

# 1 Implementation

In this section we show the reader all the details about the implementation starting from the employed software.

## 1.1 Tools

We quickly list The most important tools and libraries adopted in this project,

- ROS Melodic
- Numpy
- Matplotlib
- Pandas
- OpenCV
- PyTorch
- FastAI
- imgaug
- Blender

The framework was entirely developed on Ubuntu 18.10 with Python 3.6.

### 1.1.1 ROS Melodic

The Robot Operating System (ROS) [**ROS**] is a flexible framework for writing robot software. It is *de facto* the industry and research standard framework for robotics due to its simple yet effective interface that facilitates the task of creating a robust and complex robot behavior regardless of the platforms. ROS works by establishing a peer-to-peer connection where each *node* is to communicate between the others by exposing sockets endpoints, called *topics*, to stream data or send *messages*.

Each *node* can subscribe to a *topic* to receive or publish new messages. In our case, *Krock* exposes different topics on which we can subscribe in order to get real-time information about the state of the robot. Unfortunately, ROS does not natively support Python3, so we had to compile it by hand. Because it was a difficult and time-consuming operation, we decided to share the ready-to-go binaries as a docker image.

### 1.1.2 Numpy

Numpy is a fundamental package for any scientific use. It allows to express efficiently any matrix operation using its broadcasting functions. Numpy is used across the whole pipeline to manipulate matrices.

### 1.1.3 Matplotlib

To create almost all the plots in this report we used Matplotlib, is a Python 2D plotting library. It provides a similar functional interface to MATLAB and a deep ability to customize every region of the figure. All It is worth citing *seaborn* a data visualization library that we inglobate in our work-flow to create the heatmaps. It is based on Matplotlib and it provides an high-level interface.

### 1.1.4 Pandas

To process the data from the simulations we rely on Pandas, a Python library providing fast, flexible, and expressive data structures in a tabular form. Pandas is well suited for many different kinds of data such as handle tabular data with heterogeneously-typed columns, similar to SQL table or Excel spreadsheet, time series and matrices. We take advantages of the relationala data structure to perform custom manipulation on the rows by removing the outliears and computing the advancement.

Generally, pandas does not scale well and it is mostly used to handle small dataset while relegating big data to other frameworks such as Spark or Hadoop. We used Pandas to store the results from the simulator and inside a Thread Queue to parse each .csv file efficiently.

### 1.1.5 OpenCV

Open Source Computer Vision Library, OpenCV, is an open source computer vision library with a rich collection of highly optimized algorithms. It includes classic and state-of-the-art computer vision and machine learning methods applied in a wide array of tasks, such as object detection and face recognition. We adopt this library to handle image data, mostly to pre and post-process the heightmaps and the patches.

### 1.1.6 PyTorch

PyTorch is a Python open source deep learning framework. It allows Tensor computation (like NumPy) with strong GPU acceleration and Deep neural networks built on a tape-based auto grad system. Due to its Python-first philosophy, it is easy to use, expressive and predictable it is widely used amoung researches and enthusiast. Moreover, its main advantages over other mainstream frameworks such as TensorFlow [**tensorflow**] are a cleaner API structure, better debugging, code shareability and an enormous number of high-quality third-party packages.

### 1.1.7 FastAI

FastAI is library based on PyTorch that simplifies fast and accurate neural nets training using modern best practices. It provides a high-level API to train,

evaluate and test deep learning models on any type of dataset. We used it to train, test, and evaluate our models.

### 1.1.8 imgaug

Image augmentation (imgaug) is a python library to perform image augmenting operations on images. It provides a variety of methodologies, such as affine transformations, perspective transformations, contrast changes and Gaussian noise, to build sophisticated pipelines. It supports images, heatmaps, segmentation maps, masks, key points/landmarks, bounding boxes, polygons, and line strings. We used it to augment the heightmap, details are in section 1.4.7

### 1.1.9 Blender

Blender is the free and open source 3D creation suite. It supports the entirety of the 3D pipeline modeling, rigging, animation, simulation, rendering, compositing and motion tracking, even video editing and game creation. We used Blender to render some of the 3D terrain utilized to evaluate the trained model.

