

## ***Abstract***

*The defects of a product or materials in the manufacturing industries are one of the problems that cause headaches and costs for manufacturers as well as companies around the world so far. It is costly in recalling defective products, costly in refurbishing or remanufacturing them, and also in human resources in monitoring and evaluating products before they are put on the market. Because of that, the application of research in artificial intelligence as well as computer vision for defect detection in the manufacturing industries was started a few years ago, with the aim of reducing production, defective products as well as minimizing costs such as recalling products, remaking them, or consuming human resources for product quality assurance. More than that, it is aimed at helping the product have a better and perfect level when it reaches the user. Therefore, in this thesis, I have researched and applied the research results of the forerunners in the field of Deep learning as well as computer vision, although it is not perfect, it also achieved a possible result to contribute in the research and application of artificial intelligence, it plays a small part in the field of IoTs for the manufacturing industry. Specifically, I conducted research on the topic of defect detection on steel surfaces, it was taken from Kaggle's contest titled "Severstal: Steel Defect Detection", which was a quite big contest held in 2019. In this thesis, I went to solve the requirements of detecting and classifying defects on the steel surface through the image data set provided by the Russian corporation Severstal. More specifically, I have implemented mainly semantic segmentation and accompanied by some other techniques in deep learning to solve this problem.*

## Table of Content:

Abstract	
Table of Content	
i.Introduction	3
ii.Computer Vision	7
iii.Background Knowledge	10
iv.Implementation	70
v.Result	82
vi.Future Work Scope	94
vii.Conclusion	95

## **I. Introduction:**

In the era of industry 4.0, with technological achievements that the world has achieved through research on artificial intelligence (AI) from the last century, accompanied by well-developed infrastructure, machinery and equipment applying cutting-edge technology, which help IoT systems are gradually growing stronger and more popular, they have such a high applicability in most aspects of life. Currently, there is a rapid development of smart health care systems, smart agriculture, smart cities, smart homes, smart industries, and so on... all these "smart things" are built on top of the distributed systems and physical infrastructures of each field. If the distributed system, which is considered as the "backbone", is the main architecture building on the internet of things (IOTs) system, help the IOTs system operates smoothly, then artificial intelligence acts as a "cognitive area of the brain" in the IOTs system. It helps the system to collect and process information to make decisions, more specifically here is computer vision of deep learning, which is one of the indispensable components in the IOTs system. Using cameras in tandem with sensors such as: temperature, light, sound, smoke, or GPS in most popular IOTs, they are one of the important devices used to capture collect data about the system's cloud and use trained deep learning models to make decisions on these data, or use these data to train and improve the "brain" of the system. IoTs are getting better and better for specific purposes of the system.

With the birth and development since the 1950s of the last century, computer vision, which, is the one of the sub-disciplines of AI, has made progress as well as achieving outstanding achievements in the present. Its applications play a very important role in these systems. Looking at

deeply, Smart Homes often use them in facial recognition, identification tasks such as: people, objects, pets, or even identifying and alerting thieves to protect you and your family. Smart cities with applications for developing autonomous vehicles, or traffic monitoring camera systems, smart parking lots that let you know how many spaces are available, or smart airports that help for identification identity. Smart heaths with supports for doctors in diagnosing diseases by predicting phenomenon in images from x-ray, endoscopy, ultrasound or tomography. Smart agriculture with monitoring and diagnosing the health of crops and livestock. Smart industries with camera systems for production or monitoring of the production process of products to ensure quality, to reduce and to avoid costly human resources in tasks that can replace by using machines. Specifically, defect detection on steel surfaces based on computer vision is one of the very useful applications in the industry of manufacturing steel materials or products made from steel, it helps enterprise can identify product defects and evaluate product quality, or help reduce costs and improve product quality.

In the next parts of the thesis, I will briefly cover computer vision and its common applications in technology as well as in life and production. Then I will go to write about the knowledge in deep learning is also such as artificial intelligence is used for this project, in detail, I will clarify the basics of Artificial Neural Network (ANN) and important components of the network such as layers and types of layer, activation functions, Convolutional Neural Network (CNN) and some state-of-the-art architectures of CNN like ResNet being built from Residual Blocks, EfficientNet with applying 3D compound scaling method, U-Net with U-shaped structure including encoding part and decoding part for image

segmentation work; or about optimization methods that are commonly used in model training work.

Continuously, it is the implementation of the solution that I came up with as well as the obtaining results, particularly, I went to find a solution to the defect detection problem on steel surface through the requirements of Kaggle's contest with the accompanying dataset of Severstal corporation by applying the state of the arts models with a combination of them to give positive results. It is the method of semantic segmentation on the image dataset by a combination of U-Net as a master architecture and versions of ResNet as the backbone architecture for U-Net, or another combination of U-Net and versions of Efficient, is respectively the overall architecture and the corresponding backbone in the model. Furthermore, given the data set was provided with an imbalance in defect type, so I went to make them a balancing in the data type by augmentation method and using K-fold method to evenly divide the types of data and the amount of data for the training data set and validating data set to help the learning model give an objective result and improve the prediction accuracy at the testing data set. In addition, applied an optimization method as Adam with default values for parameters to help the model learn faster and more accurately. Using 5 folds and going through over 20 - 30 epochs for each fold in the training model, I got good results with accuracy ranging from 90% - 94% with validating data set and around 68% - 71% with testing data set through model evaluation. The following step was to make an evaluation for the model's performance with techniques of Accuracy Classification Score and Multi-label Confusion Matrix. Accuracy Classification Score method gave results ranging from 72% - 82% in total prediction of defect

types, and results from 90% - 95% in predicting each defect type. Then, the Confusion Matrix method also gave quite optimistic results of classifying labels between ground-truth masks and predicted masks of test data set. At final step, I went into re-checking the results with the naked eye by rendering the original image with ground-truth masks and the original image supplying predictive masks based on the threshold number of pixels that the model makes predictions for each defect type. To understand the solution in more detail, I will cover them in the implementation section.

Final section is conclusion with an overview of the thesis as well as the implementation and results achieved.

## **II.Computer Vision (CV)**

Computer Vision is an Artificial Intelligence (AI) area which deals with computational approaches for assisting computers in understanding and interpreting the content of digital images [38]. As a result, computer vision seeks to enable computers to view and comprehend visual data received through cameras or sensors to identify, and understand the visual world automatically by emulating vision of human being. In particular, the goal of computer vision is to artificially mimic human vision by allowing computers to comprehend visual stimuli meaningfully [39]. As a result, it is also known as machine perception. CV remains one of the most difficult subjects in computer science because of owing to the vast complexity of the varying physical world. In fact, human sight is the result of a lifetime of learning with context to train how to detect certain things or human faces or individuals in visual settings.

CV is not a newly discovered technology nowadays because the earliest attempts with it began in the 1950s, and it was used to analyze typewritten and handwritten text after that [39]. At that period of time, computer vision analysis processes were very straightforward, but they needed a significant amount of work from human operators who had to manually input data samples for analysis. So, it was difficult to supply a large amount of data manually. Furthermore, the computer capacity was insufficient making the error margin for this analysis was quite large. There are numerous supports for CV nowadays. Cloud computing, when combined with powerful algorithms, has the potential to help us tackle even the most difficult issues. The combination of new hardware, clever algorithms, and the massive amount of publicly available visual data that we generate on a daily basis is responsible for propelling computer

vision technology ahead. Modern artificial vision technology uses machine learning and deep learning approaches to train robots to recognize objects, faces, or people in visual settings based on a large pool of picture and video data. As a result, computer vision systems employ image processing methods to enable computers to locate, identify, and evaluate objects and their surroundings using data from a camera.

Moreover, CV systems are trained to examine products, monitor infrastructure, or a manufacturing asset in order to detect faults or difficulties in thousands of products or processes in real time. It can swiftly outperform human capabilities because to its speed, continuity, objectivity, correctness, and scalability. Applications of CV are employed in numerous industries, including security and medical imaging, as well as manufacturing, automotive, agricultural, construction, transportation, smart cities, and many more. More use cases become conceivable as technology progresses and becomes more versatile and scalable. In addition, the most recent deep learning models outperform humans in real-world image recognition duties such as facial identification, image classification, object detection, and image segmentation.

**Image classification:** Accepting input is an image and outputs the labels for duty of classification along with several measures (loss, accuracy, probability, etc.). For example, a picture contains a duck may be classed as a label of "duck" with some probability, whereas an image contains a chicken being made a define as a label of "chicken" with several probabilities [43].

**Object Localization:** Helping detect the present of an object in an image and marks it as a bounding box, which receives an input of image and returns the bounding box location in the form of (width, height, and location) [40].

**Object Detection:** Doing both classification and localization which receive an input of image and generates one or several bounding boxes, each box contains the attached label of class. These methods are powerful enough to handle localization and classification for multi-class and objects containing numerous occurrences as well [43]. However, the bounding boxes in object detection are always a form of rectangle. As a result, if the item contains the curvature part, it does not aid in establishing its shape, and object detection does not have an ability to estimate some measurements in high accuracy such as an object's area or perimeter from an image.

**Image Segmentation:** Being an extra extension in detection of object. It helps us designate the existence of an object in the image using pixel-wise masks built for every single object. The method is more detailed than producing of bounding box since it allows us to determine the geometry of every single object appearing on the image. This granularity is useful in a variety of applications, including medical image processing, satellite imaging, and so forth [43]. Many pictures segmentation approaches have lately been proposed. Mask R-CNN, published by K He et al. in 2017, is one of the most popular. The segmentation has two essential forms:

- **Instance Segmentation:** Making an identification for the borders of an item and labeling its pixels by distinct colors [41].

- **Semantic segmentation:** The process of marking each pixel as a label in an image including the background, the pixel is with a distinct color based on its label of class, or its category [41].

## Background Knowledge

### 1. Overview of Machine Learning (ML)

Machine learning is an artificial intelligence (AI) technology allowing computers to make an automatic learning and improvement based on experience without a flagrant programming. ML is focused on the development of computer programs being able to make an access and utilization of data to learn themselves. Its progress of learning starts from data or observations, for example: a direction of finding out data patterns to make decisions with a better quality relied on instances being supplied in the future [37]. The fundamental purpose of helping computers have an ability to learn autonomously, which does not contain any participation or assistance of people and then adjust activities of them suitably. Then, ML allows for the examination of enormous amounts of data, although it generally produces quicker and more exacting results in identifying valuable possibilities or risky hazards, it probably necessitates more time and resources to train it precisely. In addition, the combination of ML with AI and technologies of cognition is able to increase its efficiency in processing massive amounts of data. Algorithms of ML are usually classified to supervised or unsupervised [37].

- **Supervised** ML algorithms have an ability to predict events of the future by applying what they learnt previous on new data by utilizing labeled examples. The learning method generates a function of inference to predict output results based on a study of a known set of training data. After going through an adequate process of learning, the system may supply goals for any input of new data. Moreover, the algorithm also has the

function of making a comparison between its own output and output of expectation or correct output, and gives out an error detection to make a suitable adjustment for the model [37].

- On the other hand, techniques of **unsupervised** are utilized when the material required to train is both unclassed and unlabeled. These algorithms learn about mechanism of computers for producing function of inference from data without being labeled, which is to express a unseeable construction behind. Although system does not determine the correct output, it investigates the data and is able to make conclusions from sets of data to characterize concealed constitutions of unlabeled data [37].
- **Semi-supervised** methods sit midway between supervised and unsupervised learning because they train with both unlabeled and labeled data, which is often a little of labeled data amount and a huge amount of unlabeled one. This strategy allows systems to significantly enhance learning accuracy. Semi-supervised is particularly utilized when the collected labeled data requires experienced and agreeable resources for training from. Besides, obtaining unlabeled data usually does not necessitate the use of supplemental resources [37].
- **Reinforcement** ML algorithms are a type of learning system which is several interactions with its around environment by performing executions and exploring rewards or errors. The most essential aspects of these algorithms are making several searches of trials and errors and postponed rewards, which enable devices and agents of software to find an optimal behavior at a specific scenario to making an improvement for

their performance to maximum point. For the agent to know what action is the best, there is a requirement of reward's simple feedback, which is called as the signal of reinforcement [37].

## **2. Overview of Deep Learning (DL)**

Deep learning (DL) is a small field belong to ML and AI which utilizes a bunch of layers of ANNs to reach achievements of "state-of-the-art" accuracy in tasks like speech recognition, language translation, object detection, and others. There is a difference between techniques of DL and typical techniques of ML, DL is able to learn automatically from data such as text, image or video without the utilization of hand-coded rules or knowledge domain of human. When given more data, their architectures with high adaptability are able to learn directly from raw data to achieve an improvement of prediction's accuracy [42].

DL is widely employed in applications such as computer vision (CV), systems of recommendation, and AI in conversation. DL is applied for CV software to learn from digital productions such as photos and videos. The applications of Conversational AI assist devices in understanding and communicating by natural language. Besides, language, user's behaviors and digital productions are used by recommendation systems to provide searching results and services which are suitable and significant [42].

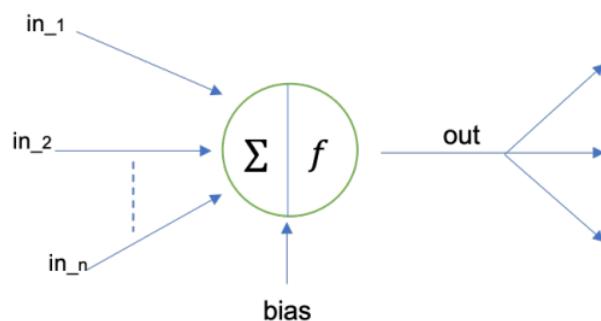
## **3. Artificial Neural Networks (ANNs)**

The concept of ANNs was built on the inspire of human brain biology characteristics, Cybenko [1] it gave out a substantiation which there is a hidden layer including neurons with a number of specific

limitations, having ability to make an approximation for any unstopping function to achieve any desire of precision in a neural network. It included a group of ML models having ability to “learn” to execute or perform some specific tasks based on supplied examples. In general, these networks were created by a group of units or nodes connected each other like a simulating version of neurons in human brain, which called artificial neurons.

### 3.1. Neuron

Analog electrical signals flow inside biological neurons, while digital values present as the signals for artificial neurons. Each artificial neuron has an output being computed by using a function of the sum of the input signal values with generally added bias, the connection between the neurons as known as edges, in each connection contains an adjustable model parameter named weight that can be “learned” through a process called “model training”. The strength of the synaptic will determines the interaction’s level between the neurons, this function combines the neurons, bias and weights in the connections called “activation function”, and there are several activation functions that mean there is not only one fixed activation function for ANN, and we can choose the activation function for the network’s output. These functions will be mentioned later in this thesis.



*Figure 1: Activation Function.*

### 3.2. Activation Functions

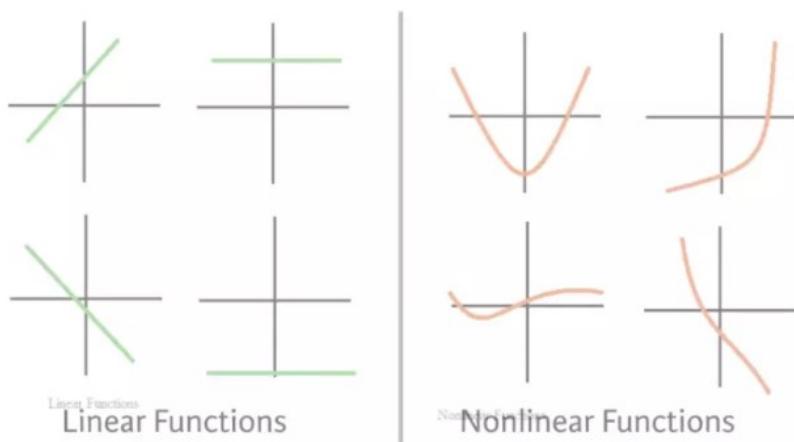
Activation Functions (AFs) in Neural Network are a component which is very essential in Deep Learning. AFs help determine the output of a Deep Learning model, the accuracy and computational efficiency of the training model – which will make a decision on the success or failure of a Neural Network system. AFs also have a great influence on the convergence and convergence speed of neural networks, or in some cases, AFs can prevent neural networks from converging in the first place.

AFs are presented in mathematical equations and this function is attached to each neuron in the network and determines whether it should be activated, based on whether the input data fits the model prediction, which also help normalize the output of each neuron to a range of 1 to 0 or between -1 and 1, depending on the type of activation function being used in the layer or network. Additionally, another aspect of AFs is that they have to have efficiency in computation because of being computed over thousands or even millions of neurons for each data sample. And modern neural networks use backpropagation technique to train the model.

The main reason that neural network models stand out from machine learning models is their ability to solve non-linear data problems through AFs of hidden layers, at which execute the task of solving problems for complex nonlinear relationships between data features and model outputs.

As mentioned earlier, in a neural network, the input data, is fed to the neurons in the input layer. Each neuron has a weight and multiplying the input by the weight gives the output of the neuron, which is then passed on to the next layer. So, AFs are applied in neural networks to be on duty of computation for the weight input's sum and biases, being utilized to make a decision on whether or not a neuron can be made activated. It manipulates the supplied data via several gradient processing, or gradient descent in typical and afterwards give out a production for a neural network's output, including the parameters. These AFs are usually known as a reference for a transfer function in some documentation or research papers. Furthermore, depending on the function they deputize, AFs can be either linear or non-linear, and they are used to modulate the outputs of neural networks in a variety of areas ranging such as: object recognition, classification, segmentation, speech recognition, cancer detection systems, finger print detection, and so on...

