

Predicting Frictional Properties of Graphene Kirigami Using Molecular Dynamics and Neural Networks

Designs for a negative friction coefficient.

Mikkel Metzsch Jensen



Thesis submitted for the degree of
Master in Computational Science: Materials Science
60 credits

Department of Physics
Faculty of mathematics and natural sciences

UNIVERSITY OF OSLO

Spring 2023

Predicting Frictional Properties of Graphene Kirigami Using Molecular Dynamics and Neural Networks

Designs for a negative friction coefficient.

Mikkel Metzsch Jensen



© 2023 Mikkel Metzsch Jensen

Predicting Frictional Properties of Graphene Kirigami Using Molecular Dynamics and Neural Networks

<http://www.duo.uio.no/>

Printed: Reprocentralen, University of Oslo

Abstract

Abstract.

Acknowledgments

Acknowledgments.

List of Symbols

F_N Normal force (normal load)

Acronyms

CNN Convolutional Neural Network. 8, 9, 10

EMA Exponetial Moving Average. 6

GA Genetic Algorithm. 14

MD Molecular Dynamics. 3, 19

MSE Mean Squared Error. 4

RMSProp Root Mean Square Propagation. 6

SGD Stochastic gradient descent. 5

Contents

| | |
|--|-----------|
| I Background Theory | 1 |
| 1 Machine Learning | 3 |
| 1.1 Neural network | 3 |
| 1.1.1 Optimizers | 5 |
| 1.1.2 Weight decay | 6 |
| 1.1.3 Parameter distributions | 7 |
| 1.1.4 Learning rate decay strategies | 7 |
| 1.2 Convolutional Neural Network | 8 |
| 1.2.1 Training, validation and test data | 10 |
| 1.3 Overfitting and underfitting | 11 |
| 1.4 Hypertuning | 11 |
| 1.5 Prediction explanation | 13 |
| 1.6 Accelerated search using genetic algorithm | 14 |
| 1.6.0.1 Repair function | 16 |
| II Simulations | 17 |
| 2 Summary | 19 |
| 2.1 Summary and conclusions | 19 |
| 2.1.1 Design MD simulations | 19 |
| 2.1.2 Design Kirigami framework | 19 |
| 2.1.3 Control friction using Kirigami | 19 |
| 2.1.4 Capture trends with ML | 20 |
| 2.1.5 Accelerated search | 20 |
| 2.1.6 Negative friction coefficient | 20 |
| 2.2 Outlook / Perspective | 20 |
| Appendices | 23 |
| A Appendix A | 25 |
| A Appendix B | 27 |
| B Appendix C | 29 |

Part I

Background Theory

Chapter 1

Machine Learning

We will use machine learning to predict the friction resulting from the stretching and loading of a given Kirigami pattern. To this end, we will generate data through MD simulations that will serve as the ground truth for training the machine learning model. The advantage of using machine learning is that it can significantly speed up the exploration of new configurations compared to full MD simulations. However, there is no guarantee that the machine learning model can accurately capture the physical mechanisms governing our system. Hence, a key objective is to assess the viability of this approach for further studies of Kirigami friction. It is not guaranteed that the machine learning model can accurately capture the physical mechanisms of our system. Hence, one of our objectives is to evaluate the applicability of this approach in the study of Kirigami friction, which we will pursue using a rather traditional machine-learning approach. In this chapter, we introduce the key concept behind machine learning and some of the concepts and techniques relevant to our implementation.

1.1 Neural network

The neural network, or more precisely the *feed forward dense neural network*, is one of the original concepts in machine learning arising from the attempt of mimicking the way neurons work in the brain [1, 2] brain. The neural network can be considered in terms of three major parts: The input layer, the so-called *hidden layers* and finally the output layer as shown in Fig. 1.1. The input is described as a vector $\mathbf{x} = x_0, x_1, \dots, x_{n_x}$ where each input x_i is usually denoted as a *feature*. The input features are densely connected to each of the *nodes* in the first hidden layer as indicated by the straight lines in Fig. 1.1. Each line represents a weighted connection that can be adjusted to configure the importance of that feature. Similar dense connections are present throughout the hidden layers to the final output layer. For a given note $a_j^{[l]}$ in layer l the input from all the nodes in the previous layer $l - 1$ are processed as

$$a_j^{[l]} = f \left(\sum_i w_{ij}^{[l]} a_i^{[l-1]} + b_j^{[l]} \right),$$

where $w_{ij}^{[l]}$ is the weight connection node $a_i^{[l-1]}$ of the previous layer to the node $a_j^{[l]}$ in the current layer. Note that having the weight belong to layer l as opposed $l - 1$ is simply a notation choice. $b_j^{[l]}$ denotes a bias and $f(\cdot)$ the *activation function*. The activation function provides a non-linear mapping of the input to each node. Without this, the network will only be capable of approximate linear functions [1]. Two common activation functions are the *sigmoid*, mapping the input to the range $(0, 1)$, and the *ReLU* which cuts off negative contributions

$$\text{Sigmoid: } f(z) = \frac{1}{1 + e^{-z}}, \quad \text{ReLU: } f(z) = \begin{cases} z & z > 0 \\ 0 & z \leq 0. \end{cases}$$

Often the same activation function is used throughout the network, except for the output layer where the activation function is usually omitted or the sigmoid is used for classification tasks. The whole process of sending data through the model is called *forward propagation* and constitutes the mechanism for mapping an input \mathbf{x} to the model output $\hat{\mathbf{y}}$. In order to get useful predictions we must *train* the model which involves tuning the model parameters, i.e. the weight and biasses.

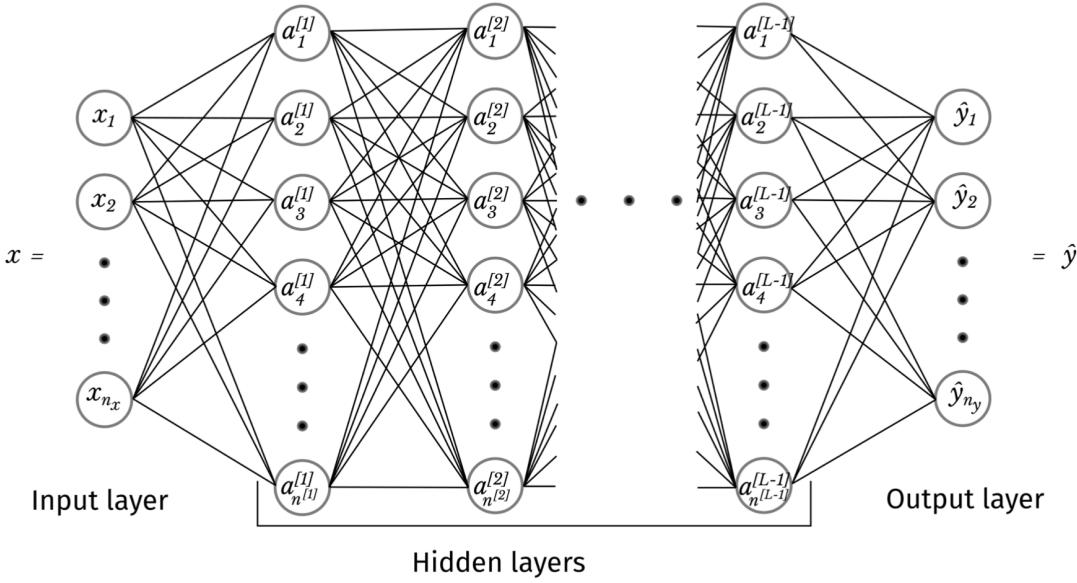


Figure 1.1: From overleaf IN5400

The model training relies on two core concepts: *backpropagation* and *gradient descent* optimization. First, we define the error associated with a model prediction, otherwise known as the *loss*, through the *loss function* L that evaluates the model output $\hat{\mathbf{y}}$ against the ground truth \mathbf{y} . For a continuous scalar output, we might simply use the mean squared error (MSE)

$$L_{\text{MSE}} = \frac{1}{N_y} \sum_{i=1}^N (y_i - \hat{y}_i)^2.$$

For a binary classification problem, meaning that the true output is True or False (1 or 0), a common choice is binary cross entropy (BCE)

$$L_{\text{BCE}} = - \sum_{i=1}^n \left[y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \right] = \sum_{i=1}^n \begin{cases} -\log(\hat{y}_i), & y_i = 1 \\ -\log(1 - \hat{y}_i), & y_i = 0. \end{cases}$$

The cross-entropy loss can be derived from a maximum likelihood estimation [SOURCE](#). Without going into details with the derivation we can convince ourselves that the error is minimized for the correct prediction and maximized for the worst prediction. When $y_i = 1$ we get the negative term $-\log(\hat{y}_i)$ where a correct prediction $\hat{y}_i \rightarrow 1$ yields a loss contribution $L_i \rightarrow 0$. For a wrong prediction $\hat{y}_i \rightarrow 0$ the loss contribution will diverge $L_i \rightarrow \infty$. Similar applies to the case of $y_i = 0$ with opposite directions.

