

*Department of physics, École Normale Supérieure, Paris*

GENERATING REALISTIC MATTER FIELDS WITH MACHINE LEARNING

---

## Machine Learning Project Report

---

Andrea Combette, Madeline Casas, Jean Goudot



**Cautionary note :** This paper is a report on numerical methods for the shallow water equations and gravity waves. It is not intended to be a complete and rigorous study of the subject. The reader is invited to refer to the references for further details. It has been made by a Master Student, with some background in physics and mathematics, but no prior knowledge of the subject. It is therefore not intended to be a reference for experts in the field.

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
I	Context . . . . .	2
1	Generating realistic matter field . . . . .	2
2	Machine Learning tackling the problem . . . . .	2
a	GANs . . . . .	2
b	DDPM . . . . .	2
c	VAE . . . . .	3
3	Datasets specifications . . . . .	3
a	Quijote Simulations . . . . .	3
b	Datasets dimensions and transformations . . . . .	3
<b>2</b>	<b>DDPM neural Network</b>	<b>4</b>
I	Theoretical background . . . . .	4
II	Architecture . . . . .	4
III	Testing . . . . .	5
1	Loss function . . . . .	5
2	Physics likelihood . . . . .	6
a	Histograms . . . . .	6
b	Power spectrum . . . . .	6
3	Third order image analysis . . . . .	7
4	Non linear contrasting . . . . .	7
a	order of magnitude for contrasting parameter . . . . .	7
b	Contrast trained model . . . . .	8
b.1	Loss function . . . . .	8
b.2	1st order distribution of pixels . . . . .	8
b.3	2nd order spectrum analysis . . . . .	9
b.4	3rd order clustering analysis . . . . .	9
IV	Conclusion . . . . .	9

# Chapter 1

## Introduction

This project aims at discovering new structures of neural networks and their applications to cosmology. In particular, we will focus on two types of neural networks: Generative Adversarial Networks (GANs) and Diffusion Deep Probabilistic Models (DDPMs). We will study their theoretical background, their architecture and their applications to cosmology. We will also compare their performances and discuss their advantages and disadvantages.

### I Context

#### 1 Generating realistic matter field

Generating density matter fields is crucial for several reasons. First, it allows us to simulate the distribution of matter in the universe, which is fundamental to cosmology. By understanding how matter is distributed, we can gain insights into the structure and evolution of the universe. Second, these simulations can help us test theories of cosmology and astrophysics. For instance, they can be used to predict the distribution of galaxies or the cosmic microwave background radiation. Lastly, generating realistic matter fields can also be a stepping stone towards more complex simulations, such as those involving dark matter or the formation of galaxies.

Previously generating matter fields was done using N-body simulations. However, these simulations are computationally expensive and time-consuming. Therefore, it is necessary to find new methods to generate matter fields. This is where machine learning comes in, because it allows us to generate matter fields in a much faster way.

#### 2 Machine Learning tackling the problem

To tackle this generation tasks, 3 types of neural networks could be used: Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs) and Diffusion Deep Probabilistic Models (DDPMs). These neural networks are all based on Gaussian approximations and are trained on a dataset of matter fields, and then used to generate new matter fields.

##### a GANs

A generative adversarial network (GAN) is a class of machine learning framework and a prominent framework for approaching generative AI. The concept was initially developed by Ian Goodfellow and his colleagues in June 2014. In a GAN, two neural networks contest with each other in the form of a zero-sum game, where one agent's gain is another agent's loss.

Given a training set, this technique learns to generate new data with the same statistics as the training set. For example, a GAN trained on photographs can generate new photographs that look at least superficially authentic to human observers, having many realistic characteristics. Though originally proposed as a form of generative model for unsupervised learning, GANs have also proved useful for semi-supervised learning, fully supervised learning, and reinforcement learning.

The core idea of a GAN is based on the "indirect" training through the discriminator, another neural network that can tell how "realistic" the input seems, which itself is also being updated dynamically. This means that the generator is not trained to minimize the distance to a specific image, but rather to fool the discriminator. This enables the model to learn in an unsupervised manner.

GANs are similar to mimicry in evolutionary biology, with an evolutionary arms race between both networks.

