

# Learning Affine Capsules by Contrasting Coordinates



Dario E. Shetni Abbaszadeh

Layout: typeset by the author using L<sup>A</sup>T<sub>E</sub>X.  
Cover illustration: Unknown artist

# Learning Affine Capsules by Contrasting Coordinates

Dario E. Shehni Abbaszadeh  
10791655

Bachelor thesis  
Credits: 18 EC

Bachelor *Kunstmatige Intelligentie*



University of Amsterdam  
Faculty of Science  
Science Park 904  
1098 XH Amsterdam

*Supervisor*  
I. Sosnovik MSc

UvA-Bosch  
Delta Lab  
University of Amsterdam  
Science Park 904, Room C3.201  
1098 XH Amsterdam

25 Jun, 2021

## **Abstract**

Convolutional neural networks (CNNs) demonstrate state-of-the-art results in many image and video analysis tasks. They gradually transform raw pixels into a meaningful representation layer by layer. While CNNs process all images in the same way, Capsule Networks use an adaptive algorithm. Raw pixels are first transformed into capsules that encode different features, their poses, and probabilities. Capsules that find an agreement are then transformed into the output of the network. Such an approach demonstrates significant improvement in tasks where the model needs to adapt to significant pose variations of the object, such as rotations or scale changes. In this project, we aim to demonstrate that by using capsule networks it is possible to build an extractor of the pose from an image without any supervision.

# Contents

1	Introduction . . . . .	1
2	Neural Networks . . . . .	3
2.1	Types of Neural Networks . . . . .	3
2.2	Building Blocks . . . . .	6
2.3	Loss Functions . . . . .	10
2.4	Backpropagation . . . . .	10
3	Capsule Networks . . . . .	12
3.1	Capsules . . . . .	12
3.2	Equivariance . . . . .	13
3.3	Dynamic Routing . . . . .	14
4	Our Caps . . . . .	15
5	Experiments . . . . .	16
5.1	Fully Connected Neural Network . . . . .	16
5.2	Convolutional Neural Network . . . . .	16
5.3	Autoencoder . . . . .	16
5.4	Transforming Autoencoder . . . . .	17
5.5	Group Equivariant Capsule Network with Dynamic Routing (CIFAR-10) . . . . .	19
5.6	Our Caps . . . . .	20
6	Discussion . . . . .	22
7	Conclusion . . . . .	23

# 1 Introduction

Artificial intelligence, or AI, can be defined as the pursuit of using machines to create intelligence. However, in practice, this is not what most AI research is focused on. Current day AI is used to perform much more specific tasks ranging from image detection to finding trends in data. One method of achieving this is by providing large amounts of data and letting the machine learn for itself. This is called machine learning, an umbrella term for various techniques that allow computers to solve problems using data. Machine learning is applied across various industries for different tasks. Common problems where machine learning is applied include fraud detection, image recognition, spam detection, and product recommendations.

Another prominent field in AI is computer vision. Unlike machine learning, computer vision is not defined by the methods used, but by the problems that it tries to solve. Computer vision deals with digital images and videos. This often involves the use of machine learning. Such as using data to train neural networks for image recognition. A prominent example of computer vision is computational photography. This involves the use of digital computation instead of optical processes to process photographs. This allows smaller devices such as smartphones to have improved camera capabilities despite hardware limitations.

The problems that AI tries to solve tend to be problems that the human brain is already capable of solving. Therefore, when attempting to find solutions to these problems it can be valuable to understand how the human brain approaches these problems. The brain consists of billions of interconnected neurons. Each neuron is a cell that receives, processes, and transmits information. Each neuron is connected to other neurons. Through these connections, a neuron outputs an electrical signal. If the strength of this input signal reaches a certain threshold the neuron is 'activated' and sends out an output signal (Abraham, 2005).

Artificial neural networks are modeled after this same architecture. They consist of layers of neurons. The neurons in a layer receive various inputs from the previous layer. The signal sent by a neuron is the result of the weighted sum of these inputs. The network improves its performance by adjusting these weights. The input is received in the input layer and moves through the network in one direction. The output layer then outputs a final result. Based on this output the weights of the model are adjusted using the backpropagation algorithm. For backpropagation to work it is necessary to have a numerical representation of this accuracy. This is where loss functions are used. A loss function takes the target output and the final output of the model. It then gives a value representing the discrepancy between the two. Backpropagation works by taking the partial deriva-

tive of the loss, with respect to each weight. This partial derivative is then used to increase or decrease the respective weight accordingly (Abraham, 2005).

Traditional neural networks are built around the neuron, which produces an output based on the weighted sum of all neurons in the previous layers. However, this is not suitable for all tasks. Image processing tasks tend to favor convolutional neural networks (CNNs). Unlike fully connected layers, convolutional layers do not receive the input from all neurons in the previous layer. Instead, they rely on convolutions taking the weighted sum of a local area. This reduces computation time and leverages the fact that pixels that are close to each other are more likely to be related or to be part of the same structure. This is why CNNs are favored in computer vision tasks, as images tend to be 2 dimensional and therefore require more information to be processed. However, although CNNs have been successful in image recognition they also have limitations. One such limitation is the difficulty in understanding the spatial relationship between these features. Max-pooling is a technique that attempts to resolve this by downsampling the input, but it provides only limited spatial invariance as the pooling mechanism only samples pre-defined areas (Jaderberg et al., 2015). CNNs are therefore limited in their ability to recognize an object when the spatial relationships between its features change.

Sabour et al. (2017) demonstrated that capsule networks can achieve higher classification accuracy than conventional convolutional models on the MNIST dataset of handwritten digits. They also propose a technique called dynamic routing that can significantly improve these results. Capsule networks perform worse on more complex data that might be found on datasets such as CIFAR-10. This is because of the higher volume of information which increases the complexity of the problem. Capsule nets are still in a research and development phase and not reliable when compared to more traditional methods such as CNNs. However, the concept is promising and further research could lead to capsule networks being more commonly used in image recognition tasks

Current experiments focus on reconstructing the affine parameters of transformations applied to individual objects. While promising, this is not sufficient for dealing with complex scenes where parts of the scene can have inconsistent spatial positions relative to the canonical image. This study will evaluate whether a capsule network trained to estimate the parameters of an affine transformation of a small patch, can be applied patch-wise to a complex scene to estimate its geometry.