Given a loss function, we can calculate the loss gradient $\nabla_{\theta} L$ with respect to each of the weights and biases in the model. This is called *backpropagation* since we follow the propagation of the errors as we go backward back through the model layers calculating the gradient using the chain rule. These gradients express how each parameter is connected to the loss and the overall idea is then to “nudge” each parameter in the right direction for reducing the loss. We usually denote a full cycle of forward-, backpropagation and an update of all model parameters as an *epoch*. We calculate the updated parameter θ_t for epoch t using the *gradient descent* method

$$\theta_t = \theta_{t-1} - \eta \nabla_{\theta} L(\theta_t).$$

Gradient descent is analog to taking a step in parameter space in the direction that yields the biggest decrease in the loss. If we imagine a simplified case with only two parameters θ_1 and θ_2 we can think of these as directions on a map and the loss being the terrain height. The gradient descent steps in the direction perpendicular to the contour lines shaped by loss function terrain as shown in Fig. 1.2a. Notice, however, that state-of-the-art models

in general contain on the order of $10^6 - 10^9$ parameters [3] which poses some challenges for the visualization. The length of each step is proportional to the gradient norm and the learning rate η . There are three main flavors to the gradient descent: Batch, stochastic and mini-batch gradient descent. In *batch gradient descent* we simply calculate the gradient based on the entire dataset by averaging the contribution from each data point before updating the parameters. This gives the most robust estimate of the gradient and thus the most direct path through parameter space in terms of minimizing the loss function as indicated in Fig. 1.2a. However, for big datasets, this calculation can be computationally heavy as it must carry the entire dataset in memory at once. A solution to this issue is provided by *stochastic gradient descent* (SGD) which considers only one data point at a time. Each data point is chosen randomly and the parameters are updated based on the corresponding gradient. This leads to more frequent updates of the parameters and a more “noisy” path through parameter space as shown in Fig. 1.2b. Under some circumstances, this might compromise the precision. However, the presence of noise can actually increase the chances of avoiding local minima in parameter loss space. The *mini-batch gradient descent* serves as a middle ground between the above methods by dividing the full dataset into a subset of mini-batches. Each parameter update is then based on the gradient within a mini-batch. By choosing a suitable batch size we get the robustness of the (full) batch gradient descent and the computational efficiency and resistance to local minima of the SGD method.

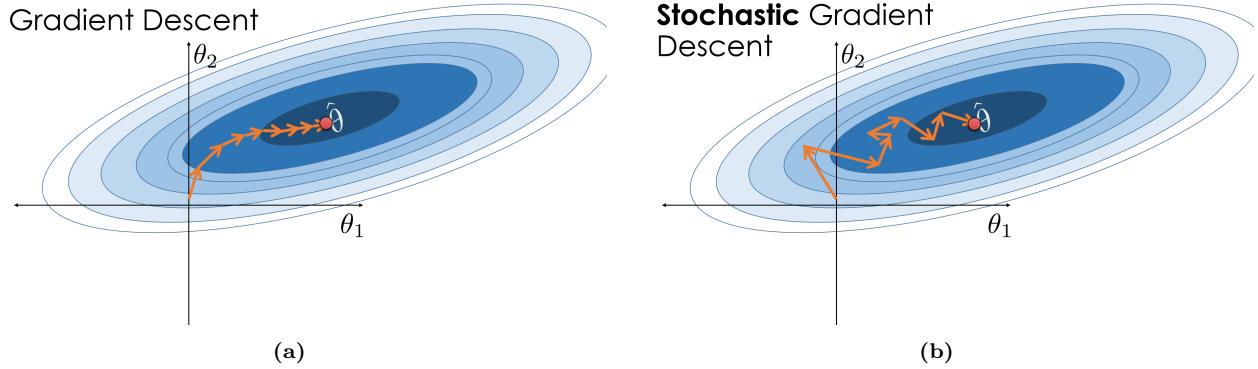


Figure 1.2: TMP

1.1.1 Optimizers

The name *optimizers* covers a variety of gradient descent methods. In our study, we will use the ADAM (adaptive moment estimation). ADAM combines several “tricks in the book” which we will introduce in the following.

One considerable extension of the gradient descent scheme is by the introduction of a momentum term m_t such that we get

$$\theta_t = \theta_{t-1} - m_t, \quad m_t = \alpha m_{t-1} + \eta \nabla_\theta L(\theta_t) \quad (1.1)$$

with $m_0 = 0$. If we introduce the shorthand $g_t = \nabla_\theta L(\theta_t)$ we find

$$\begin{aligned} m_1 &= \alpha m_0 + \eta g_1 = \eta g_1 \\ m_2 &= \alpha m_1 + \eta g_2 = \alpha^1 \eta g_1 + \eta g_2 \\ m_3 &= \alpha m_2 + \eta g_3 = \alpha^2 \eta g_1 + \alpha \eta g_2 + \eta g_3 \\ &\vdots \\ m_t &= \eta \left(\sum_{k=1}^t \alpha^{t-k} g_k \right). \end{aligned} \quad (1.2)$$

Hence m_t is a weighted average of the gradients with an exponentially decreasing weight. This act as a memory of the previous gradients and aid to pass local minima and to some degree plateaus in the loss space. A variation of momentum can be achieved with the introduction of the exponential moving average (EMA) which builds on the recursion

$$\begin{aligned}
\text{EMA}(g_1) &= \alpha \overbrace{\text{EMA}(g_0)}^{\equiv 0} + (1 - \alpha)g_1 \\
\text{EMA}(g_2) &= \alpha \text{EMA}(g_1) + (1 - \alpha)g_2 \\
&\vdots \\
\text{EMA}(g_t) &= \alpha \text{EMA}(g_{t-1}) + (1 - \alpha)g_t = \sum_{k=0}^t \alpha^{t-k} (1 - \alpha) g_t,
\end{aligned}$$

which is similar to that of momentum Eq. (1.2), but with the explicit weighting by $(1 - \alpha)$.

The second moment of the exponential moving average is utilized in the root mean square propagation (RMSProp) method which is motivated by the issue of passing long plateaus in the loss space. Since the size of the updates are proportional to the norm of the gradient

$$\theta_{t+1} = \theta_t - \eta g_t \implies \|\theta_{t+1} - \theta_t\| = \eta \|g_t\|,$$

we might get the idea of normalizing the gradient step by dividing with the norm $\|g_t\|$. However, this does not immediately solve the problem of long plateaus as we need to consider multiple past gradients which is then accommodated by the use of the EMA. When reentering a steep region again we need to “quickly” downscale the gradient steps again which can be achieved by using the squared norm $\|g_t\|^2$ for the EMA which makes it more sensitive to outliers. The RMSProp update scheme then becomes

$$\theta_t = \theta_{t-1} - \eta \frac{g_t}{\sqrt{\text{EMA}(\|g_t\|^2)} + \epsilon}, \quad (1.3)$$

where ϵ is simply a small number to avoid division by zero issues.