##### b DDPM

In machine learning, diffusion models, also known as diffusion probabilistic models or score-based generative models, are a class of generative models. The goal of diffusion models is to learn a diffusion process that generates the probability distribution of a given dataset. It mainly consists of three major components: the forward pro-



cess, the reverse process, and the sampling procedure. Three examples of generic diffusion modeling frameworks used in computer vision are denoising diffusion probabilistic models, noise conditioned score networks, and stochastic differential equations.

Diffusion models can be applied to a variety of tasks, including image denoising, inpainting, super-resolution, and image generation. For example, in image generation, a neural network is trained to denoise images with added gaussian noise by learning to remove the noise. After the training is complete, it can then be used for image generation by supplying an image composed of random noise for the network to denoise.

Diffusion models have been applied to generate many kinds of real-world data, the most famous of which are text-conditional image generators like DALL-E and Stable Diffusion. More examples are in a later section in the article.

### c VAE

In machine learning, a variational autoencoder (VAE) is an artificial neural network architecture introduced by Diederik P. Kingma and Max Welling. It is part of the families of probabilistic graphical models and variational Bayesian methods.

Variational autoencoders are often associated with the autoencoder model because of its architectural affinity, but with significant differences in the goal and mathematical formulation. Variational autoencoders are probabilistic generative models that require neural networks as only a part of their overall structure. The neural network components are typically referred to as the encoder and decoder for the first and second component respectively. The first neural network maps the input variable to a latent space that corresponds to the parameters of a variational distribution. In this way, the encoder can produce multiple different samples that all come from the same distribution. The decoder has the opposite function, which is to map from the latent space to the input space, in order to produce or generate data points. Both networks are typically trained together with the usage of the reparameterization trick, although the variance of the noise model can be learned separately.

Although this type of model was initially designed for unsupervised learning, its effectiveness has been proven for semi-supervised learning and supervised learning.

## 3 Datasets specifications

### a Quijote Simulations

Quijote provides not only thousands of simulations on different latin-hypercubes, but the a total number of 44,100 N-body simulations, with billion of halos, galaxies, voids and millions of summary statistics such as power spectra, bispectra... et, to train machine learning algorithms. In our case we will use the fiducial simulations. Those are simulations with a fiducial cosmology consistent with Planck. They only vary the initial random seed. We will only use  $N^3 = 256^3$  cubes with a redshift parameter  $z = 0$ , to get the filamentary structures of density fields (Gaussian  $\rightarrow$  filamentary structures). These cubes are avarage over 2 voxels in z-direction, to get a more isotropic dataset. The size of the cubes being  $L = 1Gpc/h$ . Then each 256 voxels width squares are sliced in 16 voxels width squares. This finally leads to the following nyquist wave vector

$$k_{nyquist} = \pi \frac{N}{L} \approx 0.8 Mpc/h^{-1}$$

### b Datasets dimensions and transformations

The raw datasets is composed of 50000 density fields of size  $64^3$  voxels. And each voxel as a value between  $[-1, 235.2]$ . In order to make the DDPM work we must normalize it between  $[-1, 1]$  to be in the range of a Gaussian distribution.

## Chapter 2

# DDPM neural Network

### I Theoretical background

A diffusion model has two main process a forward process and a reverse process. The forward process consists of a series of steps, where each datasets image is transformed into a more noisy image. Given these noisy images, it's not possible directly to learn how to clean them, the backward process consists in learning this cleaning step for each image. This give the following mathematical scheme :

Given a general distribution of the input image :  $P$  we want to reach a  $\mathcal{N}(0, I)$  distribution given  $n$  steps of diffusion. This is done by applying the following transformation to the input image  $x_0$  :

$$\begin{aligned} x_1 &= x_0\sqrt{1 - \beta_0} + \beta_0\epsilon \\ &\vdots \\ x_n &= x_{n-1}\sqrt{1 - \beta_{n-1}} + \beta_{n-1}\epsilon \end{aligned}$$

Where  $\epsilon \sim \mathcal{N}(0, I)$  and  $\beta_i$  is in the range  $[0.0001, 0.02]$ . Given this scheme it's possible to know exactly  $x_t$  with  $x_0$ . Defining :  $\alpha_i = 1 - \beta_i$  and  $\bar{\alpha}_i = \prod_{j=0}^i \alpha_j$  we have :

$$x_i = \sqrt{\bar{\alpha}_i}x_0 + \sqrt{1 - \bar{\alpha}_i}\epsilon$$

Then the learning process consists of finding the parameters of the backward process  $p(x_i, t)$ , which has the same fonctionnal form that the forward process transformation (William Feller, *Diffusion process in one dimension*). The transformation is the following :