## 1.2 Data Gathering

### 1.2.1 Heightmap generation

A heightmap is a 2D array, an image, where each pixel's value represents the terrain height. The following figure shows an heightmap and the relative terrain.

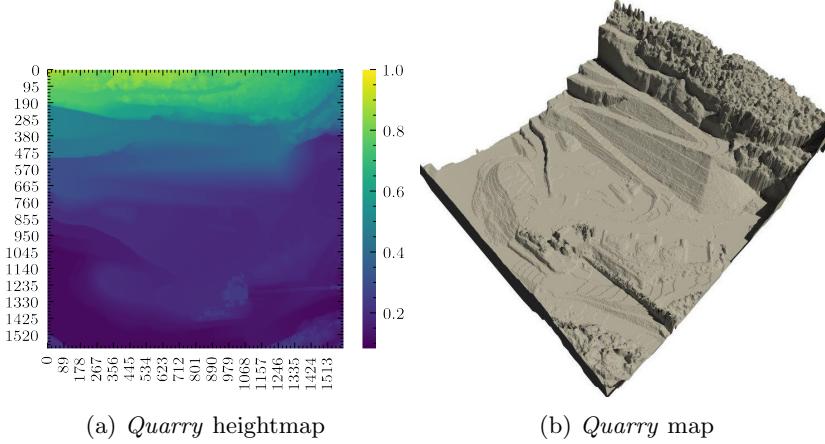


Figure 1: A heightmap

We generate fifteen terrain using random 2D simplex noise [**simplex**], a variant of Perlin noise [**perlin**], a widely used technique in terrain generation litterature.

