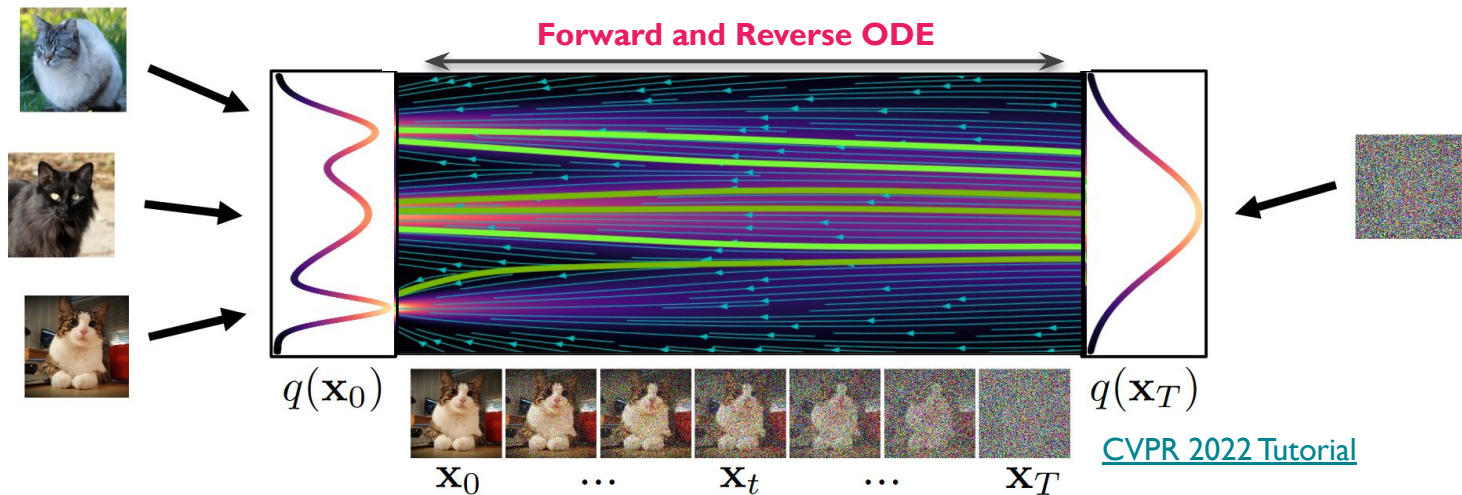


Generating Particle Cloud Jets with Denoising Diffusion

ML4Jets, Hamburg, 2023

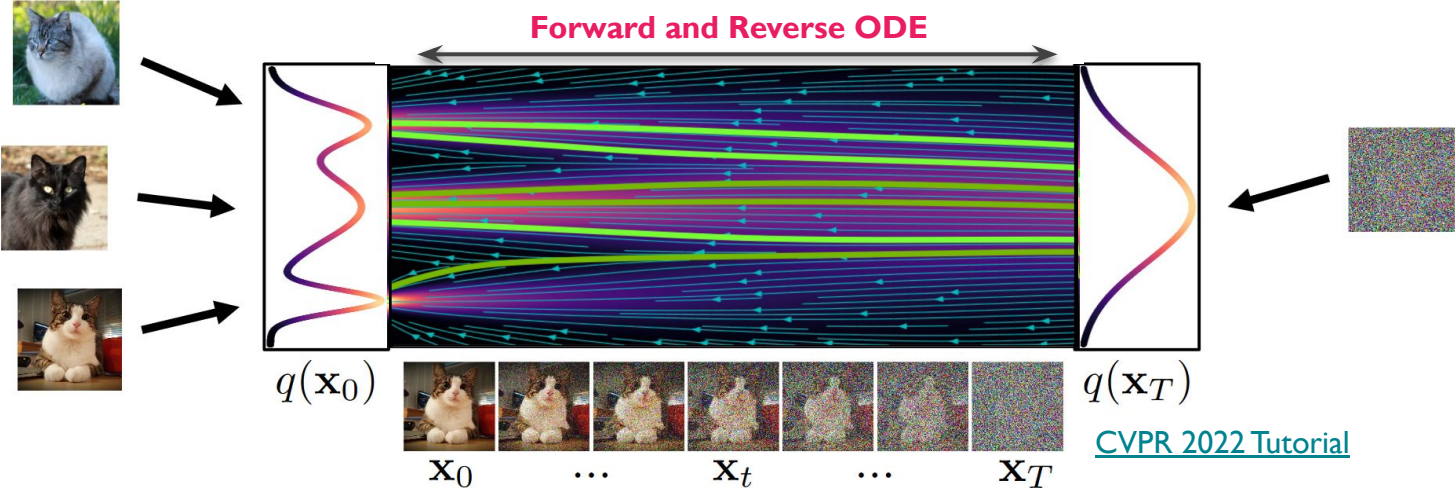
Matthew Leigh, Debajyoti Sengupta, Johnny Raine, Guillaume Quetant, Tobias Golling
UNIVERSITY OF GENEVA

Diffusion as differential equations



$$d\mathbf{x}_t = \left[f(\mathbf{x}_t, t) - \frac{1}{2}g(t)^2 \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}_{\text{score function}} \right] dt$$

Diffusion as differential equations



$$d\mathbf{x}_t = \left[f(\mathbf{x}_t, t) - \frac{1}{2}g(t)^2 \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)}_{\text{score function}} \right] dt$$

score function →

CAN LEARN THIS with a NN

Sample Generation

— — —

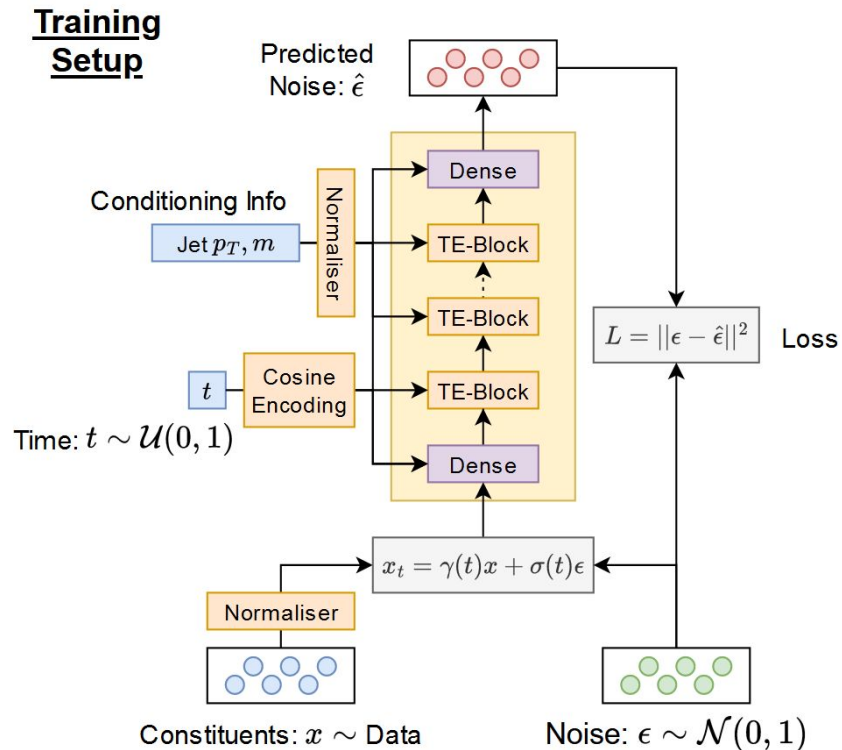
- Numerically **integrate** reverse process
 - Which process (SDE or ODE) and which integration method is **flexible**
- Each step **requires a forward pass** of the network
 - Generation needs **more computation** than GANs and Flows
 - Main **detriment** to using **diffusion** models