## 2 Neural Networks

Capsule networks are a type of neural network. To understand capsule networks it is necessary to have a general understanding of how neural networks function, and what their building blocks are. The following section will explain the following types of network architectures; fully connected neural networks, convolutional neural networks, and autoencoders. We will also discuss the process through which these networks improve themselves and how their components work.

### 2.1 Types of Neural Networks

In fully connected neural networks a neuron receives input from all neurons in the previous layer and sends its output to all neurons in the next layer (Figure 1). Due to the high number of connections this architecture is computationally intense and prone to overfitting. A fully connected architecture makes no assumptions about the structure of the data. This means that they can be used for a broad range of applications, but are typically outperformed by more specialized architectures. Fully connected networks are typically not used for image processing.

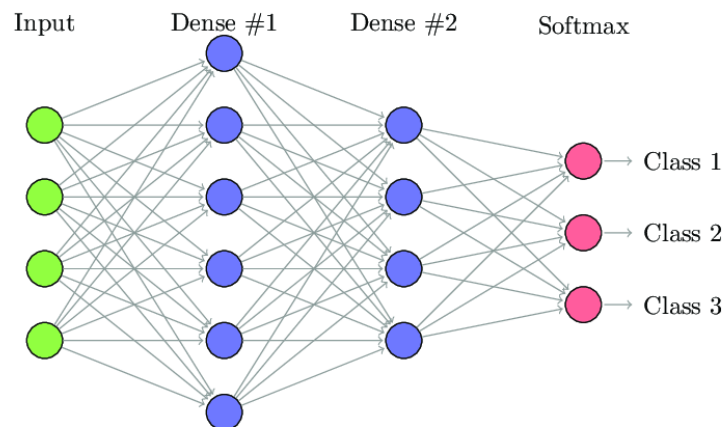


Figure 1: A fully connected neural network (Pelletier et al., 2019).

An architecture more suited for computer vision tasks is the convolutional neural network (CNN). CNNs are based on the assumption that the input is an image. Pixels in an image that are in close proximity to each other are more likely to be part of the same entity. Pixels on opposite sides of the image are less likely to be related. CNNs take advantage of this property of images by connecting each neuron only to a local region of the input. This local region is the receptive field



of the neuron. The receptive field is defined by the kernel (also called a filter) that is used to sample the input, as can be seen in figure 2.

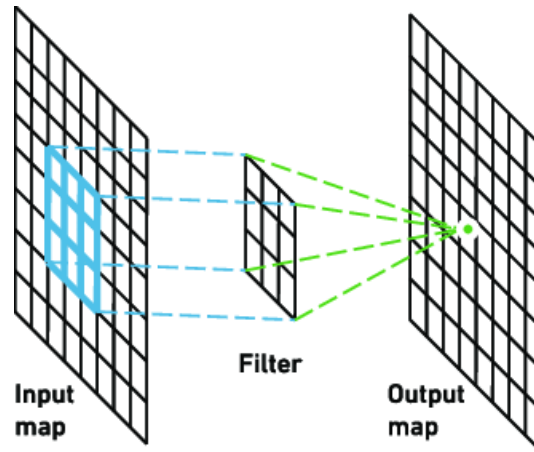


Figure 2: A representation of a kernel being applied to the input of a convolutional layer (Yakura et al., 2018).

Another difference between CNNs and regular neural networks is that the shape of the input of a CNN has three dimensions. Each pixel has a position along the width and height of the input and also a value. The dimensionality of this value is the depth of the image. A kernel has a limited receptive field with respect to the width and height of the input. But will always extend through the entire depth. The depth of the input is also called the number of channels. The kernel is applied to all channels and extracts a feature over the input, producing a 2-dimensional output map (Figure 3). As each kernel extracts a given feature, multiple kernels can be used to extract multiple features. These 2-dimensional outputs are then stacked, resulting in an output that has not only a width and height but also a depth based on the number of kernels. A convolutional layer can have multiple channels where each channel can have its own kernel. Each kernel is trained to detect a certain feature and produces a 2-dimensional output map. These outputs are then stacked. The depth of the output is defined by the number of kernels that are used.

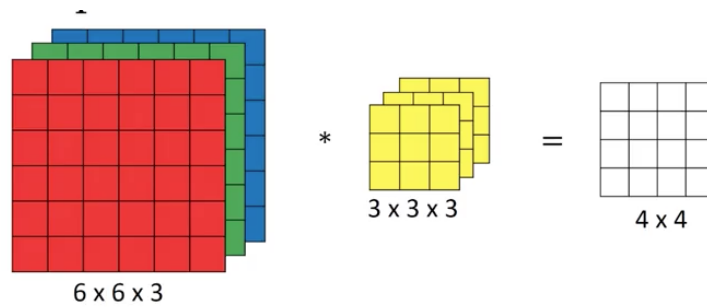


Figure 3: A kernel is applied to all channels and outputs a 2-dimensional map (Andrew Ng, 2020). The depth of the output is defined by the number of kernels that are used.

Another type of neural network architecture is an autoencoder. Autoencoders are trained to compress (encode) and decompress (decode) data. The encoder compresses the data to a latent space. This is a lower-dimensional representation of the input. The decoder then reconstructs the input from this latent space (Figure 4). Since much of the data is lost during this compression, the encoder is forced to learn how to select the features that contain the most information.

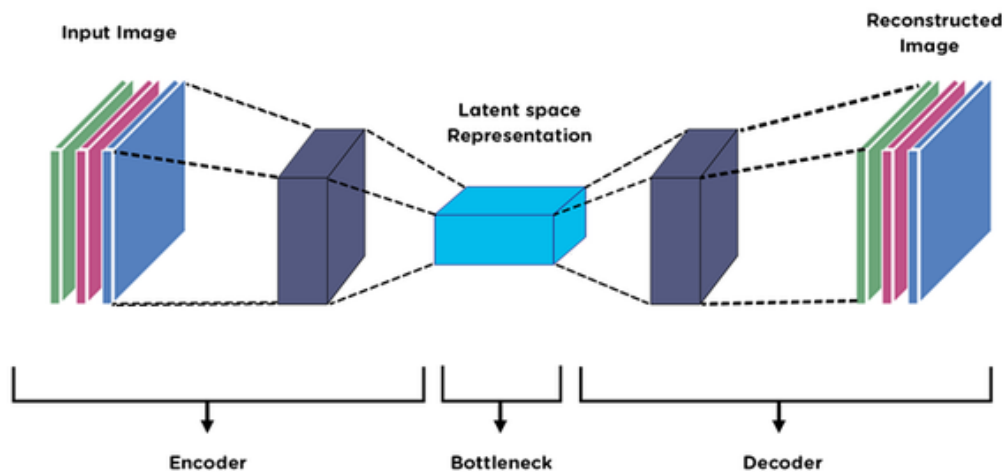


Figure 4: An autoencoder architecture using convolutional layers (Birla, 2019).