3d rendered heightmap from Blender should visualize better the idea

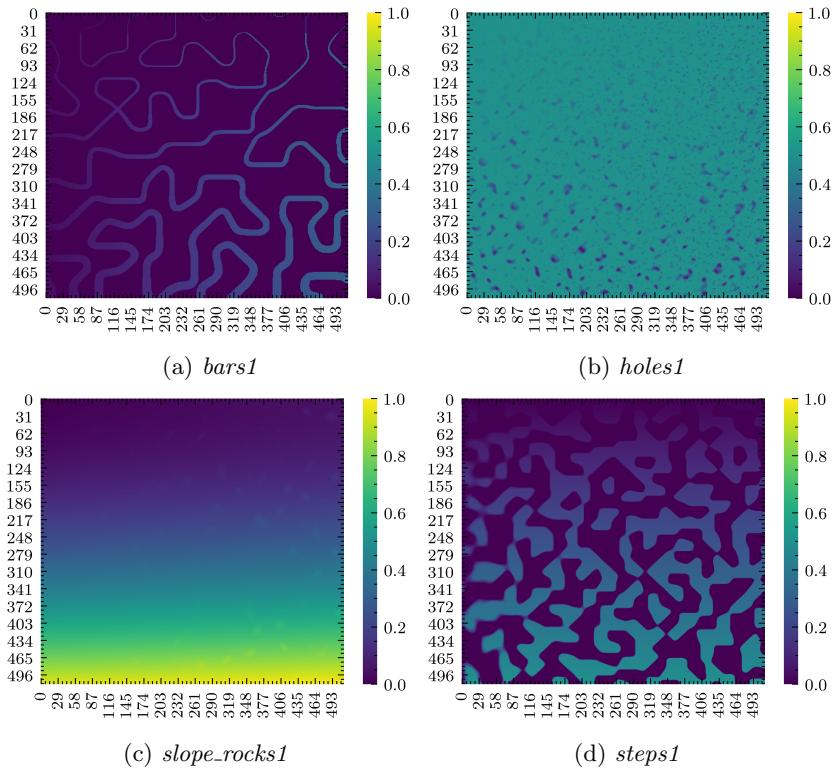


Figure 2: Some of the synthetic heightmaps

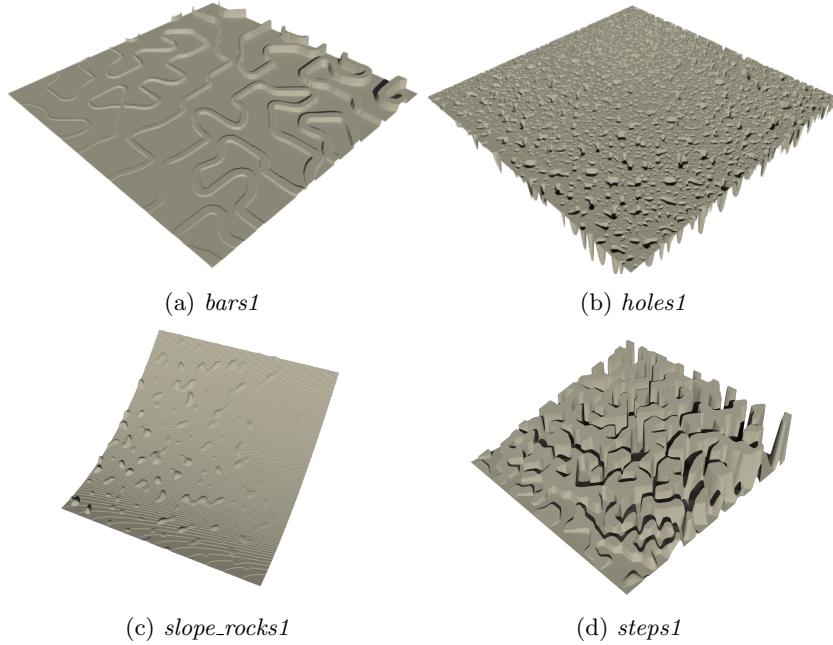


Figure 3: Some of the synthetic heightmaps

### 1.2.2 Simulator

We used Webots to move *Krock* on the generated terrain. The robot controlled was implemented by EPFL and handed to IDSIA . The controller implements a ROS' node to publish *Krock* status including its pose at a rate of 250hz. We decide to reduce it 50hz by using ROS build it throttle command. To load the map into the simulator we had to convert it to Webots's .wbt file. Unfortunately, the simulator lacks support for heightmap so we had to use a script to read the image and perform the conversion.

cite them?

cite?

To communicate with the simulator, Webots exposes a wide number of ROS services, similar to HTTP endpoints, with which we can communicate. The client can use the services to get the value of a field of a Webots' Node, for example, if we want to get the terrain height, we have to ask for the field value `height` from TERRAIN node. In addition, to call one service, we first have to get the correct type of message we wish to send and then we can call it. We decided to implement a little library called `webots2ros` to hide all the complexity needed to fetch a field value from a node.

We also implement one additional library called `Agent` to create reusable robot's interfaces independent from the simulator. The package supports callbacks that can be attached to each agent adding additional features. Finally, we used *Gym* [gym], a toolkit to develop and evaluate reinforcement learning algorithms, to define our environment. Due to the library's popularity, the code can easily be shared with other researches or we may directly experiment with

already made RL algorithm in the future without changing the code.

### 1.2.3 Simulation

To collect *Krock*'s interaction with the environment, we spawn it on the ground and let it move forward for 20 seconds. We repeat this process fifty times per each map. Unfortunately, spawning the robot is not a trivial task. In certain maps, for instance, *bars1*, we must avoid spawning on an obstacle otherwise the run will be ruined because *Krock* will be stuck form the start and we will introduce unwanted noise in the dataset. To solve the problem, we define a random spawn strategy used in most of the maps without big obstacles such as *slope\_rocks*, and a flat ground spawn strategy for the others. The random spawn just selects a random position and rotation for the robot. On the other hand, the flat ground strategy first selects suitable spawn positions by using a sliding window on the heightmap of size equal *Krock*'s footprint and check if the mean pixel value is lower than a small threshold. If so, we store the center coordinates of the patch. Intuitively, if a patch is flat the mean value will be close to zero.

We clustered those points with K-Means in  $k$  clusters where  $k$  is the number of spawning points we need, in our case fifty. By clustering, we guarantee to cover all region of the map. The following picture shows this strategy on *bars1*.