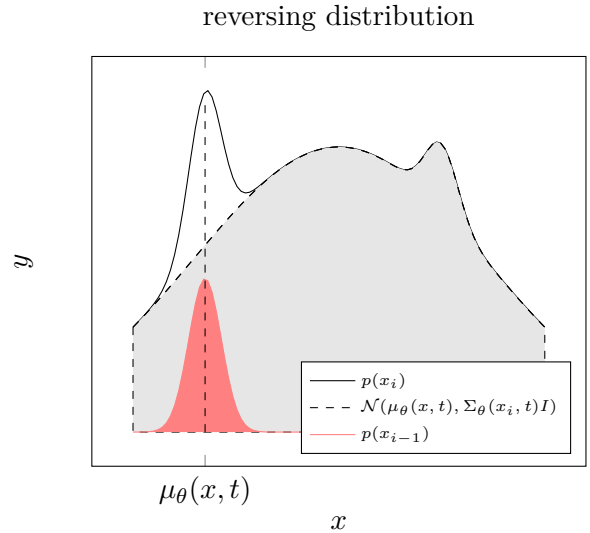
$$\bar{x}_{i-1} = \mu_\theta(x, t) + \sqrt{\Sigma_\theta(x_i, t)}\epsilon$$

. Then we can show that the parameter

$$\mu_\theta(x, t) = \frac{1}{\sqrt{\bar{\alpha}_t}}[x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}}\epsilon_\theta(x, t)]$$

and expressing  $\Sigma_\theta(x_i, t)$  as a combination of  $\alpha_i$  and  $\beta_i$  :

$$\Sigma_\theta(x_i, t) = I\frac{1 - \bar{\alpha}_{i-1}}{1 - \bar{\alpha}_i}\beta_i$$



we can find the parameters of the backward process just by knowing the noise level of the image, this is learned with a U-net neural network composed of a encoder and a decoder. By minimization of the given loss :

$$\mathcal{L} = \mathbb{E}\|\epsilon - \epsilon_\theta(x_i, t)\|^2 \quad (2.1)$$

Then we can generate a new image by applying the backward process to a noise image  $\epsilon \sim \mathcal{N}(0, I)$ . The network will then estimate the noise level and apply the desired transformation to the image. which results normally in a new image with a realistic distribution of density.

### II Architecture

The architecture used in this diffusion model is a simple encoder-decoder,

Table 2.1: layer dimensions

Layer	shape
Decoder Input	(1, 64, 64), 2 x (10, 64, 64)
Bottleneck	(160, 1, 1), 2 x (160, 1, 1)
Decoder Output	(20, 64,64), 2 x (10, 64,64)

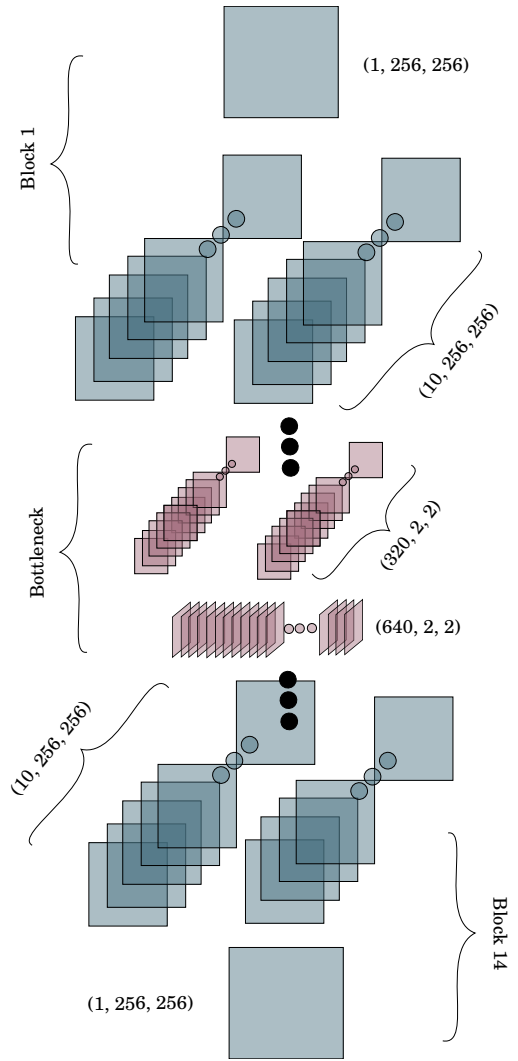
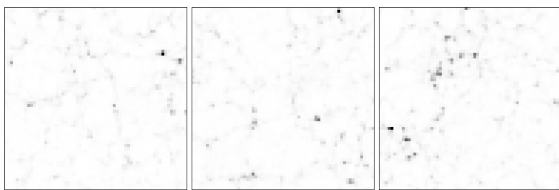
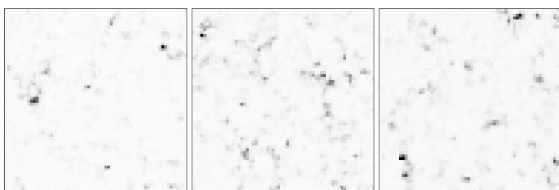