ADAM merges the idea of first order EMA for the momentum m_t , and the second order EMA v_t , as used in the root mean square propagation technique in Eq. (1.3)

$$\begin{aligned}
m_t &= \beta_1 m_{t-1} + (1 - \beta_1)g_t, \\
v_t &= \beta_2 v_{t-1} + (1 - \beta_2)g_t^2.
\end{aligned}$$

Since these are initially set to zero we can correct a bias towards zero by a scaling term $(1 - \beta_1^t)$ and $(1 - \beta_2^t)$ respectively such that the ADAM scheme becomes

$$\theta_{t+1} = \theta_t - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}, \quad \hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}. \quad (1.4)$$

1.1.2 Weight decay

By adding a so-called *regularization* to the loss function we can penalize high magnitudes of the model parameters. This is motivated by the idea of preventing overfitting during training. The most common way to do this is by the use of L2 regularization, adding the squared l^2 norm $\|\theta\|_2^2$, where $\|\theta\|_2 = \sqrt{\theta_1^2 + \theta_2^2 + \dots}$, to the model. The new loss and gradient then become

$$\begin{aligned}
L_{l^2}(\theta) &= L(\theta) + \frac{1}{2} \lambda \|\theta\|_2^2 \\
\nabla_\theta L_{l^2}(\theta) &= \nabla_\theta L(\theta) + \lambda \theta,
\end{aligned} \quad (1.5)$$

where λ is the weight decay parameter $\in [0, 1]$. The name *weight decay* relates to the fact that some practitioners only apply this penalty to the weights in the model, but we will include the biases as well (standard in PyTorch). Following the original gradient descent scheme Eq. (1.5) we get

$$\begin{aligned}
\theta_{t+1} &= \theta_t - \eta g_t - \eta \lambda \theta_t \\
&= \theta_t \underbrace{(1 - \eta \lambda)}_{\text{Weight decay}} - \eta g_t.
\end{aligned}$$

Thus we notice that choosing a high weight decay (towards 1) will downscale the model parameters while choosing a low weight decay (towards 0) yields the original gradient descent scheme. Note that this is used in combination with ADAM, but it is easier to show the consequences for the original gradient descent scheme.

1.1.3 Parameter distributions

In order to get optimal training conditions it has been found that the initial state of the weight and biases are important **SOURCE**. First of all, we must initialize the weight by sampling from some distribution. If the weights are set to equal values the gradient across a layer would be the same. This results in a complexity reduction as the model can only encode the same values across the layer **SOURCE**. Further, we want to consider the gradient flow during training. Especially for deep networks, networks with many layers, we must pay attention to the problem of vanishing or exploding gradients. If we for instance consider the sigmoid activation function and its derivative

$$f(z) = \frac{1}{1 + e^{-z}}, \quad f'(z) = \frac{df(z)}{dz} = \frac{e^{-z}}{(1 + e^{-z})^2} = \frac{e^z}{(1 + e^z)^2},$$

we notice that for large and small values we get $f(z \rightarrow \pm\infty) \rightarrow 0$. However, even a small gradient will vanish throughout a deep network as the calculation of the gradient involves the chain rule. A similar problem can be found with the ReLU activation function which gets a zero gradient for inputs $z < 0$ which can be mitigated by the so-called leaky ReLU which gives the $z < 0$ a small slope. On the other hand, we have exploding gradients, which are simply a result of the chain rule gradient calculation. For a sufficiently deep network, the gradient can grow exponentially and sometimes result in numerical overflow. While there exist techniques to accommodate this, like for instance the leaky ReLU for the vanishing gradients and so-called gradient clipping, cutting off the gradient at a maximum, they both benefit from a properly initialized set of weights **SOURCE**. That is, we want the gradients across a given layer to have a zero mean while the variance is similar to other layers. This balanced gradient flow is more likely to happen if we initialize the weight by the same set of criteria. The specific actions to achieve this depend on model architecture, including the choice of activation functions. For instance, using the ReLU activation functions it was found that the single standard deviation will depend on the number of input nodes from the previous layer $N^{[l-1]}$ as $\sim \sqrt{N^{[l-1]}}/\sqrt{2}$ and thus we simply scale a zero mean uniform distribution to match this. This is part of the Kaiming initialization scheme which is standard in Pytorch **SOURCE**. **Mention choices for bias initialization.**

Batch normalization is another technique that might also help reduce the issue of poor gradient flow. Furthermore, can benefit by speeding up convergence making it more stable **SOURCE**. In general, model parameters are modified throughout training meaning that the range of values coming from a previous layer will shift, even though the same training data is fed through the network repeatedly. By scaling the input for a given layer, for each mini-batch, we can mitigate this problem and make for a more standardized input range. This often result in a faster training convergence. For layer l we calculate the mean $\mu^{[l]}$ and variance $\sigma^{2[l]}$ across the layer with nodes $x_1^{[l]}, x_2^{[l]}, \dots, x_d^{[l]}$ for each mini-batch of size m as

$$\mu^{[l]} = \frac{1}{m} \sum_i z^i, \quad \sigma^{2[l]} = \frac{1}{m} \sum_i^d (x^i - \mu)^2.$$

We then perform normal scaling of the inputs within the batch

$$\hat{x}_i^{[l]} = \frac{x_i^{[l]} - \mu^{[l]}}{\sqrt{\sigma^{2[l]}} + \epsilon},$$

where ϵ is a small number to ensure numerical stability (similar to what we used for RMSProp gradient descent). In the final step the input values are rescaled as

$$\tilde{x}_i = \gamma^{[l]} \hat{x}_i^{[l]} + \beta^{[l]}$$

with trainable parameters γ and β . **Comment about the reason for the final step.**

1.1.4 Learning rate decay strategies

Until now we have assumed a constant learning rate, but variations use a changing learning rate beyond the adaptiveness included in the optimizers as described earlier. It can be beneficial to start with a higher learning rate to speed up the initial part of training and then lower the learning rate for the final gradient descent. One straightforward strategy is a step-wise learning rate decay where the learning rate is reduced by a factor $\gamma \in (0, 1)$

every K steps. A more smooth change can be achieved by for instance a polynomial decay $\eta_t = \eta_0/t^\alpha$ for $\alpha > 0$. More advanced approaches use multiple cycles of increasing and decreasing cycles. We will mainly concern ourselves with a one-cycle policy for which we start at an intermediate value, increase toward a maximum bound and then decrease toward a final lower learning rate bound. This can simply be done by a linear increase and decrease, but we choose a cosine function that is shifted to peak for the first 30% of the training length.

1.2 Convolutional Neural Network

Convolutional Neural Networks (CNNs) build upon some of the same concepts as introduced with the feed-forward neural network. The difference lies in its specialization for a spatially correlated input, such as pixels in an image. In a dense neural network, every node is connected to each of the nodes from the previous layers. Thus, it does not matter how the input is arranged initially, but it cannot be changed at a later stage. This is impractical if we want the model to recognize images of animals as the animal-related pixels will be in different regions of each image. The CNN can be motivated by the idea of capturing spatial relations, but without being sensitive to the relative placement within the image, i.e. being translational invariant. This is achieved by having a so-called *kernel* or *filter* which slides over the images¹ as it processes the input. A convolutional layer contains multiple kernels, each consisting of a set of trainable weights and a bias. Each kernel will produce a separate output channel to the resulting *feature map*. The kernel has a spatial size, specific to the model architecture, and a depth that matches the number of input channels to the layer. For instance, a typical RGB will have three channels. The kernel lines up with the image and calculates the feature map output as a dot product between the weights in the kernel and the aligning subset of the input. This is done for each input channel and summed up with the addition of a bias as illustrated in Fig. 1.3b. The kernel then slides over by a step size given by the *stride* model parameter and repeat the calculation. Choosing a stride of 2 or higher results in a reduction of spatial size. If we want to preserve the spatial size we must keep a stride of one and additionally apply *padding* to the input images, such that we can achieve one kernel position for each input “pixel”. The spatial size of the feature map is given as

$$N_d^{[l+1]} = \left\lfloor \frac{N_d^{[l]} - F_d + 2P}{S} + 1 \right\rfloor, \quad (1.6)$$

for padding P , stride S , spatial size of the kernel filter F_d , spatial size of input $N_d^{[l]}$, for dimension $d = x, y$ and layer l . The *down-sampling* is often done through a pooling layer. A pooling layer is reminiscent of a kernel, but instead of calculating the output as a dot product, it calculate the mean (mean pooling) or the max value (max pooling) of the values within its range. For instance, by using a max pooling of size 2×2 and stride two we essentially half the dimensions of the image as dictated by Eq. (1.6). CNNs will often use repeating series of convolution, pooling and then an activation function. Most architectures will down-sample the spatial input while increasing the number of channels in the model.