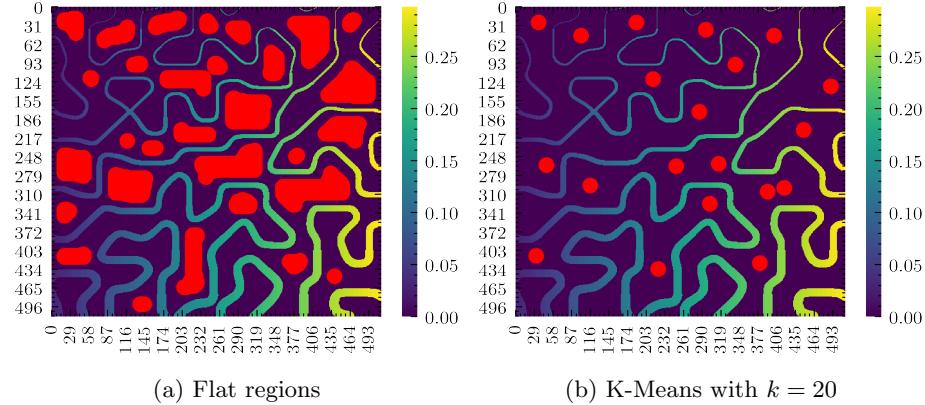


Figure 4: Examples of the spawning selection process (marked as red blobs) for the map *bars1*

## 1.3 Postprocessing

We now need to extract the patches for each pose  $p_t$  of *Krock* and compute the advancement for a given time window. To create the post-processing pipeline, we create an easy to use API called `pipeline` to define a cascade stream of function that is applied one after the other using a multi-thread queue to speed-up the process

### 1.3.1 Parse simulation data

First, we turn each `.bag` file to a pandas dataframe and cache them into `.csv` files. We used `rosbag_pandas`, an open source library we ported to python3, to perform the conversion. Then, we load the dataframes with the respective heightmaps and start the data cleaning process. We remove the rows corresponding to the first second of the simulation time to account for the robot spawning time. Then we eliminate all the entries where the *Krock* pose was near the edges of a map, we used a threshold of 22 pixels since we notice *Krock* getting stuck in the borders of the terrain during a simulation. After cleaning the data, we convert *Krock* quaternion rotation to Euler notation using the `tf` package from ROS. Then, we extract the sin and cos from the Euler orientation last component and store them in a column. Before caching again the resulting dataframes into `.csv` files, we convert the robot's position into heightmap's coordinates used later to crop the correct region of the map.

Finally, we load again the last stage and compute the advancement by projecting the pose's position,  $x$  and  $y$ , in the current line. Then, we select a time window according to the store rate, for if we select two seconds we need to multiply the rate by two, so  $50 * 2 = 100$  since *Krock* published with a  $50\text{hz}$  frequency. Once the time window is defined, we project  $x$  and  $y$  on the current line used the sin and cos values calculated before to get the advancement.

### 1.3.2 Extract patches

To create the patches, we first discover the maximum advancement for one second by running some simulations of *Krock* on flat ground and averaging the advancement. For our robot, the maximum speed is  $33\text{cm/s}$  Then, we multiply the maximum displacement by the number of seconds we are interested in. We can now crop the corresponding region in the heightmap by including the whole *Krock*'s footprint and the maximum advancement. The following figure visualizes the patch extraction process.

figure that shows *Krock* somewhere with the patch bounding boxed

. Lastly, we create a final dataframe containing the map coordinates, the advancement, and the patches paths for each simulation and store them to disk as `.csv` files.

The whole pipeline takes less than one hour to run the first time with 16 threads, and, once it is cached, less than fifteen minutes to extract the patches.

Once we extract the patches, we can always re-compute the advancement without re-running the whole pipeline. The next figure show proposed pipeline.



Figure 5: Postprocessing pipeline flow from top to bottom

**add arrows!**

## 1.4 Estimator

In this section we shows the adopted models' architecture.

