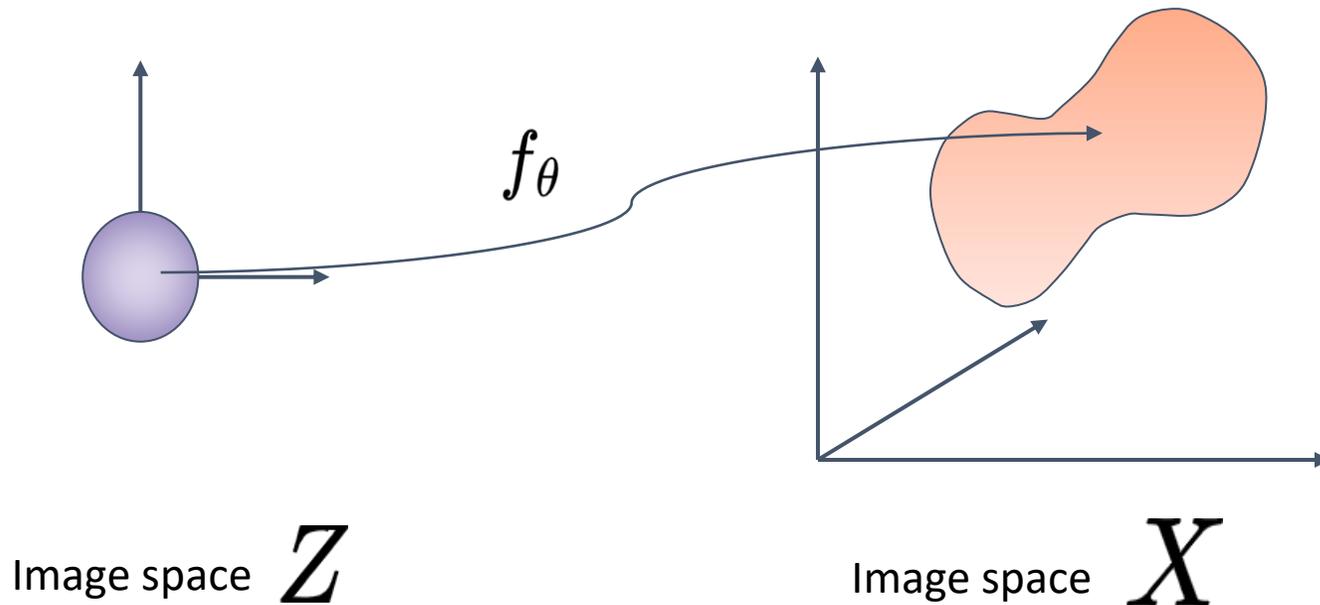
- Proof of concept ML application

# How to Generate Images

- Our data **x** is sampled from some p(**x**)
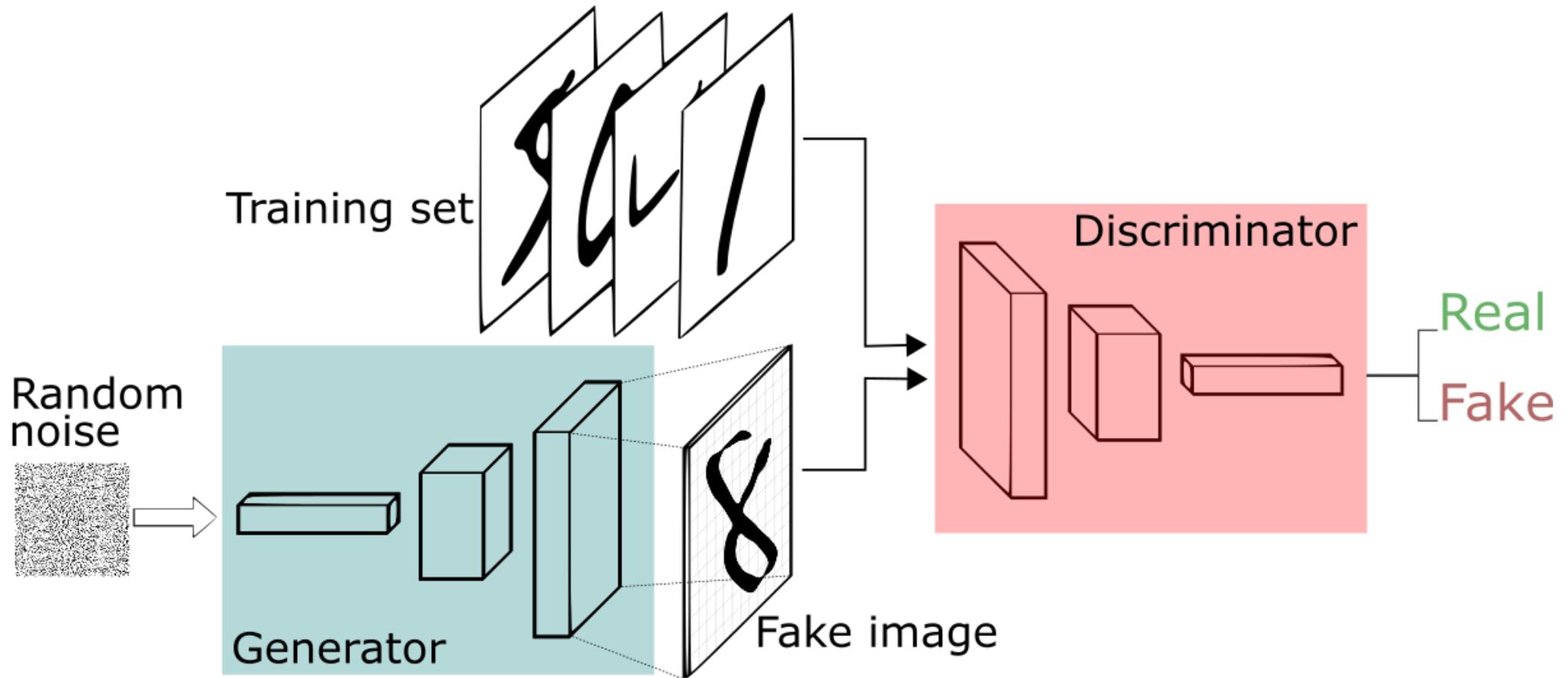
- We don't know p(**x**) directly