Figure 2.1: Unet convolutional structure

### III Testing



(a) Training image



(b) Generated image

Figure 2.2: Overview of the NN data

The training images present some specificities, we can clearly notice the presence of node-like structures, which results from a large aggregation of galaxies. And as the primordial universe converged to a filamentary-like structure, the training image presents clear and noticeable filaments, so called "cosmic-web". In the datasets, these types of structures are also recognizable, however one clear observation is the mean value of the field which is not correct at all, 2.2b, therefore we can just conclude as first sight that the distribution is range-relevant. We tried to understand why these distributions are shifted, the analysis of the forward process of the noise diffusion model is correct and leads as expected to a  $\mathcal{N}(0, I)$  distribution. The only explanation we found is a bias in the backward process.

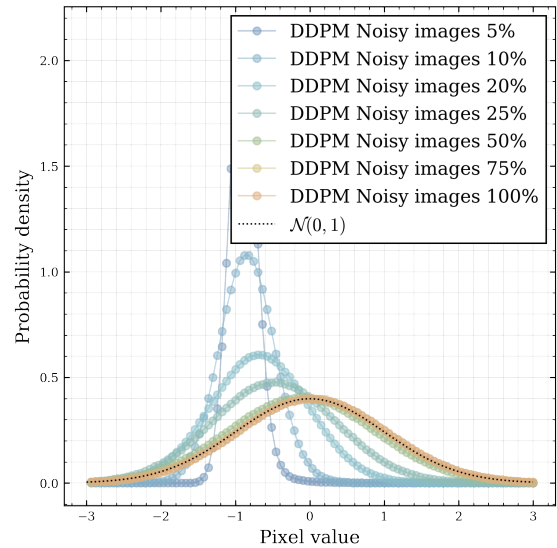


Figure 2.3: The forward process is correct; the noise level is progressively increasing, leading to the desired distribution

#### 1 Loss function

In section 2. If we have defined the loss function as a simple difference between the real noise value and the predicted one (2.1), we can then study the convergence of this loss function during the training process. For each epoch, this gives the following plot :

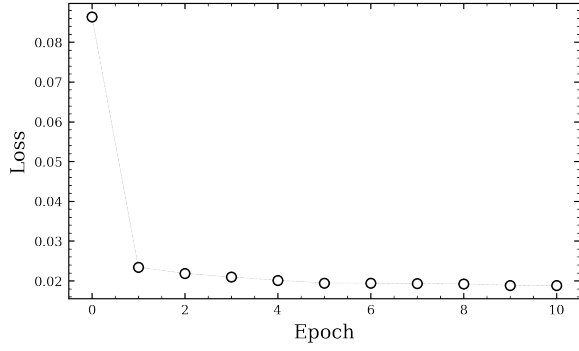


Figure 2.4: The model has been trained for 20 epochs, here only progressing epochs has been plot. This gives the loss of the Neural Network training process, it's exponentially decreasing as expected. Which means that the model is learning the noise level of the images.