Always a trade-off between time and fidelity

- Using the [JetNet](#) dataset and metrics
- Large radius **point clouds** jets
 - Gluon, Quark, Top, W, Z
 - Up to 150 constituents
 - $(\Delta\eta, \Delta\phi, p_T)$

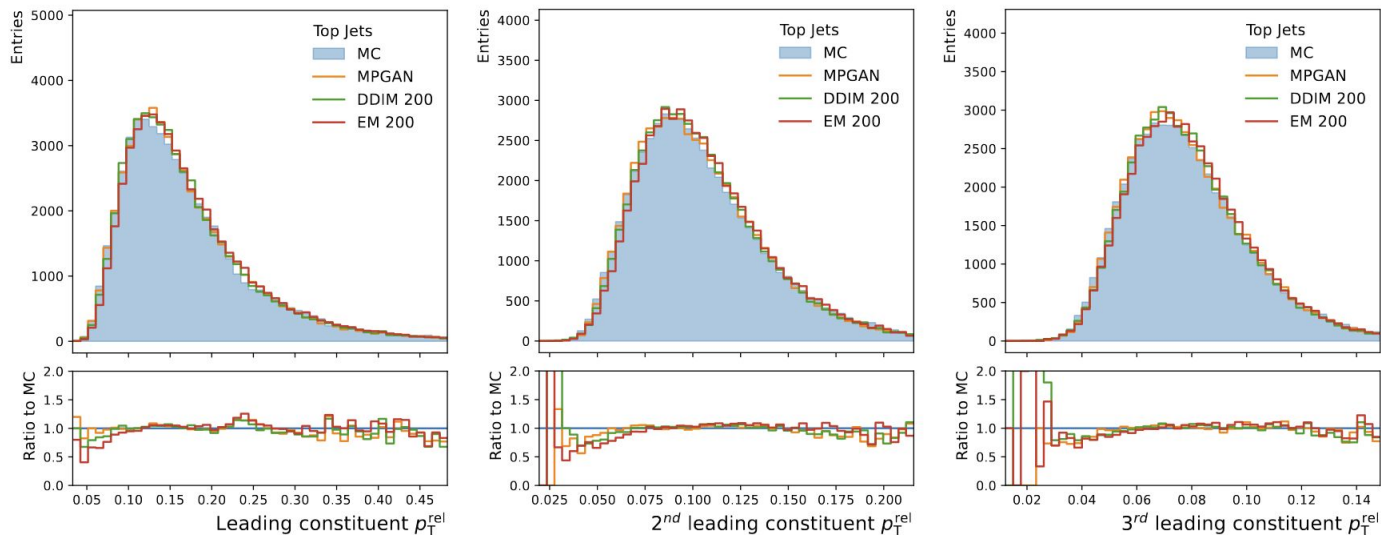
PC-Jedi Setup

- **First attempt** at particle cloud diffusion
- Based on a transformer
- Trained **separate models** for gluon and top
- Denoising objective



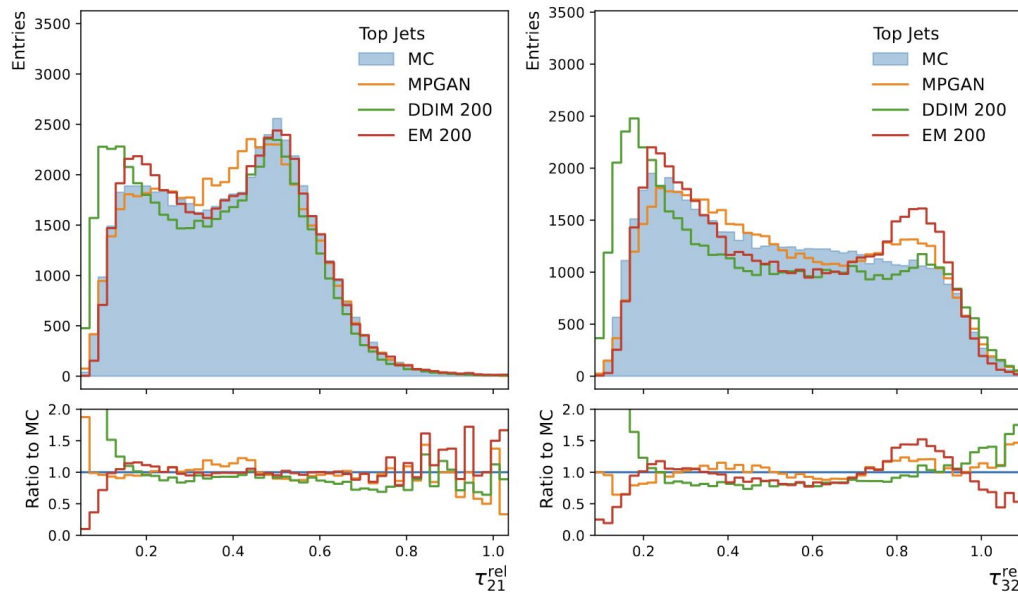
PC-Jedi Results

- Model was competitive to SOTA [MPGAN](#)