Image space $X$

# How to Generate Images

- Instead, we sample from a known distribution $z \sim \mathcal{N}(0, 1)$

- Learn a mapping $x = f_\theta(z)$

$f_\theta$
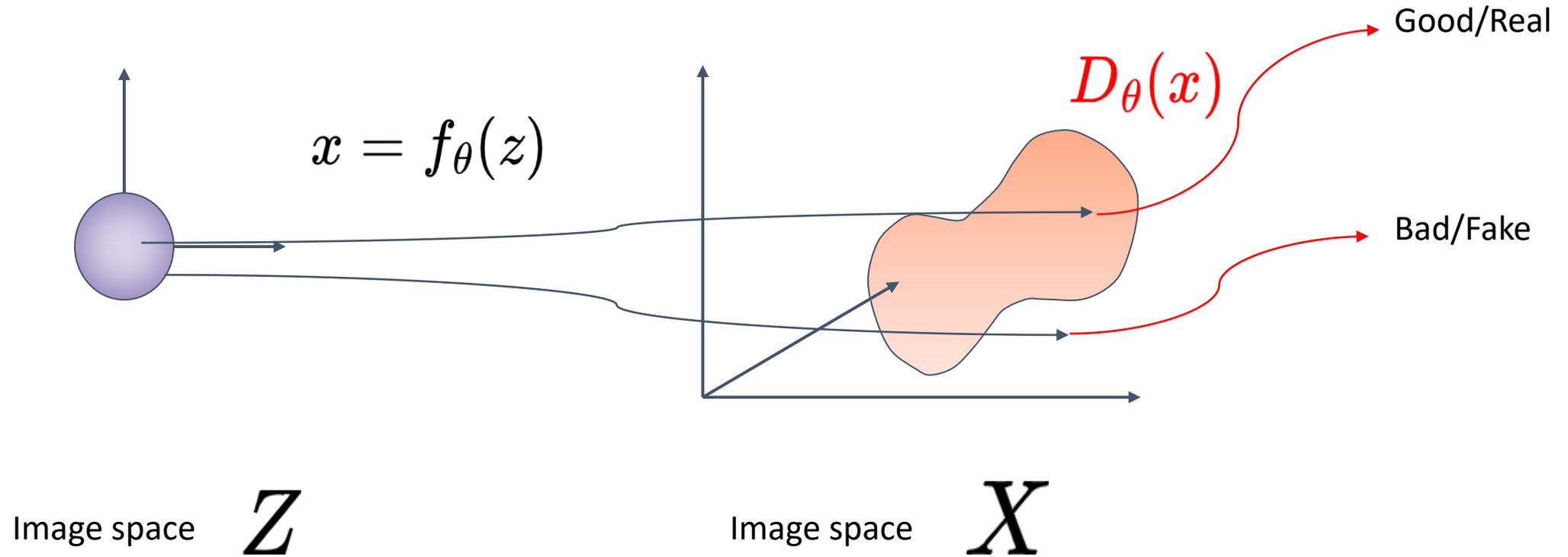
$$p(x) = p(f_\theta(z))$$

Image space $Z$

Image space $X$