[12]



*Figure 11: Linear and Non-linear functions.*

Although in fact that there are many types of AFs, in this article I will only mention some of the AFs that are commonly used in current

CNNs architectures such as: sigmoid, softmax, hyperbolic tangent (Tanh), rectified linear unit (ReLU) and Leaky rectified linear unit (LReLU).

### 3.2.1. Sigmoid function

The Sigmoid Function is also known as the Sigmoid curve. This is a mathematical function characterized by an S-shaped curve, which is one of the non-linear activation functions most widely used. If being familiar with some machine learning models, probably still remembering Logistic Regression - as one of the simple algorithms for binary classification problems, which is quite effective. The "soul" of Regression is this Sigmoid function. The sigmoid is a continuous nonlinear function, allows to pass real numbers as its input and gives a production of results in the range 0 to 1, is considered probabilistic in some problems. Suppose you trained a Neural Network to images classification of birds and dogs, what a classic problem, where bird is 1 and dog is 0. Basically, when your Model returns a value greater than 0.5 meaning the image is of a bird, or the value less than 0.5 means the image is of a dog.

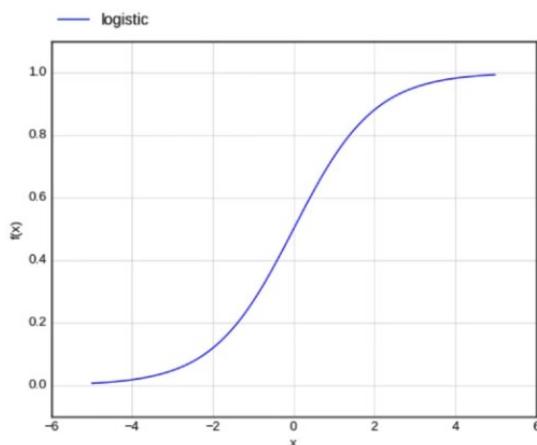


Figure 12: Sigmoid function.

In the Sigmoid function, a small adjustment from the input results in a no relative change for output. It gives a smoother and more continuous output than input, so, it has derivative in everywhere [13]. The Figure 13 shows the formular of sigmoid function:

$$f(x) = \left( \frac{1}{(1 + exp^{-x})} \right)$$

*Figure 13: Sigmoid function equation.*

In the formula,  $x$  represents for the network's output, and  $f(x)$  is its function, and  $exp$  stands for the Number of Euler. As a result, if the value of  $x$  approaches to the infinity of positive, the expected value of  $f(x)$  becomes 1. In contrast, if  $x$  approaches to the infinity of negative, the predicted value of  $f(x)$  becomes 0. And if the result of the sigmoid function is greater than 0.5, that label is classified as class 1 or positive side, and if it is less than 0.5, classified as class 0 or negative side.

On the other hand, the sigmoid has significant shortcomings such as sharp damp gradients while doing backpropagation starting from deeper hidden layers back to input layers, gradient saturation, sluggish convergence, and output of having non-zero centered, which causes gradient updates to propagate in various passageways. Other types of AF, such as the hyperbolic tangent function (Tanh), have been recommended to address some of the shortcomings of the Sigmoid AF.

### 3.2.2. SoftMax function

SoftMax Function (SF) is an exponential averaging function that calculates the probability of an event occurring. In general, the SF calculates the probability of a class occurring out of the total of all possible classes. This probability will then be used to determine the

target class for the inputs. Specifically, the SF turns a k - dimensional vector of any real values into a real-valued k-dimensional vector that sums to 1. The input value can be positive, negative, zero, or greater than 1, but the SF will always turn them into a value in the range (0:1]. The SF is calculated following this formula:

$$f(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$$

*Figure 16: SoftMax function equation.*

Where ( $x_i$  are vector of input values for the function from  $x_0$  to  $x_n$ , in which all the values can be any real number, positive number, negative number or zero. Next,  $\exp(x_i)$  is the standard exponentiation function is applied to each input value returning a positive value greater than 0. This value will be very small if the input is negative, and very large if the input is positive and it is an uncertain number in the range (0:1] as the requirement of a probability.

As such, they can be called “probabilities”. If one of the input values is negative or very small, the SF turns them into a small probability. If an input is large, it will be converted to a large probability. But the probability is always greater than 0 and less than 1, or equal to 1. So, the sum of all probabilities is always 1 and the SF brings a lot of benefits such as:

- Optimized when calculating the maximum probability in the model parameter.
- The property of the SF makes it satisfactory for interpretation of probability, which is very beneficial in ML.

- Normalization of SoftMax is a method to minimize the effect of extreme values or outliers in data with no having to modify the original data.

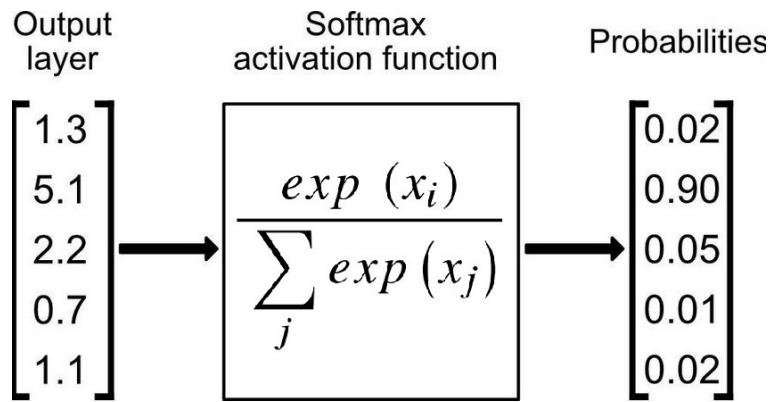


Figure 17: Softmax function applied in output layer of a network for classification task.

For the reasons mentioned above, the SF returns the value of probabilities for each class in multi-class models, with the target class having the highest likelihood. The SF may be found in practically all of the deep learning's output layers architectures where it is used. Adding, the Sigmoid and SoftMax vary in that the Sigmoid is used for binary classification and the SoftMax is utilized for the tasks of multivariate classification.

### 3.2.3. Rectified Linear Unit (ReLU) Function

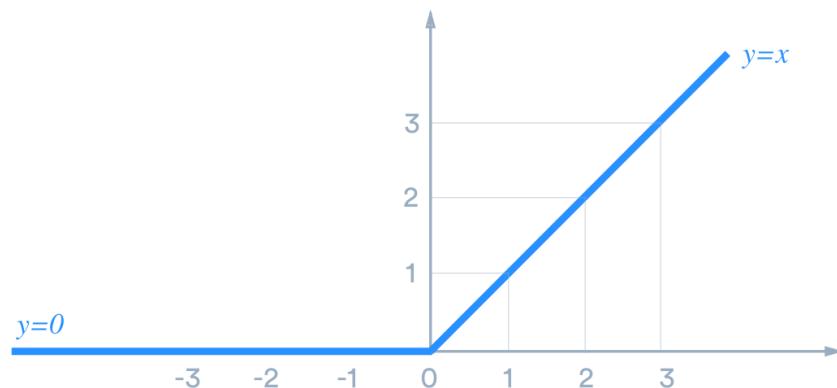
In 2010, Nair and Hinton introduced the rectified linear unit (ReLU) activation function, which has been the most extensively utilized activation function for applications in the DL field, bringing a state-of-the-art result till currently [17]. In general, The ReLU AF has an expressive speed of learning AF proven to be the most extensively and successfully applied function [14]. In DL, it outperforms and generalizes

the Sigmoid and Tanh activation functions [18]. In particular, there is a representation of a roughly linear function in ReLU function and hence retains the features of linear models that make them easier to optimize using the approaches of gradient decent. The mechanism of the ReLU AF, which applies a functioning of threshold to each element of input having values smaller than zero to be replaced to zero, and the formula of the ReLU AF as follows:

$$f(x) = \max(0, x) = \begin{cases} x_i, & \text{if } x_i \geq 0 \\ 0, & \text{if } x_i < 0 \end{cases}$$

*Figure 18: Rectified Linear Unit function equation.*

The ReLU will correct the inputs' values which are smaller than zero, driving them to zero and removing the issue about vanishing gradient that was noticed in the previous types of AFs. Mainly, the ReLU function is put to utilize into the hidden units of deep neural networks (DNNs). Besides, another AF as sigmoid or SoftMax is often applied in the network's output layers with popular purpose of classification in image [12] and speech recognition.

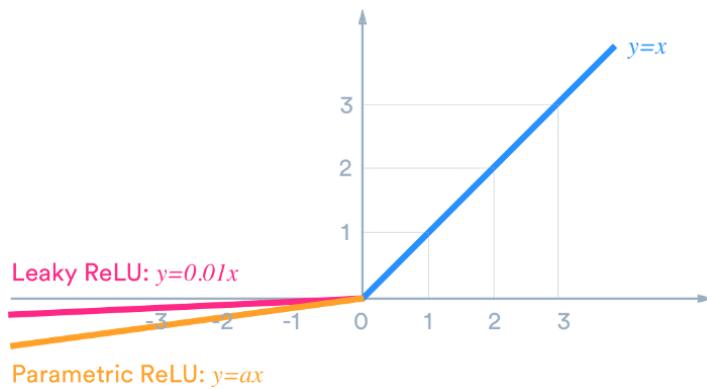


*Figure 19: Rectified Linear Unit Function [33].*

There are many advantages when using ReLU. In detail, the fundamental advantage of ReLU in computation is that they bring a more expeditious computation because they do not need to quantify exponentials and divisions, resulting in increased overall computation performance [18]. creation of sparsity in the hidden units is another feature of the ReLU, which squishes the values from zero to maximum. The ReLU has a problem in that it readily overfits when compared to the sigmoid function, despite the fact that the approach of dropout has been used to lessen the overfitting's influence of ReLUs and the rectified networks enhanced deep neural network (DNN) performances [19].

There is a notable weakness in the ReLU, which it is occasionally brittle during training, leading to the death of the several gradients. This will make neurons to be dead and direct weight updates to inactivate data points in the future, hampered learning because of dead neurons result in zero activation [3]. So that, the leaky ReLU was offered as a solution to the dead neuron problem.

### 3.2.4. Leaky ReLU



*Figure 20: Leaky Rectified Linear Unit Function [33].*

The leaky ReLU was released in 2013 that give out several small negative slope to the ReLU to maintain and stay survived the weight

updates throughout the phase of propagation. The building of the alpha parameter serves for the solution to issue of dead neuron of the ReLUs, which ensure that the gradients are not fixed zero anymore during training. The computation of the LReLU bases on the gradient with a tiny constant value for the negative gradient  $\alpha$  in the range of 0.01; hence, the LReLU AF is computed following this formula:

$$f(x) = \alpha x + x = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{if } x \leq 0 \end{cases}$$

*Figure 21: Leaky Rectified Linear Unit function equation.*

The comparison of the ReLU to an exception of no having zero gradients over the entire of time, which report a similar result. Therefore, excepting of distribution and sparsity in the comparison of ReLU to Tanh, no having any noteworthy enhancement of result. As a result, Leaky ReLU offers two advantages:

- Because it lacks zero-slope portions, it solves the "dying ReLU" issue, so training process is expedited. And, the "mean activation" near to zero is evidence for speeding up training. (It aids in keeping the diagonal entries of the Fisher information matrix short, but it can disregard judiciously.) Unlike ReLU, leaky ReLU is more "balanced," and as a result, it may learn faster.
- Keep in mind that the outcome is not always consistent. Leaky ReLU isn't necessarily better than regular ReLU, and it should only be used as a last resort.

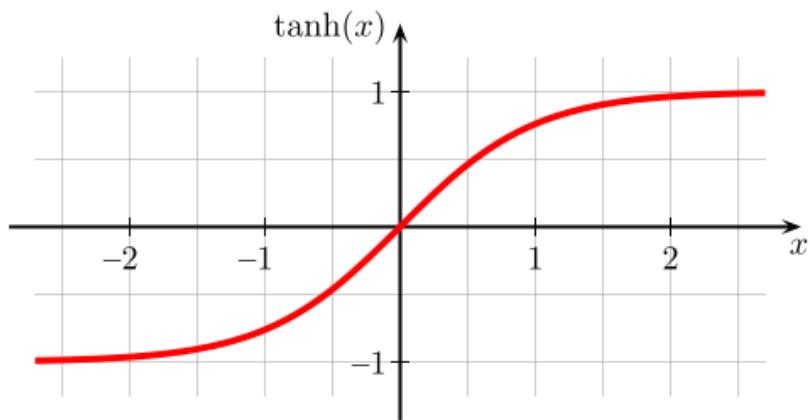
### 3.2.5. Hyperbolic tangent function (Tanh)

Another form of AF utilized in DL is the hyperbolic tangent function, which has several versions employed in DL applications [14]. The hyperbolic tangent function, often known as the tanh function, is a smoother zero-centered function with a range of -1 to 1, and its output is given by:

$$f(x) = \left( \frac{e^x - e^{-x}}{e^x + e^{-x}} \right)$$

*Figure 14: Tanh equation.*

Because of supplying higher training performance for multilayer neural networks than the sigmoid function, the tanh function has become the favored function over the sigmoid function [15]. Otherwise, the tanh function, like the sigmoid functions, was unable to overcome the vanishing gradient problem. The function's key advantage is that it generates zero-centered output, which aids in the operation of back-propagation.

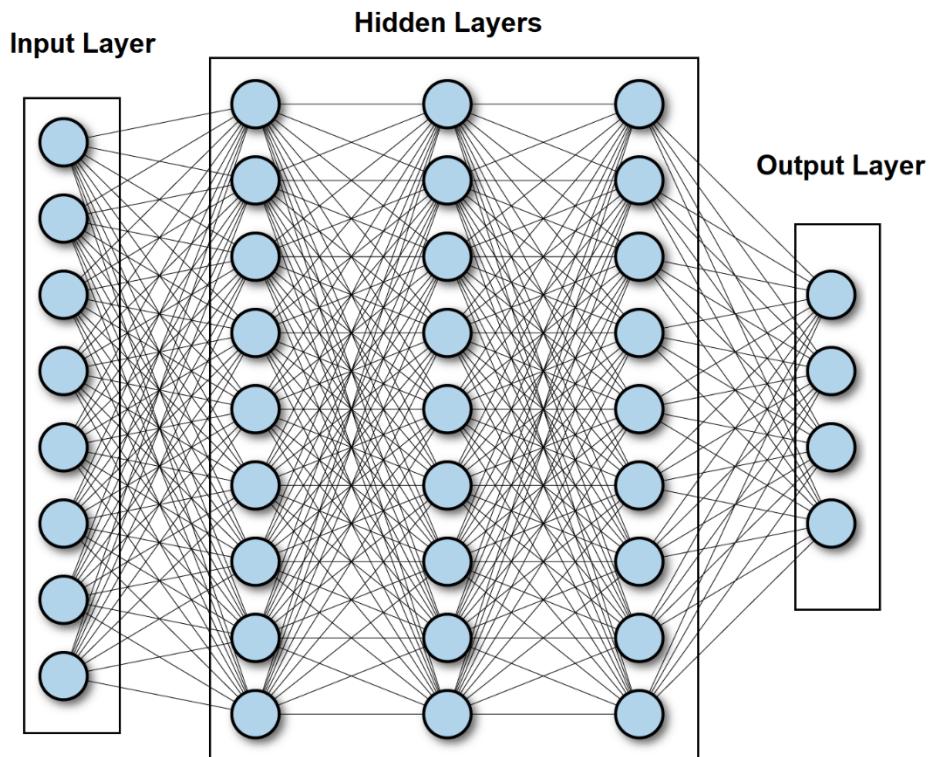


*Figure 15: Tanh function with values moving in range [-1: 1].*

The tanh function has the peculiarity of only being able to achieve a gradient of 1 only when the input's value is equal to 0, that equivalent to  $x$  equal zero. As a result, the tanh function generates some inactivated neurons during computation. The unworked neuron is a condition in which the activation weight rarely employed as a result of a zero gradient. In addition, this obstruction of the tanh function prompted additional study in activation functions to remedy the problem, giving rise to the ReLU activation function. Furthermore, the tanh function has mostly been employed in recurrent neural networks for natural language processing and speech recognition [16].

### 3.3. Layer

All neurons are linked directly in network, on the other hand, they are grouped into subsets of neurons known as layers, which are generally connected in a sequential way based on the activation function in every node to compute nodes' outputs from node's inputs. In detail, the neurons in a layer are made a connection to single set of neurons in both previous layer and to another one in next layer. Because of the modeling of ANNs as a connection of a neurons' structure in layer, neurons' outputs at a certain layer become neurons' inputs at the following layer, which networks are understood as feedforward neural networks. There are no loops in present layer, therefore the inputs are always sent forward. As a result, the neurons into a layer cannot be connected with each other. When all the neurons in a layer are created a connection to all other ones in prior layer and following layer, being called "Fully Connected Layers", this state of the connection is kept in whole network as a series of fully connected layers that will be called "Fully Connected Neural Network".



*Figure 2: Fully Connected Neural Network.*

As Figure 2 shows above, there is an existing of 3 kinds of layer inside a network, they are input layer, hidden layers and output layer.

The input layer represents for the input of the network, the number of nodes as known as neurons are size of the input that we call “input dimension”. The hidden layers lay on the middle of input layer and output layer, there can be so many hidden layers in a network, however, there is one input layer and one output layer in general, the output layer represent for the network’s output. In addition, the bias contains in input layer, not in output layer. Although simple decisions are made at the first layer relied on the input, sophisticated judgments are made at the second layer depending on determinations established at the first layer. As one moves into the network more deeply, more difficult and judgments of

abstraction emerge. The phrase deep neural networks refer to a network with numerous hidden layers.