## 2 Physics likelihood

### a Histograms

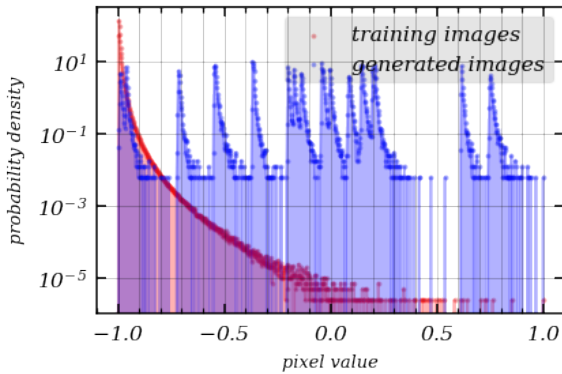


Figure 2.5: Histogram of intensities for both training and generated images, in red line the training images and in blue the generated images

With 2.6, we can lay the emphasis on the distribution similarity between generated and training images, however we can notice that the mean value of generated images is random in  $[-1, 1]$ , as we discussed in **2.III**. I've not find the cause of this issue . . . . To get rid of this issue, we can shift the generated images to the mean value of the training images, this gives the following histograms :

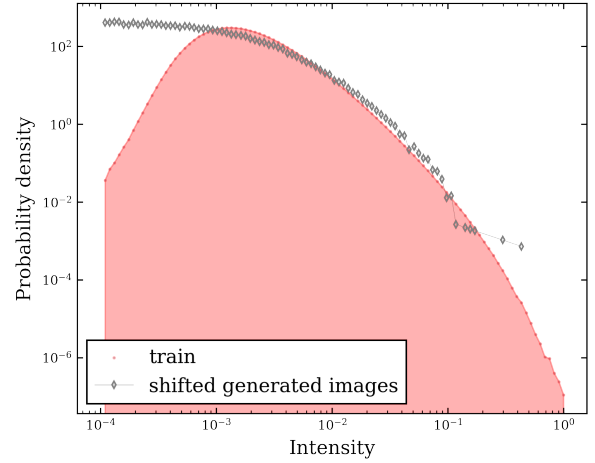


Figure 2.6: Shifted images to the training images mean value, the generated images present quite the same distribution for reasonable pixel value. However, the model tends to increase the void around filaments (excess of small pixel value), and to reduce the density of nodes (lack of high pixel value). The model is clearly exaggerating the filamentary structure making big voids around them and reducing the density of nodes.

### b Power spectrum

As we discussed before, the generated images are not perfect, however besides the mean value issue, the range of the distribution seems correct. To study more precisely the characteristics of the images, one common way in astronomy is to study the power spectrum of the images. We can suppose the power spectrum isotropic in space, meaned over each y-voxel. The power spectrum is defined as the fourier transform of  $I^2$  with  $I$  the pixel intensity. It describes the repartition in wave vector of the power, which is not related to with the mean value of the given image, which is a good point to study the generated images. Here we will study the normalized power spectrum, which is defined as : [cite]

$$P(k) = k^3 p(k)$$

with  $p(k)$  the power spectrum and  $k$  the wave vector.