## 2.2 Building Blocks

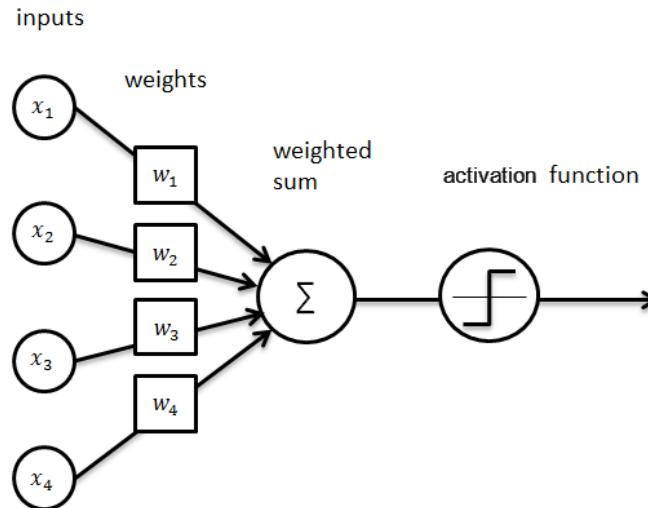


Figure 5: A single neuron in a neural net (Ognjanovski, 2019).

A neuron is the basic computational unit of a neural network. A neural net consists of a large number of processing neurons that can be fully or partially interconnected. They are organized into layers of neurons. The data then moves through the network in one direction. An individual neuron is connected to several neurons in the previous layer, from which it receives data, and several neurons in the following layer, to which it sends data (Figure 5). Each of these incoming inputs is multiplied by a weight, allowing the network to prioritize certain connections over others (Equation 1). The weights decide how strong these connections are.

$$z = b + \sum_i W_i x_i \quad (1)$$

CNNs have a specific structure called a convolutional layer. As shown in figure 2, convolutional layers perform a convolution on the input using an array of learnable weights called a kernel. The kernel is applied systematically across the input and transforms each patch into a scalar value. Similar to fully connected layers, the result is a weighted sum over the input as seen in figure 6.

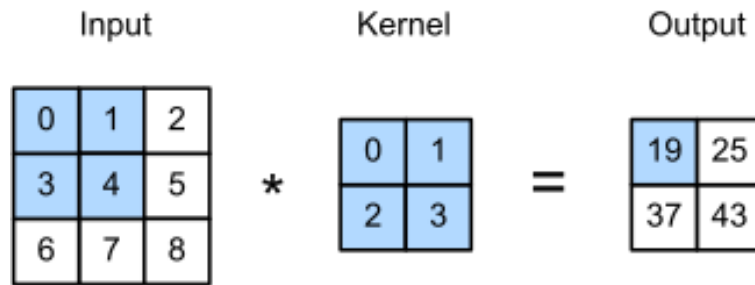


Figure 6: A kernel taking the weighted sum over a local patch of the input (Zhang et al., 2020).

This weighted sum is then transformed by a nonlinear activation function. An activation function in a neural network defines how the weighted sum of the input is transformed into the output of a neuron. One of the purposes of this activation function is to add nonlinearity to the network. Linear equations are less complex, and therefore less able to learn complex mappings. Without this nonlinearity, the neural network would behave as a single-layer network regardless of how many layers there are. The output of this activation function decides whether the neuron gets activated or not. There are various types of activation functions. The most commonly used ones are the rectified linear unit (ReLU) function, the sigmoid function, and the softmax function.

The ReLU function is commonly used in the hidden layers. It is calculated by taking  $\text{MAX}(0, x)$ . This means that it will not change the input unless it is below 0, in which case ReLU will output 0 (Figure 7). One advantage of the ReLU function is that it is relatively fast to compute compared to other activation functions. A disadvantage of ReLU is that it has a gradient of zero for inputs below or equal to 0. This causes multiplication by zero which results in gradient descent not altering the weights. This is known as the dying ReLU problem (Lu et al., 2020).

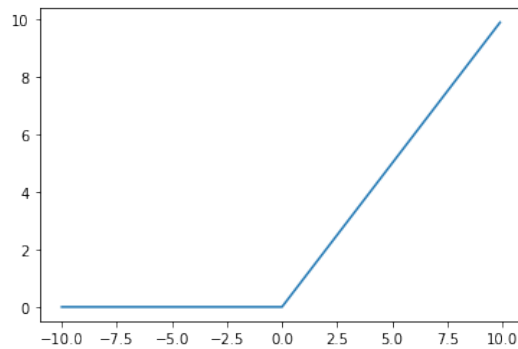


Figure 7: The ReLU function

The sigmoid function, also known as the logistic function, maps a number to a value between 0 and 1 (Equation 2). This makes the sigmoid function useful for two-class classification, such as deciding whether a neuron should be activated or not. It can also be used for computing a probability. Like ReLU it is commonly used in the hidden layers. A disadvantage of the sigmoid is that an input near 0 or 1 will have a low gradient approaching zero. Consecutive multiplication of these gradients will therefore quickly result in a value approaching zero.

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

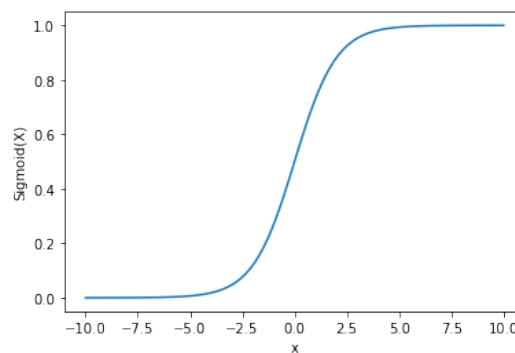


Figure 8: The sigmoid function

Softmax is used to compute probability distributions from a vector of scalars. It converts each of the scalars to a value between 0 and 1, with the sum of the vector being 1. This can be useful when trying to predict multiple classes, as softmax will return the relative probabilities of all the classes. The difference with the sigmoid function is that softmax takes into account the probabilities of the

other classes. This is why it is used in multiclass classification where there are more than two classes. Softmax is most commonly used in the output layer of the network.

$$f(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}} \quad (3)$$

A common problem when training neural networks is overfitting. This occurs when a model fits too closely to the training data. Larger neural networks are more likely to overfit. Especially when the training dataset is relatively small. One method that helps to reduce overfitting is the use of dropout layers. Dropout layers can be used to randomly discard some of the outputs during the training. By randomly disabling some percentage of neurons, the network learns not to rely too much on any specific connection. Dropout layers are typically used after fully connected layers in the network (Srivastava et al., 2014).