¹Note, that we will be using the word “image” as a reference for a spatially dependent input, but in reality, it does not have to be an actual image in the classical sense.

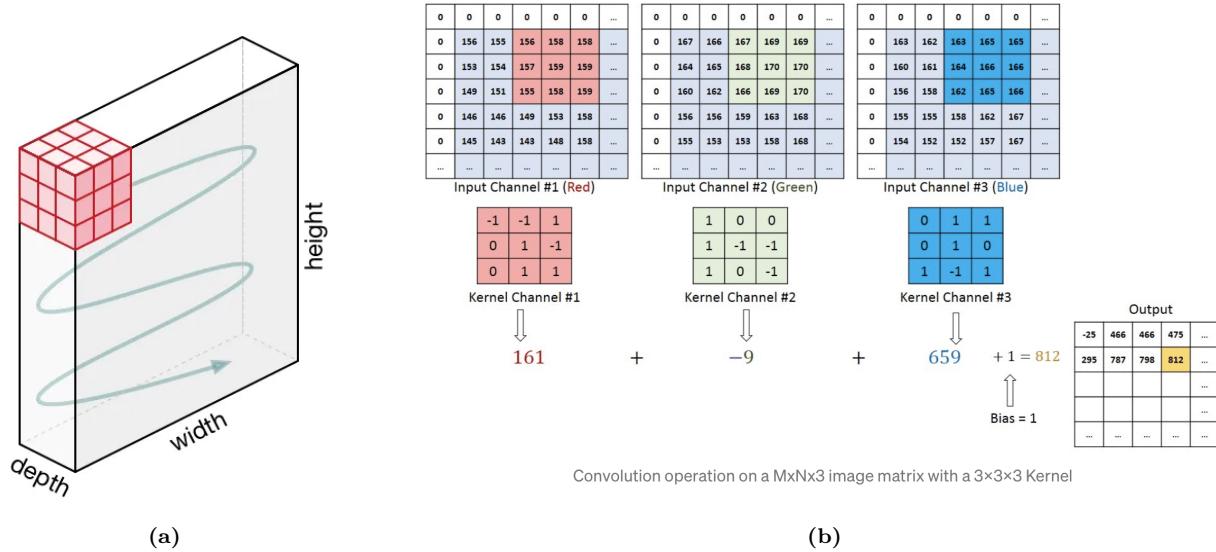


Figure 1.3: TMP

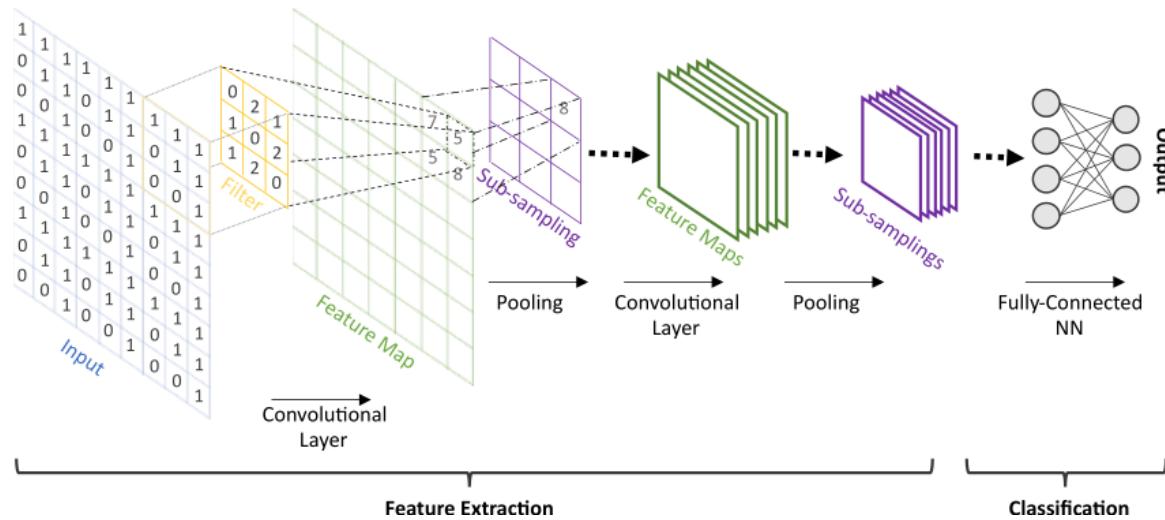


Figure 1.4: TMP Not used at the moment...

For a CNN, we often consider the *receptive field*. The receptive field relates to the spatial size of the input that affects a given node in the feature map at a given layer of the model. In Fig. 1.5 the receptive field is illustrated for a 1D representation of a CNN with repetitive use of a kernel of size 2 and stride 1. Going from the output and backward, we see that the output layers are connected to two nodes in the previous layer. Each of these nodes is connected to two nodes in the layer before that, however with one of them being the same due to the stride of 1. By back-tracking to the input we see that this corresponds to a receptive field of $D = 5$. By increasing the filter size and the stride the 2D receptive field will grow a lot faster than shown in this simple 1D example. For a receptive field D_d , with respect to the spatial dimension d , a spatial size of the filter F_d , stride S_l (from layer $l-1$ to l) we have

$$D_l = D_{l-1} + \left[(F_l - 1) \cdot \prod_{i=1}^{l-1} S_i \right],$$

with $D_0 = 1$ and $l = 0$ as the input layer. Note that by convention, the product of zero elements is 1, such that for the first layer, the product is 1. The receptive field is important in understanding the connectivity in the

model. The model output will be completely independent of the inputs and feature maps outside the receptive field. Furthermore, we differentiate between the theoretical receptive field and the effective receptive field. The effective receptive field will have a Gaussian distribution within the theoretical receptive field due to the fact that these nodes in the center of the receptive field will have more connections leading to the output, as seen in Fig. 1.5. Thus, in practice the effective receptive field will be smaller than the theoretical Implementations like dilated convolutions, which make the filter expand in circumference and skipping positions within the filter, which can be used to further increase the effective field.