PC-Jedi Results

- Struggled recreating substructure variables for top jets
- And was slow



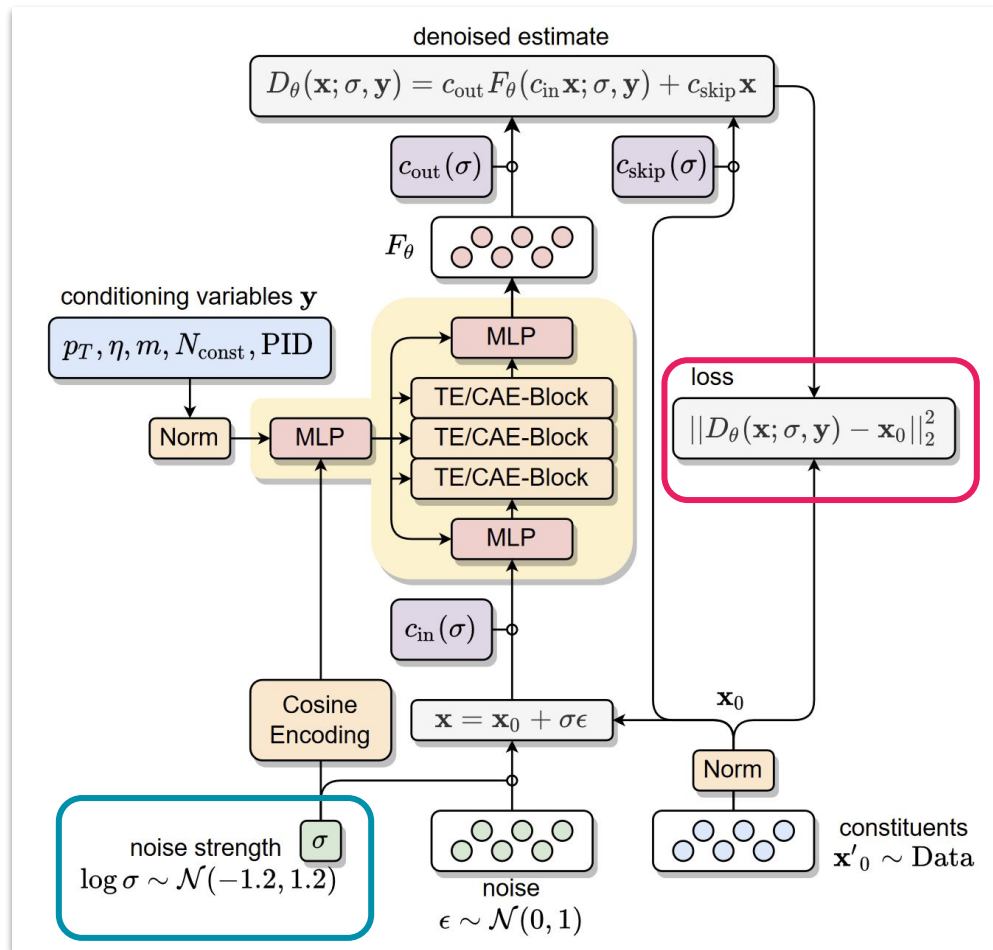
JeDi → Droid

PC-Droid

One conditional model for all jet types

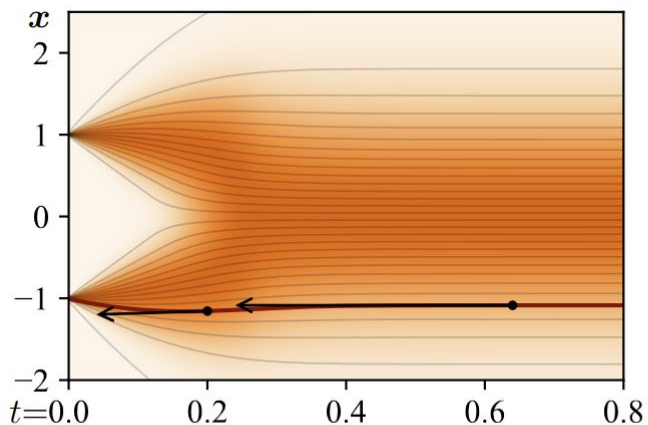
- Now **predicting denoised data** (not noise)
- **Smarter noise sampling** to focus on key areas of the trajectory during training.
- **Skip connections** for stability during training.

Compatible with SotA SDE/ODE Solvers.



Improvements with PC-Droid: ODE Trajectories

1. Change to [EDM](#) setup

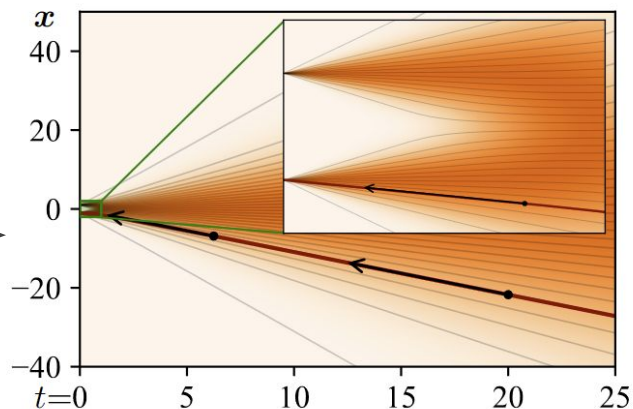


(a) Variance preserving ODE [49]

JeDi

Straighter trajectories

- Fewer truncation errors
- Easier solve



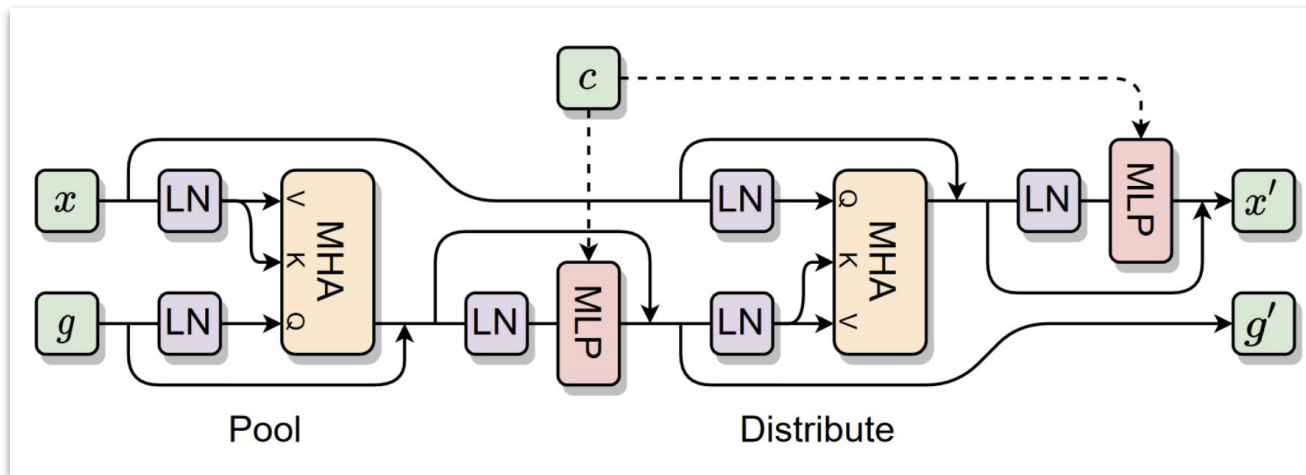
(c) DDIM [47] / Our ODE

Droid

Speed and scaling improvements: CAE Architecture

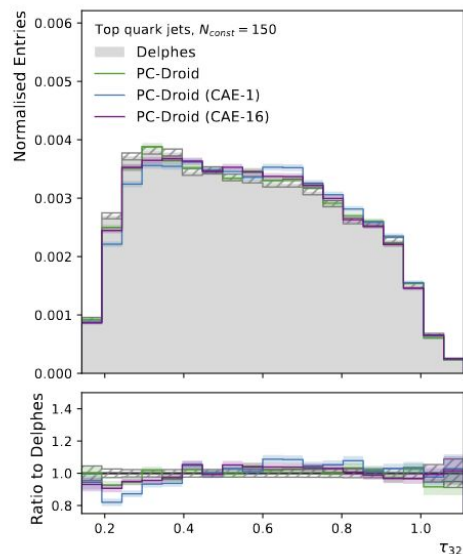
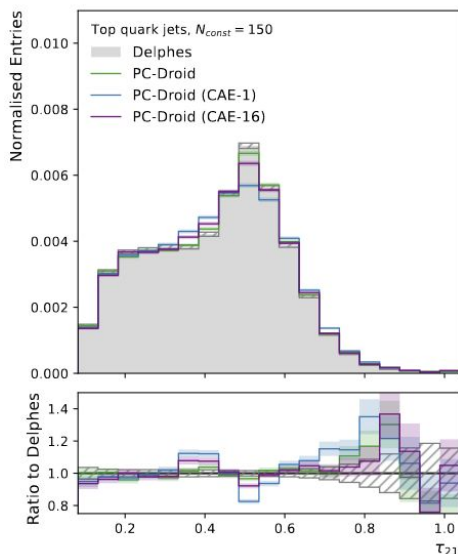
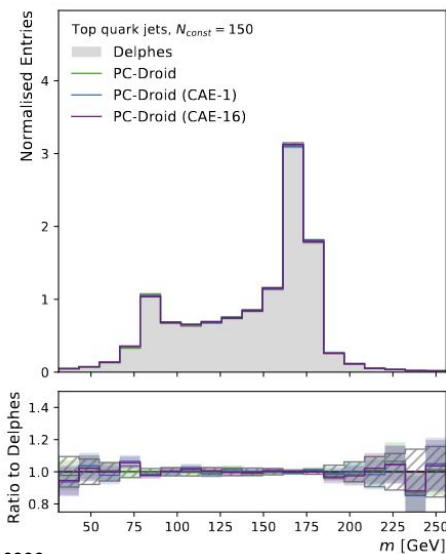
— — —
Increase number of constituents from 30 to 150