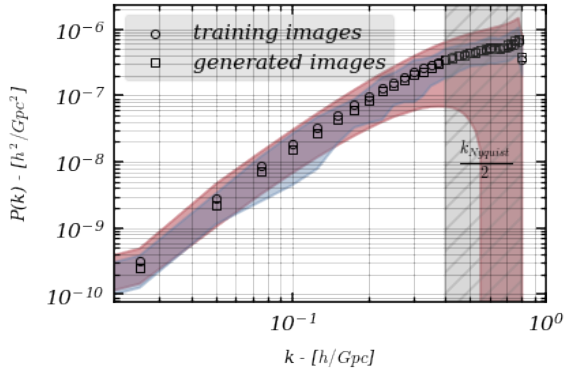


Figure 2.7: The blue fill is the standard deviation of the generated datasets, and the red one corresponds to the training image deviation. The grey shaded domain delimits the  $k_{nyquist}$  validity domain, due to the discretization of space. The psd leads to the same conclusion, the generated images and the training image have the same power spectrum, which strengthened the previous conclusion

```

1 def compute_isotropic_psd(img : Iterable)
2   -> tuple:
3     """
4     compute_isotropic_psd
5     compute_isotropic_psd compute
6     isotropic power distribution of a
7     given image
8
9     Parameters
10    -----
11    img : Iterable
12
13    _description_
14
15    Returns
16    -----
17    tuple
18    tuple of Iterable (wave vector k,
19    power vector)
20    """
21
22    delta_x = 250/64
23    k, p = welch(img, fs = 2 * np.pi /
24    delta_x, axis = 0, scaling= 'spectrum
25    ', nperseg = 64)
26
27    return k, np.mean(p, axis = 1)

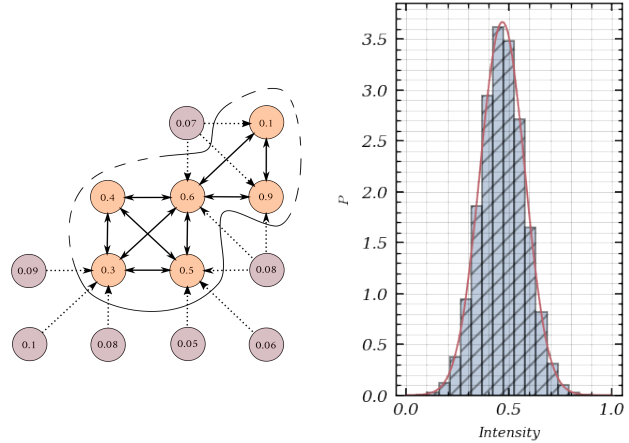
```

### 3 Third order image analysis

The goal of this analysis is to isolate complex structures such as filaments or nodes and to study their statistical characteristics. For isolating such structures, we used a graph approach, consider-

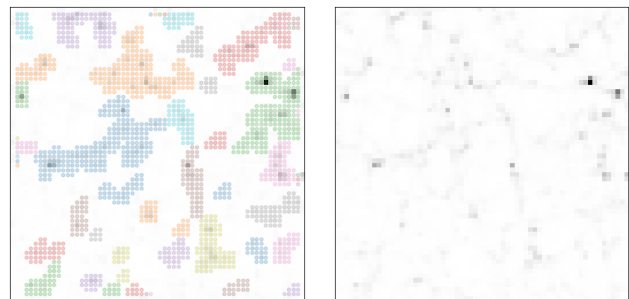
ing every pixel as a graph node connected to its nearest neighbors, this connection is removed if the pixel values is smaller than a typical threshold value fixed by the typical deviation of the gaussian noise distribution of the primary universe. To find this typical length we use *Fiducial* simulations at  $z = 127$ , which gives after renormalization process :

$$\sigma_0 = 0.107$$



(a) Graph construction of the density field (b) Density distribution for  $z = 127$

Figure 2.8



(a) Analyzed image with many filament cluster (b) Generated image

Figure 2.9

## 4 Non linear contrasting

### a order of magnitude for contrasting parameter

To improve the contrast of previous images, a solution has been proposed by [cite], it consists of the following transformation :

$$\sigma : x \mapsto \frac{x}{x + a}$$

Let's consider the number  $\mathcal{C}_{nl} = \frac{\sigma_0}{(\sqrt{a}-a)}$ . The interesting point about this typical number is the





following. If  $C_{nl} \ll 1$ , the contrasting function is too weak for the filament structure, whereas for  $C_{nl} \gg 1$  the contrasting function is too strong for the noise. Therefore, we must choose  $a$ , such that  $C_{nl} \sim 1$ . The idea behind this is the following if  $C_{nl} \ll 1$  the derivative of the contrasting function at the typical length  $\sigma_0$  is too big, it means that we do not isolate the filamentary structure. On the other hand if  $C_{nl} \gg 1$  the derivative of the contrasting function at the typical length  $\sigma_0$  is too small, it means that we isolate the noise.

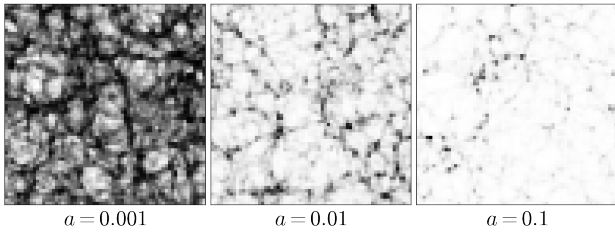
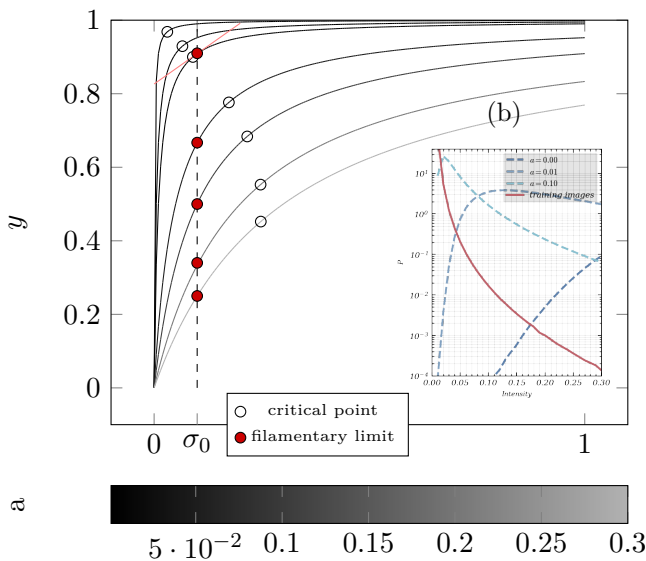