# Attempt 1: Generative Adversarial Network

# GAN Mapping



$$x = f_\theta(z)$$

$D_\theta(x)$

Good/Real

Bad/Fake

Image space $Z$

Image space $X$
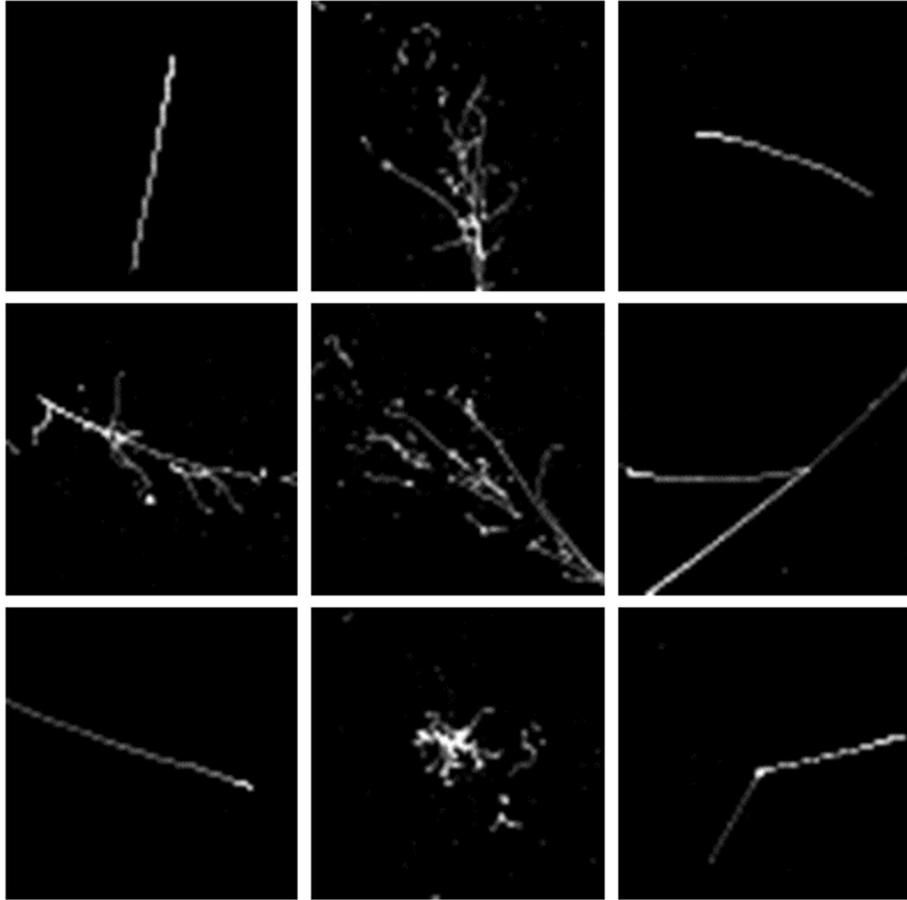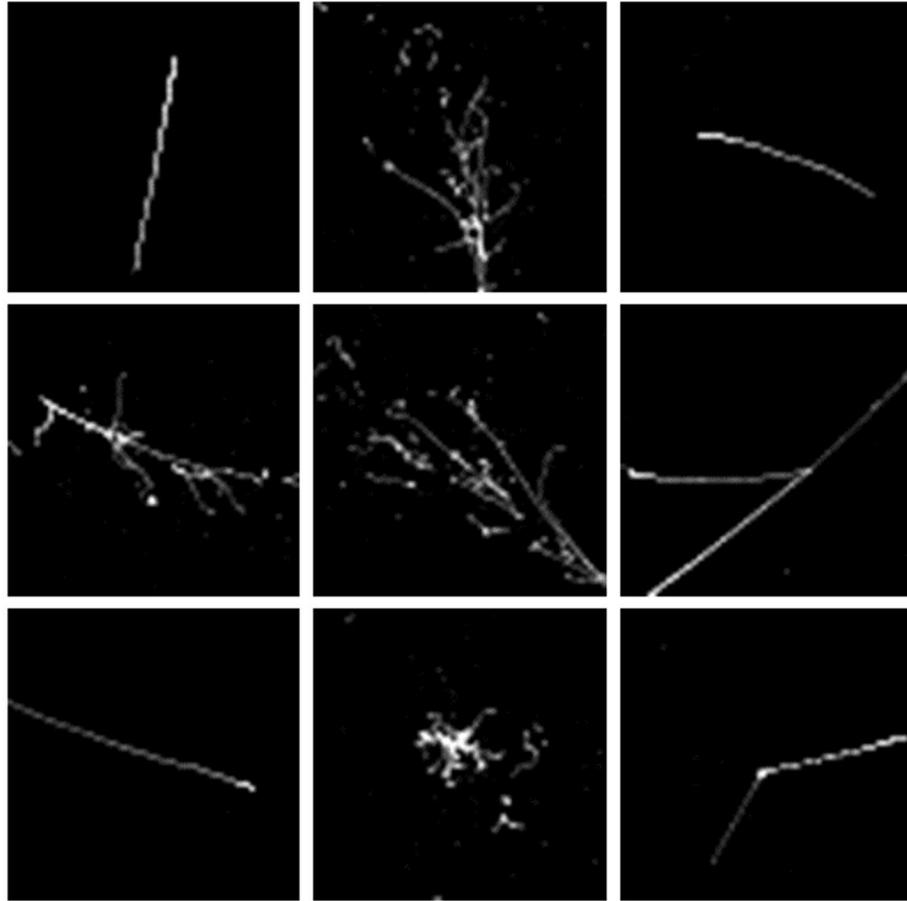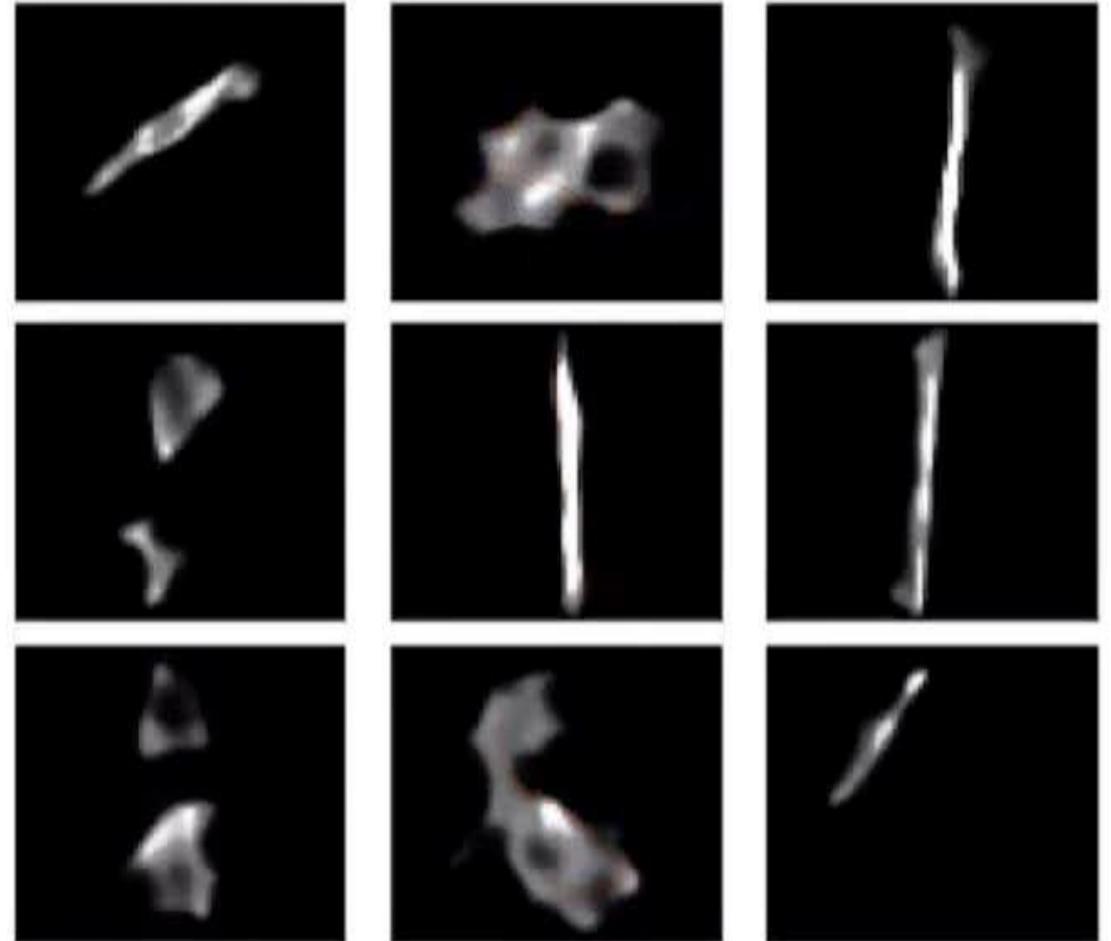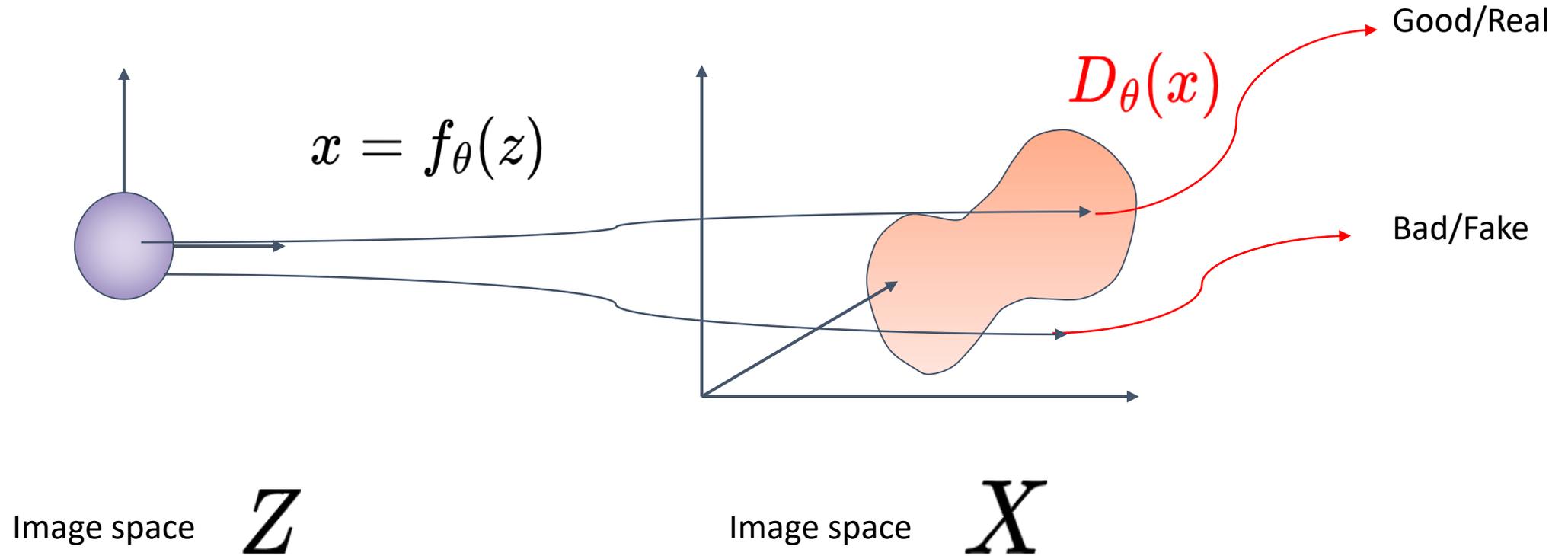
# LArTPC GAN