- Introduced new network type: Cross Attention Encoder
- Bipartite graph between point cloud and collection of global tokens
 - Number of global tokens is a hyperparameter (M)
 - $O(NM)$ computations compared to $O(N^2)$ of standard transformer



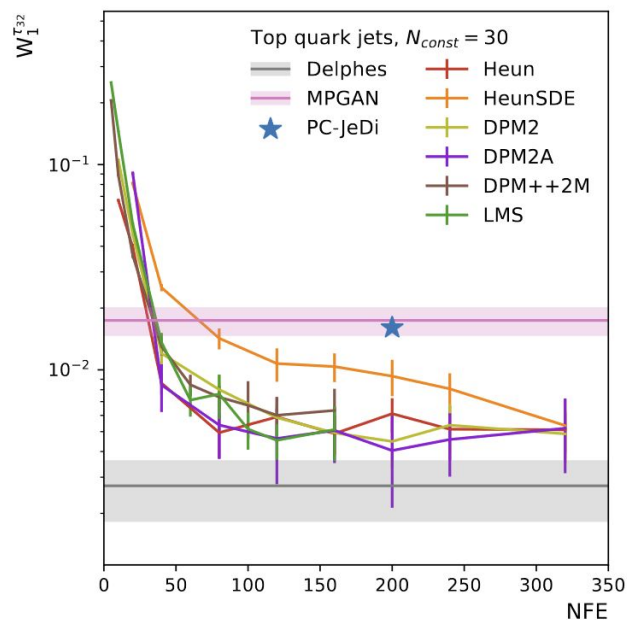
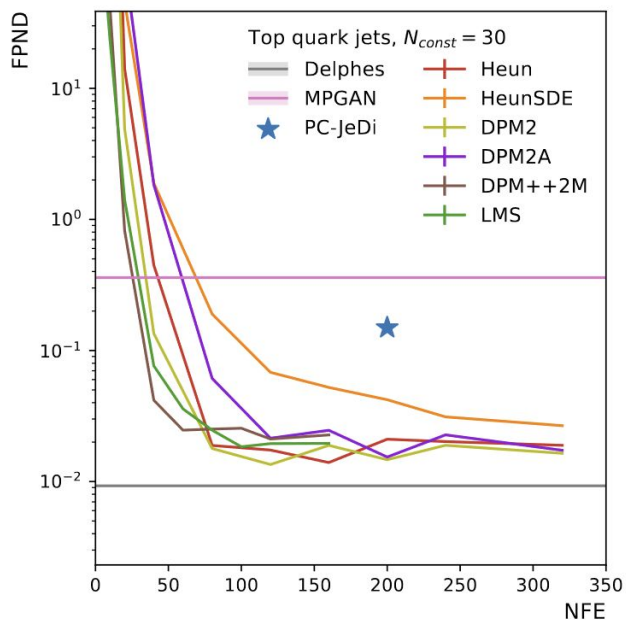
PC-Droid Results

- Great performance on 150 dataset
- New CAE network performs similarly with a big increase in generation speed



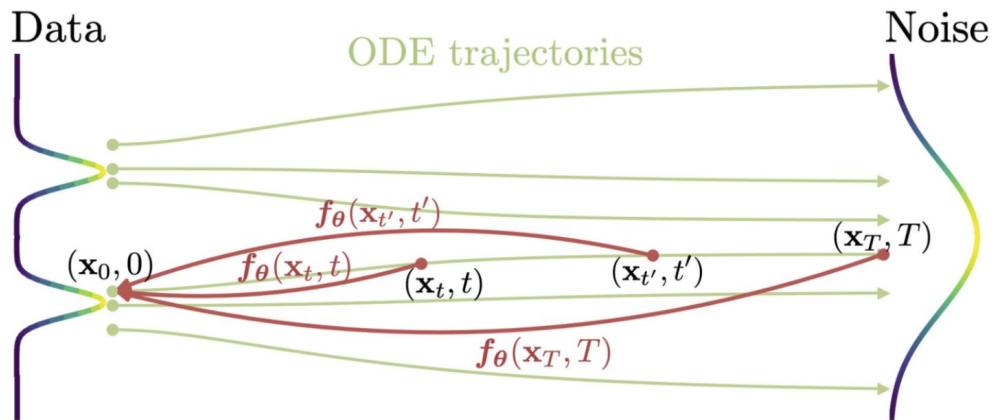
PC-Droid Results

- Massive improvements over our older diffusion model and MPGAN on 30 constituent dataset
- Significantly overtaking SOTA models



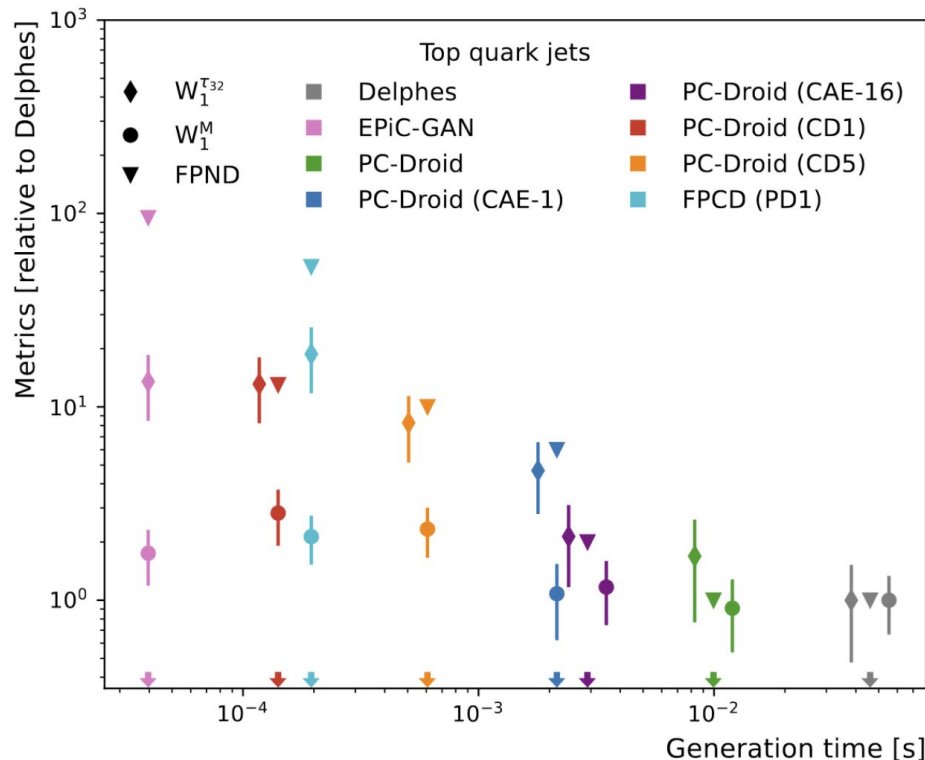
Further speed improvements: Consistency Distillation