CNNs can use a pooling layer to downsample the input data. As it reduces the number of parameters, it makes the model more robust to variations in the positioning of features. Pooling is applied separately to each depth dimension of an input. The most commonly used operation for pooling is the MAX operation, which takes the local maximum, using a  $2 \times 2$  filter and a stride of 2 (figure 17). The stride parameter indicates the step size when sliding the kernel across the input map. This reduces the number of parameters by 75%, thus also reducing the computational cost. Pooling decreases the width and height of the input, whilst retaining the depth. The kernel size for the pooling layer tends to be small, as to not lose too much of the information (Gilon, 2021).

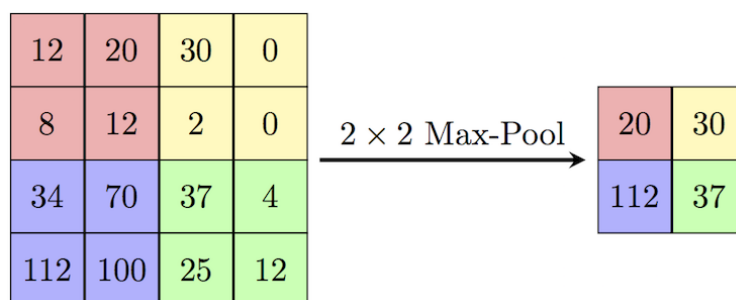


Figure 9:  $2 \times 2$  max-pooling with a stride of 2 on a  $4 \times 4$  input map (Andrew Ng, 2020).

Batch normalization is a technique proposed by Ioffe and Szegedy (2015), that can be used to reduce the number of iterations required for training. As the parameters of the network change, the distribution of the input also changes for each

layer of the network. This is referred to as internal covariate shift. Batch normalization normalizes the input of each layer so that the distribution remains equal. This speeds up the network as it requires less time to initialize the parameters.

### 2.3 Loss Functions

A loss function quantifies the discrepancy between the output and the desired output. A lower value for the loss represents an output that is close to the desired output. Different type of problems require their own loss function to give a meaningful representation of the performance of the network. A loss function that can be used to compare images is the mean squared error (MSE) (figure 4. MSE squares the error for each predicted value. This has the effect of penalizing a single large error more severely than multiple smaller errors. The MSE is then defined as the mean over all squared errors.

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \hat{\mathbf{y}}_i)^2 \quad (4)$$

Cross-entropy (CSE) loss is often used for classification problems. CSE is used after softmax, which outputs a prediction vector of probabilities between 0 and 1. It measures the entropy by comparing the distribution of the prediction vector with the distribution of the desired output. The probability for each class is compared to the desired output, which is either 0 or 1. Figure 5 shows how the prediction for a class is multiplied by the logarithm of the desired output. The CSE loss is then defined as the mean error of all class predictions (Gilon, 2021).

$$-\frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i \cdot \log(\hat{\mathbf{y}}_i)) \quad (5)$$

### 2.4 Backpropagation

Backpropagation adjusts the weights in the network to minimize the loss. All weights in a neural network are initialized randomly at first. When the network outputs a prediction, a loss function is used to calculate the loss. The loss for a given input is a function of the weights of the network. The relationship between the loss and the weights is visualized in figure 10. This means that to minimize the loss it is necessary to find the global minimum of the loss function. Since there is no straightforward way to find the global minimum of a complex function such as a neural network, backpropagation uses stochastic gradient descent (SGD) to find a local minimum by making incremental changes to the weights based on the gradient.

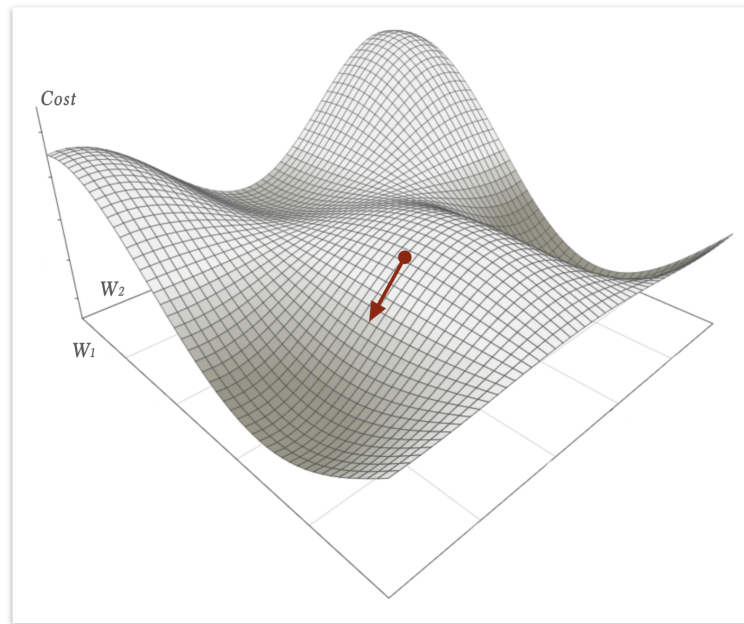


Figure 10: The loss function plotted for a network with two weights (Hill, 2021).