### **3.4. Forward and Backward (back propagation) steps**

As mentioned before, there is a journey of the signals with an order starting from input layer, hidden layers and till output layer; which operations performed over the values received from the nodes in previous layers (in all layers but in the input layer) and whole activation function results are sent forwards to the next layer till the output layer. However, it is not the end of process. After getting the result of output layer, there is a backwards process named back-propagation applying in here as known as training step for ANN. The algorithm is used to effectively train a neural network through a technique called chain rule. In simple terms, after each forward pass through a network, backpropagation performs a backward pass while adjusting the model's parameters (weights and biases), which helps minimize the error for network's result, it is called "loss function" or "cost function" which is the evaluation between predicted output and expected output at the final point of forward step. Basically, this function calculates the gradient of the loss function to each parameter in network respectively that means it works from the top layer (output layer) to the bottom layer (input layer) to calculate the contribution of each parameter in that value [2], which updates the weights until reaching the minimum value. Furthermore, the key point of calculating the minimum value for loss function is gradient descent, which is a process of determining the modifications needed to minimize the loss and optimize the model to a specific evaluation measure by calculating the derivative of the loss function [3].

## **4. Gradient Descent Optimization Algorithms (Optimizers)**

During the training phase, we frequently modify and change our model's parameters (weights) in order to try to minimize its loss function and help produce predictions as optimal and accurate as possible. As a result, optimizers were developed to assist in resolving this issue. They make a combination between the loss function and model parameters by making updates of the model in response to the loss function's output. In other terms, optimizers are techniques or approaches which are utilized to help error function (loss function) reaching a minimum value or to make production efficiency peaking a maximum value [28]. The optimizers are built from mathematical functions that are influenced by learnable parameters of the model such as Weights and Biases. Additionally, the optimizers assist in determining how to adjust the weights and learning rate of a neural network in order to minimize losses. In a nutshell, the optimizers will mould your model into its most accurate possible form by tinkering with the weights [28]. And the loss function serves as a guide to the terrain, informing the optimizer whether it is travelling in the correct or incorrect path. There are several common optimizers such as Batch Gradient Decent (BGD), Stochastic Gradient Decent (SGD), Mini-Batch Gradient Decent, SGD with Momentum, Adaptive Learning Rate methods (AdaGrad, RMSprop, Adam).

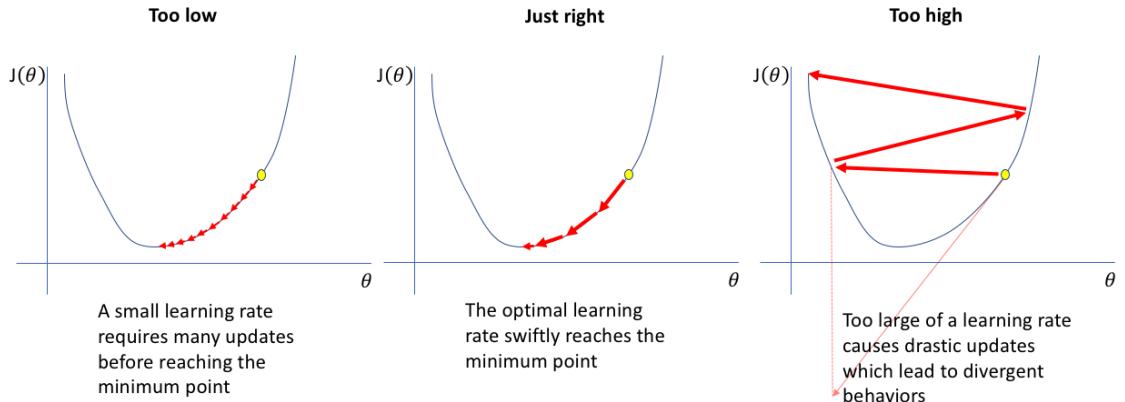
### **4.1. Gradient Descent (GD):**

Gradient descent is an optimization algorithm which is commonly-used to train machine learning models and neural networks [4]. It is built on a convex function and iteratively changes its parameters to minimize a given function to its local minimum. In the other words, GD begins by defining the initialized parameter values, and then iteratively adjusts the

values to minimize the given cost-function. Deeply, a gradient is a function's derivative with several input variables. In terms of mathematics, a gradient is known as the function's slope, and it simply measures the change in all weights in relation to the change in error. The steeper the slope and the faster a model can learn, the higher the gradient. However, if the slope is zero, the model will stop learning. A gradient is a partial derivative with regard to its inputs.

Imaging that a blindfolded man wants to hike down to the bottom of a valley from the top of a hill with as few steps as possible. He might begin going down the hill by taking large steps in the sharpest direction, which he can do as long as he is not close to the bottom. However, as he gets closer to the bottom, his steps will become smaller and smaller in order to prevent overshooting it. The gradient can be used to mathematically characterize this process. The learning rate  $\eta$  determines the size of the steps we take to reach a (local) minimum. In other words, we follow the direction of the slope of the surface created by the objective function downhill until we reach a valley [5].

To reach the local minimum for GD, we should choose a relevant value for the learning rate, which is not both too low and too high. This is very significant because if the steps are too large, it could fail to reach the local minimum since it bounces between the convex function of GD. Besides, GD will eventually reach the local minimum if we set the learning rate to a very low value, but it may take a lot of time [6]. So, the key point is choosing an optimal learning rate that helps model converge to the minimum value without spending much time.



*Figure 3: Learning rate choices that decide gradient descent result.*

GD has three forms that differ in the amount of data used to compute the gradient of the objective function. We choose between the accuracy of the parameter update and the time it takes to complete an update based on the amount of data.

#### 4.1.1. Batch Gradient Descent (BGD)

Batch gradient descent calculates the sum of the error for each point in a training set and updates the model only after all training instances have been evaluated. This procedure is known as a training epoch. Though this batching supplies a computation efficiency, it can still take a long time to execute big training datasets because it must still keep all of the data in memory [5]. Batch gradient descent also frequently yields a stable error gradient and convergence, although that convergence point isn't always the best, locating the local minimum rather than the global minimum.

#### 4.1.2. Stochastic Gradient Descent

Stochastic gradient descent (SGD) performs a training epoch for each example in the dataset and updates the parameters of each training

example at a time. They are easier to remember because you only need to hold one training example [5]. In spite of the fact that these frequent updates provide more detail and speed, they can lead in computational efficiency losses when it be compared to batch gradient descent. Despite having ability of noisy gradients production, these frequent updates can also aid in escape the local minimum and locating the global one.

#### **4.1.3. Mini-batch Gradient Descent**

Mini-batch gradient descent combines batch gradient descent with stochastic gradient descent ideas. It divides the training dataset into small batches and updates each of those batches. This method strikes a balance between batch gradient descent's computing efficiency and stochastic gradient descent's speed [5].

Besides, gradient descent has its own set of challenges such as: Local minima and saddle points, vanishing and exploding gradients.

- Local minima resemble global minima in shape with the slope of the cost function increasing on each side of the current location.
- Saddle points happen when the negative gradient exists just on one side of the point, reaching a local maximum on one side and a local minimum on the other. Its name was inspired by horse's saddle;
- Vanishing gradients happens when the gradient is too small. The gradient continues to shrink as we move backwards during backpropagation, leading earlier levels of the network to learn more slowly than later layers. When this occurs, the weight

parameters are updated until they become inconsequential, i.e., 0, resulting in an algorithm that is no longer learning.

- Exploding gradients occurs when the gradient becomes too big, resulting in an unstable model. In this instance, the model weights will get excessively huge and will be represented as “NaN” finally. One approach to this problem is to use a dimensionality reduction technique, which can help to reduce model complexity.

## 4.2. Stochastic Gradient Descent with Momentum

SGD has difficulty navigating ravines, which are widespread around local optima and are regions where the surface curves considerably more steeply in one dimension than another. In these settings, SGD oscillates across the ravine's slopes, making only sluggish progress along the bottom toward the local optimum [28]. That is the reason why momentum was established to reduce extreme volatility in SGD and soften the convergence. It promotes convergence in the relevant direction while decreasing fluctuation in the irrelevant direction. This approach employs another hyperparameter known as momentum, denoted by the symbol " $\gamma$ ".

$$\begin{aligned}v_t &= \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta) \\ \theta &= \theta - v_t\end{aligned}$$

Figure 40: Momentum Formula [29]

Some implementations swap the equations' signs. Typically, the momentum term  $\gamma$  is set to 0.9 or a similar number. We basically use momentum to drive a ball down a hill. As it rolls downhill, the ball gains momentum, increasing faster and faster (until it achieves its terminal velocity at the point being resisted by the air, i.e.  $\gamma < 1$ ) [29]. The term of momentum increment for dimensions has the point of gradients in the same directions and decreases for dimensions have change directions in gradients. Consequently, there is an achievement of quicker convergence and less oscillation. Although momentum helps eliminate oscillations and high variance in the parameters and converges faster than gradient descent, there is one more hyper-parameter which must be added and adjusted manually and precisely [28].

### **4.3. Adaptive Learning Rates Methods:**

The difficulty with employing learning rate schedules is that their hyperparameters must be established in advance and are strongly dependent on the type of model and issue. Another issue is that all parameter updates use the same value of learning rate [29]. If the data is sparse, we desire to update the parameters to a different extent. With this cause, AdaGrad, Adadelta, RMSprop, and Adam are adaptive learning rate algorithms that offer an alternative solution for traditional SGD. These per-parameter learning rate approaches give a heuristic approach without the need for costly manual tuning of hyperparameters for the learning rate schedule.

AdaGrad, in summary, conducts greater changes for more sparse parameters and smaller changes for less sparse ones. It performs well with sparse data and in training large-scale neural networks. However, during training deep neural networks, its monotone learning rate is

frequently excessively aggressive and stops learning too soon, that means AdaGrad makes learning rate changing adaptively in iterations and has an ability to train spare data as well, but it will make dead neuron issue when the neural network is deep the learning rate becomes very small number [29].

So, AdaDelta is an Adagrad addition that aims to minimize its belligerent, monotonically declining learning rate, when using AdaDelta, it is not essential to set a default learning rate, but the cost for computation is expensive [29].

In addition, RMSprop makes a very modest adjustment to the Adagrad method in an attempt to minimize its aggressive, monotonically falling learning rate. Specifically, the RMSprop helps learning rate automatically be tweaked and picks a different learning rate for every single parameter, however, it makes the training model very slow [29].

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

*Figure 41: RMSprop function [29]*

The learning rate is also divided by an exponentially decaying average of squared gradients by RMSprop. It is recommended that  $\gamma$  be set to 0.9, with 0.001 as an acceptable default value for the learning rate  $\eta$  [29].

Adam stands for Adaptive Moment Estimation, is one of the most well-known and often used gradient descent optimization methods, which is a technique for calculating adaptive learning rates for every

single parameter, stores the decaying average of previous gradients, comparable to momentum, as well as the decaying average of previous squared gradients, similar to RMS-Prop and AdaDelta. It incorporates the benefits of both strategies. Furthermore, Adam brings some benefits that are easy to implement, computationally efficient, using little memory requirements, converging rapidly, rectifying vanishing learning rate and high variance, but its computation is costly [28].

$$\begin{aligned}
 m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\
 v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \\
 \hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \\
 \theta_{t+1} &= \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t.
 \end{aligned}$$

*Figure 42: Adam function [28]*

In figure 42,  $m(t)$  and  $v(t)$  are values of the first moment which is the Mean and the second moment which is the uncentered variance of the gradients respectively. There is a set of default values for  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , and  $\epsilon = 10^{-8}$ .

## 5. Convolutional Neural Network

Convolutional Neural Networks (CNNs), like ANNs, being inspired by the structure of the human visual cortex [8]. CNNs, which are a type of neural network typically utilized for computer vision tasks in deep learning, include convolutional layers, pooling layers and fully connected layers. These networks are identical to traditional neural networks, except that instead of generic matrix multiplication, a

convolution kernel is used in one or more of their layers [3]. These networks have the role of reduction the images into a form which is easier to process, without losing features which are critical for getting a good prediction. Moreover, images are tensors in 3 dimensions of width, height, channels, and CNNs have ability to take the tensor as an input, identify which image characteristics or features are essential for classification or differentiation, and output these classifications [1, 7].

Therefore, the basic structure of a Convolutional Network in order starting from a Convolutional Layer or so-called convolution process (do not count the initial input layer), has the effect of picking out the 3-dimensional characteristics of the image input. Specifically, the features of the input original image are passed through the “convolutional layer”, at which there is a mechanism built up by matrix multiplication between the input image known as a matrix and one or more matrices with a certain size (usually smaller than the input matrix) being called “filters”, this matrix multiplication will be shifted by rows and columns with a certain rule that will give generate one or more new matrices based on the number of filter layers, the product from this process is called feature maps which contain the attributes of the original input image.

The pooling layer helps the network synthesize the features of the image in a deeply specific way (max pooling) or more objectively general (average pooling or mean pooling) from the output of the previous layer. Next, the size of the image has been significantly reduced that is very meaningful for memory in process of training model. Although the dimensions are decreased, the number of channels or feature maps are still kept unchanged. So, this process is as a composition of a reduction for the size of the input image feature in

length and width as well as a reduction of the number of learnable parameters. At last, when the input size has been reduced to a reasonable value, the 3-dimensional image will be converted to dimensional vector or known as "flatten layer", at which this process is generally implemented as a Fully Connected Layer, the number of units in this layer will correspond to the number of classes that defined for the network before, relied on specific activation function, the final output of the network will produce a specific final result in range of the classes.

In general, using CNNs because of their capacity to learn translation-invariant patterns as well as spatial hierarchies utilizing appropriate filters. Learning translation invariant patterns means that once a pattern is learned in one area of a given image, the network will be able to recognize it anywhere else in the image. This improves image processing efficiency by requiring fewer training samples to learn generalizable representations [1].

Spatial hierarchies show how earlier layers in the network can learn tiny local patterns and grow into bigger patterns composed of these tiny patterns in the next layers, which also improves the network's ability to learn complicated and abstract concepts of visualization. In addition, this also allows these networks to deliver class predictions with far higher accuracy than vectorizing a complex image with pixel dependencies throughout. Therefore, CNNs have to be used for any image classification task so as to track these crucial properties and correlations between features [1].

## 5.1. Convolution Layer

Convolution is an operation that produces a new function by combining two functions of a real-valued argument. Let  $s(t)$  is the output estimate function based on time  $t$ ,  $x(a)$  is the age-based input position function, and  $w(a)$  is the weighting function that prioritizes recent measurements. This is the general formula for convolution utilizing this defined function:

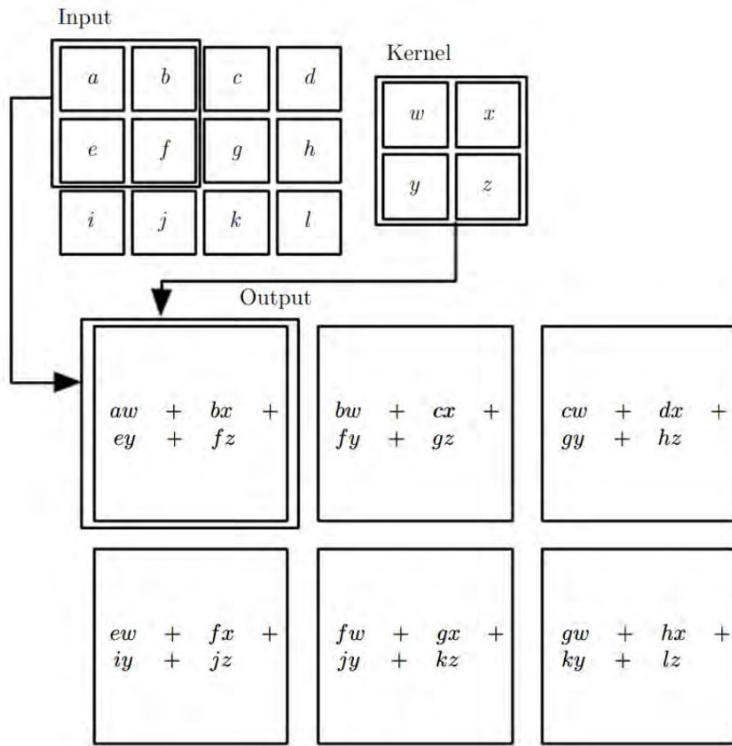
$$s(t) = \int x(a)w(t - a)da.$$

*Figure 4: Convolution Function.*

In detail, the input is the function  $x(a)$ , the kernel is the weighting function  $w(a)$ , and the output  $s(t)$  is the feature map in the CNN. the CNNs use convolutions that are performed over two-dimensional tensors that are made up of the input images' height, width, and channels of color. These techniques take out patches and alter the input to build a feature map with varied depth. One of the two critical criteria for defining these convolutions is the size of these extracted patches. The depth of the output feature map is the second parameter to consider. This is the number of filters built by the layer that can encode diverse features of the input data.