- One of many diffusion **distillation** methods
- Use a **pretrained model** to train a **student model** to perform diffusion in less steps
- In some cases even allowing generation in **1 step**



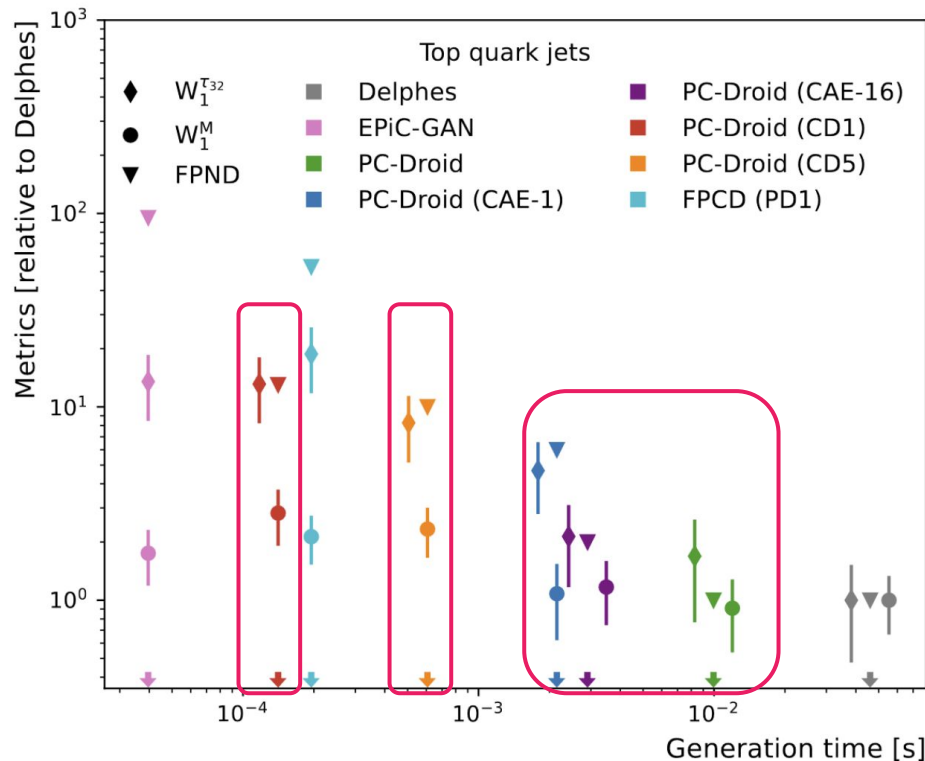
Time vs Fidelity Trade-Off

- Comparison with other generative models on 150 dataset
 - [FPCD](#)
 - [EPiC-GAN](#)



Time vs Fidelity Trade-Off

- Comparison with other generative models on 150 dataset
 - [FPCD](#)
 - [EPIC-GAN](#)
- PC-Droid performance on higher end is now close to ideal and 5 times faster
- Can sacrifice fidelity to get up to 100 times faster with distillation



Outlook

— — —

- Introduced diffusion models into HEP for point cloud generation with PC-JeDI
- Significantly improved quality with PC-Droid
- Looked at all models in terms of time-vs-quality trade off
- We are now looking at new ways to use such models beyond fast-sim (Next talk!)

Current work

— — —

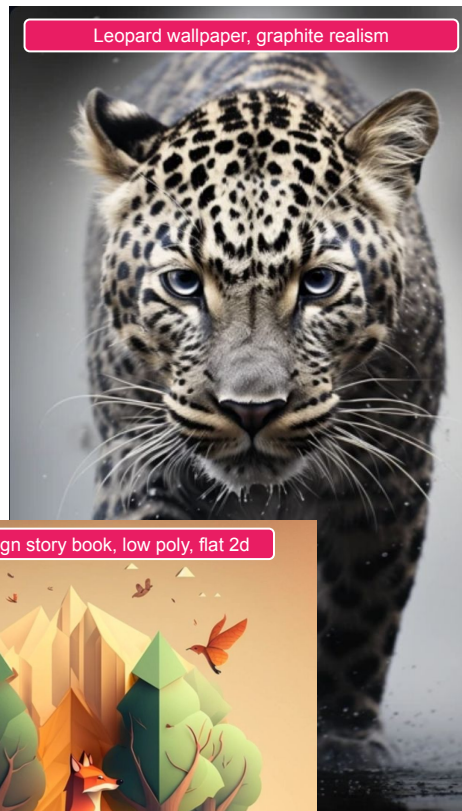
- [PC-JeDi: Diffusion for Particle Cloud Generation in High Energy Physics](#)
 - March 2023
 - Theory based on [Score-Based Generative Modeling through Stochastic Differential Equations](#)
- [PC-Droid: Faster diffusion and improved quality for particle cloud generation](#)
 - July 2023
 - Theory based on [Elucidating the Design Space of Diffusion-Based Generative Models](#) and [Consistency Models](#)
- [EPiC-Iy Fast Particle Cloud Generation with Flow-Matching and Diffusion](#)
 - September 2023

Thank You

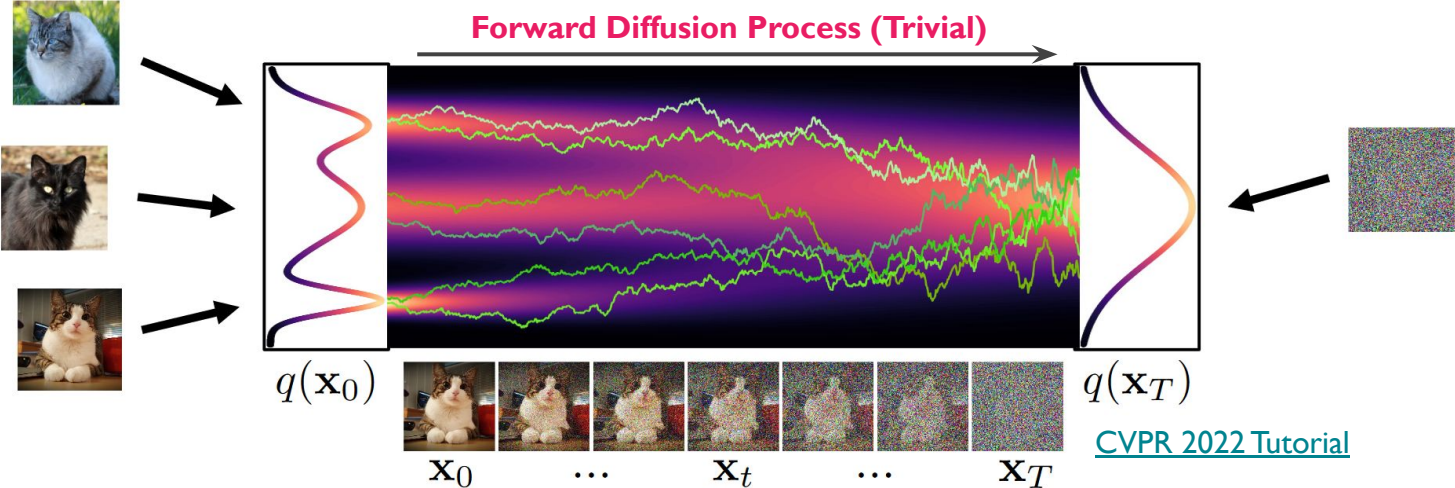
Backup

Proposal

- Use **Diffusion**
 - Generation = iterative **denoising steps**
- **Point clouds**
 - Replace the typical UNet with a message passing network
- Can use the **conditional generation**
 - Generate jets with desired high-level features
 - Momentum, mass, signal type
 - Required for Fast-Sim and template building