The gradient of the loss returns a vector pointing in the direction where the loss increases most rapidly. By taking the negative gradient it is possible to calculate the direction where the loss decreases the most. The network can be seen as a complex function where each weight is a parameter in that function. Given a loss function  $\mathcal{E}$ , the partial derivative can be computed for a single weight or bias as seen in equation 6. This allows the network to optimize each weight individually.

$$\frac{\partial \mathbf{E}}{\partial w} \quad \text{and} \quad \frac{\partial \mathbf{E}}{\partial b} \quad (6)$$

The partial derivative is used to describe the local curvature. This allows SGD to understand which weights to increase or decrease. A learning rate parameter  $\alpha$  decides the increment by which the weights are then adjusted. A lower value for  $\alpha$  means that the weights can be adjusted more precisely. It will however take more iterations to find a minimum. Another downside is that SGD can find itself in a local minimum from which it can not escape, because the step size is too small. This can be avoided by choosing a higher value for  $\alpha$ . This will help SGD to explore a larger area of the problem space, but has the downside of potentially skipping over a local minimum.



## 3 Capsule Networks

### 3.1 Capsules

A key feature of capsule networks is that they are equivariant. Information about the location and pose of an object is preserved throughout the network. A capsule network consists of layers of capsules. A capsule is a group of neurons encapsulating the properties of a single entity. Each of these capsules is its own neural network. Capsules correspond to an entity or sub-entity in the input image. Capsules in lower layers will detect simple features such as edges or gradients. Higher-level capsules will combine these simple features into more complex features.

Not only can a capsule detect the presence of such an entity, it is also able to detect properties of that entity. Each capsule in the layer outputs a vector that encodes this information and passes it on to the next layer. A squash function is used to ensure that the vector length is between 0 and 1. The capsules in this next layer also correspond to an entity. These entities are typically composite entities of those in the previous layer. Each of these capsules will then integrate the vectors of the previous layer and come up with their own output vector encoding the information for their corresponding entity.

In a typical neural network, the outputs of a layer would be multiplied with a weight matrix to determine the inputs for the subsequent layer. In a capsule network, however, the capsules can decide where to direct their output. A capsule will route its output to the appropriate capsule in the following layer depending on its output. This is accomplished through an iterative procedure via the routing mechanism.

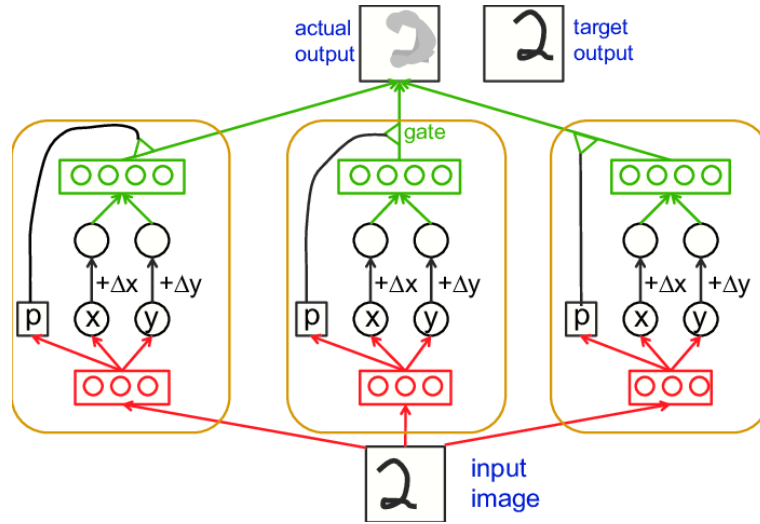


Figure 11: Three capsules that model translation (Hinton et al., 2011).

An implementation of a capsule layer can be seen in figure 11. The transforming autoencoder proposed by Hinton et al. (2011) contains capsules that extract the probability  $p$  that a given feature is present and a pose  $\vec{x}\vec{y}$  from an input image. Each capsule is trained to detect a certain feature. The capsules take an image and a desired translation  $\Delta\vec{x}\vec{y}$ . This translation is added to the outputted pose vector  $\vec{x}\vec{y}$  of the encoder. A decoder then reconstructs the image and the difference between the output image and the translated image is backpropagated through the network. By adding the translation to the output  $\vec{x}\vec{y}$ , from which the image is reconstructed, the network is trained such that the output  $\vec{x}\vec{y}$  must represent the pose of the input image.

### 3.2 Equivariance

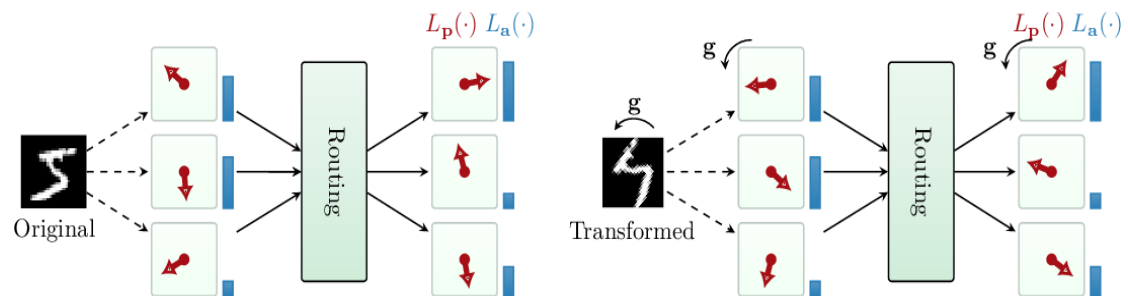


Figure 12: Equivariance of the pose vectors (Lenssen et al., 2018).

Equivariance is a form of symmetry. A function is said to be equivariant with respect to a given transformation if applying that transformation and then computing the function produces the same result as computing the function and then applying the transformation (Equation 7). This means that if you transform the input. The output will be transformed accordingly. Figure 12 visualizes how rotating the original images causes the pose vectors to be rotated in the exact same way. Thus showing that it is equivariant with respect to rotation.

$$f(L(x)) = L'(f(x)), \quad \forall x \quad (7)$$

### 3.3 Dynamic Routing

Dynamic routing is done through a method called routing by agreement. The lower-level capsules all output a prediction vector. This vector contains the pose of the feature that the capsule is trained to detect, and also the probability that the feature is present in the input. This vector is only routed to the higher-level capsule if the other low-level capsules in that layer agree with each other. Agreement is measured by comparing the output vectors of the capsules in a given layer. If the capsules agree, then the output is routed to the next capsule. The agreement between the vectors is dependent on which capsule in the following layer the information is being routed to. Capsules can be in agreement with each other when routing to a given capsule in the next layer but disagree when being routed to another capsule. This allows for different combinations of capsules to be routed to each capsule in the next layer. Over each iteration of the routing algorithm, the network will increasingly settle on a routing path for the capsules.