### 1.4.1 Vanilla Model

### 1.4.2 ResNet

We decide to use a Residual Network, ResNet [he2015deep], variant. Residual networks are deep convolutional networks consisting of many stacked Residual Units : Intuitively, the residual unit allows the input of a layer to contribute to the next layer's input by being added to the current layer's output. Due to possible different features dimension, the input must go through and identify map to make the addition possible. This allows a stronger gradient flows and mitigates the degradation problem. A Residual Units is composed by a two  $3 \times 3$

*Convolution*, *Batchnorm* [**ioffe2015batch**] and a *Relu* blocks. Formally, it is defined as:

$$y = \mathcal{F}(\mathbf{x}, \{W_i\}) + h(\mathbf{x}) \quad (1)$$

Where,  $x$  and  $y$  are the input and output vector of the layers considered. The function  $\mathcal{F}(\mathbf{x}, \{W_i\})$  is the residual mapping to be learned and  $h$  is the identity mapping. The next figure visualises the equation.

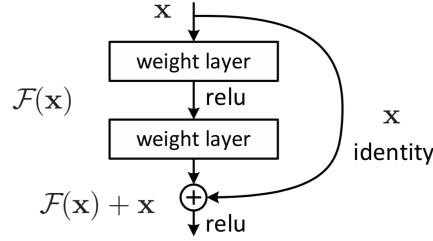


Figure 6: *Resnet* block [**he2015deep**]

When the input and output shapes mismatch, the *identity map* is applied to the input as a  $3 \times 3$  Convolution with a stride of 2 to mimic the pooling operator. A single block is composed by a  $3 \times 3$  *Convolution*, *Batchnorm* and a *ReLU* activation function.

#### 1.4.3 Preactivation

Following the recent work of He et al. [**he2015identity**] we adopt *pre-activation* in each block. *Pre-activation* works by just reverse the order of the operations in a block.

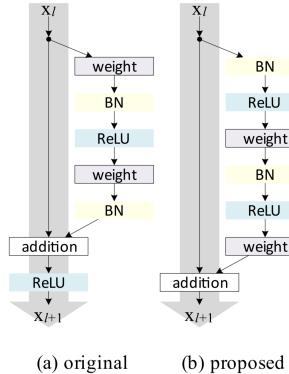


Figure 7: *Preactivation* [**he2015identity**]

#### 1.4.4 Squeeze and Excitation

Finally, we also used the *Squeeze and Excitation* (SE) module [hu2017squeeze]. It is a form of attention that weights the channel of each convolutional operation by learnable scaling factors. Formally, for a given transformation, e.g. Convolution, defined as  $\mathbf{F}_{tr} : \mathbf{X} \mapsto \mathbf{U}$ ,  $\mathbf{X} \in \mathbb{R}^{H' \times W' \times C'}$ ,  $\mathbf{U} \in \mathbb{R}^{H \times W \times C}$ , the SE module first squeeze the information by using average pooling,  $\mathbf{F}_{sq}$ , the it excites them using learnable weights,  $\mathbf{F}_{ex}$  and finally, adaptive recalibration is performed,  $\mathbf{F}_{scale}$ .

explain notation

The next figure visualises the SE module.

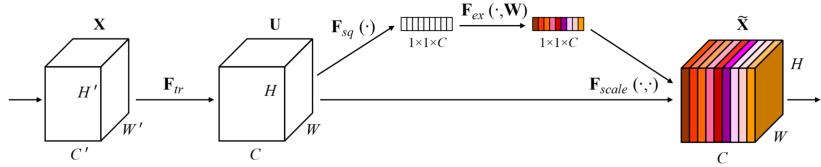


Figure 8: *Squeeze and Excitation* [hu2017squeeze]

#### 1.4.5 Micro resnet

Our network differs is composed by  $n$  ResNet blocks, a depth of  $d$  and a channel incrementing factor of 2. We select  $n = 1$  and  $d = 3$  with a starting channel size of 16, we called this model architecture *micro-resnet*.

micro resnet
depth = 3
n = 1
3 x 3, 16 stride 1
2 x 2 max-pool
$\begin{bmatrix} 3 \times 3, & 16 \\ 3 \times 3, & 16 \end{bmatrix} \times 1$
SE-module
$\begin{bmatrix} 3 \times 3, & 32 \\ 3 \times 3, & 32 \end{bmatrix} \times 1$
SE-module
$\begin{bmatrix} 3 \times 3, & 64 \\ 3 \times 3, & 64 \end{bmatrix} \times 1$
SE-module
average pool, 1-d fc, softmax

Table 1: *micro-resnet* architecture

add model picture

#### 1.4.6 Normalization

Before feeding the data to the models, we need to make the patches height invariant. This must be done to correctly normalize different patches taken from different maps with different height scaling factor. We subtract the height of the map corresponding *Krock*'s position from the patch to correctly center it. The following figure shows the normalization process on the patch with the square in the middle.

