

<u>DeGeSim – Pileup Modeling using</u> <u>Conditional Denoising Probabilistic</u> <u>Models</u>

<u>14-06-24</u>

Stephen Jiggins – DESY Judith Katzy - DESY



HELMHOLTZ AI SUPPORTED PROJECT







Contents





Problem:

The Large Hadron Collider (LHC) smashes bunches of protons (~10¹¹) approximately every 25ns. At present this yields ~40 collisions per crossing, in the High Luminosity LHC this will increase to ~200 \rightarrow This is pileup ~ noise





Problem:

The Large Hadron Collider (LHC) smashes bunches of protons (~10¹¹) approximately every 25ns. At present this yields ~40 collisions per crossing, in the High Luminosity LHC this will increase to ~200 \rightarrow This is pileup ~ noise





Problem:

The Large Hadron Collider (LHC) smashes bunches of protons (~10¹¹) approximately every 25ns. At present this yields ~40 collisions per crossing, in the High Luminosity LHC this will increase to ~200 \rightarrow This is pileup ~ noise





Problem:

The Large Hadron Collider (LHC) smashes bunches of protons (~10¹¹) approximately every 25ns. At present this yields ~40 collisions per crossing, in the High Luminosity LHC this will increase to ~200 \rightarrow This is pileup ~ noise



DeGeSim





DeGeSim





Project Aim: Generative modeling of ATLAS calorimeter response due to pileup interactions using ML image synthesis techniques, but map from one image class to another



Transform from a image X_{I} sampled from a distribution $P_{I}(\bar{x}|\bar{\theta})$ to a target image X_{I} sampled from $P_{I}(\bar{x}|\bar{\theta})$:





Model Overview

Calorimeter Image Model



Latent Variable 3

Calorimeter Image Model



Latent Variable 3







Reverse Denoising Process





Reverse Denoising Process



 \in



Reverse Denoising Process

Latent DDPM Calo. Imaging









Partially Diffused DDPM Generation

Train a conditional DDPM and partially diffuse a sample to a *release time* \mathbf{t}_{r} . Then switch class label (0 \rightarrow 1) and denoise.

Partially Diffused DDPM Operation



 \rightarrow How do we determine the release time t_r ?





Neural Likelihood Ratio Estimation

For a neural network with configurable set $\Phi = \{w_i\}_{i=1}^{N}$ parameters and a given loss functional $\mathcal{L}(s)$ of the form:

$$\mathcal{L}(s) = -\int_{\mathbb{R}^n} d^n x \left(p(x|t,\theta_0) \cdot \ln(s(x|t)) + p(x|t,\theta_1) \cdot \ln(1-s(x|t)) \right)$$

The extrema ($\delta \mathcal{L}(sh)/\delta s = 0$) of this general loss yields:

$$\frac{1 - s(x|t)}{s(x|t)} = \frac{p(x|t, \theta_1)}{p(x|t, \theta_0)} = r(x|t, \theta_0, \theta_1)$$





Prototype





 \rightarrow In particle physics we identify these particles by their mass:

Top Mass $(m_t) \sim 175 \text{ GeV}$ W-Boson Mass $(m_W) \sim 81 \text{ GeV}$

^[1] JetNet: A Python package for accessing open datasets and benchmarking machine learning methods in high energy physics, https://joss.theoj.org/papers/10.21105/joss.05789

Prototype



^[1] JetNet: A Python package for accessing open datasets and benchmarking machine learning methods in high energy physics, https://joss.theoj.org/papers/10.21105/joss.05789

Prototype





<u>Key</u>

- → **Top Pane:** Gaussian generated (Gen.) sampling with forward diffused *real data*
- → **Middle Pane:** Time evolution of generative data $p_{\theta}(x_{t-1} | y, x_t)$ for time each 10th time step:
 - Black contours = equal probability White line = random event path
- → Bottom Pane: Real data vs diffusion generated data from a guassian prior (top pane).
 - → **Metrics:** Wasserstein distance in 1D projection between total real and total generated data