## 4 Our Caps

In this research, we aim to demonstrate that by using capsule networks it is possible to build an extractor of the pose from an image without any supervision. This was done by extracting the translation from an image, which is simpler to implement compared to affine pose extraction. This is equivalent because an affine transformation is a composition of one-parameter groups. Each of those groups can then be described as a translation by reparameterizing the coordinate system, e.g. a rotation in cartesian coordinates is the same as translating in polar coordinates (Hel-Or & Teo, 1998).

The architecture will be based on a transforming autoencoder as described by (Hinton et al., 2011). To extract the pose, the model consisted of a primary capsule layer with multiple capsules that each detect a different feature. A randomly translated version of the input image was created and passed to each of these capsules, together with the original input image. The capsules predicted a probability  $p$  and a vector  $\vec{x}\vec{y}$  for each image.  $p$  is the probability that the feature corresponding to the capsule was present. The vector  $\vec{x}\vec{y}$  is the pose that represents the position of the feature with respect to its canonical view. The canonical view being the feature, as learned by a specific capsule, where  $\vec{x}\vec{y}$  is equal to zero. Capsules are trained to detect only one specific feature. Therefore, the learned feature can be reconstructed with any given pose from just the translation vector  $\vec{x}\vec{y}$ , as the feature itself is inherent to the capsule. The reconstructed image was then multiplied by the probability that the feature is present. The output of each capsule is then summed up, resulting in the final reconstructed image.

Each capsule also received a translated version of the input image and the parameters of its translation  $\Delta\vec{x}\vec{y}$ . The translation vector  $\Delta\vec{x}\vec{y}$  was added to the pose  $\vec{x}\vec{y}$  extracted from the original image. The translated image was then reconstructed by a decoder from only the translation vector  $\Delta\vec{x}\vec{y}$  and the pose extracted from the original image  $\vec{x}\vec{y}$ .

The loss was calculated by summing up the mean squared error of both the reconstruction of the image and the reconstruction of the translated image. Additionally, the pose vector extracted from the original images was subtracted from the pose vector extracted from the translated image. The difference of these vectors was expected to be equal to the translation vector  $\Delta\vec{x}\vec{y}$ . This vector difference was then squared and multiplied by the probability  $p$ , as to not penalize predictions for features that are not present.

## 5 Experiments

Capsule networks use capsules to encode some geometric transformation. The goal of the experiments is to show that it is possible to train the network to detect objects from different viewpoints by training the network to extract the pose. This will be done in an unsupervised fashion as described in section 3.3. Multiple experiments will be conducted. Initial experiments will attempt to correctly classify the images. Then the images will be compressed to a lower-dimensional space and reconstructed. Finally, the pose will be extracted from the images.

All models were made in PyTorch using Python 3.6. The models were trained on the MNIST dataset. A labeled 28x28 grayscale image dataset of handwritten images. The labels are the digits from 0 to 9. All experiments use a batch size of 64. The model with dynamic routing was trained using the CIFAR-10 dataset. This is a 32x32 RGB color image dataset with 10 classes.

### 5.1 Fully Connected Neural Network

First, a fully connected neural network was created to predict the class of the handwritten digits. The network was built with five fully connected layers, each followed by a ReLU activation function. The final layer outputted a 10-dimensional vector. One prediction for each class. A softmax activation was used to normalize the output. The output was then evaluated with a cross-entropy loss function. This model was trained for 8 epochs and reached an accuracy of 91,9%.

### 5.2 Convolutional Neural Network

A convolutional neural network was created with 3 convolutional layers to predict classes of MNIST images. Each convolutional layer was followed by a batch normalization, a ReLU activation, and max-pooling. The network ran for 8 epochs and reached 98,1% accuracy.

### 5.3 Autoencoder

For this experiment, an autoencoder was made for prediction and reconstruction on MNIST. the model consisted of an encoder and a decoder. The encoder outputted a prediction of the class. The decoder then used this prediction to reconstruct the image. The model used three convolutional layers in the encoder, and a linear layer to transform the data to a 10-dimensional vector containing a probability for each class. One scalar for the probability of each class. Another linear layer then upsampled the data from those 10 values. The result was reshaped to 2

dimensions. A decoder with three convolutional layers was used to reconstruct the input image. The model used two loss functions. Cross entropy loss was used on the output of the predictor for the classification. Mean squared error was used to evaluate the reconstructed image. The loss was then defined as the sum of the output of these functions. The reconstruction loss was multiplied with a factor  $\alpha$ . The model ran for 20 epochs. The highest accuracy achieved for classification was 99,2% with a loss of 0.021 for the reconstruction.

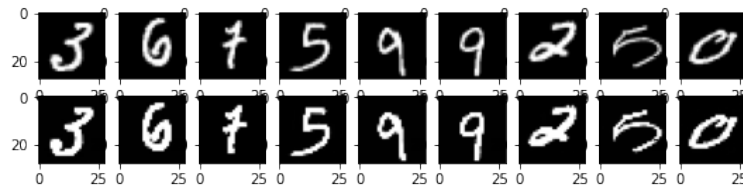


Figure 14: A plot showing the target images above and the reconstructions below.

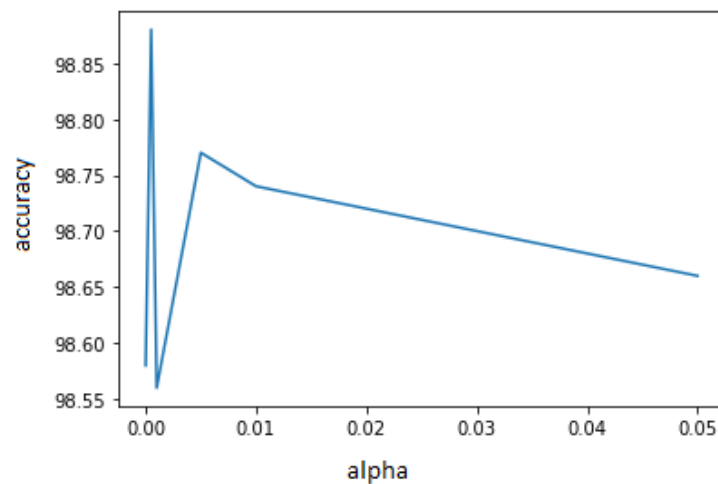


Figure 15: The accuracy plotted against different values for  $\alpha$

## 5.4 Transforming Autoencoder

A transforming autoencoder (TAE) was trained for reconstruction on the MNIST dataset. For this experiment, the target reconstructions were randomly translated. The parameters for this random translation (x and y) were also passed on to the network. The TAE consists of an encoder, which outputs 3 values; A probability p and two parameters for the translation, x and y. The translation on the target output was then also applied to the x and y output of the encoder. The image

was then reconstructed solely from these translated x and y values. Each capsule was not trained to reconstruct the whole image. Instead, each capsule was trained to reconstruct a specific feature with a specific translation. Each capsule also outputted a probability representing how likely it is that that given feature is present. The reconstruction of each capsule was then multiplied with its respective probability. The final output of the model was obtained by summing up the output of each capsule and can be seen in below in figure 16.