Diffusion as an SDE

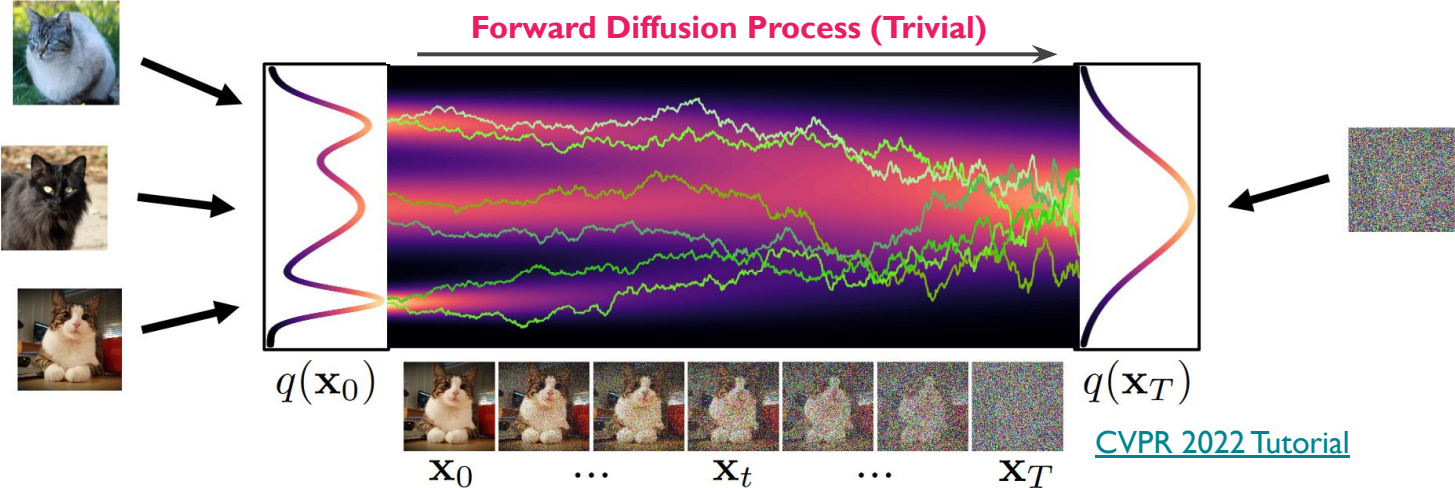


[CVPR 2022 Tutorial](#)

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

f, g are hyperparameters

Diffusion as an SDE



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

If f is affine

$$p(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \gamma(t)\mathbf{x}_0, \sigma(t)^2 \mathbf{I})$$

Denoising Learning Objective

Approximating the **score function** with a network is **impossible**

$$\min_{\theta} \underbrace{\mathbb{E}_{t \sim U(0,1)}}_{\text{time}} \underbrace{\mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)}}_{\substack{\text{diffused} \\ \text{data}}} \underbrace{\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)\|}_{\substack{\text{neural network} \\ \text{score of diffused data}}}^2$$

Denoising Learning Objective

— — —

Approximating the **score function** with a network is **impossible**

$$\min_{\theta} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} \left\| \mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \right\|^2$$

**marginal diffused densities are
intractable**

Denoising Learning Objective

Approximating the **score function** with a network is **impossible**

$$\min_{\theta} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} \|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)\|^2$$

Instead we look at the diffusion process of a **single sample \mathbf{x}_0**

$$\min_{\theta} \underbrace{\mathbb{E}_{t \sim U(0,1)}}_{\text{time}} \underbrace{\mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0)}}_{\text{data sample}} \underbrace{\mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_0)}}_{\text{diffused sample}} \underbrace{\|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0)\|^2}_{\text{score of diffused sample}}$$

Denoising Learning Objective

Approximating the **score function** with a network is **impossible**

$$\min_{\theta} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} \|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)\|^2$$

Instead we look at the diffusion process of a **single sample \mathbf{x}_0**

$$\min_{\theta} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_0)} \|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log \underline{p(\mathbf{x}_t | \mathbf{x}_0)}\|^2$$

This change is allowed because after expectations

$$\mathbf{s}_{\theta}(\mathbf{x}_t, t) \sim \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

Denoising Learning Objective

Approximating the **score function** with a network is **impossible**

$$\min_{\theta} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t)} \|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)\|^2$$

Instead we look at the diffusion process of a **single sample \mathbf{x}_0**

$$\min_{\theta} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_0)} \|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log \underline{p(\mathbf{x}_t | \mathbf{x}_0)}\|^2$$

This change is allowed because after expectations

$$\mathbf{s}_{\theta}(\mathbf{x}_t, t) \sim \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

This change is useful because the conditional density is tractable

$$p(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \gamma(t)\mathbf{x}_0, \sigma(t)^2 \mathbf{I})$$

Denoising Learning Objective

— — —
Old Learning Objective

$$\min_{\theta} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0)} \mathbb{E}_{\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{x}_0)} \|\mathbf{s}_{\theta}(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0)\|^2$$

Conditional Density:

$$p(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \gamma(t)\mathbf{x}_0, \sigma(t)^2 \mathbf{I})$$

Diffused Sample Score:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{x}_0) = \nabla_{\mathbf{x}_t} \frac{(\mathbf{x}_t - \gamma(t)\mathbf{x}_0)^2}{2\sigma(t)^2} = -\frac{\mathbf{x}_t - \gamma(t)\mathbf{x}_0}{\sigma(t)^2} = -\frac{\gamma(t)\mathbf{x}_0 + \sigma(t)\boldsymbol{\epsilon} - \gamma(t)\mathbf{x}_0}{\sigma(t)^2} = -\frac{\boldsymbol{\epsilon}}{\sigma(t)}$$

Neural Network Parameterisation:

$$\hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t) = -\sigma(t)s_{\theta}(\mathbf{x}_t, t)$$

New Learning objective:

$$\min_{\theta} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0)} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \frac{1}{\sigma(t)^2} \|\hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t) - \boldsymbol{\epsilon}\|^2$$

Denoising Learning Objective

$$\min_{\theta} \mathbb{E}_{t \sim U(0,1)} \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0)} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \frac{1}{\sigma(t)^2} \|\hat{\boldsymbol{\epsilon}}_{\theta}(\mathbf{x}_t, t) - \boldsymbol{\epsilon}\|^2$$