→ When the densities are equivalent we define the *release time* of each event/data point

$$r(x|t, \theta_0, \theta_1) = \frac{\mathcal{L}(\theta_1|x, t)}{\mathcal{L}(\theta_0|x, t)} = 1$$
27





→ When the densities are equivalent we define the *release time* of each event/data point

$$r(x|t,\theta_0,\theta_1) = \frac{\mathcal{L}(\theta_1|x,t)}{\mathcal{L}(\theta_0|x,t)} = 1$$
28





→ When the densities are equivalent we define the *release time* of each event/data point

$$r(x|t,\theta_0,\theta_1) = \frac{\mathcal{L}(\theta_1|x,t)}{\mathcal{L}(\theta_0|x,t)} = 1$$
29



→ When the densities are equivalent we define the *release time* of each event/data point

$$r(x|t,\theta_0,\theta_1) = \frac{\mathcal{L}(\theta_1|x,t)}{\mathcal{L}(\theta_0|x,t)} = 1$$
30



→ When the densities are equivalent we define the *release time* of each event/data point

$$r(x|t,\theta_0,\theta_1) = \frac{\mathcal{L}(\theta_1|x,t)}{\mathcal{L}(\theta_0|x,t)} = 1$$
31



→ When the densities are equivalent we define the *release time* of each event/data point

$$r(x|t,\theta_0,\theta_1) = \frac{\mathcal{L}(\theta_1|x,t)}{\mathcal{L}(\theta_0|x,t)} = 1$$
32





DeGeSim – JetNet preliminary





Calorimeter Image Model





Calorimeter Image Model





Summary

→ **Goal:** Emulate calorimeter *images* of pileup seeded from MC based simulations

→ Status: Partial diffusion using self-conditioning and classifier free guidance to preserve class-specific features prior to image denoising → Image-to-image (I2I) translation using partial information loss

→ **Problems:** Null value noise injection into calorimeter images:

 \rightarrow CNNs struggle with noisy images and null noise injection

 \rightarrow Point clouds introduce cell \leftrightarrow point assignment ambiguity

→ **Further work** will include:

- \rightarrow Consistency models for faster inferencing
- → **Cross-domain** conditioning to enhance domain translation duing diffusion

Project Aim:

Helmholtz AI funded project for simulating detector simulations in proton-proton colliders using deep generative ML models as joint project between ATLAS + CMS + Juelich + TRIUMF

HELMHOLTZAI ARTIFICIAL INTELLIGENCE

ATLAS (DESY)



CMS (DESY)



Juelich

TRIUMF





Backup



Detector Image - Point Cloud



 \rightarrow **Point cloud** is actually the only way forward due to *simple geometry* used in the ATLAS simulation chain due to the irregular shape of the ECAL/HCAL

→ Artifacts induced due to projection will be problematic therefore simplest approach to solve the **data volume & geometry issues**





Problem:

Pileup from multiple proton interactions currently miss-model data. This trend worsens with the number of interactions per bunch crossing leading to an explosive issue moving towards the HL-LHC:



Self-Conditioning

Self-Conditioned DDPM



Neural Network Architecture











DDPMs - Recipe

\rightarrow Noise Scheduler:

Cosine β -scheduling noise injection:

$$\beta_t = 1 - \frac{\overline{\alpha_t}}{\overline{\alpha_{t-1}}} \qquad \overline{\alpha_t} = f(t)/f(0) \qquad f(t) = \cos\left(\frac{(t/T+s)\pi}{2.(1+s)}\right)$$

- \rightarrow Denoising Score Matching
 - $L(\theta) = \mathbb{E}_{t,x_0,\epsilon}[\left|\epsilon \epsilon_{\theta}(x_t(x_0,\epsilon),t)^2\right|]$
- \rightarrow Density Ratio Estimating Classifer

 $L_{DRE} = \sum \left[p_c(t, x_0) \log \left(\widehat{p}_c(t, x_0) \right) \right]$





\rightarrow Noise Scheduler:

Cosine β *-scheduling noise injection:*

$$\beta_t = 1 - \frac{\overline{\alpha_t}}{\overline{\alpha_{t-1}}} \qquad \overline{\alpha_t} = f(t) / f(0) \qquad f(t) = \cos\left(\frac{(t/T+s)\pi}{2.(1+s)}\right)$$

 \rightarrow Denoising Score Matching

$$L(\theta) = \mathbb{E}_{t,x_0,\epsilon}[\left|\epsilon - \epsilon_{\theta}(x_t(x_0,\epsilon),t)^2\right|]$$

 \rightarrow Density Ratio Estimating Classifer

 $L_{DRE} = \sum \left[p_c(t, x_0) \log \left(\widehat{p}_c(t, x_0) \right) \right]$



 \rightarrow The *origin* conditional probability is tractable:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t,\mathbf{x}_0) = q(\mathbf{x}_t|\mathbf{x}_{t-1},\mathbf{x}_0) rac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)}$$

$$\rightarrow$$
 Network needs to learn the reverse process:
 $p_{ heta}(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_{ heta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$

 \rightarrow Lower variational bound objective:

$$L_{ ext{VLB}} = \mathbb{E}_{q(\mathbf{x}_{0:T})} \Big[\log rac{q(\mathbf{x}_{1:T} | \mathbf{x}_0)}{p_{ heta}(\mathbf{x}_{0:T})} \Big]$$

 \rightarrow DSM needs to learn the term:

$$L_{\text{VLB}} = L_T + \underbrace{L_{T-1} + \cdots + L_0}_{T-1}$$

 $-\sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}}$

DDPMs - Recipe

\rightarrow Noise Scheduler:

Cosine β -scheduling noise injection:

$$\beta_t = 1 - \frac{\overline{\alpha_t}}{\overline{\alpha_{t-1}}} \qquad \overline{\alpha_t} = f(t)/f(0) \qquad f(t) = \cos\left(\frac{(t/T+s)\pi}{2.(1+s)}\right)$$

- \rightarrow Denoising Score Matching
 - $L(\theta) = \mathbb{E}_{t,x_0,\epsilon}[\left|\epsilon \epsilon_{\theta}(x_t(x_0,\epsilon),t)^2\right|]$
- \rightarrow Density Ratio Estimating Classifer

 $L_{DRE} = \sum \left[p_c(t, x_0) \log \left(\widehat{p}_c(t, x_0) \right) \right]$

- $\rightarrow \text{Conditional Embedding}$ $L(\theta) = \mathbb{E}_{t, x_0, \epsilon, y}[\left|\epsilon \epsilon_{\theta}(x_t(x_0, \epsilon), y, t)^2\right|]$
- \rightarrow Classifier Free Guidance

 $\overline{\epsilon_{\theta}}(x_t, \mathbf{y}, t) = (w+1) \epsilon_{\theta}(x_t, \mathbf{y}, t) - w. \epsilon_{\theta}(x_t, t)$



 \rightarrow Conditional dependency can be integrated:

$$q(\mathbf{x}_{t-1}|x_{t}, x_{0,y}) = q(\mathbf{x}_{t}|x_{t-1}, x_{0,y}) \cdot \frac{q(x_{t-1}|x_{0,y})}{q(x_{t}|x_{0,y})}$$
$$p_{\theta}(\mathbf{x}_{0:T}|y) = p_{\theta}(\mathbf{x}_{T}) \prod_{t=1}^{T} p_{\theta}(\mathbf{x}_{t-1}|x_{t}, y)$$

 \rightarrow Score estimation naturally extends with this conditionality:

$$\nabla_{\mathbf{x}_{t}}^{w}\log(p_{\theta}(\mathbf{x}_{t}|\mathbf{y})) = (1 - w)\nabla_{\mathbf{x}_{t}}\log(p_{\theta}(\mathbf{x}_{t})) + w\nabla_{\mathbf{x}_{t}}\log(p_{\theta}(\mathbf{x}_{t}|\mathbf{y})) - (1 - w)\nabla_{\mathbf{x}_{t}}\log(p_{\theta}(\mathbf{x}_{t}|\mathbf{y})) - (1 - w)\nabla_{\mathbf{x}_{t}}\log(p_{\theta}(\mathbf{x}_{$$