2D Convolutions are computed by moving a square window of defined height and width across all pixels of the input feature map, which connect to each local region or receptive field of the image corresponding to the size of the applied filters. For images, the input feature map has three dimensions: height, width and color bands, which commonly correspond to red, green, and blue. When these feature maps perform 2D convolution, two dimensional patches of surrounding

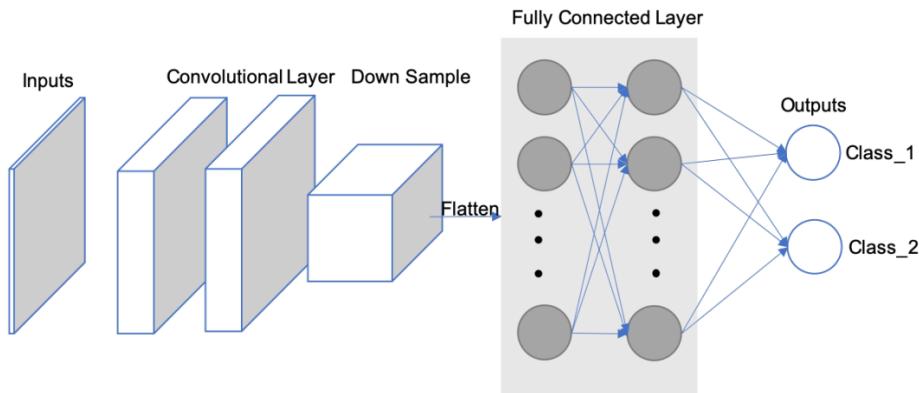
features are created. These patches are then turned into a one-dimensional vector reflecting the output depth via a convolution kernel. After that, the vectors are reassembled spatially into a two-dimensional output map that corresponds to all points in the input map. The convolution process is depicted in the Figure 5 below.



*Figure 5: Feature maps are computed by stride mechanism of convolution process.*

In other words, a feature map is the output of a convolutional layer, which is typically a three-dimensional tensor with dimensions width, height, and depth. To process inputs or feature maps, the convolutional layer employs a filter, similar to the filtering procedure in picture pre-processing. The primary distinction is that in standard image pre-processing, the filter weights (values) are hand-crafted, whereas in CNNs, during training, these weights are learnt. Alternatively, the receptive field in the input feature map corresponds to the output of a

convolutional layer, and the size of the receptive field is governed by the kernel size and dilation. It is not essential to flatten images while using CNNs. However, flattening is required to feed the features maps, being the outputs of convolutional layers, to the fully connected layers. Fully connected layers are the same as the ANNs mentioned above, however ANN is a catch-all term for all artificial neural networks, including CNNs. Fully connected layers understand the linear and non-linear correlations between extracting features and classifying each sample into many abstract classes. Figure 6 shows a simple CNN architecture.

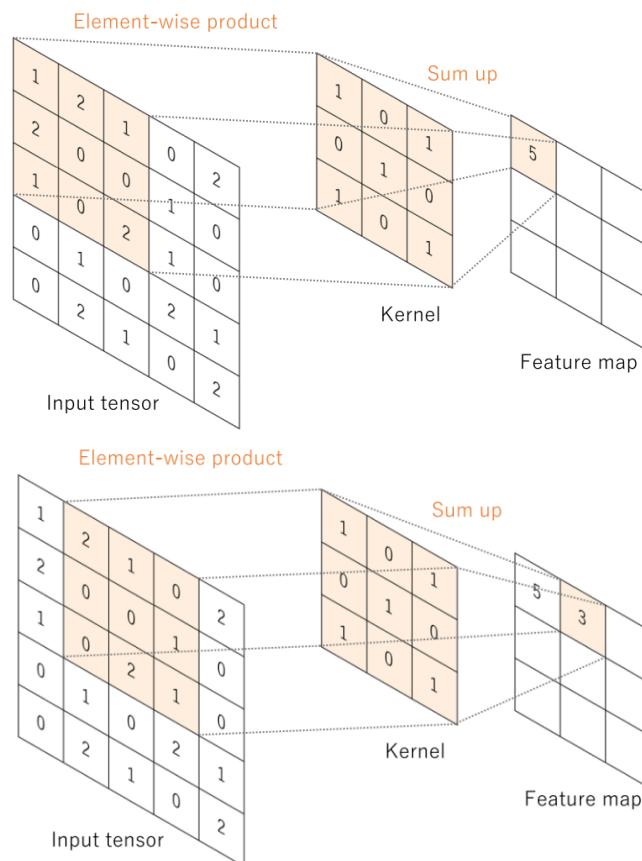


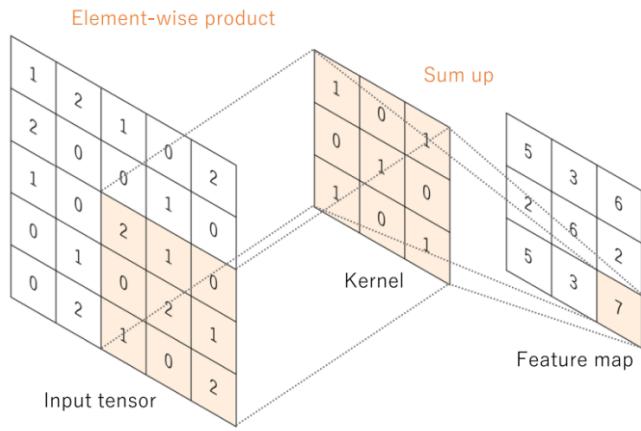
*Figure 6: CNN Architecture.*

The size of the output volume is controlled by hyperparameters which are depth, stride, and zero-padding.

The first component of the output volume decision is depth, which is the number of filters to utilize, each learning to search for something different in the input. For example, if the first Convolutional Layer receives the raw image as an input, distinct neurons along the depth dimension may fire in the presence of different oriented edges or color blobs. A column of depth is a group of neurons being all staring at the same region of the input [9].

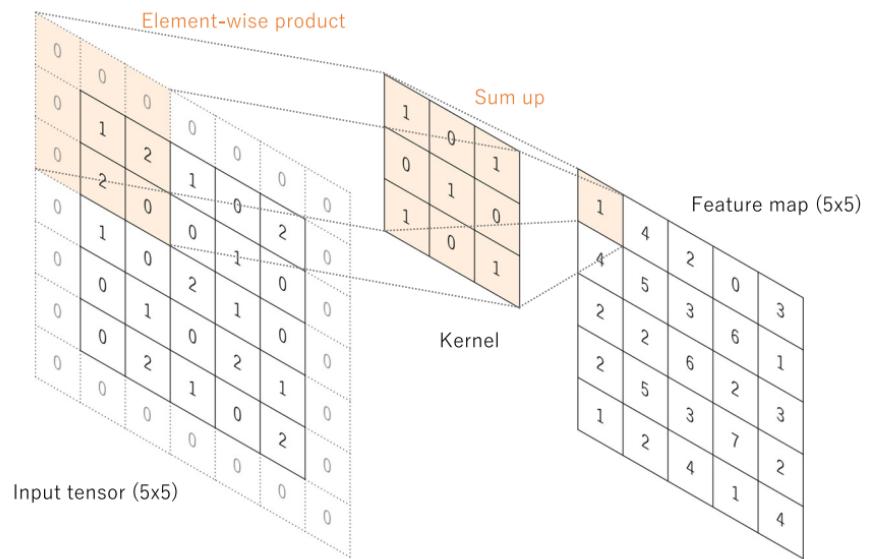
Stride aids in picture and video data compression tuning. The stride of the convolution is a characteristic which can contribute to a different output size in CNNs, which is a convolution operation parameter that defines the distance between patches derived from the input feature map. With a stride of 2, the output feature map's width and height are down-sampled by a factor of two with no padding. For example, if the stride of a neural network is set to 1, step by step, the filter will make a reflection onto the matrix of input, which move one pixel or one unit at a time from left to right in the row till the end of horizontal point and top to bottom in the column at the end of the matrix. Because the size of the filter influences the encoded output volume, stride is frequently a positive integer rather than a fraction or decimal.





*Figure 7: Generating feature map process.*

It is sometimes useful to pad the input volume with zeros around the border. This zero-size padding's is a hyperparameter. The great thing about zero padding is that it allows to tweak the spatial size of the output volumes, which most typically using it to precisely preserve the spatial size of the input volume so as to be same of the input and output width and height [9].



*Figure 8: Generating feature map process with zero-padding.*

Combination of these components will give out the formula for calculating the size of Convolutional layer:

$$[(N + 2P - F) / S + 1] \times [(N + 2P - F) / S + 1] \times C'$$

*Figure 9: The formula of output size of Convolutional layer.*

The spatial dimension of the output volume may be calculated as a function of the input volume size (N) which represents for the dimension by width and height of input image, the receptive field size of the Conv Layer neurons (F) as the dimension of filter by width and height, the amount of zero padding placed the border for input image (P), the positive integer number of strides with which they are applied (S), and the number of filter (C').

The results of a linear operation, such as convolution, are subsequently processed by a nonlinear activation function. Although smooth nonlinear functions such as the sigmoid or hyperbolic tangent (tanh) function were previously used because they are mathematical representations of biological neuron behavior, the rectified linear unit (ReLU) function is now the most commonly used nonlinear activation function, which simply computes the function:  $f(x) = \max(0, x)$

In convolutional neural networks, rectified linear units (ReLUs) provide three key advantages:

- The ReLUs efficiently transmit the gradient, reducing the risk of a vanishing gradient problem, which is prevalent in deep neural architectures.

- The ReLUs set negative values to zero, so solving the problem of cancellation and resulting in a much sparser activation volume at its output. Sparsity is helpful for a variety of reasons; However, it is most commonly used to supply robustness to tiny adjustments in input such as noise.
- The ReLUs that including just simple computations (primarily comparisons) and are thus substantially more efficient to implement in convolutional neural networks.

## 5.2. Pooling Layer

There are also pooling layers between convolutional layers. They are used to down sample the feature maps and reduce the number of parameters. The later layers after feature extraction will need a large number of parameters because the depth is determined by the number of channels in the later layers often increasing exponentially. That increases the number of parameters and the amount of computation in the neural network. Therefore, to reduce the computational load we will need to reduce the dimensions of the input matrix block or reduce the number of layer units. Since each unit would be a representative result of applying a filter to find a specific feature, reducing the number of units would not be feasible. Reducing the input matrix block size by finding a representative value for each spatial region that the filter passes through will not change the main contours of the image but reduce the size of the image. Hence the matrix reduction process is applied.

Average pooling and max pooling are the two most commonly used types of pooling layers helping make this reduction purpose. As the suggestion of their names, average pooling takes out the average (mean)

of the values on the feature map corresponding to a kernel size. Besides, max pooling calculates the maximum value in the neighbourhood corresponding to the kernel size position on the feature map. Both of them travel through the whole feature map by a stride size without overlapping local region position. In addition, the pooling layers do not have any parameters which need to be learned during training and the stride in the convolutional layer is able to also be used to replace the pooling layers. Figure below shows the difference between these two pooling types.

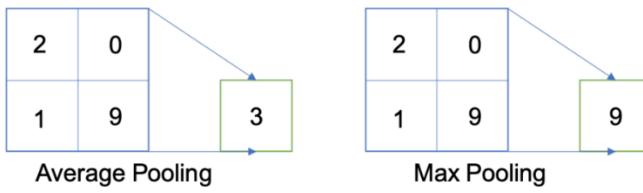


Figure 10: Average (Mean) and Max Pooling Layer.

### 5.3. Batch Normalization:

Another critical component of a CNN architecture is batch normalization (BN). Although this is not a required component in CNN, it is a relatively important component with many benefits for training the model to reach the best results. BN allows faster training and stabilization of deep neural networks by stabilizing the distribution of layer inputs during training. This approach mainly involves Internal Covariate Shift (ICS). To improve model training, it is important to reduce ICS by controlling the means and variances of the input data layers. In detail, BN normalizes the feature maps at the output of a convolutional layer in the same way which it normalizes the input pictures. BN makes subtracts the batch mean and divides the feature map values by the batch standard deviation to ensure that the distribution of

feature maps is the same. As a result, BN lowers value shift in hidden layers and BN allows to obtain an additional flexibility in the initialization of kernel weights. Furthermore, because of guarantee of BN which activation does not diverge to very large values, greater learning rates can be employed. In addition, BN is quite similar to a dropout layer which also has some regularization effects, because of adding noise to the feature maps. During training and evaluation, both the BN and dropout layers will behave differently.

When we would like the “function of probability density” (FPD) of the model to be closest the true real distribution of data as possible in deep learning. Unfortunately, in fact that we only have a narrow amount of training data at our disposal. As a result, the goal is to get the FPD of the model as near to the FPD of the train as possible by training on a dataset that replicates the domain of application in the real-world. A sufficiently big training dataset is required to get improved representation. When the training dataset is insufficiently large, there will be a significant difference in the error between the FPD of real and the FPD of train.

In other expressions, even if we produce a model FPD that is very close to the FPD of the train, it will not perform with the FPD of test skillfully, being a sample of the FPD of the real. This is referred to as the overfitting problem. Overfitting occurs when a model produces a good performance in train data set and a very bad performance in test data set. However, the model is not overfitted when performing impressive in training and a bit worse in testing. Dropout layer is applied to stay away from overfitting by randomly selecting some neurons to brush aside throughout the training process. A dropped-out node's incoming and

outgoing edges are also eliminated. Thus, dropout not only reduces the problem of overfitting, but it also simultaneously speeds up training by eliminating training with all nodes.

#### 5.4. Fully Connected Layer

Typically, the activation function applied to the last connected layer differs from the others. Each job necessitates the selection of an appropriate activation function. A soft-max function is used as an activation function in the multiclass classification job to normalize output real values from the final fully connected layer to target class probabilities, where each value ranges between 0 and 1 and all values will be made a sum to 1. Typical final layer activation function selections for various sorts of jobs.

In addition, the majority of the trainable parameters are located in fully connected layers in CNN. Fully connected layers can give out outputs that are categories for image classification or localization via bounding boxes [11]. Convolutional layer work can also be used as fully connected layers to execute the same duty more efficiently. Instead of a class, the result of an image segmentation job is a mask with the same resolution as the input pictures. We may use a 1x1 kernel as the output layer to reduce the depth to the same level as the target classes [10], giving the model pixel-wise classification capability.

Training a CNN consists of two steps: forward and backward. We feed the CNN with input information flow in the forward. To acquire the input for the following layer, feature maps are generated in each layer with predefined weights and biases. A loss value will be computed at the output layer by applying a loss function to the difference between the

calculated output and the ground truth. The chain rule is utilized in the backward step to find out the gradients of loss value for each trainable parameter. After that, the gradients are utilized to adjust each parameter for the following iteration. This type of update is what allows CNNs to “learn”.

After processing enough iterations, at which the loss value converges to an acceptable small number that means it is approximately close to zero, the training process can be terminated. The loss function computes the difference or error between the outputs and the ground truth. Minimizing loss implies minimizing the difference between model’s function of probability density and train one when the loss function is applied as a guide. Moreover, the loss function will change depending on the application and objective. Because of its ability of quantify the similarity/difference between two distributions, cross-entropy is a widely used loss function.

Otherwise, there are several optimizer options to help you decide how to update the network parameters. Adam and Stochastic Gradient Descent (SGD) are two options. In general, SGD will produce better results while learning at a slow pace. It will take several iterations to achieve a satisfactory outcome. SGD has a fixed learning rate that does not vary during training. Adam, on the other hand, converges quicker since it employs a vector of learning rates that are adjusted as learning continues during training progresses. They will be mentioned more detail in the next sections of Optimizer.

## 6. Residual Network (ResNet)

ResNet, an abbreviation for Residual Network, is a form of NN introduced in 2015 by Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun in their paper "Deep Residual Learning for Image Recognition". The ResNet models were immensely successful with winning first place in the competition of the ILSVRC 2015 for classification with a top-5 error proportion of 3.57 percentage (An ensemble model), as well as first place in the competitions of ILSVRC and COCO 2015 in ImageNet both detection and localization, Coco detection, and Coco segmentation [20]. Furthermore, ResNet-101 was used to replace VGG-16 layers in Faster R-CNN. They discovered a relative improvement with a proportion of 28 percentage. ResNet was developed because of the vanishing or exploding gradient issue on some state of the arts model architectures at that time.

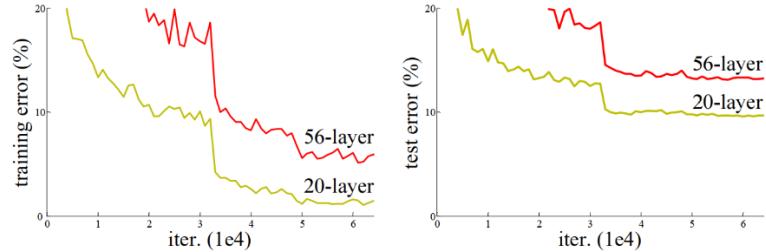


Figure 22: Performance of 20 layers and 56 layers plain networks on CIFAR-10 [20].

Before the ResNet, with the advent of neural networks deeper and deeper by a common argument, which is "go deeper", growing up the depth of the network will be an improvement for the quality of the network. Based on that, several new DNN architectures have been born such as AlexNet, GoogleNet or VGG, which apply the mechanism to increase the number of layers or number of layer's group. Because these added layers are able to solve complicated issues more efficiently as the non-identical layers could be trained for various duties to get high results

of accuracy. Additionally, the number of stacked layers can enrich the features of the model. However, with an increasing number of layers of complexity in deeper CNNs, a new trouble has also emerged called vanishing gradients. In specific, as deeper networks may begin to converge, the accuracy of the network gradually saturates to a certain threshold and then slowly moves down during backpropagation step, and higher error increases may occur. This issue goes against the purpose of reducing error and increasing accuracy during network model training. In other words, the values go through the gradient step that are chain-rule by the calculation of partial derivatives method, which will quickly return a very small result that approaches 0 after several chain-rule applications [22]. It will lead the learning of the network stopped soon, so this backpropagation will not be executed in previous layers, in this case, it may be stopped at an unknown hidden layer. This problem of training very deep networks has been alleviated with these ResNets which are made up from Residual Blocks (RBs).