Validation LArTPC Data

# LArTPC GAN

Validation LArTPC Data

GAN Generated

# GAN Mapping

# GAN Mapping

- LArTPC images exist as thin manifold in image space



$$x = f_\theta(z)$$

$$D_\theta(x)$$

Good/Real

Bad/Fake

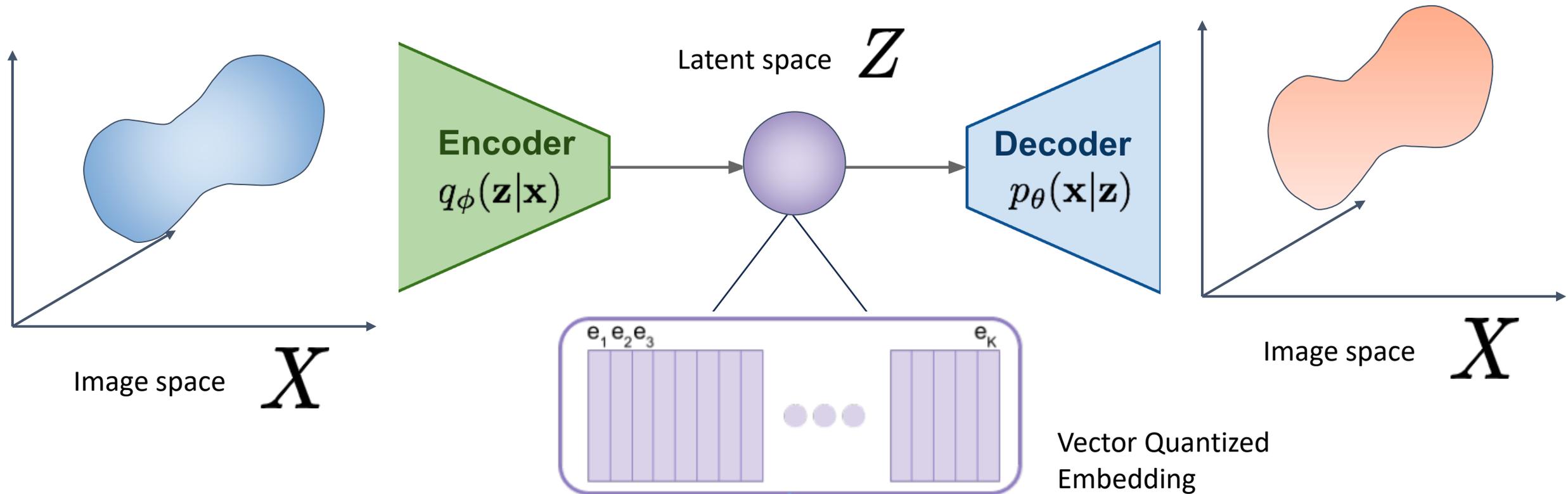Image space $Z$

Image space $X$

# Attempt 2: VQ-VAE

- Vector Quantized Variational Autoencoder

# Attempt 2: VQ-VAE

- Vector Quantized Variational Autoencoder

# LArTPC VQ-VAE

Validation LArTPC Data

Lutkus, Wongjirad, & Aeron; arXiv:2204.02496

# LArTPC VQ-VAE

Validation LArTPC Data

VQ-VAE Generated

# What is Good Enough?

- No standard quality tests for LArTPC images

- 64x64 are too small for traditional physics analysis

- We developed several options

# Semantic Segmentation Network (SSNet)



Background

Track

Shower

**Tracks**

**Showers**

# Attempt 3: Diffusion



$p_t(x)$ @ t=0

$p_t(x)$ @ t=T

Image space $X$

Image space $Z$

$\mathbf{x}_0 \quad \mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \dots \quad \mathbf{z}$