→ Guidance scalar (w) used to control conditional ↔ unconditionally interpolation of score/noise estimators:

$$\overline{\epsilon}(x_t,t,y) = w \epsilon_{\theta}(x_t,t,y) - (1-w) \epsilon_{\theta}(x_t,t) \blacktriangleleft$$

DDPMs - Recipe

\rightarrow Noise Scheduler:

Cosine β -scheduling noise injection:

$$\beta_t = 1 - \frac{\overline{\alpha_t}}{\overline{\alpha_{t-1}}} \qquad \overline{\alpha_t} = f(t) / f(0) \qquad f(t) = \cos\left(\frac{(t/T+s)\pi}{2.(1+s)}\right)$$

 \rightarrow Denoising Score Matching

 $L(\theta) = \mathbb{E}_{t,x_0,\epsilon}[\left|\epsilon - \epsilon_{\theta}(x_t(x_0,\epsilon),t)^2\right|]$

 \rightarrow Density Ratio Estimating Classifer

$$L_{DRE} = \sum \left[p_c(t, x_0) \log \left(\widehat{p_c}(t, x_0) \right) \right]$$

 $\rightarrow \text{Conditional Embedding}$ $L(\theta) = \mathbb{E}_{t, x_0, \epsilon, y}[\left|\epsilon - \epsilon_{\theta}(x_t(x_0, \epsilon), y, t)^2\right|]$

$$\overline{\epsilon_{\theta}}(x_t, \mathbf{y}, t) = (w+1) \epsilon_{\theta}(x_t, \mathbf{y}, t) - w. \epsilon_{\theta}(x_t, t)$$

 \rightarrow Self-Conditioning

 $L(\theta) = \mathbb{E}_{t, x_0, \widehat{x}_0, \epsilon, \mathbf{y}}[\left|\epsilon - \epsilon_{\theta}(x_t(x_0, \epsilon), \mathbf{y}, t, \widehat{x}_0)^2\right|]$



 \rightarrow Self-conditioned dependency can be integrated:

$$p_{\theta}(\mathbf{x}_{0:T}|y,\widetilde{x}_{0}) = p_{\theta}(\mathbf{x}_{T}) \prod_{t=1}^{T} p_{\theta}(\mathbf{x}_{t-1}|x_{t}, y, \widetilde{x}_{0})$$

 \rightarrow Provided that there is a stop gradient on the estimate of the original image $\widetilde{x_{o}}$:

 $\nabla_{x_t}^w \log(p_{\theta}(x_t|y,\widetilde{x}_0)) = (1 - w) \nabla_{x_t} \log(p_{\theta}(x_t)) + w \nabla_{x_t} \log(p_{\theta}(x_t|y,\widetilde{x}_0))$

 \rightarrow Scalar based guidance becomes:

$$\overline{\epsilon}(x_t, t, y) = w \,\epsilon_{\theta}(x_t, t, y, \widetilde{x}_0) - (1 - w) \,\epsilon_{\theta}(x_t, t)$$
52

Delphes Pile-up Simulation



Problem:

Pileup from multiple proton interactions currently miss-model data. This trend worsens with the number of interactions per bunch crossing leading to an explosive issue moving towards the HL-LHC:



Denoise Score Matching Neural Network

Latent Diffusion Residual Network



\rightarrow Residual Neural Network

Deep embedding blocks to prevent mode collapse. Therefore need residual connections to overcome lost gradients

\rightarrow Simple Sinusoidal Embedding

 $\begin{aligned} PE_{t,2i} = \sin(t/\alpha^{2i/d_{emb}}) \\ PE_{t,2i+1} = \cos(t/\alpha^{(2i+1)/d_{emb}}) \\ Natural normalisation of conditional variables \end{aligned}$

that scales dimensionally with embedding space

\rightarrow FFN Blocks

 \bigcirc

Simple dense feed-forward neural network blocks for dimensionality scaling and expressibility





Latent Diffusion Residual Network









DeGeSim – **JetNet** preliminary