Sample time: $t \sim U[0,1]$

Sample data: $\mathbf{x}_0 \sim \{\text{Training set}\}$

Sample noise: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})_d$

Corrupt data: $\mathbf{x}_t = \gamma(t) * \mathbf{x}_0 + \sigma(t) * \boldsymbol{\epsilon}$

Get loss: $L = c(t) * [\text{NN}(\mathbf{x}_t, t) - \text{eps}]^2$

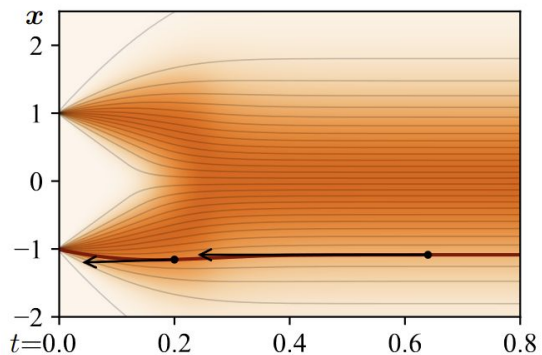
PC-Jedi: Paper 1

Improvements with PC-Droid

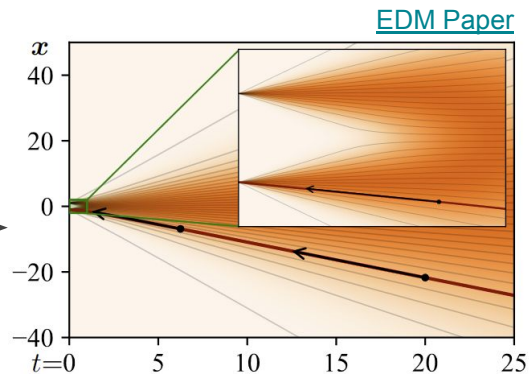
1. Change to EDM setup with preprocessing and sigma sampling

$$\text{SDE: } dx_t = \sqrt{2t} dw$$

$$\text{ODE: } dx_t = -t \nabla_x \log p(x; t) dt$$



(a) Variance preserving ODE [49]

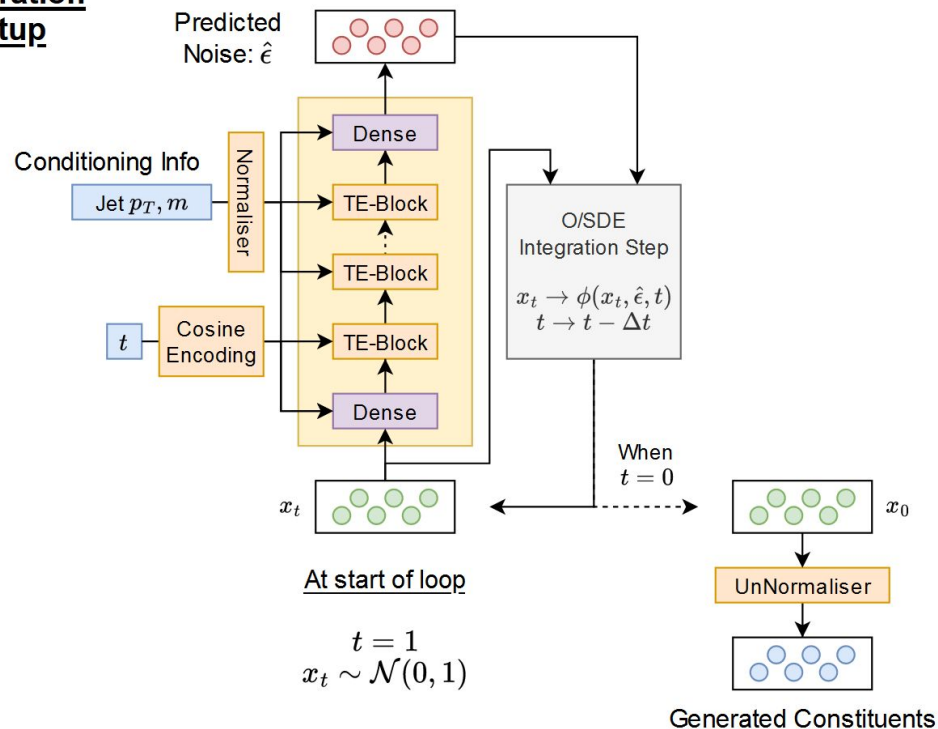


(c) DDIM [47] / Our ODE

PC-Jedi Setup

- For generation we tested:
 - Euler
 - Euler-Maruyama (SDE)
 - RK4
 - DDIM

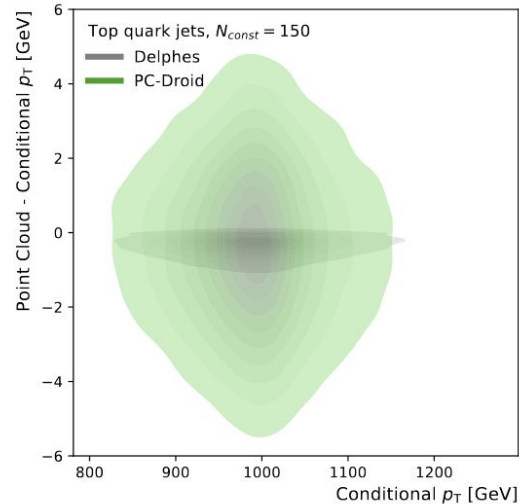
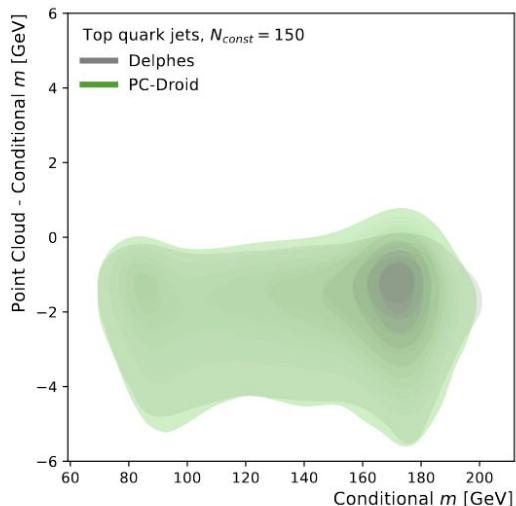
Generation Setup