### 6.1. Residual Blocks (RBs)

ResNet takes advantages of using RBs to increase model accuracy. The strength of this form of neural network is the concept of "skip connections" or skip paths, which is at the heart of the RBs. There are two ways for the skip connections to work. At first, they address the issue of vanishing gradients by creating an alternate path as a shortcut for the gradient to go through. Second, they also allow the model to learn a function of identity, which assures that the model's upper levels will have a performance not being worse than the layers lays underneath.

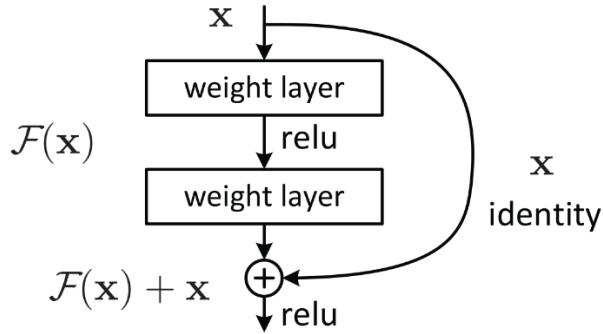


Figure 23: Residual Block Unit with skip connection [20].

The construction of an RB is depicted in the Figure 23. The ResNet is developed to create its stack of 2 layers being able to fit to  $F(x) = H(x) - x$  rather than training only stacked convolutional kernels have an ability to be suitable to the intended mapping. Supposing that, the original mapping has been recast as  $H(x) = F(x) + x$  [20]. They claim that optimizing the residual mapping substantially easier, which is  $F(x)$ . In spite of being the optimal option of the identity in an extreme circumstance, model will give out a result in  $F(x)$  being zero. Using such a framework, networks with more than 100 layers might be trained. Besides, this method appears to have a minor flaw when there is a difference between input's dimensions and those output, which have ability to occur with pooling and convolutional layers. When the dimensions of  $F(x)$  differ from those of  $x$ , we have two options:

- To up-sample its dimensions, extra zero entries will be padded to the skip connection.
- To make the dimensions equal, the linear projection approach is utilized, which is accomplished by attaching more  $1 \times 1$  convolutional layers to the input. In this example, the output is as follows:  $H(x) = F(x) + W(x)$ . In this case, adding an extra parameter  $W$ .

In short, the RBs help the layers acquire a knowledge of identity functions substantially easier than plain networks. So that, ResNet raises up the performance of DNNs with additional neural layers while reducing the proportion of error. In other words, the skip paths combine the outputs of prior layers with the outputs of stacked layers, allowing for considerably deeper networks to be trained than previously feasible.

## 6.2. ResNet Architecture

### 6.2.1. ResNet – 34 Architecture

The ResNet-34 was the first architecture built, which included inserting shortcut connections into a plain network to transform it into counterpart of its residual network. The plain network in this case was influenced by VGG neural networks (VGG-16, VGG-19), while the convolutional networks had 3 x 3 filters. On the other hand, ResNets have fewer filters and are less sophisticated than VGGNets. There is a comparison between the 34-layer ResNet attains 3.6 billion FLOPs of performance and larger VGG-19 with 19.6 billion FLOPs [20]. ResNet 34 takes a proportion of 18 percentage of VGG-19, thus, the VGG-19 has more filters and higher complexity than ResNet 34.

The ResNet-34 also adhered to two simple design principles that the layers must to have the same filter amounts for the same size of feature map's output, and the filter amounts is twofold if the size of feature map is cut off a half to maintain the time complexity for each layer [21]. And one more thing is 34 weighted layers being involved in it. In addition, this plain network plus shortcut connections. Despite being identical in the input and output dimensions, the identity shortcuts were employed by a direct route. There were two choices for considering when the dimensions raised up. The first one was the shortcut would

continue to execute identity mapping while padding extra zero entries for an increment of dimensions. The projection shortcut could also be implemented to fit dimensions.

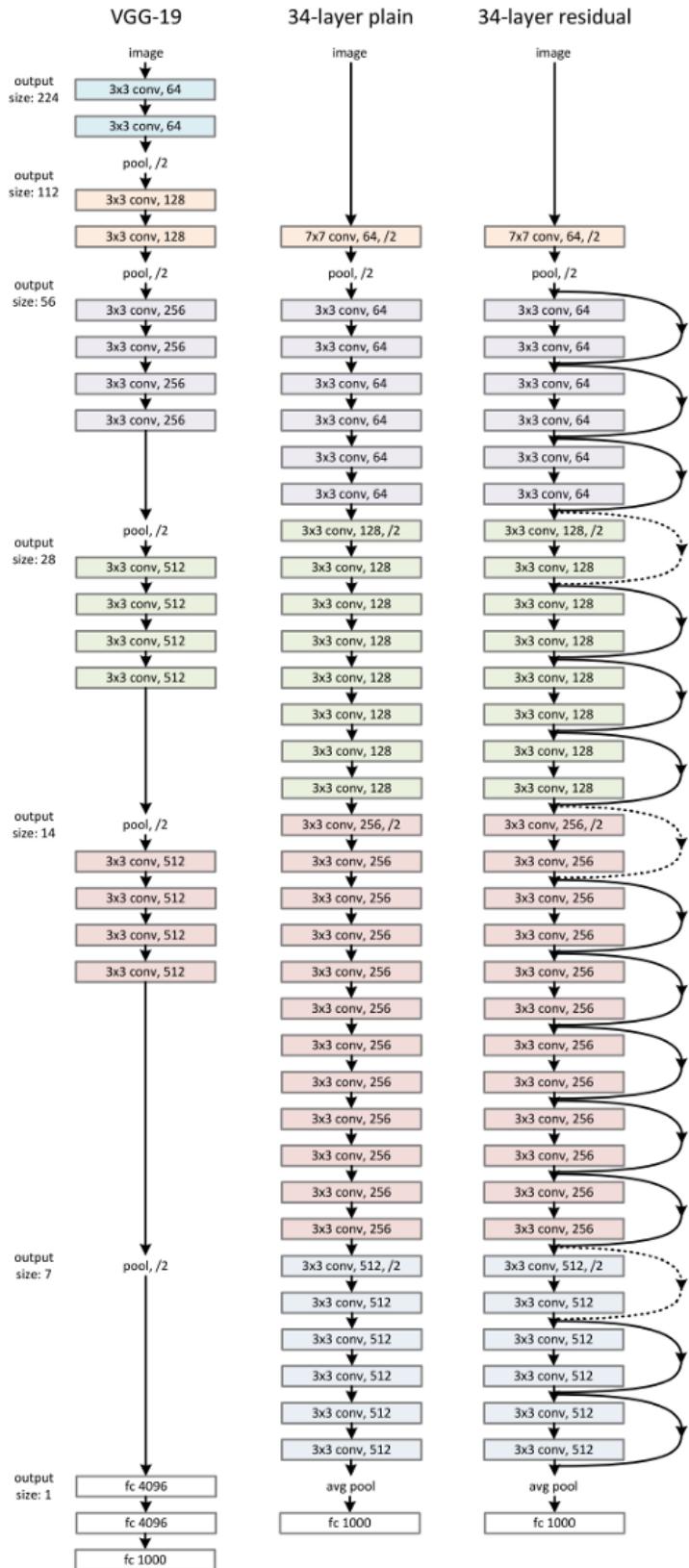


Figure 24: VGG-19 model, 34 Plain Network, ResNet 34

[20].

#### **4.2.2. ResNet-50, ResNet-101, and ResNet-152 Architecture**

Although the ResNet-50 architecture is based on the above model, there is one significant variation. In this situation, the RB was re-created as a design of bottleneck appearance cause of concerning about the requirement of time for training the layers. Instead of using two layers, a three-layer stack was used for building of RB. As a result, each of the ResNet-34's 2-layer blocks was taken a place of a 3-layer bottleneck block for the construction of ResNet-50, which is known as "Deeper Bottleneck Architecture". This model is generally more accurate than the 34-layer ResNet model with performing of 3.8 billion FLOPS compared with 3.6 billion FLOPS [20].

Furthermore, Larger Residual Networks are the 101-layer ResNet-101 and ResNet-152, which are built with additional 3-layer blocks as well as ResNet-50 [21]. Even with an increment of network depth, the 152-layer ResNet reaches to 11.3 billion FLOPS, and has a specifically lower complexity than VGG-16 or VGG-19 nets with the number of FLOPS about 15.3 and 19.6 billion in respective.

## **7. EfficientNet**

Following the previous success of AlexNets, GoogleNets, VGGs, ResNets, Gpipe, RCNN, or Mask RCNN with outstanding results in image recognition contests as well as the impressive accuracy in image recognition tasks bringing in reality. EfficientNet was born in 2019 by Quoc V. Le et al from GoogleBrain. Based on inspiration from the process of scaling up ConvNets, these researchers systematically studied it and came up with a completely new solution compared to the previous state of the arts solutions, which is scaling up the main component of the

network model in 3 dimensions: in depth - scaling to add more layers to the model, in width - scaling to add more channels for layers, in resolution - scaling in resolution being also known as the size of the layers, feature maps, they called it the "compound scaling method"[25].

Before going into more detail about EfficientNet, we will learn about the reason why the Model Scaling definition of the previous network models are the fundamental foundation to construct the EfficientNet.

## 7.1. Model Scaling

In fact, the state-of-the-arts models is built by Model scaling method, and "go deeper" is also one of this method. ResNet is able to be scaled up or down by modifying the layers of network - in depth, whereas WideResNet and MobileNets have ability to scale based on adjustment of the network's channels - in width. Besides, the network model is also widely acknowledged that larger input image sizes improve accuracy while incurring the burden of additional FLOPS [25]. Although several previous researches have proven that network depth and width are both important for the power expression of the ConvNets. It is still not easy to productively scale a ConvNet for a better achievement of accuracy and efficiency.

**In Depth:** Many ConvNets, including ResNet and Inception Net, scale network depth in this manner. The assumption is that deeper ConvNet is capable to capture more affluent and more complicated characteristics of image while also generalizing successfully on new tasks. However, because of the vanishing gradient problem, deeper networks are harder to train. In spite of having numerous approaches as

ResNet-skip connections and batch normalization to mitigate the training issue, the accuracy gain of very deep networks will drop, for example: Despite having so many more layers than ResNet-101, ResNet-1000 has equal accuracy [23].

**In Width:** For small scale models like MobileNet, MobileNetv2, and MnasNet, scaling network width is widely applied. Wide residual networks of Zagoruyko and Komodakis in 2016 explored how wider networks can take more fine-grained features and are easier to train [25]. On the other hand, outstandingly wide but shallow networks, having a difficulty in catching the features in higher level. Therefore, when networks expand substantially wider with bigger channels, the accuracy quickly saturates [23].

**In Resolution:** ConvNets may be able to catch more fine-grained patterns with higher resolution input photos. Early beginning with 224x224 with ConvNets, later ConvNets often employ 299x299 or 331x331. And in 480x480 resolution [23], GPipe just achieved state-of-the-art accuracy of ImageNet. In object detection ConvNets, higher resolutions, such as 600x600, are also commonly employed [23]. The results of scaling network resolutions are shown in the figure below, where greater resolutions indeed enhance accuracy, but the accuracy gain reduces for very high resolutions ( $r = 1:0$  means resolution 224x224 and  $r = 2:5$  denotes resolution 560x560).

Summary of this thing, although the accuracy of the model will be enhanced when scaling up any single dimension of network width, depth, or resolution, the accuracy gain reduces with larger models.

## 7.2. Compound Scaling Method – EfficientNet

Based on the above 3 options of Model scaling method, the EfficientNet's researchers have come up with the compound scaling method which is a uniform combination of 3-dimensional expansion in the width, depth and resolution of the network, with a set of the ratio is fixed, this aims to create a balance in the dimensions of all 3 dimensions [25].

$$N = F_k \odot F_{k-1} \odot \dots \odot F_1(X_1) = \bigodot_{j=1 \dots k} F_j(X_1)$$

$$N = \bigodot_{j=1 \dots s} F_i^{L_i} X_{\langle H_i, W_i, C_i \rangle}$$

*Figure 25: Common ConvNets equation [25].*

The formula in Figure 25 presents for commonly structure of network model base on ConvNets layer, the ConvNets layer  $i$  was defined as an equation  $Y_i = F_i(X_i)$ , where  $F_i$  is an operator,  $Y_i$  is an output of tensor,  $F_i^{L_i}$  presents  $F_i$  layer iterating  $L_i$  times at  $i$  state,  $X_i$  is tensor input, spatial dimensions are  $H_i$  and  $W_i$  as width and height, and the  $C_i$  represents for the number of channels, all of these things make a combination for a list of layers. In fact, the Convnet layers are divided into a plenty of stages and all the layers in each state have the same architecture such as ResNet [25]. EfficientNet expands three dimensions of width ( $C_i$ ), length ( $L_i$ ), and/or resolution ( $W_i, H_i$ ), this action makes no change for  $F_i$ . Based on fixing  $F_i$ , issue of the design for limitations of new resource will be simplified by model scaling, however, there is an existing of a huge space of design for discovering various of these 3 dimensions in every single layer. So, there is a limitation that all layers have to be scaled consistently with unchanged ratio to get a further

reduction of design space. It serves for purpose of maximizing the model accuracy in any case of having any limitation of provided resource. [25]. Formula below shows the optimal solution of the model scaling, where  $(w, d, r)$  are coefficients for scaling network width, depth, and resolution,  $F^i, L^i, H^i, W^i, C^i$  are predefined parameters in baseline network.

The model can achieve better accuracy and efficiency, more specifically this ratio is followed by this principle:

$$\begin{aligned} \text{depth: } d &= \alpha^\phi \\ \text{width: } w &= \beta^\phi \\ \text{resolution: } r &= \gamma^\phi \\ \text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 &\approx 2 \\ \alpha \geq 1, \beta \geq 1, \gamma \geq 1 & \end{aligned}$$

*Figure 26: The setting of coefficient values in compound scaling method [25].*

Where  $\alpha, \beta, \gamma$  are constants that can be affirmative by a tiny grid search. Visionally,  $\phi$  (Phi) is a coefficient able to adjust and increase the number of available resources for model scaling by user's designation, while  $\alpha, \beta, \gamma$  are respective to the size of width, depth and resolution that are additional to the network [25]. The FLOPS of a regular convolution operation is proportional to  $d, w^2, r^2$ , i.e., increasing network depth twice, will make a double increment of FLOPS too, however, if there is a double gaining of network width or resolution, it will be a cause of a quad-increment for FLOPS [25]. That said, convolution operations often prevail over the cost of reckoning in ConvNets, when making a scale for a ConvNet with the above equation, it will nearly improve overall FLOPS by  $(\alpha \cdot \beta^2 \cdot \gamma^2) \phi$ . And the default value of equation was set:  $\alpha \cdot$

$\beta_2 \cdot \gamma_2 \approx 2$  in order to be with any new  $\varphi$ , the total FLOPS will approximately increase by  $2\varphi$ .

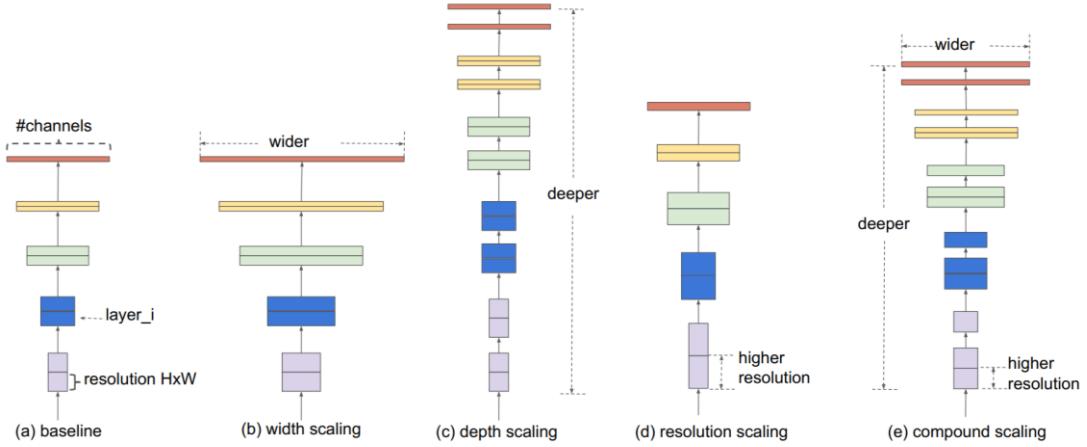


Figure 27: Model Scaling: (a) baseline network; (b) – (d): the networks only scale up one dimension of width, depth, or resolution; (e): compound scaling network has a scaling uniform of 3 dimensions with a fixed ratio [25].

### 7.3. EfficientNet Architecture

Because model scaling did not modify layer operators  $F^i$  of baseline network and had a good baseline network is also critical. The EfficientNet was built based on existing ConvNets with scaling method, in specific, they are MobileNet and ResNet, the mechanism is with a manipulation of a multi-objective neural architecture search, it made an optimization in both accuracy and FLOPS.

Indeed, manipulating  $ACC(m) \times [F \text{ LOPS}(m)=T]$  w as a goal of optimization, where  $ACC(m)$  plays the accuracy and  $LOPS(m)$  plays FLOPS of model m, T is the target FLOPS and w = -0.07 is a

hyperparameter for adjusting the trade-off between accuracy and FLOPS [25]. This compound scaling method was applied with two steps:

- Firstly, there is a fixed  $\varphi = 1$ , assumption of having an availability of double resources and more, then making a small grid search of  $\alpha, \beta, \gamma$  based on Equation 2 and 3. In particular, we find the best values for EfficientNet-B0 are  $\alpha = 1:2, \beta = 1:1, \gamma = 1:15$ , under constraint of  $\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$ .
- Secondly,  $\alpha, \beta, \gamma$  will be fixed as constants and expanded baseline network with different  $\varphi$  using Equation 3, to obtain EfficientNet-B1 to B7.