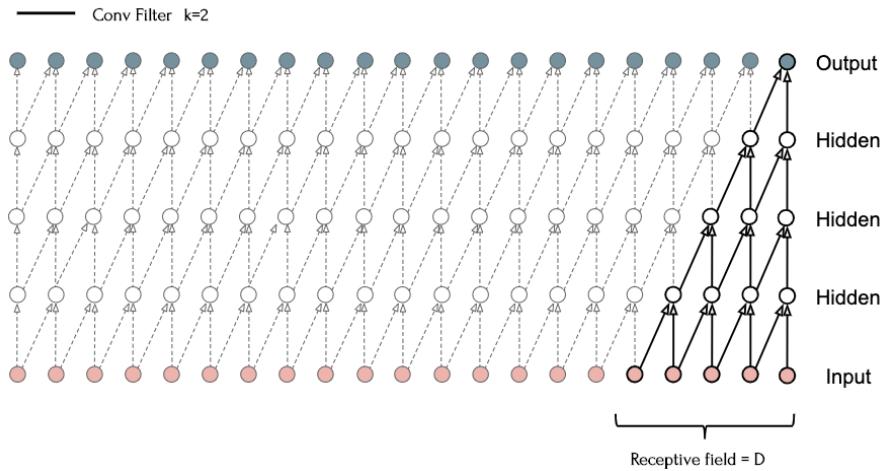


Figure 1.5: TMP

On a final note regarding the CNN we point out that convolution is often used in combination with a dense network, or *fully connected*, at the end. We can then think of the convolution part to handle the translation from a spatial import to some internal features. For the animal detection network, we would perhaps suggest features like, number of legs, size, color and so on. In practice, the network will not create easily interpreted features for the processing of the fully connected layer. We discuss one approach for interpreting of the model internals in Sec. 1.5

1.2.1 Training, validation and test data

So far, we have simply considered the concept of training data as a means to update the model parameters. Yet, we want to evaluate the model performance as it improves. If a model has enough complexity, i.e. enough layers with a sufficient amount of nodes, it can be proven that it can fit any function [SOURCE](#). Thus for a complex model, it is just a matter of time before the model eventually will learn to give matching predictions for the training data. However, we want the model to learn general trends and not to “memorize” all the data points which is known as overfitting. While the predictions can get arbitrarily good the performance on unseen data within the same task will yield awful results if the model is overfitted. Hence, we put aside some subset of the data which we use as a *validation* throughout and after training. By keeping this *validation* set separate from the training data we can get a more reliable performance estimate for the model. This splitting of data should be done randomly in order to ensure a similar distribution of data features in each set. For the size ratio we want to strike a balance between training quality and validation quality, usually a 20%–80% split favoring the training set is a good figure. Other techniques exist which aim to optimize the data use for sparse data situations, like cross-validation and bootstrap [fact check](#), but we will not consider such methods for this thesis. A third training set which is often forgotten is the *test* set. While the validation set should be kept unseen from the model training, the test set should be kept unseen from the model developer. As we choose the model architecture and hyper-parameters like learning rate, momentum and weight decay, we will use the performance on the validation set as a metric to steer by. This represents a slower pace parameter fitting which can also lead to overfitting on a higher level. Hence, we should denote a test set for the final evaluation of our model which have not been considered before the end. Formally, this is the only reliable performance metric for the model.

1.3 Overfitting and underfitting

Underfitting and overfitting represent a crucial balance going on when training a model. This concept is highly related to the model complexity and the chosen hyper-parameters such as learning rate, momentum and weight decay. The textbook visualization of underfitting and overfitting is shown in Fig. 1.6a. As we begin to train or model both the training and validation loss is decreasing. At some point, the model will start to pick up to specific trends which marks the transition into overfitting where the validation loss will start to increase again. Notice that the performance on the training data will keep increasing. *Early stopping* can be utilized to detect this transition and stop the training. We will use a variation of this which is based on storing the model parameters at the best validation performance but letting the model training go on. We can imagine a similar curve as seen in Fig. 1.6a with the replacement of *model complexity* on the x-axis. For a certain amount of epochs a simple model will yield underfitting and an overly complex model will yield overfitting. In figure Fig. 1.6b we see how under- and overfitting can be represented by a fitting to a 2D curve. We see how a simple, not complex, underfitted model will make a crude approximation for the true curve. An overly complex overfitted model will follow pick up the noise in the training data and miss the general trend. In practice the diagnosing of under- and overfitting is not as simple as the figures in Fig. 1.6 imply. First of all, the Fig. 1.6a is simplifying the idea, and the training and validation loss will rarely showcase such a clear indication of the transition to overfitting. The problem with Fig. 1.6b lies in the fact that we do not know the true curve. If we did, we would not need machine learning to estimate it in the first place. Without having additional insight into the governing source of the data the overfitting case seems to produce the most confident fit.

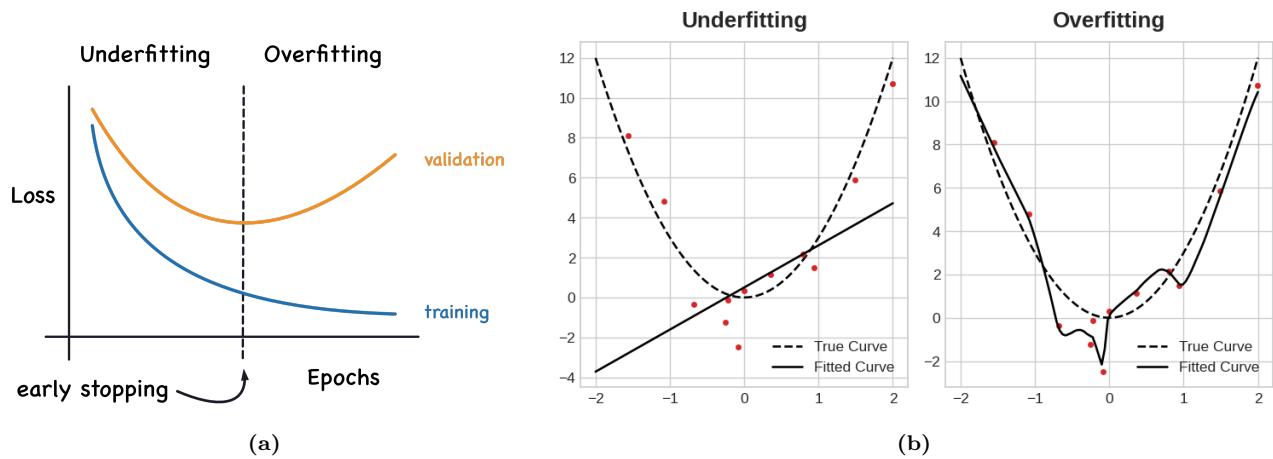


Figure 1.6: TMP

1.4 Hypertuning

Define hyper-parameter properly like done in [4]

Training a machine learning model revolves around tuning the model parameters such as weights and biases. However, a handful of so-called *hyper-parameters* remains for us to decide. First of all, we need to choose an architecture for the model. This includes high-level considerations, for instance, whether to use a neural network or a convolutional network, but also lower-level considerations, such as the depth and the width of the model, i.e. how many layers and how many nodes/channels. In addition, we have to define and consider the loss function and the optimizer which come with hyper-parameters such as learning rate, momentum and weight decay. This extensive list of choices makes the designing of an a functional model more complicated than simply starting the training. As N. Smith [5] puts it: “Setting the hyper-parameters remains a black art that requires years of experience to acquire”. In the following we will review a general approach for choosing the learning rate, momentum and weight decay hyper-parameter based mainly on the findings of [5]. The traditional approach is to perform a grid-search of hyper-parameter combinations for multiple training sessions, but this might rather quickly become computationally expensive and ineffective. In addition, hyper-parameters will depend on the training data, the model architecture and not at least each other.

The learning rate is often regarded as the most important hyper-parameter to tune [4]. Typical values are in the range $[10^{-6}, 1]$. Instead of simply running a grid search, i.e. multiple full training sessions using different learning rates, we can perform a so-called *learning rate range test* (LR range test). One specifies the minimum and maximum learning rate boundaries and a learning rate step size. A minimum and maximum bound of 10^{-7} to 10 will most likely cover an appropriate range, but the test will reveal this immediately. The idea is then to vary the learning rate throughout the given range in small steps during a short pre-training. We vary the learning rate for each iteration, i.e. each parameter update following a mini-batch, and thus we can run this test for a few epochs, or even a single one, depending on the number of mini-batches. The learning rate can be varied in a linear increasing or decreasing manner which is found to produce similar results [6]. We choose the linear increasing version for simplicity. For small learning rates, the model will converge slowly. As the learning rate approaches an appropriate value the convergence will accelerate which we see as a drop in the validation loss. Eventually the convergence will stop and the validation loss will pass a minimum for which it will begin to diverge. This general behaviour can be understood for the simplified 1D example of finding the minima of a second-order polynomial as shown in Fig. 1.7. Small learning rates will step in the right direction, but for overly small values this will simply result in a slow convergence. If the learning rate becomes too large, we will effectively step past the minimum. For each step we will overshoot the minimum more and more leading to a diverging trend. The point of divergence can be used as an upper bound for the learning rates when considering a cyclic learning rate scheme. The steepest decline of the validation loss can be used as an estimate for the best constant learning rate choice.