# Attempt 3: Diffusion

Validation LArTPC Data

Imani, Aeron, & Wongjirad; PhysRevD.109.072011

# Attempt 3: Diffusion

Validation LArTPC Data

Diffusion Generated

Imani, Aeron, & Wongjirad; PhysRevD.109.072011

# Tracks

# Showers

# Physics Quality Tests: Showers



Shower Charge Distribution

Legend:
- LArTPC Val (N=3102)
- VQ-VAE (N=8667)
- Diffusion (N=15791)

# Physics Quality Tests: Tracks



Track Length Distribution

LArTPC Tracks Val (N=6898)
VQ-VAE (N=41333)
Diffusion (N=34209)

Track Width Distribution

LArTPC Tracks Val (N=6898)
VQ-VAE (N=41333)
Diffusion (N=34209)

# Additional Quality Tests

- High dimensional goodness of fit tests
  - Maximum Mean Discrepancy (MMD)
  - Sinkhorn divergence
  - Wasserstein-1 (EMD)

- SSNet-FID

- Turing test survey

# Next Steps

- Scale up to larger images
  - Goal of 512x512 image size to do physics analyses
  - Use latent diffusion to overcome scaling issue

- Conditional generation on energy and particle type

- Improve generation speed and efficiency

# Visualizing Distributions

- T-distributed Stochastic Neighbor Embedding (T-SNE)

- Nonlinear dimensionality reduction, maintains relative distance



MNIST T-SNE

# T-SNE on LArTPC

- Pretty, but no clear structure



LArTPC Val + Gen Data

Legend: Val Showers, Val Tracks, Gen Showers, Gen Tracks

# T-SNE on LArTPC

- Darker points = longer/more charge



Euclidean T-SNE

Legend:
- Val Tracks (blue circle)
- Gen Tracks (green plus)
- Val Showers (gray square)
- Gen Showers (red triangle)

# Digression: Distance Metrics

- Euclidian distance (L2 norm)    $\|\boldsymbol{x}\|_2 := \sqrt{x_1^2 + \cdots + x_n^2}$

# Digression: Distance Metrics

- Euclidian distance (L2 norm)

$$\|\boldsymbol{x}\|_2 := \sqrt{x_1^2 + \cdots + x_n^2}$$

- Earth Mover's Distance (EMD)

$$\mathrm{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim P}[f(x)] - \mathbb{E}_{y \sim Q}[f(y)]$$

  - Wasserstein-1 distance

  - 'Natural' metric for particle physics

$$\min_F \sum_{i=1}^{m} \sum_{j=1}^{n} f_{i,j} d_{i,j}$$

# Digression: Distance Metrics

- Euclidian distance (L2 norm)

$$\|\boldsymbol{x}\|_2 := \sqrt{x_1^2 + \cdots + x_n^2}$$

- Earth Mover's Distance (EMD)
  - Wasserstein-1 distance
  - 'Natural' metric for particle physics

- red distribution: "dirt"
- blue distribution: "holes"

Komiske, Metodiev, & Thaler; arXiv:1902.02346

# T-SNE EMD

- Separation of track and shower events

- Ongoing exploration of this data representation



EMD T-SNE

Legend:
- Val Tracks (blue circle)
- Gen Tracks (green plus)
- Val Showers (gray square)
- Gen Showers (red triangle)

# Key Takeaways

1. LArTPC data differs from natural images

   - Globally sparse, but locally dense

2. Diffusion is a versatile method of data generation

   - Can handle our LArTPC data

3. Development of some quality metrics for LArTPC images

4. Earth Mover's Distance is a useful metric for particle event data

*Score-based Diffusion Models for Generating Liquid Argon Time Projection Chamber Images*