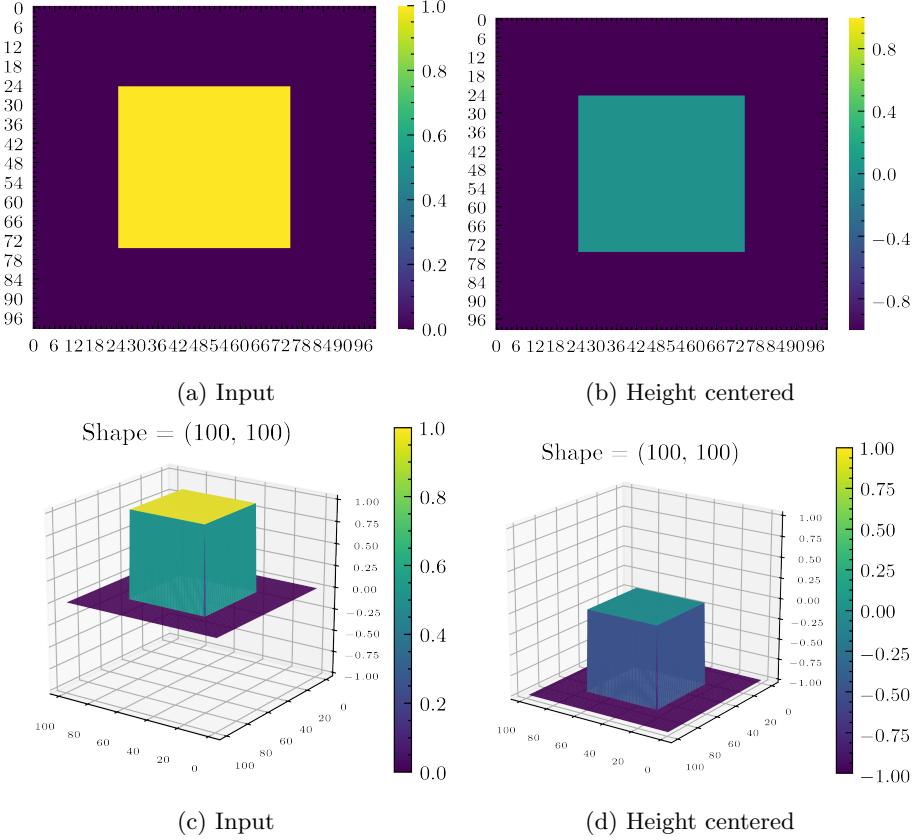


Figure 9: Normalization process

#### 1.4.7 Data Augmentation

Data augmentation is used to change the input of a model using different techniques to change it in order to produce more training examples. Since our inputs are heightmaps we cannot utilize the classic image manipulations such as shifts, flips, and zooms. Imagine that we have a patch with a wall in front of it, if we random rotate the image the wall may go in a position where the patch it is now traversable but its label is still not traversable, we have to be more creative. We decided to apply dropout, coarse dropout, and random simplex noise since they are traversability invariant. To illustrate those techniques we are going to use the following example patch of size  $100 \times 100$ .