Conditional Adherence

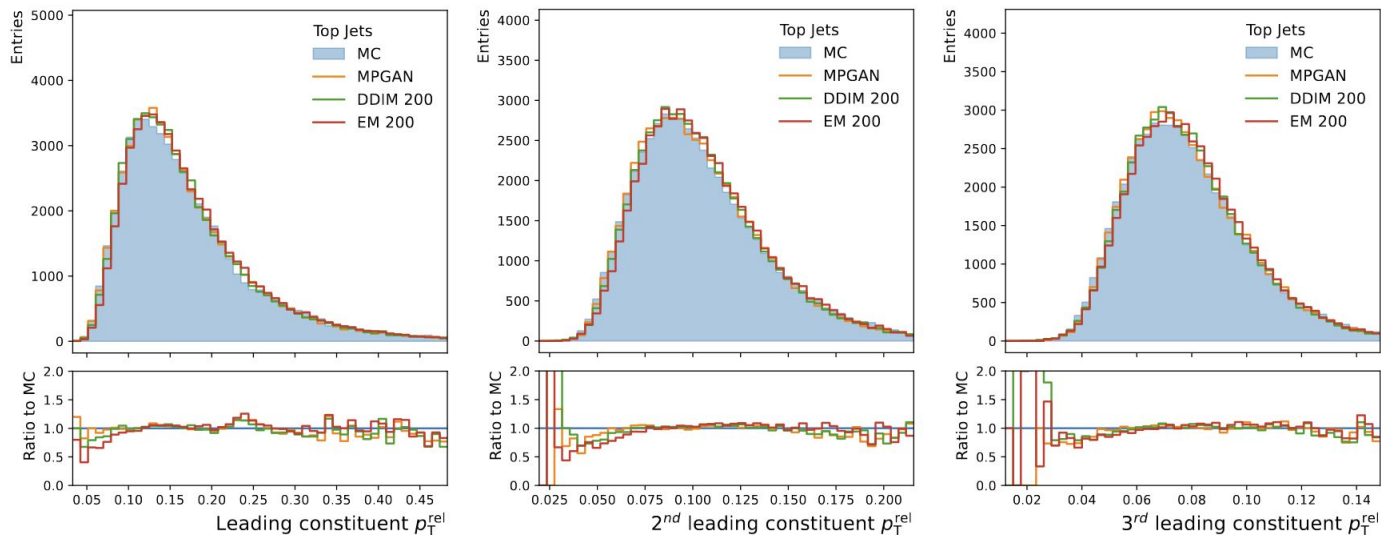
Is our conditional model actually obeying its conditions?

- Natural difference between conditional and point cloud variables in the data
- Slightly larger spread in \mathbf{p}_T
 - Majority within 0.3%



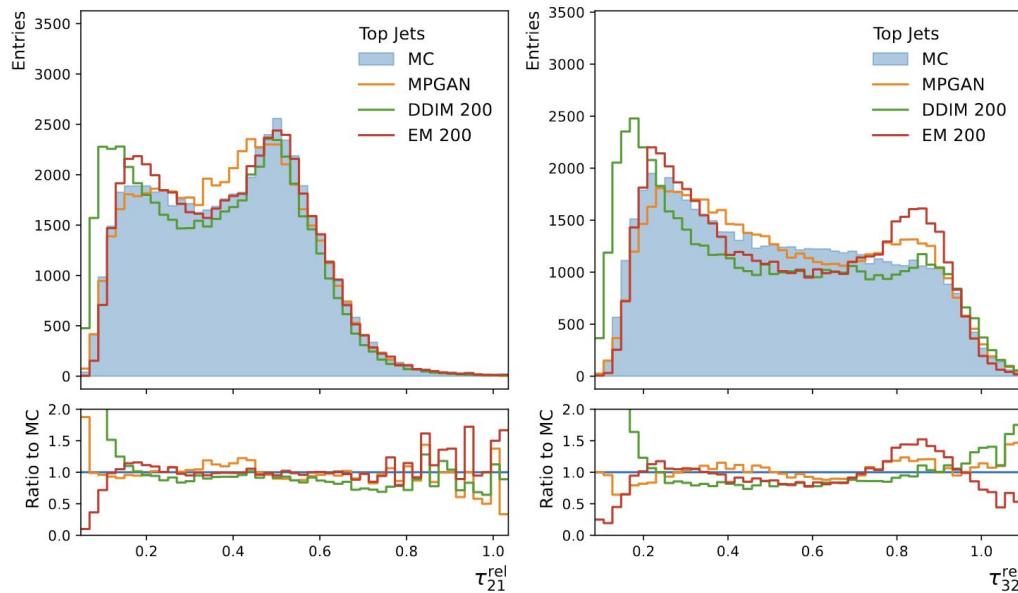
PC-Jedi Results

- Model was competitive to SOTA [MPGAN](#)



PC-Jedi Results

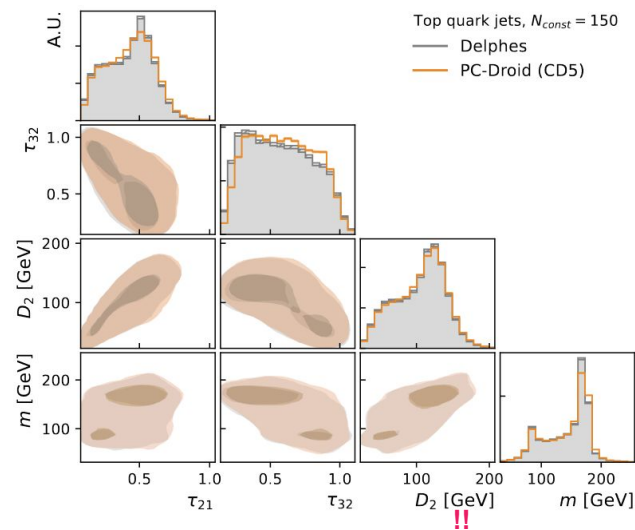
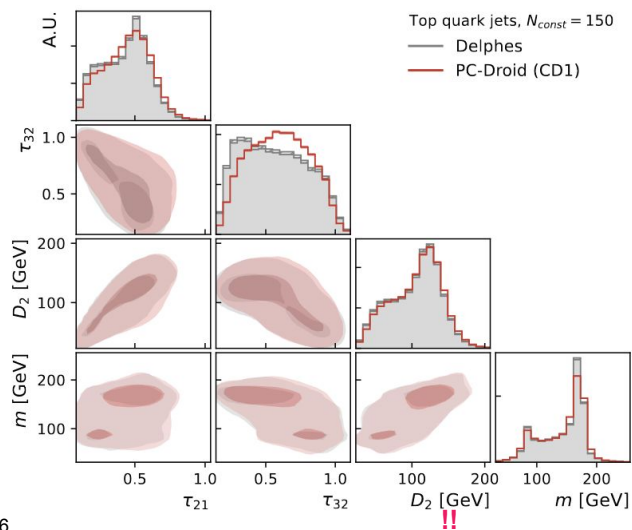
- Struggled recreating substructure variables for top jets



PC-Droid: Paper 2

CD Models Results

- Tested CD model with **1 and 5** step generation
- **Significantly faster** than base model (100x) but lower quality



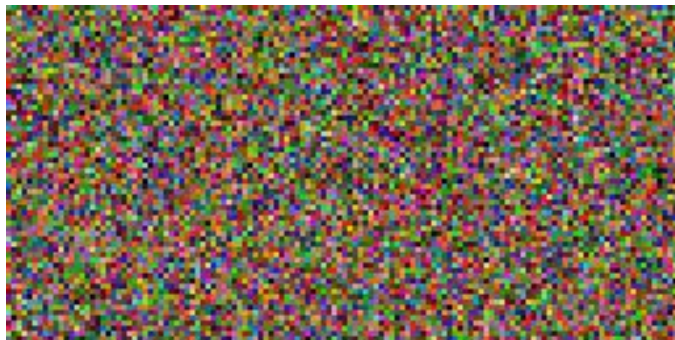
Sample Generation

— — —

- Numerically **integrate** reverse process
 - Which process (SDE or ODE) and which integration method is flexible
- Each step **requires a forward pass** of the network
 - Generation needs **more computation** than GANs and Flows
 - Main **detriment** to using **diffusion** models

Always a trade-off between time and fidelity

SDE with Euler-Maruyama



ODE with DDIM