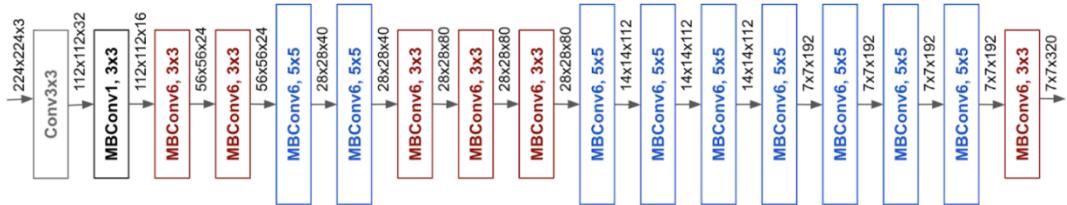
$$\begin{aligned} & \max_{d,w,r} \quad \text{Accuracy}(\mathcal{N}(d, w, r)) \\ \text{s.t.} \quad & \mathcal{N}(d, w, r) = \bigodot_{i=1 \dots s} \hat{\mathcal{F}}_i^{d \cdot \hat{L}_i} (X_{\langle r \cdot \hat{H}_i, r \cdot \hat{W}_i, w \cdot \hat{C}_i \rangle}) \\ & \text{Memory}(\mathcal{N}) \leq \text{target\_memory} \\ & \text{FLOPS}(\mathcal{N}) \leq \text{target\_flops} \end{aligned}$$

*Figure 28: Optimal formula of scaling up method [25].*

Outstandingly, it is feasible to attain even higher performance by looking for  $\alpha; \beta; \gamma$  directly around a large model, but the search cost becomes prohibitively expensive for larger models. This method overcomes this problem by conducting a single search on the tiny baseline network at first step and then using the same coefficients of expanding for all other models at second step [25].

Stage <i>i</i>	Operator $\hat{\mathcal{F}}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels $\hat{C}_i$	#Layers $\hat{L}_i$
1	Conv3x3	$224 \times 224$	32	1
2	MBCConv1, k3x3	$112 \times 112$	16	1
3	MBCConv6, k3x3	$112 \times 112$	24	2
4	MBCConv6, k5x5	$56 \times 56$	40	2
5	MBCConv6, k3x3	$28 \times 28$	80	3
6	MBCConv6, k5x5	$14 \times 14$	112	3
7	MBCConv6, k5x5	$14 \times 14$	192	4
8	MBCConv6, k3x3	$7 \times 7$	320	1
9	Conv1x1 & Pooling & FC	$7 \times 7$	1280	1

Figure 29: Optimal formula of scaling up method [25].



*Figure 30: Efficient Net – B0 [25].*

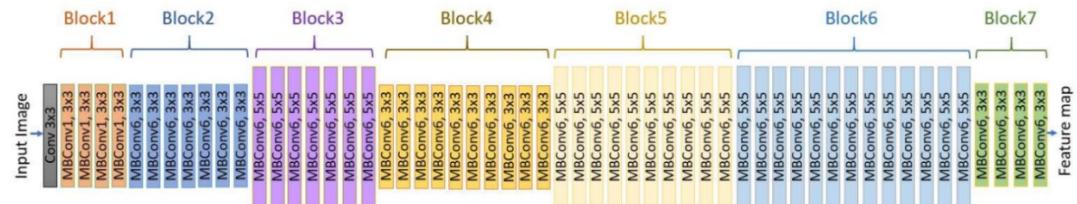


Figure 31: Efficient Net – B7 [31].

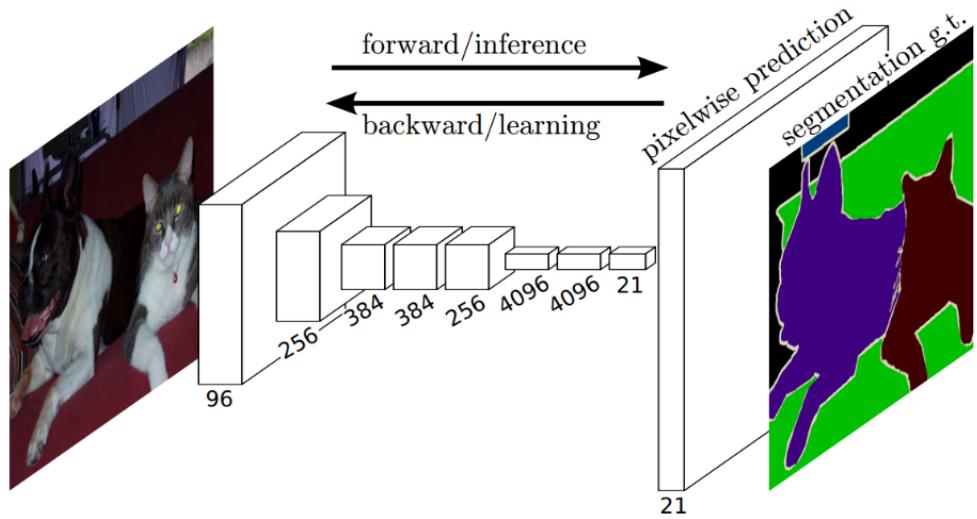
As a result, EfficientNets give out an outstanding performance to other ConvNets. Specifically, EfficientNet-B7 got over the best existing accuracy of GPipe with using 8.4x less parameters than and running 6.1x faster on inference [24], it was about 66 million parameters and 37 billion FLOPS that reach to a proportion of 84.3 percentage of the accuracy. Adding, there was a comparison of ResNet-50 and EfficientNet-B4 that made an improvement of the peak of accuracy from 76.3 percent to 83.0 percent, increasing about 6.7 percent with similar FLOPS.

Model	Top-1 Acc.	Top-5 Acc.	#Params	Ratio-to-EfficientNet	#FLOPs	Ratio-to-EfficientNet
<b>EfficientNet-B0</b>	<b>77.1%</b>	<b>93.3%</b>	<b>5.3M</b>	<b>1x</b>	<b>0.39B</b>	<b>1x</b>
ResNet-50 (He et al., 2016)	76.0%	93.0%	26M	4.9x	4.1B	11x
DenseNet-169 (Huang et al., 2017)	76.2%	93.2%	14M	2.6x	3.5B	8.9x
<b>EfficientNet-B1</b>	<b>79.1%</b>	<b>94.4%</b>	<b>7.8M</b>	<b>1x</b>	<b>0.70B</b>	<b>1x</b>
ResNet-152 (He et al., 2016)	77.8%	93.8%	60M	7.6x	11B	16x
DenseNet-264 (Huang et al., 2017)	77.9%	93.9%	34M	4.3x	6.0B	8.6x
Inception-v3 (Szegedy et al., 2016)	78.8%	94.4%	24M	3.0x	5.7B	8.1x
Xception (Chollet, 2017)	79.0%	94.5%	23M	3.0x	8.4B	12x
<b>EfficientNet-B2</b>	<b>80.1%</b>	<b>94.9%</b>	<b>9.2M</b>	<b>1x</b>	<b>1.0B</b>	<b>1x</b>
Inception-v4 (Szegedy et al., 2017)	80.0%	95.0%	48M	5.2x	13B	13x
Inception-resnet-v2 (Szegedy et al., 2017)	80.1%	95.1%	56M	6.1x	13B	13x
<b>EfficientNet-B3</b>	<b>81.6%</b>	<b>95.7%</b>	<b>12M</b>	<b>1x</b>	<b>1.8B</b>	<b>1x</b>
ResNeXt-101 (Xie et al., 2017)	80.9%	95.6%	84M	7.0x	32B	18x
PolyNet (Zhang et al., 2017)	81.3%	95.8%	92M	7.7x	35B	19x
<b>EfficientNet-B4</b>	<b>82.9%</b>	<b>96.4%</b>	<b>19M</b>	<b>1x</b>	<b>4.2B</b>	<b>1x</b>
SENet (Hu et al., 2018)	82.7%	96.2%	146M	7.7x	42B	10x
NASNet-A (Zoph et al., 2018)	82.7%	96.2%	89M	4.7x	24B	5.7x
AmoebaNet-A (Real et al., 2019)	82.8%	96.1%	87M	4.6x	23B	5.5x
PNASNet (Liu et al., 2018)	82.9%	96.2%	86M	4.5x	23B	6.0x
<b>EfficientNet-B5</b>	<b>83.6%</b>	<b>96.7%</b>	<b>30M</b>	<b>1x</b>	<b>9.9B</b>	<b>1x</b>
AmoebaNet-C (Cubuk et al., 2019)	83.5%	96.5%	155M	5.2x	41B	4.1x
<b>EfficientNet-B6</b>	<b>84.0%</b>	<b>96.8%</b>	<b>43M</b>	<b>1x</b>	<b>19B</b>	<b>1x</b>
<b>EfficientNet-B7</b>	<b>84.3%</b>	<b>97.0%</b>	<b>66M</b>	<b>1x</b>	<b>37B</b>	<b>1x</b>
GPipe (Huang et al., 2018)	84.3%	97.0%	557M	8.4x	-	-

Figure 32: EfficientNets Performance Result on ImageNet [25].

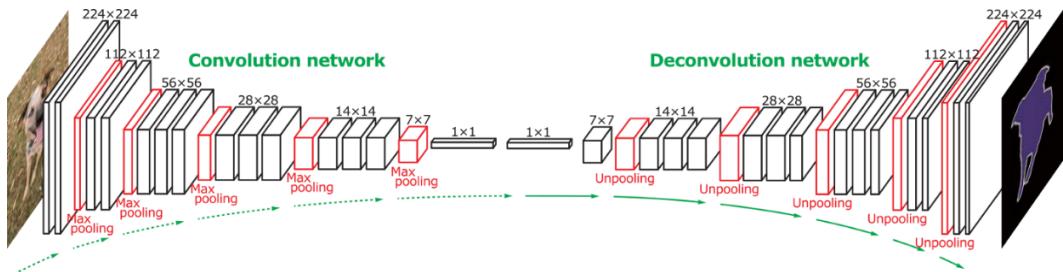
## 8. Fully Convolutional Networks

One of the turning points in the field of image detection was the introduction of the Fully Convolutional Network (FCN), which is known as a variant of CNN for image segmentation work [10]. The FCN has a difference from ordinary CNNs in which the fully connected layers are replaced by convolution layers at the end of the network. So, there is a generation of a network making a computation of a nonlinear filter for each layer's output vectors. As a result, not only can the completed network function on an input of any size and make a production of an output with respective spatial dimensions, but this also allows the classification network to generate a heatmap of the class of desire. In addition, adding layers and a spatial loss to the network results in an efficient machine for end-to-end dense learning [10]. Figure 33 shows the transformation of FCN.



*Figure 33: Fully Convolutional Network serves for semantic segmentation tasks [10].*

The use of FCNs makes introduces allow to take multiple input image sizes, and to have a better efficiency for dense predictions of learning via in-network up-sampling [10]. Furthermore, the FCN may preserve spatial information from the input that is very important for semantic segmentation task because of involving both localization and classification for the task. Although an FCN has an ability to accept any size input image, there is a decrease of the output resolution by being applied convolutions without any padding. These were added to make filters minimal and requirements of calculation manageable. As a result, the coarse output is reduced in size by a factor equal to the pixel stride of the output unit's receptive field [10].

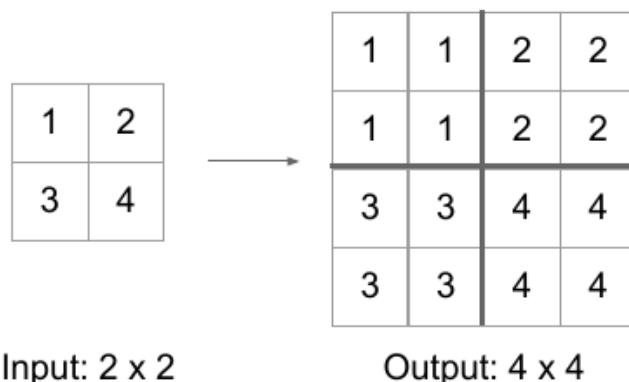


*Figure 34: Fully Convolutional Network Architecture [32]*

In other words, the FCN is divided in 2 parts that are convolution network (down-sample) as an encoder step and the following is deconvolution network (up-sample) known as decoder step, the convolution network part is very similar with casual CNN architecture so we will go deep down to up-sample part.

In Semantic Segmentation, after using the method of convoluted features extraction to decrease the input's resolution, the following step is to upscale the low resolution turning back to the original resolution of the input image. Pooling is the process of converting a bunch of values to a solitary value, whereas up-sample is the process of converting a solitary value to a bunch of values and it can be named Unpooling, like Pooling, the unpooling can be carried out in a variety of ways [26].

- **Nearest Neighbour:** Choosing a value and populating it into the surrounding cells (number of cells depending on the increment in resolution), figure below shows the copy process of all values in matrix  $2 \times 2$  to every cell with size of  $2 \times 2$  square in new matrix  $4 \times 4$  [26].



*Figure 35: Nearest Neighbour Unpooling [26].*

- **Bed of Nails:** Inserting the value in turn at a predefined cell and adding the rest with 0, the figure below is an example of filling the value from the top-left cell of each  $2 \times 2$  square in new matrix with size  $4 \times 4$  [26].

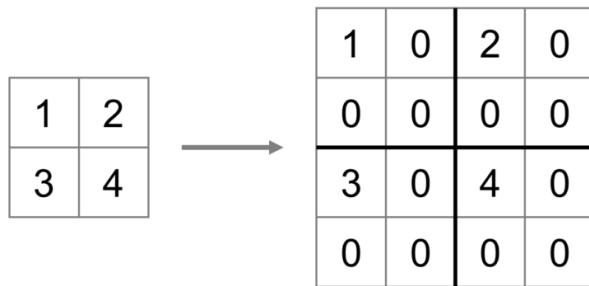


Figure 36: Bed of Nails Unpooling [26].

- **Max Unpooling:** In CNN, the max-pooling layer pulls the maximum values from the picture regions covered by the filter in order to down-sample the data, and the unpooling layer provides the value to the location from which the values were picked up in order to up-sample the data [26]. However, these approaches are not data-dependent, which means they are able to learn nothing from data; they are just programmed computations making them become as a specific task and costly in terms of calculation timing. The generic answer for these situations is transposed convolutions.

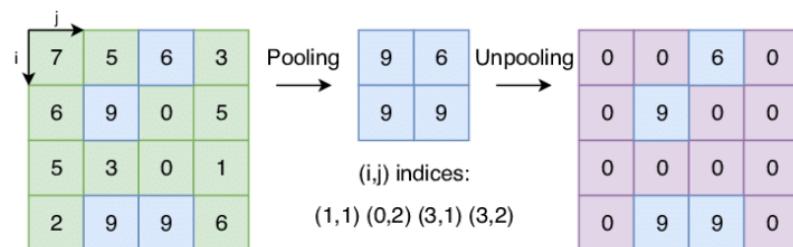
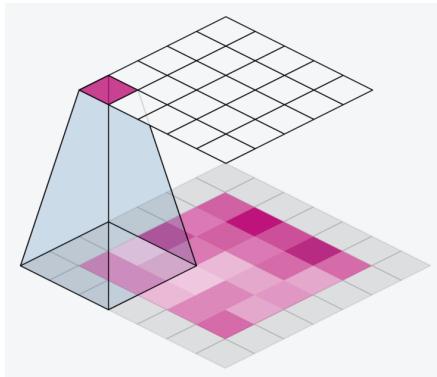


Figure 37: Max Unpooling [26].

Transposed Convolution, also known wrongly as Deconvolution, is the inverse action of Convolution. A transposed convolution layer, as opposed to convolution, is used to up-sample the decreased resolution feature back to its original resolution. To obtain the final output from the lower resolution features, a set of strides and padding values is learned. Transposed convolution is superior than other Up-sampling techniques because of being unlike earlier techniques (Nearest Neighbour, Bed of Nails, Max Unpooling), the transposed convolution is a learnable up-sampling technique being highly efficient. The image below describes the method in a very simple way [26].



*Figure 38: Transposed Convolution [26].*

Next, because of an existing main problem with in-network down-sampling in an FCN, which is the reduction of the resolution of the input by a huge proportion, making it impossible to recreate finer features during up-sampling even with complicated techniques like Transpose Convolution. As a result of getting a coarse output. One solution is to include "skip connections" in the up-sampling from previous layers stage and then making a sum of total the two feature maps. These skip connections supply enough information to following layers to produce correct borders of segmentation. This combination of fine and coarse layers yields local forecasts with virtually accurate global (spatial)

structure. The skip connection learns how to integrate coarse and fine layer information. So, there is a performance of grids with varying degrees of spatial coarseness for layers, with the intermediate convolution layers of FCN deleted for clarity [26].