By Zeviel Imani, Shuchin Aeron, & Taritree Wongjirad

PhysRevD.109.072011

# Questions?

# Backup Slides

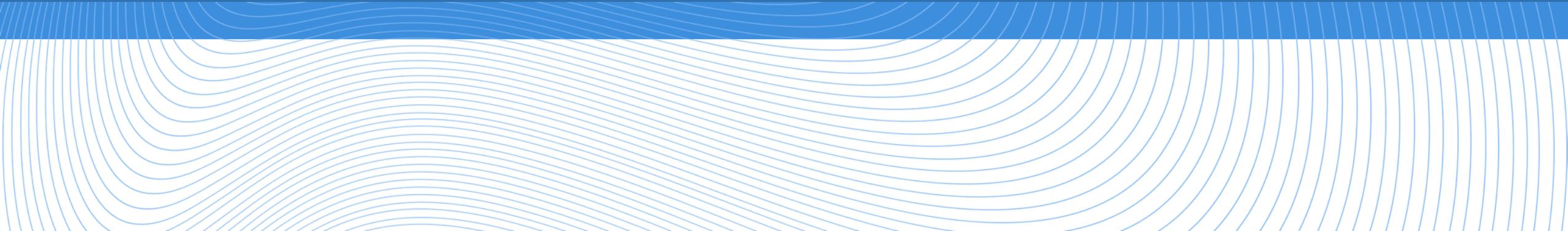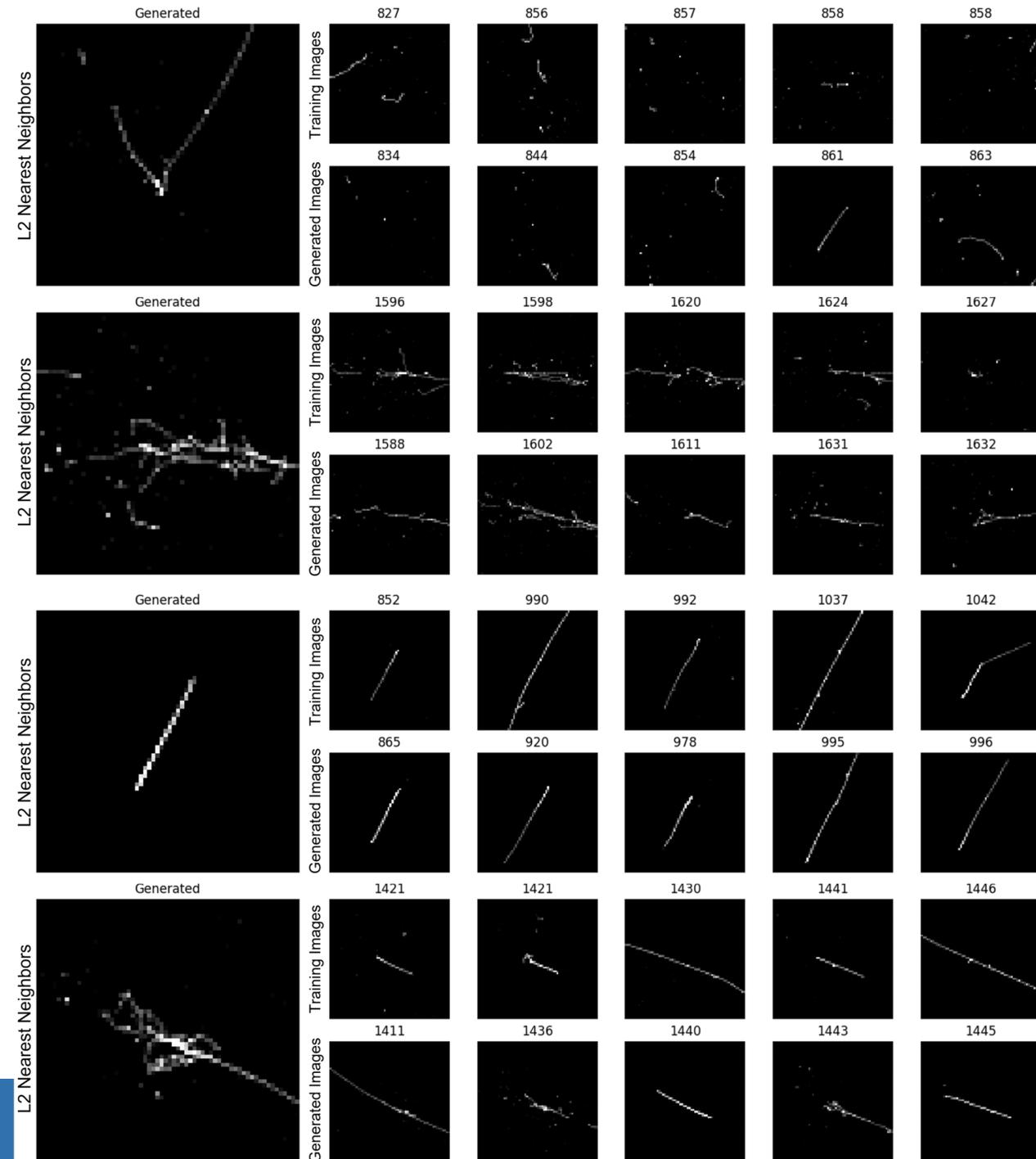**(and skipped sections)**

# Mode Collapse

- Nearest neighbors using

  L2 Euclidian Norm distance

# Mode Collapse

- Nearest neighbors using

  Earth Mover's Distance (EMD)

Shower Charge Distribution

Track Length Distribution

Track Width Distribution

# Physics Metrics: Chi-Squared

| $\chi^2$ Test | Track Length | Track Width | Shower Charge |
|---|---|---|---|
| 10 Epochs | 206 | 825 | 6458 |
| 50 Epochs | **126** | 418 | **228** |
| 150 Epochs | 130 | **175** | 382 |

Loss

$$\mathcal{L}(\boldsymbol{\theta}; t) = \frac{1}{N} \sum_i^N || \boldsymbol{s_\theta}(\vec{\boldsymbol{X}}_t^i, t) - \frac{-(\vec{\boldsymbol{X}}_t^i - \vec{\boldsymbol{\mu}}_t^i)}{\vec{\boldsymbol{\sigma^2}}_t^i} ||_2^2$$

Legend:
- Training Set (N=50000)
- Validation Set (N=10000)

# Turing Test



Turing Test Accuracy

Imani, Aeron, & Wongjirad; PhysRevD.109.072011

# High Dimensional Goodness of Fit Tests

# Fréchet Inception Distance (FID)

- Process:

  1. Get **layer activations** from classifier

     - Typically use Google's Inception v3 deepest activation layer (pool3)
       - 2048-dimensional activation vector

  2. Fit activations to multidimensional Gaussian distribution

  3. Find Wasserstein-2 distance between the Gaussians

- We can use activations from SSNet instead

# SSNet-FID

# Conditional 1: Statistical Reframe

- Given random variables **x** (LArTPC image) and **y** (energy) we want to sample from p(**x** | **y**)

- Approach 1) Extend score: $s_\theta(\mathbf{x}, t) \rightarrow s_\theta(\mathbf{x}, t, \mathbf{y})$

- Or…

# Conditional 2: Inverse Problem

- We know how to get **y** (energy) from **x** (LArTPC image)

- Bayes' Rule: $$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x})p(\mathbf{y}|\mathbf{x})}{\int p(\mathbf{x})p(\mathbf{y}|\mathbf{x})\mathrm{d}\mathbf{x}}$$

- Take gradient: $$\nabla_{\mathbf{x}} \log p(\mathbf{x} \mid \mathbf{y}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y} \mid \mathbf{x})$$

**score**        **classifier**

# Score-based Diffusion Model

# How to Generate Images



Data samples

$$\{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_N\} \overset{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

# How to Generate Images



score matching

Data samples

$$\{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_N\} \overset{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