The network was trained for 15 epochs using 10 capsules. The final loss was 0.75106

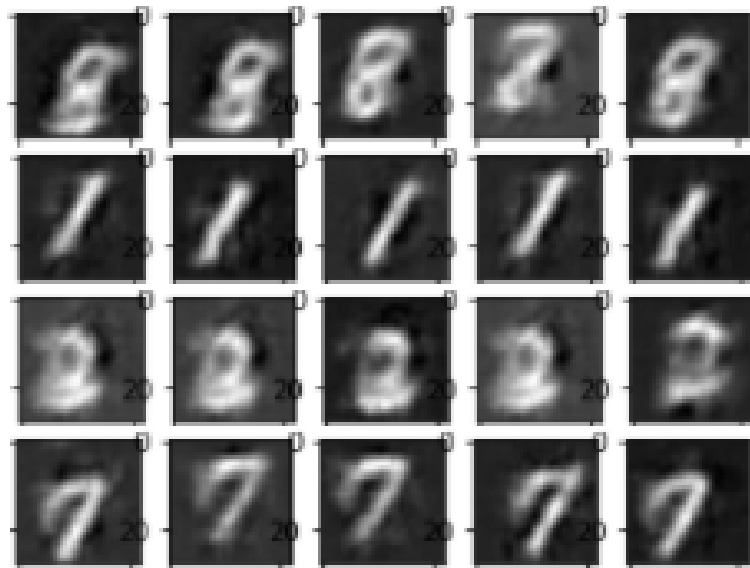


Figure 16: Reconstructed images for classes 8, 1, 2, and 7 respectively from top to bottom.

The results showed that most reconstructions are recognizable as their target class. The third row in figure 16 appeared less recognizable and exhibited features from both a 2 and a 3. This is likely caused by the capsule encoding the features for label 3, being activated by similar features in label 2.

### 5.5 Group Equivariant Capsule Network with Dynamic Routing (CIFAR-10)

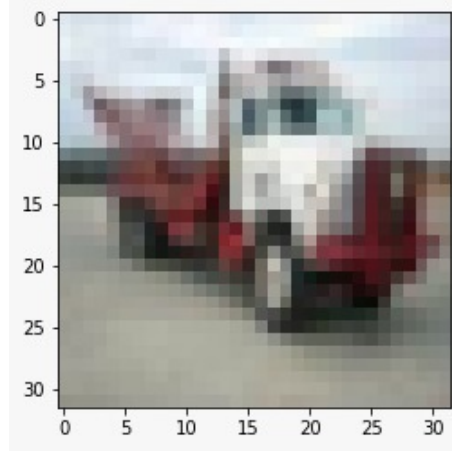


Figure 17: An image in the CIFAR-10 dataset.

For this experiment an existing model was used as described in Lenssen et al. (2018) using dynamic routing. In contrast to the method from Sabour et al. (2017), softmax was not used over the output capsule dimension but the sigmoid function for each weight individually. The sigmoid function makes it possible for the network to route information to more than one output capsule as well as to no output capsule at all.

Gradients are computed in the x and y directions by convolving the image with the kernels shown in figure 18. The resulting gradients,  $G_x$  and  $G_y$ , together form the vector  $\vec{G}_{xy}$  of the direction of the gradient (Figure 19). Taking the norm of each vector results in  $G_{xy}$ , as shown in figure 21.

-1	0	+1
-2	0	+2
-1	0	+1

$G_x$

+1	+2	+1
0	0	0
-1	-2	-1

$G_y$

Figure 18: The kernels used for computing the gradients in the x and y directions.



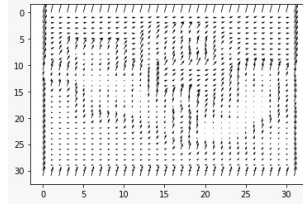


Figure 19: Vector view of the combined x and y gradients.

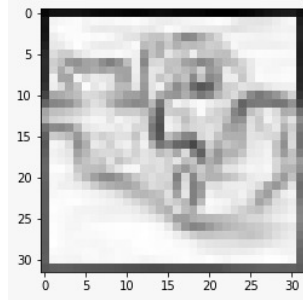


Figure 20: Visualized gradient length.

The implementation used deprecated software libraries and could not be run on newer versions of PyTorch. An attempt was made to get the model to work but the implementation did not provide enough flexibility. The experiment yielded no results.

## 5.6 Our Caps

For the final experiment the capsule network as described in the our caps section was implemented. The network has 16 capsules. Each capsule has an encoder with two convolutional and one linear with three output nodes predicting  $p$ ,  $x$ , and  $y$ . The decoder starts with a linear layer that upsamples  $x$  and  $y$ , and then reshapes the output to an image. Two more convolutional layers reconstruct the image.

Two loss functions were used. Mean squared error was used on the reconstructed images and a variation thereof for the extracted pose, where the errors were weighted by the probability. The final loss was the sum of these losses. The model used 16 capsules and ran for 20 epochs. The final loss was 0,050.