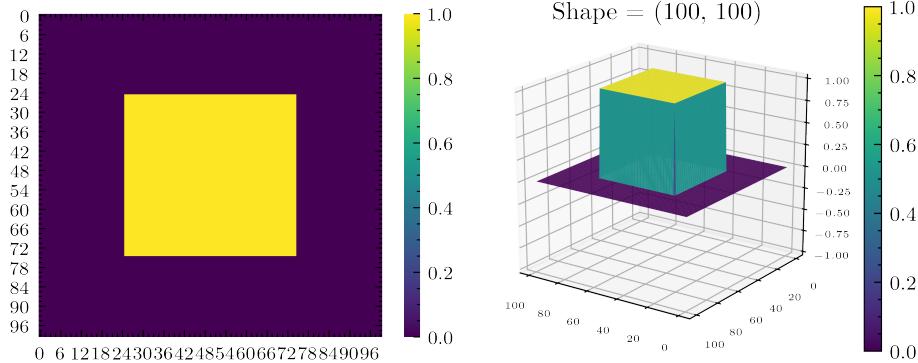


Figure 10: A patch with a square in the middle

**Dropout** is a technique to randomly set some pixels to zero, in our case we flat some random pixel in the patch.

**Coarse Dropout** similar to dropout, it sets to zero random regions of pixels.

**Simplex Noise** is a form of Perlin noise that is mostly used in ground generation. Our idea is to add some noise to make the network generalize better since lots of training maps have only obstacles in flat ground. Since it is computationally expensive, we randomly fist apply the noise to five hundred images with only zeros. Then, we randomly scaled them and add to the input image.

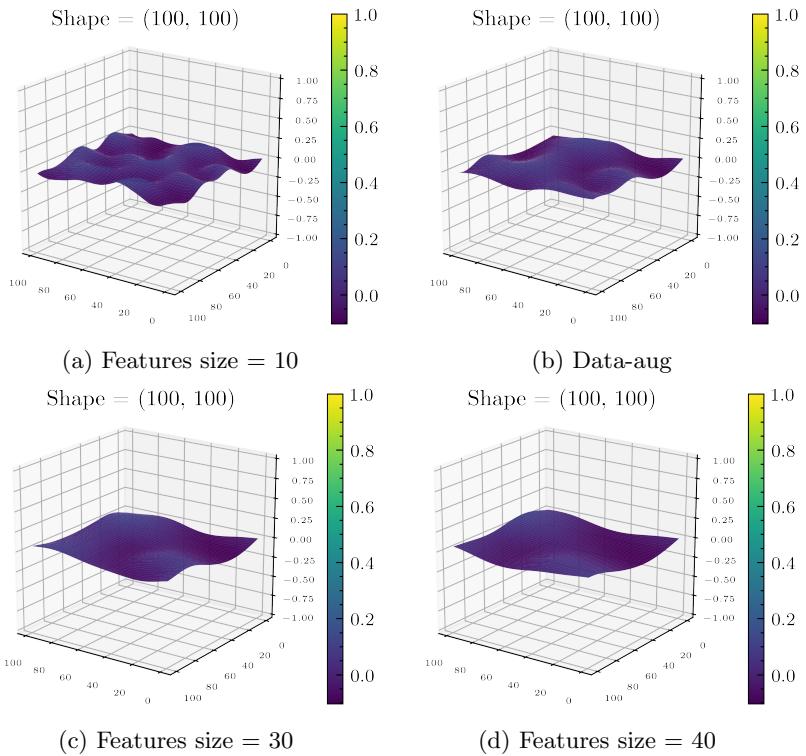


Figure 11: Simplex Noise on flat ground

The following images show the tree data augmentation techniques used applied the input image.