Scores

$$\mathbf{s}_\theta(\mathbf{x}) \approx \nabla_\mathbf{x} \log p(\mathbf{x})$$

# How to Generate Images



Data samples

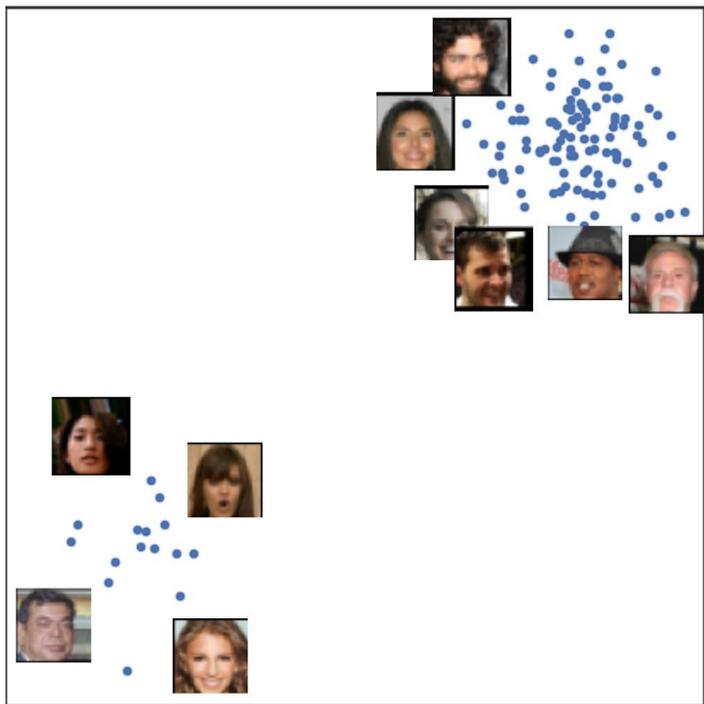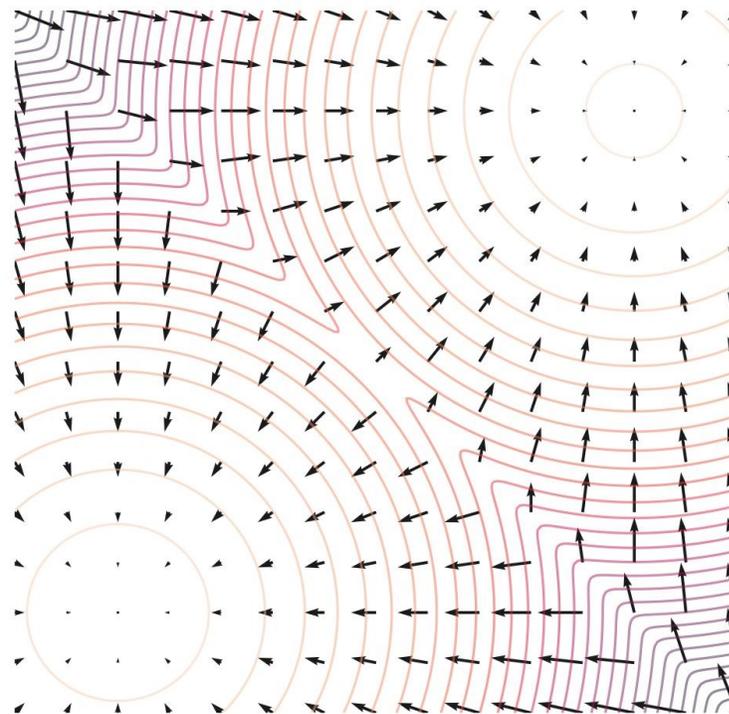$$\{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_N\} \overset{\text{i.i.d.}}{\sim} p(\mathbf{x})$$
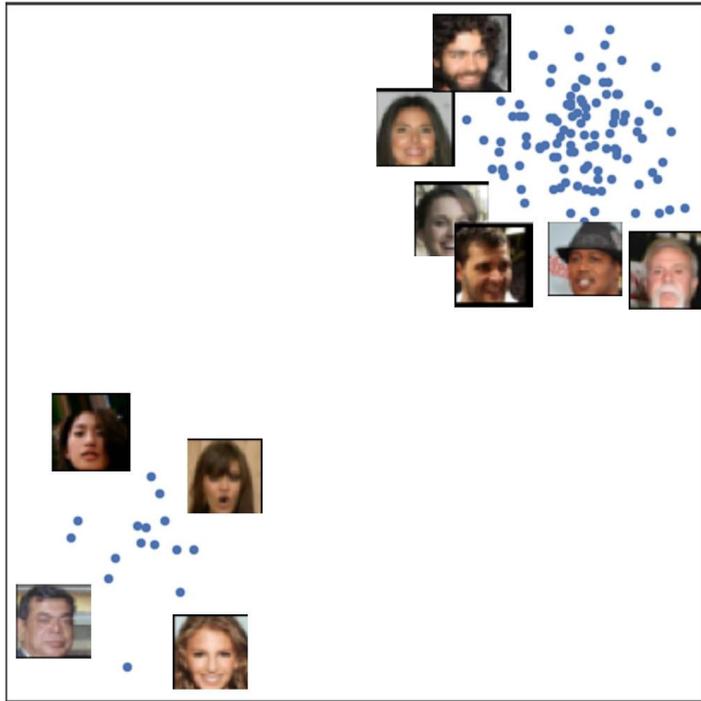
score matching