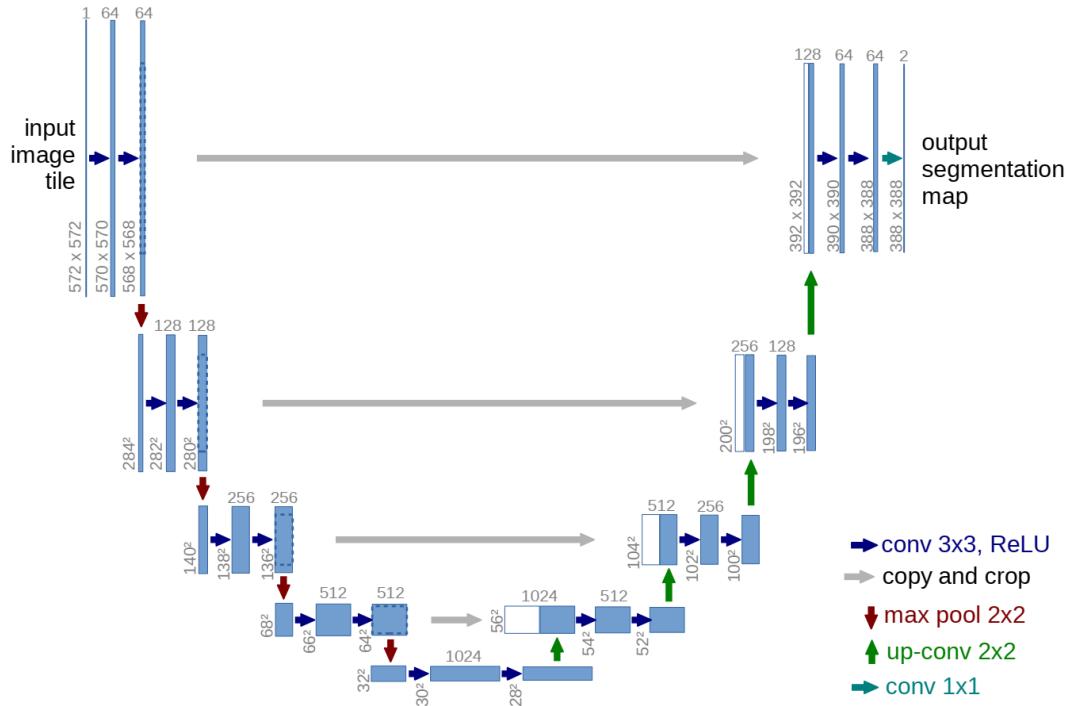
## 9. U-Net

The U-Net architecture is a variant of the FCN that is commonly used for semantic segmentation duties. Ronneberger et al. created the first U-Net architecture in 2015 to conduct picture segmentation and localization for biological applications [27]. U-Net is superior to traditional CNNs since they can give localization as well as classification in their output. In this scenario, localization implies labelling each single pixel from an image with a certain class. It is also more favorable than FCN because U-Net can function with less training photos while producing more exact segmentations. This is accomplished by constructing up-sampling layers with a huge amount of feature channels allowing context information to be propagated to layers having better resolution. Additionally, the segmentation networks of U-Net are the most appropriate for this application because they are the most easily scalable and sizeable FCN. Specially, they can avoid a compromise between localization and contextual information, and are frequently employed in cutting-edge semantic segmentation approaches [27].

### 9.1. U-Net Architecture

The U-Net architecture is made up of a path of contraction at the input and a path of expansion at the output. Ronneburger's paper introduces the basic architecture of this network, the first path is the contraction path, also known as the encoder path, which is essentially a

stack of convolution, activation, and pooling layers designed to take the context from the input image. The encoder output, which is smaller than the input, is gradually expanded in the expansion path. With transposed convolutions, the expansion path has another name is decoder path which allows for exact localization. Through a series of up-convolutions and concatenation with appropriate feature maps from the contracting path, the expansion pathway integrates high level features with spatial information. The architecture is depicted in figure 39 [27].



*Figure 39: U-Net Architecture (The representation of a multi-channel feature map is each blue box. The number of channels is indicated on the box's top. The x-y size is specified at the box's lower left side. White boxes represent feature maps that have been replicated. The arrows represent the various operations) [27].*

The encoder route of U-Net is typical CNN design, including recurring 3x3 convolutions and operations of max pooling with stride 2 for down-sampling which double amount of feature channels before appearing of up-sampling the feature map and performing a 2x2 convolution on the expanding route. In detail, these actions remove the feature channels amount in a half before concatenating them with the feature map being cropped from the route of contraction, this concatenation of the U-Net is the different point with the skip connection of FCN. After that, an extra convolution layer is added to map each feature vector to the expected number of classes [27]. The U-Net can collect image context by employing the contracting route, on the other side, the symmetric expanding route allows for exact localization. By integrating the characteristics with high resolution which come from the contracting path being up-sampled output, localization is activated. Added, the propagation of context information to the layers with higher resolution is allowed by the model's up-sampling section with a huge amount of feature channels [27].

Based on the feature maps with low level from the encoder contain more spatial information which is helpful to the analysis of complicated scenes containing multiple objects and their relative configurations, there is a combination of the feature maps at intermediate low level from EfficientNet or ResNet and the feature maps at intermediate high level from U-Net decoder. The network is able to transfer context information to higher resolution layers due to the huge amount of feature channels in the up-sampling section. Because of the symmetric characteristics in conventional U-Net. In my work, I have tried to use EfficientNets and ResNets as the encoders in contraction route instead of regular set of

convolution layers. The decoder module is similar to the primordial U-Net. The detailed implementation for the combination of these architectures will be mentioned at the following section named “implementation”.

### **III. Implementation**

#### **1. Dataset – Image Processing**

As I mentioned in the previous section, this thesis I used the dataset provided by Severstal, it is a dataset about images of steel containing defects, specifically here there are 4 types of defects, set of the data is divided into 2 main parts including images for train set and test set, and csv files containing content about pre-classification of defects on each image. Specifically, the two image folders of train set and test set are RGB images with 1600 x 256 pixels in size, they are both photos of manufactured steel and the content of the defects on them is not the same, most of them contain 1 type of defect, and a small part contains 2 types of defects. The total number of images in the data set is more than 12k5 records, while the images in the test set are more than 5k5 records. Besides, the csv file is a file that summarizes the information of the image data file, more specifically, in the data set provided by the provider, there are 2 csv files, 1 train.csv file carries the information of the train set including the train set. includes columns of information about file name (ImageId), type of defect (ClassId), pixels on image are defects (EncodedPixels), while the remaining csv file has the same structure but no information about EncodedPixels and the ClassId in the test set from the provider. Therefore, the requirement here is to build a model based on the train set for the purpose of detecting defects on the test set without given ground-truth masks to verify.

EncodedPixels is an encoder method that minimizes the amount of information and memory size of pixel information containing defects in an image. Since an RGB image is essentially formed by 3 matrices stacked on top of each other in layers or channels, and the values in each matrix represent each pixel on the image, and each value has specific index certainly, so the provider encoded defect information based on starting position of pixels and run length based on these matrices, which has format [pixel index] [run length] e.g., 10 25 100 200, it means there are 2 defects on this image, which runs from 10<sup>th</sup> pixel to 35<sup>th</sup> pixel (10 + 25) and the second one is from 100<sup>th</sup> to 300<sup>th</sup> (100 + 200).

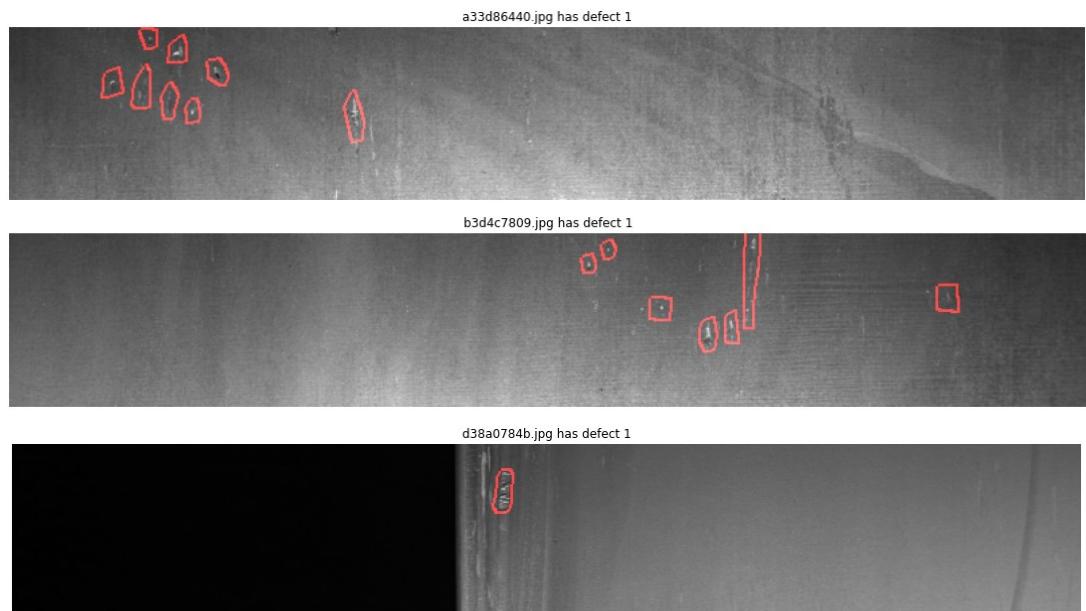


Figure 43: Defect type 1 with red border

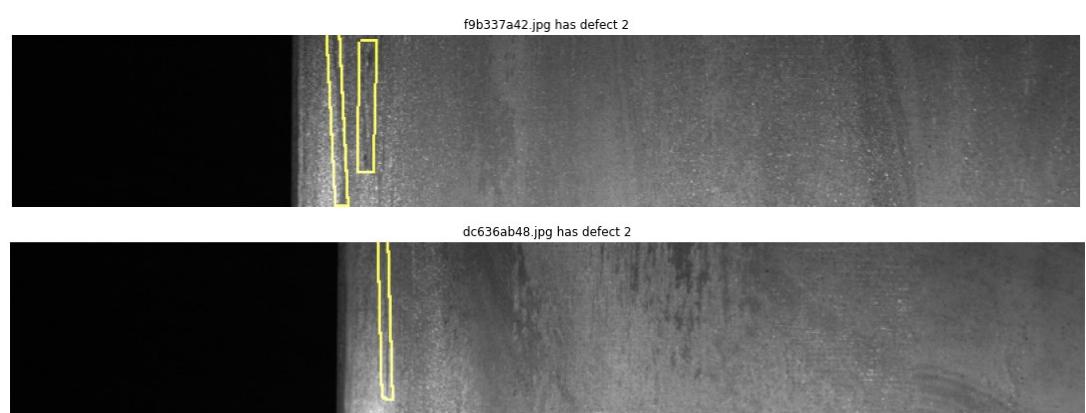




Figure 44: Defect type 2 with yellow border

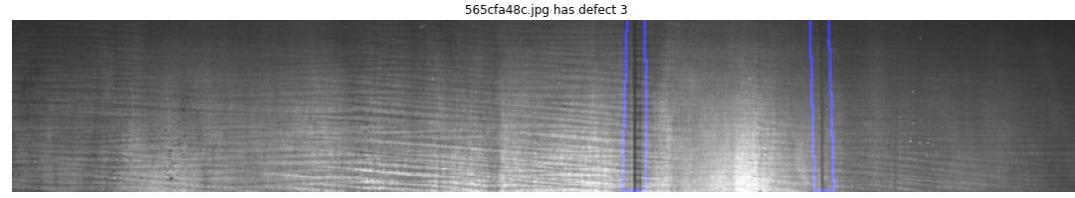
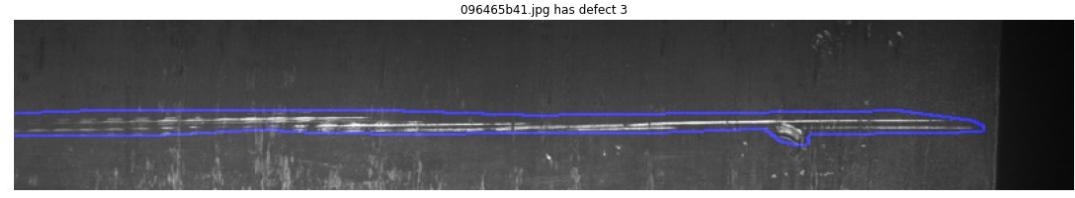


Figure 45: Defect type 3 with blue border

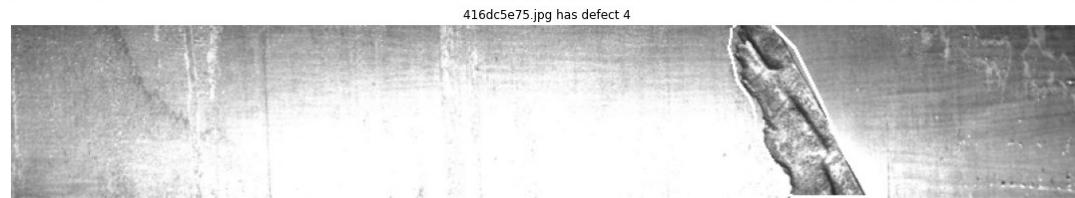
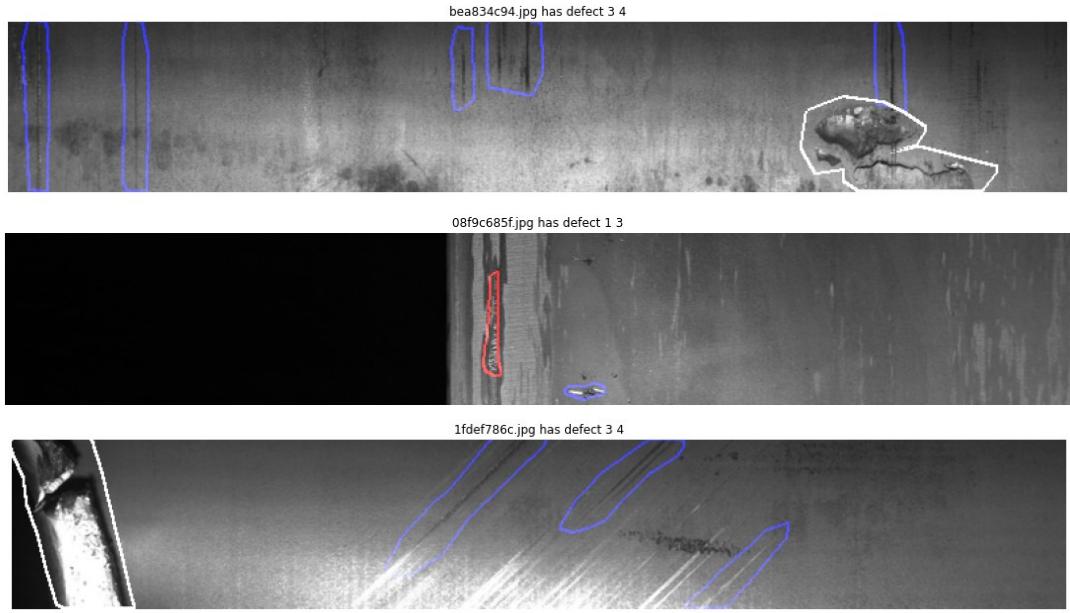


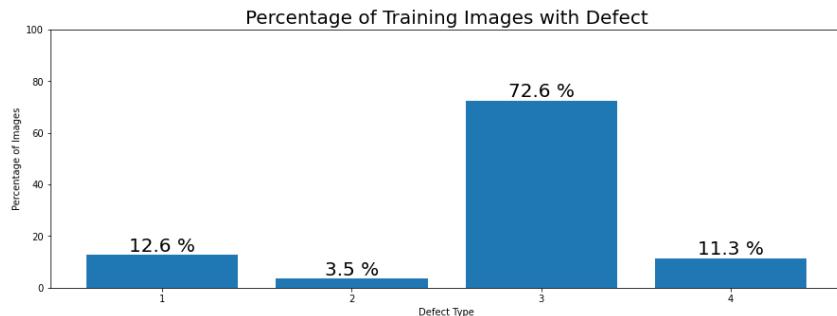
Figure 46: Defect type 4 with white border



*Figure 47: Image with multiple defect types*

### 1.1. Data Augmentation

Train set with more than 12k5 records, but in reality, the number of defects is not balanced and the number of records with defects in the csv file is just over 6k6. In this data set, the number of type 3 defects accounts for the majority compared to other types of defects, and type 2 defects account for the smallest number in the data set, specifically, in total counting of both single and multiple defects on image, type 1 defects account for 12.6% with 897 records, type 2 accounted for 3.5% with 247 records, type 3 accounted for 72.6% of 5150 records, and type 4 accounted for 11.3% of 801 records. That is reflected in Figure below.



*Figure 48: Percentage of defect types include in train data*

For that reason, I applied data augmentation method to be able to help the amount of data types of defects become an equilibrium. Because, when the labels in the data for training model are not balanced, they will have a great influence on the accuracy of the model, in this case it will detect and predict defect type 3 very well, and so bad at defect type 2, and there is a relative level in the other two categories. Therefore, balancing the data is a very necessary job to improve the quality of the model.

In detail, with a smaller number in categories 1, 2, 4 I used the method of rotating 180 degrees, 10 degrees and flipping the image to increase the number of records of these types, besides, I dropped given a certain number of type 3 defects, I ended up getting a result number of defects appearing in the train set at equilibrium with 1291 records in both single and multiple defects appearing in image.

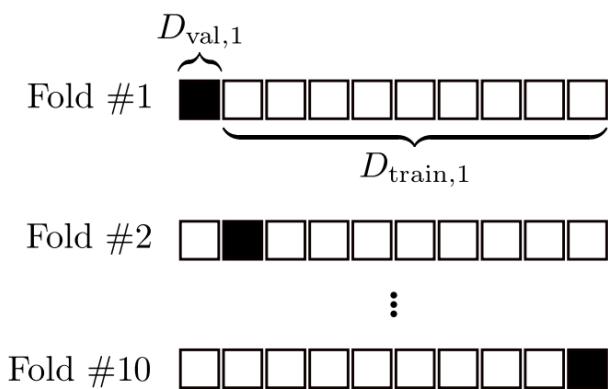
Before implementing the augmentation method, I retrieved 25 records for each single defect type and 25 image records containing multiple defect types for the test set, totaling 250 image records. The reason is that the test set provided by the supplier does not have ground-truth masks data, resulting in no data source to re-evaluate the model with the test set after conducting the training. Therefore, with these 250 records with ground-truth masks, it is easy for me to re-evaluate the model after training.