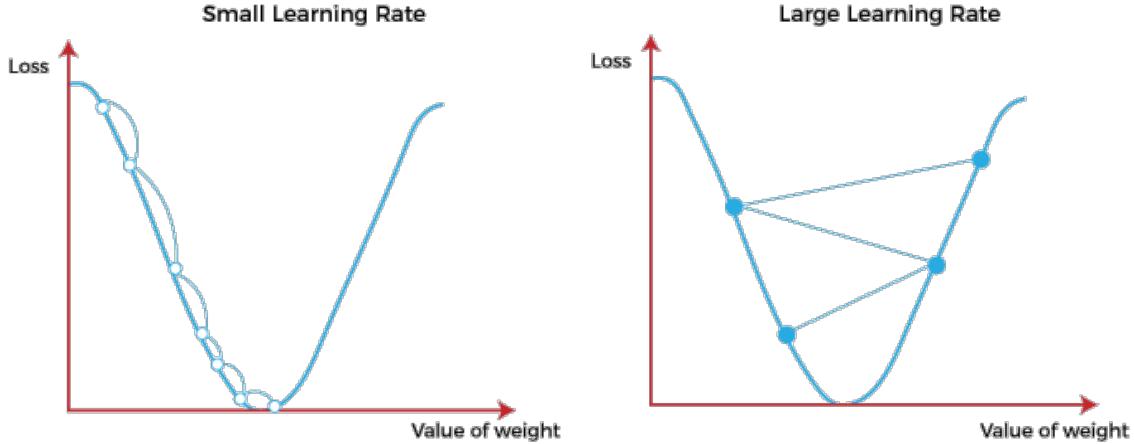


Figure 1.7: TMP

Next, we consider the choice of momentum. Momentum and learning rates are found to affect each other considerably. From the gradient descent scheme with momentum Eq. (1.1) we see that the momentum parameter α and the learning rate η have a similar effect on the parameter update

$$\theta_t = \theta_{t-1} - \eta g_t - \alpha m_{t-1},$$

as the m_t is a moving average of the gradient g_t as well. Like the learning rate, we want to set the momentum value as high as possible without causing instabilities in the training. However, it is found that these values are somewhat inversely related. Choosing a high learning rate should be coupled with a lower momentum and vice versa. N. Smith [5] reports that a momentum range test is not useful to find the right momentum. Instead, he suggests doing a few short runs with different values of momentum, such as 0.99, 0.97, 0.95, and 0.9, to determine a suitable choice. By including momentum in the LR range test we can balance the learning rate accordingly. Moreover, for a cyclic learning rate scheme he suggests using a cycling momentum scheme, reversed with respect to the learning rate. When the learning rate increase toward the upper bound the learning rate should decrease toward the lower bound and vice versa. Choosing a lower momentum of 0.80–0.85 often gives similar stable results.

Finally, we address the weight decay. N. Smith reports that weight decay is different from learning rate and momentum by the fact that weight decay is better chosen as a constant value as opposed to a cyclic scheme.

However, the weight decay is dependent on the learning rate and momentum choice and this can often be dialed in after setting those. Furthermore, the weight decay will be dependent on the model complexity. We can estimate a suitable choice by doing a rough grid search for values such as 0, 10^{-6} , 10^{-5} and 10^{-4} for complex architectures and 10^{-4} , 10^{-3} and 10^{-2} for more shallow architectures. Choosing the weight decay on the scale of exponential exponents will often provide good enough precision in practice.

1.5 Prediction explanation

With the rise of constantly improving machine learning models and computational power, the capabilities of machine learning systems are improving as well. However, with deep learning models, we have essentially no insight into the considerations that went into a prediction other than the input data. This is known as the *black box* problem. A lot of effort is currently being developed for making more transparent models, like decision trees with interpretable rules, and numerical tools for unpacking the inner workings of the model. We will consider a gradient based method called *Grad-CAM* [7]. The idea is to use the gradients for a certain feature map with respect to the loss.

First, we forward propagate the input through the model, and decide on a feature map of interest. We then calculate the gradients for the feature map with respect to the loss of a target output. For a classification task, one would often choose the predicted class as the target output, the class with the highest score. The gradients will then highlight which part of the feature maps is most important for the prediction loss. A ReLU activation layer is then applied to only highlight the positive contributions. Since the convolutional layers preserve relative spacing we can rescale the heatmap provided by the feature map gradient to make an input-sized heatmap we can overlay with the input image. This provides a visual clue of what the prediction is based on. We can do this for different depths of the model and even combine the results for multiple layers. In Fig. 1.8 we see how this can be used to reveal the basis for a prediction. In this case, the Grad-CAM analysis shows the difference between a biased and unbiased model for the predictions of professions. The biased model shows to be considering the person more than the actually objective clues given by relevant equipment and work-related clothing

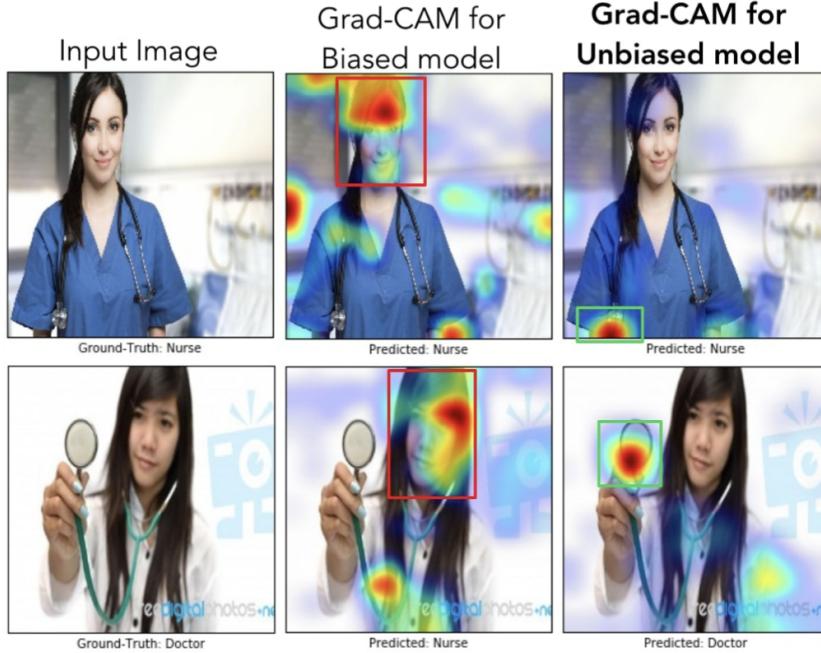


Fig. 8: In the first row, we can see that even though both models made the right decision, the biased model (model1) was looking at the face of the person to decide if the person was a nurse, whereas the unbiased model was looking at the short sleeves to make the decision. For the example image in the second row, the biased model made the wrong prediction (misclassifying a doctor as a nurse) by looking at the face and the hairstyle, whereas the unbiased model made the right prediction looking at the white coat, and the stethoscope.