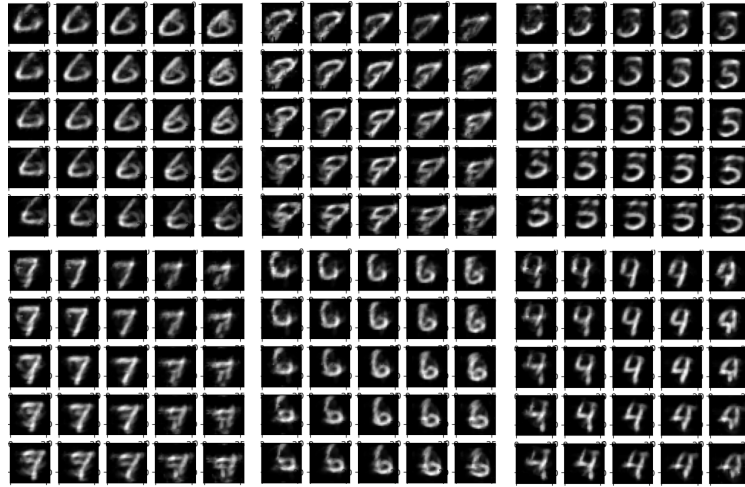


Figure 21: Capsule reconstructions MNIST images with varying degrees of translation. The center images are reconstructed from the original.

Similar to the transforming autoencoder in the previous experiment, the reconstructions are recognizable. However, some images show overlap with other features. This is demonstrated in the reconstruction of the 3, which sometimes overlaps with the features of the label 5.

## 6 Discussion

The initial experiments show that high accuracy can be achieved on MNIST classification with a relatively low number of layers. Only 8 epochs were required to reach this accuracy. The results also show that CNNs are more suited to solve this problem. This is expected as CNNs make the assumption that the input is an image by taking into account feature proximity.

Image reconstruction downsamples the data to a lower-dimensional representation. This makes it more difficult than classification. Much of the information is lost, forcing the network to select for the most important features. Successful reconstruction from lower-dimension latent space shows that much of the information about the images can be retained in the network architecture. This is useful for creating capsules because in capsule networks, images are reconstructed from a given pose. This means that feature information has to be inherent to the capsule itself and can not be derived from the input. This is further supported by the transforming autoencoder experiment. Where images are successfully reconstructed from the pose alone. Successful extraction of the translation parameters is also promising, as affine transformations can be decomposed and reparameterized as translations.

The group equivariant experiment on CIFAR-10 did not produce any results. However, this did not stand in the way of answering our research question as the final experiment did provide meaningful results. Unsupervised pose extraction was the main purpose of the experiment. Extraction of the translation parameters showed to be successful and the findings of Hel-Or and Teo (1998) can be used to modify this implementation for affine parameters in future research. Furthermore, Jaderberg et al. (2015) also proposed an improvement on the capsule architecture by creating a deep capsule network architecture and a 3D convolution-based dynamic routing algorithm that aims to improve the performance of the network. This helps improve network performance on complex images by adding additional layers. Current implementations do not scale well when adding additional layers. Future experiments can build on this research to implement a deep architecture for pose extraction on more complex scenes.

## 7 Conclusion

We aimed to demonstrate that it is possible to equip capsules with a notion of geometric transformation with no supervision at all. We did this by successfully extracting the translation parameters from the input images. Results for the final experiment demonstrated that the pose can indeed be extracted from an object in an unsupervised fashion. We have also shown that our model is equivariant with respect to translation and that this can be modified to work for all affine transformations. However, we have not successfully demonstrated that this can be generalized to more complex scenes as we were not able to implement multi-layered architecture for our capsule network with dynamic routing, which is required to train our model on more complex images.

## References

- Abraham, A. (2005). Artificial neural networks. *Handbook of measuring system design*.
- Andrew Ng, Y. B. M., Kian Katanforoosh. (2020). Deep learning specialization. <https://www.coursera.org/specializations/deep-learning>
- Birla, D. (2019). *Basics of autoencoders*. Retrieved June 10, 2021, from <https://medium.com/@birla.deepak26/autoencoders-76bb49ae6a8f>
- Gilon, Y. (2021). Cs231n convolutional neural networks for visual recognition. <http://cs231n.stanford.edu/>
- Hel-Or, Y., & Teo, P. C. (1998). Canonical decomposition of steerable functions. *Journal of Mathematical Imaging and Vision*, 9(1), 83–95. <https://doi.org/10.1023/A:1008274211102>
- Hill, T. (2021). *Part 2: Gradient descent and backpropagation*. Retrieved 2018, from <https://machinelearning.tobiashill.se/part-2-gradient-descent-and-backpropagation/>
- Hinton, G. E., Krizhevsky, A., & Wang, S. D. (2011). Transforming auto-encoders. In T. Honkela, W. Duch, M. Girolami, & S. Kaski (Eds.), *Artificial neural networks and machine learning – icann 2011* (pp. 44–51). Springer Berlin Heidelberg.
- Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In F. Bach & D. Blei (Eds.), *Proceedings of the 32nd international conference on machine learning* (pp. 448–456). PMLR. <http://proceedings.mlr.press/v37/ioffe15.html>
- Jaderberg, M., Simonyan, K., Zisserman, A., & Kavukcuoglu, K. (2015). Spatial transformer networks. *CoRR*, abs/1506.02025. <http://arxiv.org/abs/1506.02025>
- Lenssen, J. E., Fey, M., & Libuschewski, P. (2018). Group equivariant capsule networks. *CoRR*, abs/1806.05086. <http://arxiv.org/abs/1806.05086>
- Lu, L., Shin, Y., Su, Y., & Karniadakis, G. E. (2020). Dying relu and initialization: Theory and numerical examples. *Communications in Computational Physics*, 28(5), 1671–1706. <https://doi.org/10.4208/cicp.oa-2020-0165>
- Ognjanovski, G. (2019). *Everything you need to know about neural networks and backpropagation — machine learning easy and fun*. Retrieved June 10, 2021, from <https://towardsdatascience.com/everything-you-need-to-know-about-neural-networks-and-backpropagation-machine-learning-made-easy-e5285bc2be3a>
- Pelletier, C., Webb, G., & Petitjean, F. (2019). Temporal convolutional neural network for the classification of satellite image time series. *Remote Sensing*, 11, 523. <https://doi.org/10.3390/rs11050523>

- Sabour, S., Frosst, N., & Hinton, G. E. (2017). Dynamic routing between capsules. *CoRR*, *abs/1710.09829*. <http://arxiv.org/abs/1710.09829>
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.*, *15*(1), 1929–1958.
- Yakura, H., Shinozaki, S., Nishimura, R., Oyama, Y., & Sakuma, J. (2018). Malware analysis of imaged binary samples by convolutional neural network with attention mechanism, 127–134. <https://doi.org/10.1145/3176258.3176335>
- Zhang, A., Lipton, Z. C., Li, M., & Smola, A. J. (2020). *Dive into deep learning* [<https://d2l.ai>].