Scores

$$\mathbf{s}_\theta(\mathbf{x}) \approx \nabla_{\mathbf{x}} \log p(\mathbf{x})$$

Langevin dynamics

New samples

# How to Generate Images



score
matching

Langevin
dynamics

Data samples

Scores

New samples

$$\{\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_N\} \overset{\text{i.i.d.}}{\sim} p(\mathbf{x})$$

$$\mathbf{s}_\theta(\mathbf{x}) \approx \nabla_\mathbf{x} \log p(\mathbf{x})$$

# Manifold Hypothesis

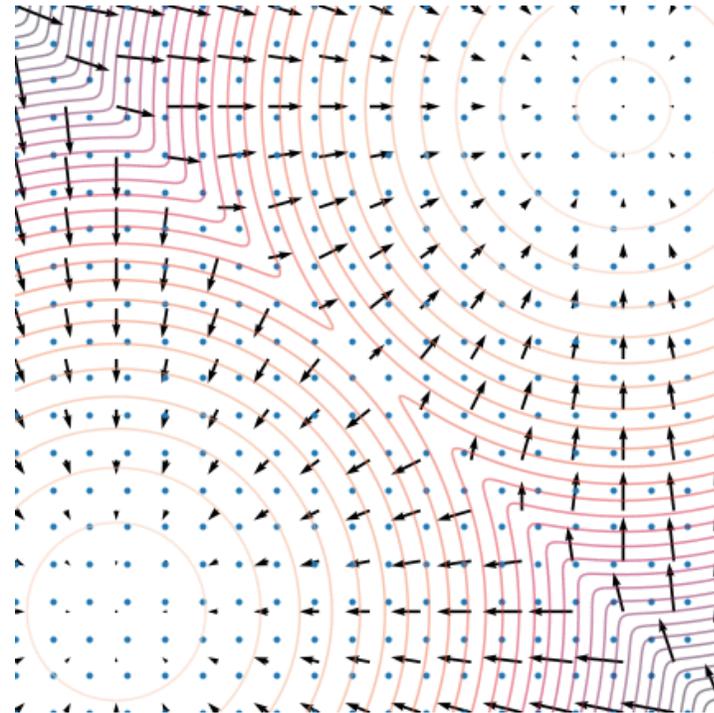# Add Diffusion

Perturbed density

Perturbed scores

Accurate

yang-song.net/blog/2021/score

# Annealed Langevin Sampling



$\sigma_1$        $\sigma_2$        $\sigma_3$

# LArTPC Image Generation

Training Images

Generated Images

# All Together Now

Imani, Aeron, & Wongjirad; PhysRevD.109.072011

# Where is the mapping?



$p_t(x)$

Image space $X$

$p_t(x)$

Image space $X$

# Forward Stochastic Differential Equation (SDE)

Forward SDE (data → noise)

$$\mathbf{dx} = \mathbf{f}(\mathbf{x}, t)\mathrm{d}t + g(t)\mathbf{dw}$$

Drift $\mathbf{f}(\mathbf{x}, t)\mathrm{d}t$

Deterministic evolution

$\mathbf{f}(\mathbf{x}, t) = -\mathbf{x}\frac{1}{2}\beta_t$

d$t$ = time increment

Diffusion $g(t)\mathbf{dw}$

Scale factor $g(t) = \sqrt{\beta_t}$

d$\mathbf{w}$ = Brownian motion
(Random walk)

# Forward SDE



Data — Forward SDE — Prior

$x(0)$ — $\mathrm{d}x = f(x,t)\mathrm{d}t + g(t)\mathrm{d}w$ → $x(T)$

$x$

$p_0(x)$ — $p_t(x)$ → $p_T(x)$

# Forward SDE

# Forward SDE

# Reverse Stochastic Differential Equations (SDE)

Drift (Reverse) $\mathbf{f}(\mathbf{x}, t)\mathrm{d}t$

Diffusion (Reverse) $g(t)\mathrm{d}\bar{\mathbf{w}}$

score function

$$g^2(t)\boxed{\nabla_{\mathbf{x}} \log p_t(\mathbf{x})}$$

Scale factor $g^2(t) = \beta_t$

Reverse SDE (noise → data)

score function

$$\mathrm{d}\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - g^2(t)\boxed{\nabla_{\mathbf{x}} \log p_t(\mathbf{x})}\right]\mathrm{d}t + g(t)\mathrm{d}\bar{\mathbf{w}}$$

# Reverse SDE



Prior — Reverse SDE — Data

$$x(T) \quad \mathrm{d}x = \left[ f(x,t) - g^2(t)\nabla_x \log p_t(x) \right] \mathrm{d}t + g(t)\mathrm{d}\bar{w} \quad x(0)$$

$$p_T(x) \quad\quad\quad p_t(x) \quad\quad\quad p_0(x)$$

# Reverse SDE



Prior           Reverse SDE           Data

$$\mathrm{d}x = \left[f(x,t) - g^2(t)\nabla_x \log p_t(x)\right]\mathrm{d}t + g(t)\mathrm{d}\bar{w}$$

$x(T)$      $x(0)$

— Reverse stochastic process

$p_T(x)$          $p_t(x)$          $p_0(x)$

# Reverse SDE



Prior $\longrightarrow$ Reverse SDE $\longrightarrow$ Data

$x(T) \longrightarrow \mathrm{d}x = \left[ f(x,t) - g^2(t)\nabla_x \log p_t(x) \right] \mathrm{d}t + g(t)\mathrm{d}\bar{w} \longrightarrow x(0)$

— Reverse stochastic process

$p_T(x) \longrightarrow p_t(x) \longrightarrow p_0(x)$

# Full Process



Data — Forward SDE — Prior — Reverse SDE — Data

$x(0)$ $\longrightarrow$ $\mathrm{d}x = f(x,t)\mathrm{d}t + g(t)\mathrm{d}w$ $\longrightarrow$ $x(T)$ $\longrightarrow$ $\mathrm{d}x = \left[f(x,t) - g^2(t)\nabla_x \log p_t(x)\right]\mathrm{d}t + g(t)\mathrm{d}\bar{w}$ $\longrightarrow$ $x(0)$

SDE — Probability Flow ODE

$p_0(x) \longrightarrow p_t(x) \longrightarrow p_T(x) \longrightarrow p_t(x) \longrightarrow p_0(x)$

# Full Process



Forward SDE (data → noise)

Reverse SDE (noise → data)

$$\mathrm{d}\mathbf{x} = \mathbf{f}(\mathbf{x}, t)\mathrm{d}t + g(t)\mathrm{d}\mathbf{w}$$

score function

$$\mathrm{d}\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - g^2(t)\boxed{\nabla_{\mathbf{x}} \log p_t(\mathbf{x})}\right]\mathrm{d}t + g(t)\mathrm{d}\bar{\mathbf{w}}$$