Figure 1.8: TMP [7]

1.6 Accelerated search using genetic algorithm

The genetic algorithm (GA) is inspired by biological evolution and mimics the Darwin theory of the survival of the fittest. [8]. GA is a population-based algorithm for which the basic elements are chromosome representation, fitness selection and biological-inspired operators. The chromosomes represent the genes for each individual in the population and typically take the form of a binary string. Each position within the chromosome is called a *locus* and has two possible values (0 or 1). A fitness function is defined to assign a score for all chromosomes based on some optimization objective. This plays a role for the biologically inspired operators for which the main ones are selection, mutation and crossover. Selection is the process of selecting chromosomes based on their fitness score for further processing. In mutation, some of the loci within a chromosome are flipped and in crossover chromosomes are merged to create offspring. GA has been implemented in many areas such as the traveling salesman problem [9], function optimization [10], adaptive agents in stock markets [11] and airport scheduling [12]. Wang et al. [13] note that a general drawback is a need for expertise when choosing parameters that match specific applications. They propose an accelerated genetic algorithm based on a Markov chain transition probability matrix to perform a guided search that reduces the amount of parameter choices one have to make. We will implement this method and thus the following introduction of the method on [13]

We define the binary population matrix $A_{ij}(t)$ at generation t , consisting of N rows denoting chromosomes $i \in \{0, 1, \dots, N\}$ and L columns denoting the loci $j \in \{0, 1, \dots, L\}$. For our application, we let the locus represent an atom in the Kirigami pattern matrix which is flattened to fit the format of the population matrix. We carry forward the binary values with 0 meaning a removed atom and 1 a present atom. By the use of a fitness function $f(t)$, we sort the population matrix row-wise in descending order by fitness score, i.e. $f_i(t) \leq f_k(t)$ for $i \geq k$.

In the spirit of Markov chains, we assume that some transitions probability exists for the transition between the current state $A(t)$ and the next state $A(t + 1)$. We assume that this transition probability only takes into account the mutation process, and thus we omit operators like crossover. For each generation, the chromosomes are sorted according to the fitness function and the chromosome at the i^{th} fittest place is assigned a ranking score $r_i(t)$ by some monotonic increasing ranking scheme. We take this to be

$$r_i(t) = \begin{cases} (i - 1)/N', & i - 1 < N' \\ 1, & \text{else} \end{cases}$$

with $N' = N/2$ from [13]. We assign a row mutation probability $a_i(t)$ meaning that the probability for a mutation will increase towards the lower fitness scores. For the considerations of mutation with respect to each locus in the columns of $A_{ij}(t)$, we define the count of 0's and 1's as $C_0(j)$ and $C_1(j)$ respectively. These are normalized as

$$n_0(j, t) = \frac{C_0(j)}{C_0(j) + C_1(j)}, \quad n_1(j, t) = \frac{C_1(j)}{C_0(j) + C_1(j)}.$$

We can thus describe the state of the j^{th} locus column as the state vector $\mathbf{n}(j, t) = (n_0(j, t), n_1(j, t))$. In order to direct the current population to a preferred state for locus j we consider the highest weight $W_i = 1 - r_i$ among the chromosomes for the case of the locus being 0 or 1 respectively. This corresponds to the targets

$$\begin{aligned} C'_0(j) &= \max\{W_i | A_{ij} = 0; i = 1, \dots, N\} \\ C'_1(j) &= \max\{W_i | A_{ij} = 1; i = 1, \dots, N\}. \end{aligned}$$

These are normalized

$$n_0(j, t + 1) = \frac{C'_0(j)}{C'_0(j) + C'_1(j)}, \quad n_1(j, t + 1) = \frac{C'_1(j)}{C'_0(j) + C'_1(j)}. \quad (1.7)$$

to produce the target state vector $\mathbf{n}(j, t + 1) = (n_0(j, t + 1), n_1(j, t + 1))$. This will serve as a direction for each locus to evolve in and thus we can formulate the Markov chain as

$$\begin{bmatrix} n_0(j, t + 1) \\ n_1(j, t + 1) \end{bmatrix} = \begin{bmatrix} P_{00}(j, t) & P_{10}(j, t) \\ P_{01}(j, t) & P_{11}(j, t) \end{bmatrix} \begin{bmatrix} n_0(j, t) \\ n_1(j, t) \end{bmatrix},$$

where the matrix represents the transition matrix. Since the probability must sum to one for the rows in the transition matrix we get

$$P_{00}(j, t) = 1 - P_{01}(j, t), \quad P_{11}(j, t) = 1 - P_{10}(j, t).$$

These conditions allow us to solve for the transition probability $P_{10}(j, t)$ in terms of the single variable $P_{00}(j, t)$

$$P_{10}(j, t) = \frac{n_0(j, t + 1) - P_{00}(j, t)n_0(j, t)}{n_1(j, t)} \quad (1.8)$$

$$P_{01}(j, t) = 1 - P_{00}(j, t) \quad (1.9)$$

$$P_{11}(j, t) = 1 - P_{10}(j, t) \quad (1.10)$$

The remaining part is to define $P_{00}(j, t)$. We adopt the choice from [13] and start from $P_{00}(j, t = 0) = 0.5$ and choose $P_{00}(j, t) = n_0(j, t)$ for the following generations. Thus for a locus $A_{ij}(t)$ we mutate it, changing the binary value, by the probability

$$p_{ij}(t) = \begin{cases} a_i(t)P_{01}(t), & A_{ij}(t) = 0 \\ a_i(t)P_{10}(t), & A_{ij}(t) = 1 \end{cases}$$

In summary, each generation update involves the following steps.

1. For generation t calculate the fitness score $f_i(t)$ of each chromosome i and sort the population matrix $A_{ij}(t)$ row-wise according to a descending score.

2. From a defined ranking scheme $r_i(t)$ set the chromosome mutation probability to $a_i(t) = r_i(t)$ and the weighting of each row $W_i(t) = 1 - r_i(t)$.
3. Calculate the target states Eq. (1.7) and the transition probabilities using Eq. (1.8) to (1.10) and $P_{00}(j, t = 0) = 0.5$, $P_{00}(j, t > 0) = n_0(j, t > 0)$.
4. Mutate (flip) each locus $A_{ij}(t)$ by the p_{ij} given by ??.

Notice that this algorithm treats every locus as an independent gene. Thus, we do not incorporate any effects from spatial dependencies in the Kirigami pattern matrix.

1.6.0.1 Repair function

A numerical scheme for repairing the Kirigami matrix to correspond to a non-detached sheet. This was implemented in order to get candidates that are not immediately marked as a rupture. But it might be better placed in the Random walk section even though we did not make this algorithm before creating the random walk configurations. .

Part II

Simulations

Chapter 2

Summary

The work presented in this thesis covers several topics (find a better opening line?). We have created an MD simulation which enabled us to study the frictional behavior of a graphene sheet sliding on a Si substrate. In addition, we have created a numerical framework for creating Kirigami design patterns and introducing these into the friction simulations. This was used to study the effects of the out-of-plane buckling induced by a selected pair of Kirigami designs in relation to a non-cut sheet under the influence of strain. Further, we have created a dataset of various Kirigami designs for the scope of investigating the possibilities with Kirigami design. We have investigated the possibility to use machine learning on this dataset and attempted an accelerated search. Finally we look into the prospects of achieving a negative friction coefficient for a system with coupled load and stretch. In this chapter we will summarize the findings and draw some final conclusions. We will also provide some topics for further research.

2.1 Summary and conclusions

2.1.1 Design MD simulations

We have designed an MD simulation for the examination of friction for a graphene sliding on a substrate. Some of the key features for the numerical procedure were that we managed the sheet through pull blocks in the ends. We could then apply load and stretch the sheet without acting directly on the inner parts. Say something about parameter dependencies from the Pilot study.