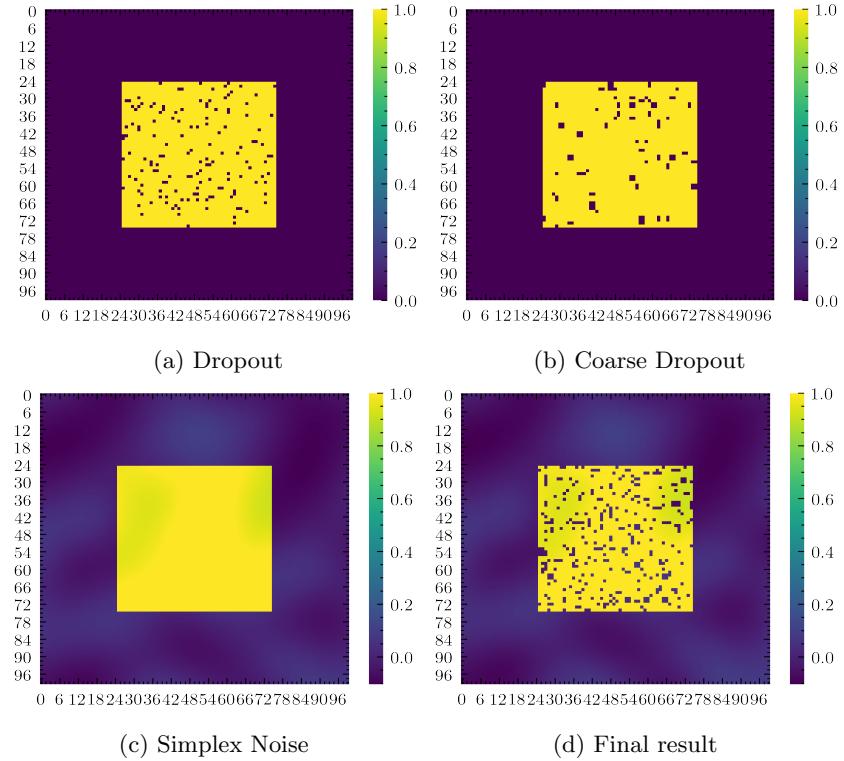


Figure 12: Data augmentation

It follows an other set of figures that shows the data augmentation we utilised on different inputs.

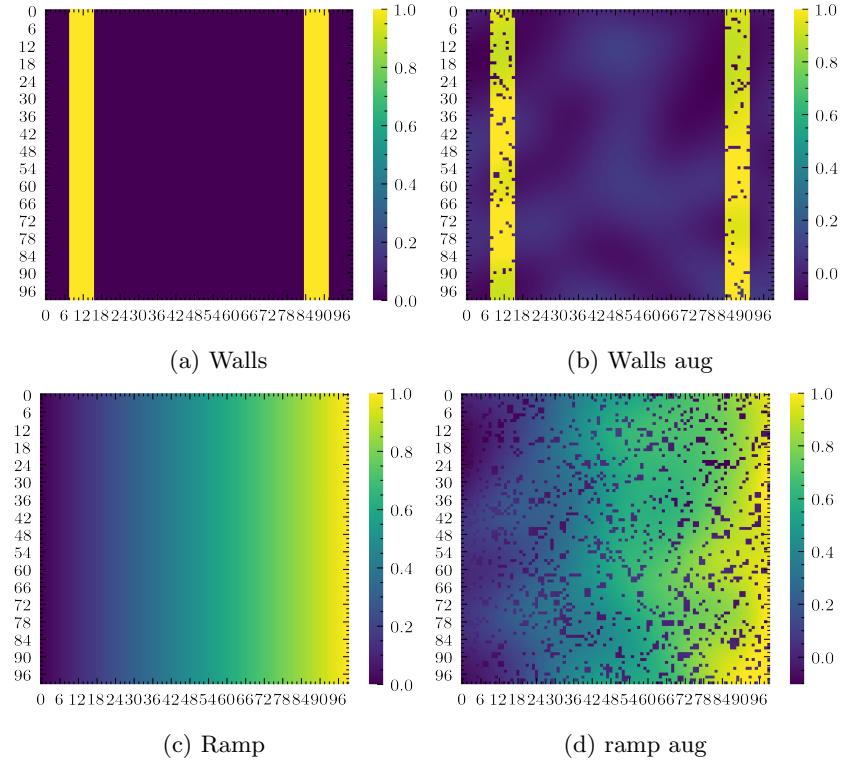


Figure 13: Wall

add the parameters that we set

In all the traning epochs, we apply data-augmentation to each input image  $x$  with a probability of 0.8. Dropout has a probability between 0.05 and 0.1. Coarse dropout with a probability of 0.02 and 0.1 with a size of the lower resolution image from which to sample the dropout between 0.6 and 0.8. Simplex noise with a feature size between 1 and 50 with a random scaling factor between 6 and 10.