Figure 2.10: Transformations for multiple values of  $a$ , the extreme loglike behaviour at low  $a$  is totally deforming the physics matter field. which should drastically lower the learning efficiency of the NN. For relatively high  $a = 0.1$  we do not increase a lot the filaments contrast as opposed to medium  $a$  values, where the filamentary structure is clearly distinguishable, without deforming the physics field 2.16.



(a) Analyzed image with many filament cluster

Figure 2.12: In the Contrast transformations where the parameter  $a$  is chosen to optimize the contrasts given the previous normalization process. These functions allow us to shrink the image domain after a critical value (2.11a), which should be chosen accordingly to threshold used in the clustering process. Indeed we want to isolate the filamentary structures from the noisy background, which is not too small, otherwise we will isolate the noise, and not too big, otherwise we will not isolate the filaments. Given 2.11b we can observe the migration of intensity distribution for low  $a$ , it means that a lot of noisy pixels have been strengthened, which is not what we want. For relatively high  $a \sim 0.3$  we do not isolate the filaments anymore, because the critical intensity is much higher

**b Contrast trained model**

**b.1 Loss function**

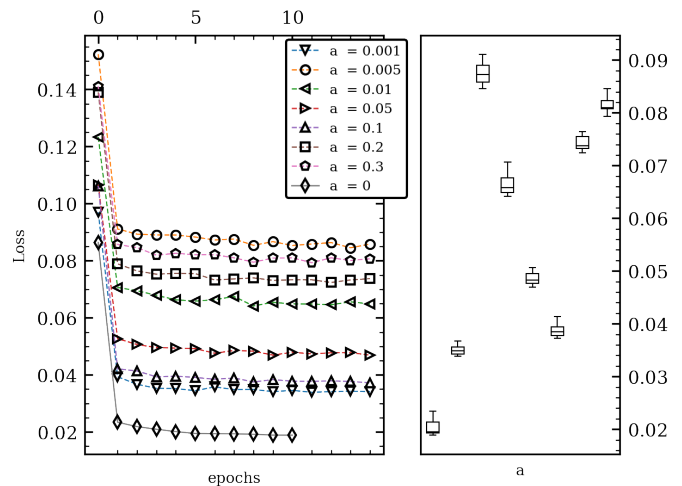


Figure 2.13: Here is plotted the different losses for the training process, for multiple values of the contrasting parameter  $a$ . We can observe that the lowest value of the loss function is obtained for  $a = 0$ , i.e. for no contrasting transformation, which can be opposed to the decision of [1]. However we will push forward our analysis, and try to see, if the contrasting function helps in creating a more realistic field

**b.2 1st order distribution of pixels** This first order analysis enables us to verify the pixel distribution range for all the models trained with different contrast parameters.

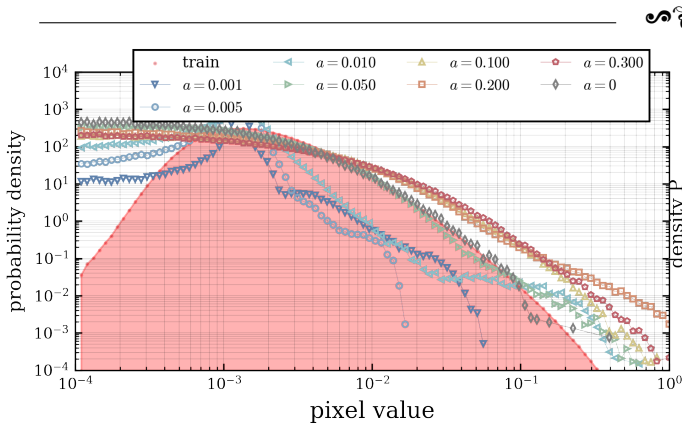


Figure 2.14: Note that as discussed in the fig.2.3 the mean value of the generated images is not the same as the training images, to study the distribution we renormalized the mean to fit the training images mean for each generated images. This leads to the former figure. We can note that, all the generated images have poor results for low pixel value, however the non contrasted generated images and the contrasted images with  $a = 0.05$  have a strong similarities with the original distribution for all pixel values  $\geq 1e^{-3}$ , which is a good point.