2.1.2 Design Kirigami framework

We have designed a numerical framework for creating Kirigami designs. By defining an indexing system for the hexagonal lattice structure we were able to define the Kirigami designs as 2D matrix for numerical implementation. We digitalized two different macroscale designs, which we named the *Tetrahedron* and *Honeycomb* pattern respectively, that successfully produced out-of-plane buckling when stretched. Through a numerical framework we could create an ensemble of perturbed variations which gave approximately 135k configurations for the Tetrahedron pattern and 2025k patterns for the size of the sheet used in our study. When considering the possibility to translate the patterns this gave roughly a factor 100 more of unique perturbations. We also created a framework for creating Kirigami designs through a random walk. This was further controlled by introducing features such as bias, avoidance of existing cuts, preference to keeping a direction and procedures to repair the sheet for simulation purposes. The capabilities of the numerical framework for generating Kirigami designs was far larger than the capabilities for producing MD designs within the time constraint of this thesis. Thus we believe that this contains the possibility to benefit more extended studies and for the creation of a larger dataset.

2.1.3 Control friction using Kirigami

We have investigated the friction behavior of the non-cut sheet and a selected Tetrahedron and Honeycomb pattern under various stretch and load. The non-cut sheet did not exhibit significant out-of-plane buckling as opposed to the Tetrahedron and Honeycomb pattern. This is even when considering that the non-cut sheet had

a yield strain of 0.35 while the Tetrahedron had a lower yield strain of 0.21 and the Honeycomb a considerable larger one at 1.27 based on a stretch in vacuum. The out-of-plane buckling resulted in a significant reduction of the contact area as the sheet were stretched, towards a minimum of X for the Tetrahedron and y for the Honeycomb pattern. However, this disagreed with the asperity theory hypothesis of a decreasing friction with decreasing contact area. We found that the strain-induced buckling was initially (at low relative strain) associated with an increase in friction. Moreover, the friction-strain curve produced a non-linear behaviour which was not compatible with the approximately monotonic decreasing contact area as strain were increased. This is shown in ???. This led us to the conclusion that the contact area cannot be attributed a dominant mechanism for friction throughout the straining of the studied Kirigami sheets. In general we found a non-existing relationship between friction and load considering the uncertainties in the simulation. This is best attributed to the superlubric state of the graphene sheet on the substrate. The slope of the friction-load curves were not significantly affected by the straining of the Kirigami sheet and thus we conclude that the load effect on friction is negligible compared to the strain effects.

2.1.4 Capture trends with ML

From the dataset we found some correlations in the data. That is, we found a positive correlation between the mean friction force and strain (0.77) and porosity (0.60) and negative correlation to contact area (-0.67). These trends aligns with initial study showing that the Kirigami effects is associated to a decreasing contact area and a non-zero porosity (atoms need to be removed to produce a Kirigami sheet). By the use of machine learning we found that the trends in the data could generally be captured to a R_2 score above 0.97 for the validation data. However, the subset of configurations exhibiting the most significant effects from Kirigami and stretching provided a bigger challenge. Say something about hyperparameter search. By choosing our hyperparameters for a selected subset of challenging configurations we got a final network with a selected set score of $R^2 = 0.89$ at the worst for the Tetrahedron pattern. However these configurations was partly included in the training data and thus these numbers should be taken carefully. By creating a training set based on best candidates in an accelerated search using random walk configurations the network failed completely. This can to some extent be attributed to overfitting of the model but we believe that this is more deeply anchored in the insufficient distribution in the dataset not matching the created test set. It produced a reasonable ranking of the properties of interest within the dataset at least.

2.1.5 Accelerated search

Using the machine learning model we performed two types of accelerated search. One by going generating an extended dataset of Tetrahedron, Honeycomb and Random walk patterns and evaluating the results using the model. Another was based on the genetic algorithm where we perturbated the candidates in order to optimize for the max drop friction property. The generative accelerated search method did produce some new candidates, but it is unsure how reliable this is due to the insufficiency in the machine learning model. We concluded that the minimization of friction was not applicable. The model revealed a sensitivity to edge effect which is likely due to a model insufficiency, but this serves as an interesting topic for future research. For the genetic algorithm search it did mainly not find any new candidates and we believe that the simplicity of the genetic algorithm was not immediately suited for creating new designs with the necessary spatial correlations presented

2.1.6 Negative friction coefficient

By enforcing a coupling between load and stretch, mimicking a nanomachine attached to the sheet, we investigated the load curves arising from load the Tetrahedron and Honeycomb pattern from the pilot study. The non-linear trend observed for increasing strain carried over to this simulation setup and produce a highly non-linear friction-load curve. This demonstrated a negative friction coefficient say something about the values.

2.2 Outlook / Perspective

Topics for further studies.

- How is this behavior effected by scaling?

- How does the distribution of normal load effect the Kirigami friction behavior?
- Things to vary: load range, scan directions, adhesive forces, longer relaxation time, different potential (AIREBO)
- Investigate if the contact area is effecting the friction non-linear by turning of friction force for atoms corresponding to those that lift off from the sheet during the out-of-plane buckling.
- Investigations of commensurability effects.
- Study dependency of translation of the patterns as suggested by the ML results.
- Investigate effects from pull blocks...
- Investigate effects from the thermostat since the top and bottom edges was shown interest by the model prediction.

Appendices

Appendix A

Appendix A

Appendix A

Appendix B

Appendix B

Appendix C

Bibliography

- ¹J. Lederer, *Activation functions in artificial neural networks: a systematic overview*, 2021.
- ²P. Shankar, “A review on artificial neural networks”, **3**, 166–169 (2022).
- ³N. C. Thompson, K. Greenewald, K. Lee, and G. F. Manso, *The computational limits of deep learning*, 2022.
- ⁴Y. Bengio, “Practical recommendations for gradient-based training of deep architectures”, in *Neural networks: tricks of the trade: second edition*, edited by G. Montavon, G. B. Orr, and K.-R. Müller (Springer Berlin Heidelberg, Berlin, Heidelberg, 2012), pp. 437–478.
- ⁵L. N. Smith, *A disciplined approach to neural network hyper-parameters: part 1 – learning rate, batch size, momentum, and weight decay*, 2018.
- ⁶L. N. Smith, *Cyclical learning rates for training neural networks*, 2017.
- ⁷R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-CAM: visual explanations from deep networks via gradient-based localization”, *International Journal of Computer Vision* **128**, 336–359 (2019).
- ⁸S. Katoch, S. S. Chauhan, and V. Kumar, “A review on genetic algorithm: past, present, and future”, *Multimedia Tools and Applications* **80**, 8091–8126 (2021).
- ⁹R Jiang, K. Szeto, Y. Luo, and D. Hu, “Distributed parallel genetic algorithm with path splitting scheme for the large traveling salesman problems”, in Proceedings of conference on intelligent information processing, 16th world computer congress (2000), pp. 21–25.
- ¹⁰K. Szeto, K. Cheung, and S. Li, “Effects of dimensionality on parallel genetic algorithms”, in Proceedings of the 4th international conference on information system, analysis and synthesis, orlando, florida, usa, Vol. 2 (1998), pp. 322–325.
- ¹¹K. Y. Szeto and L. Fong, “How adaptive agents in stock market perform in the presence of random news: a genetic algorithm approach”, in Intelligent data engineering and automated learning—ideal 2000. data mining, financial engineering, and intelligent agents: second international conference shatin, nt, hong kong, china, december 13–15, 2000 proceedings 2 (Springer, 2000), pp. 505–510.
- ¹²K. L. Shiu and K. Y. Szeto, “Self-adaptive mutation only genetic algorithm: an application on the optimization of airport capacity utilization”, in Intelligent data engineering and automated learning—ideal 2008: 9th international conference daejeon, south korea, november 2–5, 2008 proceedings 9 (Springer, 2008), pp. 428–435.
- ¹³G. Wang, C. Chen, and K. Y. Szeto, “Accelerated genetic algorithms with markov chains”, in *Nature inspired cooperative strategies for optimization (nicso 2010)*, edited by J. R. González, D. A. Pelta, C. Cruz, G. Terrazas, and N. Krasnogor (Springer Berlin Heidelberg, Berlin, Heidelberg, 2010), pp. 245–254.