**b.3 2nd order spectrum analysis** The second order analysis, is a more precise analysis based on spatial frequencies redundancy.

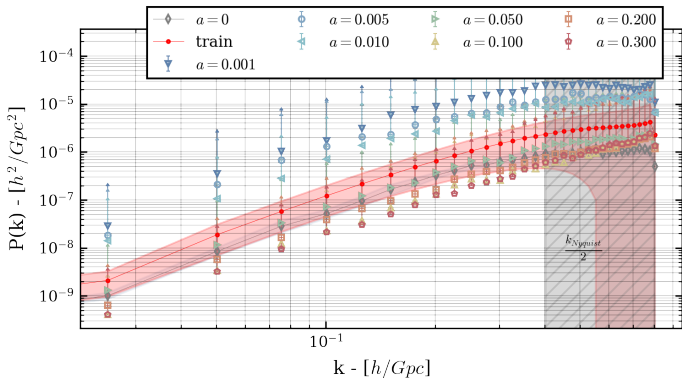


Figure 2.15: The spectrum analysis leads to the same conclusion, indeed it seems that only the non contrasted images and the contrasted images with  $a = 0.05$  have a range compatible power spectrum with training images. Note that only the high std is plotted for the generated images because the low std is too small to be plotted.

**b.4 3rd order clustering analysis** The third order analysis is the optimal way to evaluate the model, it allows us to have deeper insights on the generated images structures.

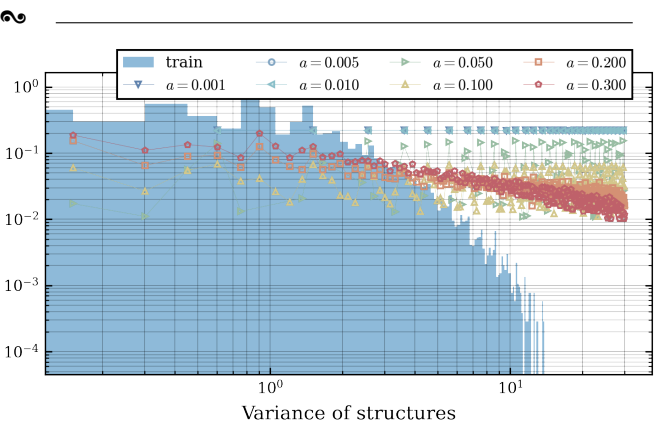


Figure 2.16: The clustering analysis on inverse transform images, dig out some interesting results, all the contrast trained models are unable to have a correct rendering of structures, indeed, it seems that contrast-based models exaggerate the filamentary structure to minimize their loss, which seems to be reasonable for contrasted images, but not for inverse transform images.

## IV Conclusion

In this study, we implemented a Diffusion Deep Probabilistic Model (DDPM) neural network for image denoising. The theoretical foundation of the model involves a forward diffusion process and a reverse process. The mathematical formulation was presented, emphasizing the transformation applied at each step to achieve a target distribution. The architecture of the neural network is a simple encoder-decoder structure. The generated images, however, exhibited some discrepancies in terms of mean value and distribution. Further analyses involved studying the power spectrum, histogram distributions, and third-order image analysis using clustering techniques. The power spectrum analysis showed that the generated images and training images had similar power spectra, strengthening the conclusion that the model learned the noise level effectively. Additionally, a contrasting transformation was introduced to improve contrast in the generated images. Different values of the contrasting parameter 'a' were explored, and it was observed that for certain values, the contrasted images exhibited similarities to the original distribution of training images. However, the clustering analysis on inverse-transformed images revealed that contrast-based models tended to exaggerate filamentary structures. In conclusion, the DDPM neural network, while showing promising results in learning matter density field distributions, still requires further improvements to achieve a more realistic rendering of structures.