## 1.2. K-fold cross-validation

With the train set balanced after the augmentation method, I applied K-fold strategy to them during the training model to increase the efficiency of the model.

K-fold cross-validation divides the available learning set into k disjoint subsets of roughly similar size. The number of resulting subsets is denoted by the term "fold", which this partitioning is accomplished by randomly picking examples from the learning set without replacing them. The model is trained using k - 1 subsets that represent the training set as a whole. The model is then applied to the remaining subset, known as validation set, and its performance is evaluated [30].

This approach is repeated until all k subsets have served as validation sets. The cross-validated performance is the average of the k performance measurements across the k validation sets. This process is depicted in Figure below for k = 10, i.e., 10 - fold cross-validation. The first subset in the first fold acts as the validation set  $D_{val,1}$ , while the following nine subsets serve as the training set  $D_{train,1}$ . The validation set is the second subset in the second fold, while the other subsets are the training set, and so on [30].



*Figure 49: K – Fold Cross-Validation [30].*

### 1.3. Data Generation

With a certain amount of data for training and validation work, I have created a data generator class for the purpose of matching and

optimizing the hardware in the training model, which has function generating an output with processed image and ground-truth masks. Specifically, the processed image, which is a reduction of the input image size by half from 1600 x 256 pixel to 800 x 128 pixel, dividing the number of input image into 16 batch sizes for each training epoch, moreover, the image has been normalized from RGB format with range 0 - 255 to grayscale format with range 0 – 1 that helps model easy to make a decision of being pixel of defect or not. Then the masks are 4 matrices with size 800 x 128 equivalent to the image, each matrix represents a defect type contained in the image, in which the pixel defined as defect having the value 1 and other ones without defect having a value of 0. These matrices are generated by decoding information from the EncodedPixels provided by the provider.

## 2. Model

Basing on the achievements of state-of-the-art networks such as ResNet, EfficientNet, and U-Net. I have used two separate combinations: the combination of ResNet and U-Net, and the combination of EfficientNet and U-Net, where U-Net was the architecture acting as the main framework of the model, and Resnet and EfficientNet played a role as the backbone of the model, they replaced the original encoder part of the U-Net architecture, with this combination there were models with a huge amount of parameters created, and with skip connections in the ResNet Blocks and compound scaling blocks in EfficientNet made it possible to train the model more deeply than the original one, which gave the model a very significant improvement.

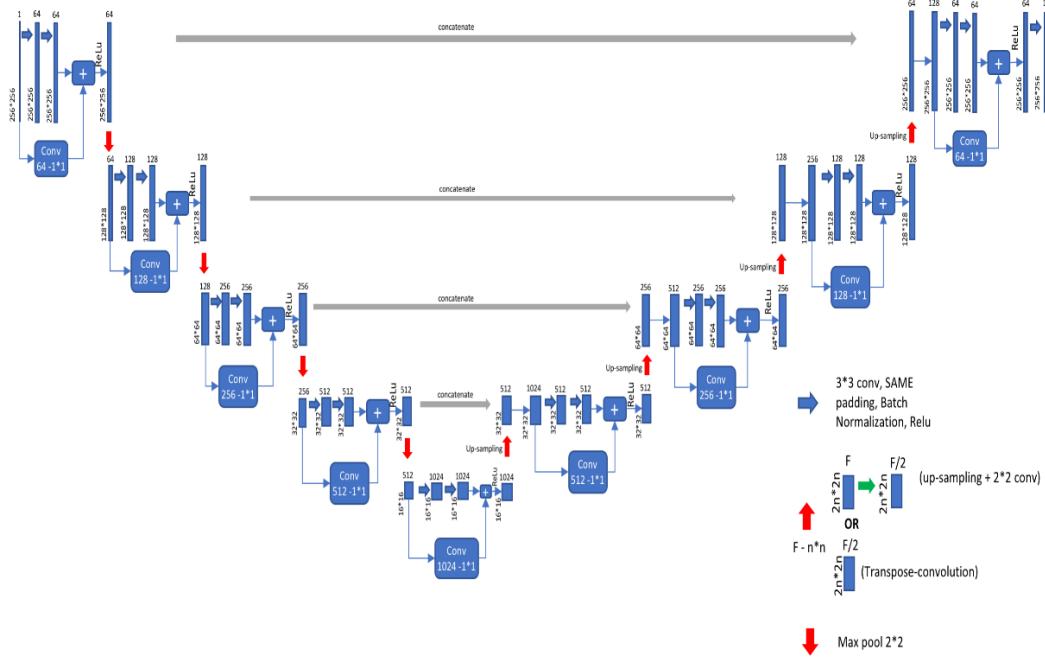


Figure 50: U-Net with ResNet as the backbone [34].

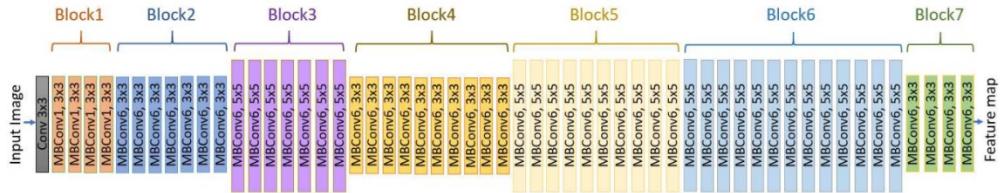


Figure 51: EfficientNet-B7 architecture [31].

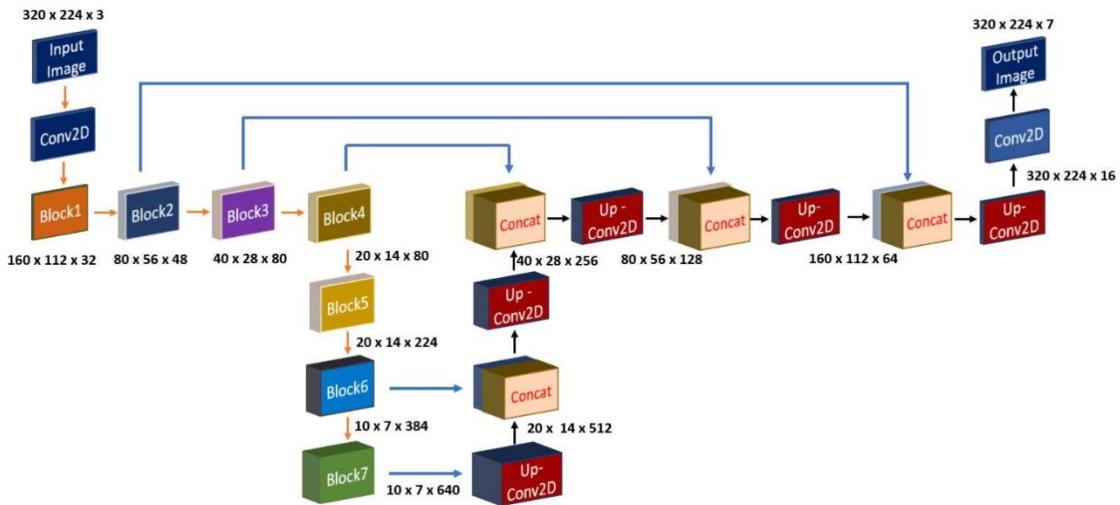


Figure 52: U-Net with EfficientNet-B7 as backbone [31].

## 2.1. Evaluation Metrics

As the evaluation matrix for this thesis, I used the mean Dice coefficient that can be used to make a comparison about the pixel-wise agreement between a predicted segmentation and its matching ground truth. The dice coefficient is calculated by counting the similar pixels (taking intersections that appear in both images) in the two images being compared, multiplying it by 2, and dividing it by the total pixels of both images. The formula is shown in figure 53:

$$\frac{2 * |X \cap Y|}{|X| + |Y|}$$

Figure 53: Dice – Coefficient [36].

where X is the predicted set of pixels and Y is the ground truth.

## 2.2. Loss function

For the purpose of detecting which pixels are defect or not on the input image, I have used the loss function of binary cross entropy for this thesis, with this loss function operating on the returned real numeric values in range 0 - 1, with values  $\geq 0.5$  will be considered as defects, and vice versa values  $< 0.5$  will be considered as background.

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

Figure 54: Binary cross entropy [35].

where y is the label (1 for pixels of defect and 0 for pixels of background) and p(y) is the predicted probability of the pixel being defect for all N pixels [35]. The formula tells that for each pixel of defect

( $y = 1$ ), it adds  $\log(p(y))$  to the loss, that is, the log probability of it being defect. Conversely, it adds  $\log(1 - p(y))$ , which is the log probability of it being background, for each pixel of background ( $y = 0$ ) [35].

### **2.3. Execution Environment**

The thesis is used Colab Pro platform of Google, it is an application of cloud computing running Linux operating system, which provides a very useful tool with code cells and run time of python with version 3.7, quite similar to Anaconda notebook running on local machine. The platform offers quite a powerful hardware for ML as well as DL research, with GPU of NVIDIA TESLA P100 16GB, 27.3 GB for memory, and 168 GB for the local storage. This specification of hardware helps save training time and supply enough space for storing the data in research.

### **2.4. Model Training Process**

Implementing this idea, I have done experiments on several versions of combination of ResNet and U-Net, and some versions of combination of EfficientNet and U-Net. I have used transfer learning for all of the above architectures based on "Segmentation Models" library version 1.0.1, specifically with ResNet, I have used ResNet-34, ResNet-50, and ResNet-101 as backbone. On the opposite side, with EfficientNet, I use instances B0, B1, B3, B5 and B7 as the backbone of the model.

All versions of the model use the same loss function of binary cross entropy, and dice-coefficient as an evaluation matrix. Since the requirement is to identify 4 different types of defects, the number of

classes needs to be segmented accordingly with the total number of defect types. The output layer of the model has a shape of 128 x 800 x 4 corresponding to the height, width, and number of defects, which is applied the sigmoid function as the activation function. Furthermore, in order to help the model quickly converge, I have used pre-trained weights of "imagenet" corresponding to each type of backbone used, and the weights of this model encoder are frozen during training. Therefore, there are a large number of parameters in the encoder part that cannot be trained and only the remaining parameters in the decoder part are trained. Adding, the optimization method I use for the model is Adam.

In the training process, I used the K - fold method, with K = 5 to create a training and validation set with a ratio of 80 - 20 respectively in the total training data. The models are trained this bunch of data in 5 folds, each fold made a duty on running through 150 epochs, however, I implemented EarlyStopping callback with monitoring the change of "val\_loss" value in range of 15 for patience value with mode is "min" aims to stop the training work when the model has reached the convergence point and there is no possibility of further improvement. With this selection of implementation, regular models are trained through a total of 90 - 120 epochs to reach convergence and for the EarlyStopping callback is executed.

And these implementations all yield positive and slightly different results.

## IV. Result

After several experiments training the models, I have obtained relatively positive results. Specifically, after training for about 90 - 120 epochs out of a total of 5 folds, I went to evaluate the model with the validation set and the test set that collected from the original training set with 250 records. The average train loss and validation loss values of binary cross entropy accounted for the range of 0.0017 - 0.0031, because using the k-fold cross validation tactic which leads to the fact that the values of loss and validation loss are approximately equal, and this also happens with the average values of the evaluation matrix of dice coefficient in training and validation, specifically they accounted from 0.904 to 0.941. It is shown in the Figures 56.

U-Net Model with backbone	Time Execution (minutes)	Average Loss in 5 - folds	Average Val Loss in 5 - folds	Average Dice in 5 - folds	Average Val Dice in 5 - folds	Test loss	Test dice
ResNet-34	<b>104.75</b>	<b>0.00311</b>	<b>0.00311</b>	<b>0.904</b>	<b>0.904</b>	<b>0.033</b>	0.692
ResNet-50	147.91	0.00194	0.00194	0.938	0.938	0.038	0.701
ResNet-101	182.82	0.00252	0.0025	0.926	0.928	0.036	0.701
EfficientNet-B0	117.58	0.00254	0.00254	0.927	0.926	<b>0.047</b>	<b>0.668</b>
EfficientNet-B1	126.39	0.0023	0.0023	0.927	0.927	0.039	0.696
EfficientNet-B3	162.6	0.0022	0.0023	0.924	0.927	0.034	0.714
EfficientNet-B5	255.26	<b>0.0017</b>	<b>0.0018</b>	<b>0.941</b>	<b>0.94</b>	0.045	0.711
EfficientNet-B7	<b>389.69</b>	0.0023	0.0024	0.926	0.925	0.034	<b>0.714</b>

*Figure 55: Result of training models with evaluation values.*

U-Net Model with backbone	Total Parameters	Trainable Parameters
ResNet-34	24,456,589	<b>3,167,495</b>
ResNet-50	32,561,549	9,059,079
ResNet-101	51,605,901	9,111,303
EfficientNet-B0	<b>10,115,936</b>	6,106,404
EfficientNet-B1	12,641,604	6,126,436
EfficientNet-B3	17,868,268	7,170,052
EfficientNet-B5	37,469,108	9,126,340
EfficientNet-B7	<b>75,048,532</b>	<b>11,259,588</b>

*Figure 56: Model Parameters.*

Based on figure 56 we can see, the combination of ResNet-34 and U-Net gives the least efficient results, with results of the average loss of training set and validation set, and loss of test set were 0.00311, 0.00311, and 0.033 respectively. In addition, the average training and validation of dice coefficient had a value of 0.904, while the test dice coefficient came in at 0.692, which reached almost the bottom of the results. In contrast, the model with the EfficientNet-B5 backbone had the best results, with the average loss of training set and validation set, and loss of test set values of 0.0017, 0.0018, and 0.045, respectively. Then, not only the accuracy of dice coefficient in both training and validation accounted for 0.941 and 0.94, but also the test dice coefficient value reached to 0.711 being almost highest result of value. Looking at the results of the test dataset with 250 records, the model using EfficientNet-B0 gave the worst result with 0.047 for the test loss and 0.668 for the test dice coefficient. Oppositely, the model with EfficientNet-B7 backbone gave the best and most optimal results for loss test and dice coefficient test, which accounted for values of 0.034 and 0.714 respectively.

Furthermore, when looking at both Figure 55 and Figure 56, we see that the number of trainable parameters between the models were not the same, with the models with the fewest number of trainable parameters giving the lowest results and the least time cost for training, specifically, ResNet-34 backbone gave us more than 3M trainable parameters and only took more than 104 minutes to perform the training task. Whereas EfficientNet-B7 backbone was the choice that consumes the most parameters and execution time and also brought the best results for the test set, with more than 11M trainable parameters they took more than 389 minutes to train the model.

The next step is to evaluate the performance of the classification model, which is to use the "accuracy\_score" and "multilabel\_confusion\_matrix" functions provided by the "Sklearn" library. With the "accuracy\_score" function, I performed the aggregation of the image's ground-truth masks as a value of the "y\_true" parameter and synthesized the predicted masks of each image as a value of the "y\_pred" parameter, then I used them to perform a general evaluation of accuracy for the test data set with multiple defect types, which is supposed to be how many defect types appear on an image, used its results to calculate an average result of accuracy score. And then it was an evaluation step for each defect type separately. After that, with the use of "multilabel\_confusion\_matrix" function to help give the number of matrices corresponding to the number of defect types to be predicted, here with 2 parameters "y\_pred" and "y\_true" used in this function returned the result as 4 matrices, in each matrix encloses the number of records of true positive (TP), true negative (TN), false positive (FP), and

false negative (FN) results on the total number of records of the test data set. Specific results are shown in Figure 57.

U-Net Model with backbone	Accuracy classification score for total defects (250 records)	Accuracy classification score in each defect (250 * 4 = 1000 records)	Confusion Matrix in defect type 1 [TN FP FN TP] [Accuracy]	Confusion Matrix in defect type 2 [TN FP FN TP] [Accuracy]	Confusion Matrix in defect type 3 [TN FP FN TP] [Accuracy]	Confusion Matrix in defect type 4 [TN FP FN TP] [Accuracy]
ResNet-34	0.7	0.915	167 15 12 56 (0.892)	189 1 13 47 (0.944)	129 28 13 80 (0.836)	168 2 1 79 (0.988)
ResNet-50	0.764	0.933	172 10 12 56 (0.912)	190 0 14 46 (0.944)	139 18 11 82 (0.884)	169 1 12 79 (0.992)
ResNet-101	0.756	0.927	173 9 11 57 (0.92)	187 3 13 47 (0.936)	138 19 15 78 (0.864)	168 2 1 79 (0.988)
EfficientNet-B0	0.748	0.928	173 9 17 51 (0.896)	189 1 9 51 (0.96)	139 18 15 78 (0.868)	168 2 1 79 (0.988)
EfficientNet-B1	0.772	0.933	177 5 15 53 (0.92)	188 2 8 52 (0.96)	141 16 17 76 (0.832)	167 3 1 79 (0.984)
EfficientNet-B3	0.804	0.943	178 4 15 53 (0.924)	188 2 8 52 (0.96)	140 17 10 83 (0.892)	169 1 0 80 (0.996)
EfficientNet-B5	0.82	0.949	177 5 13 55 (0.928)	189 1 7 53 (0.968)	144 13 10 83 (0.908)	169 1 1 79 (0.992)
EfficientNet-B7	0.804	0.943	178 4 15 53 (0.924)	188 2 8 52 (0.96)	140 17 10 83 (0.892)	169 1 0 80 (0.996)

*Figure 57: Accuracy Scores and Confusion Matrix Result of Training Models.*

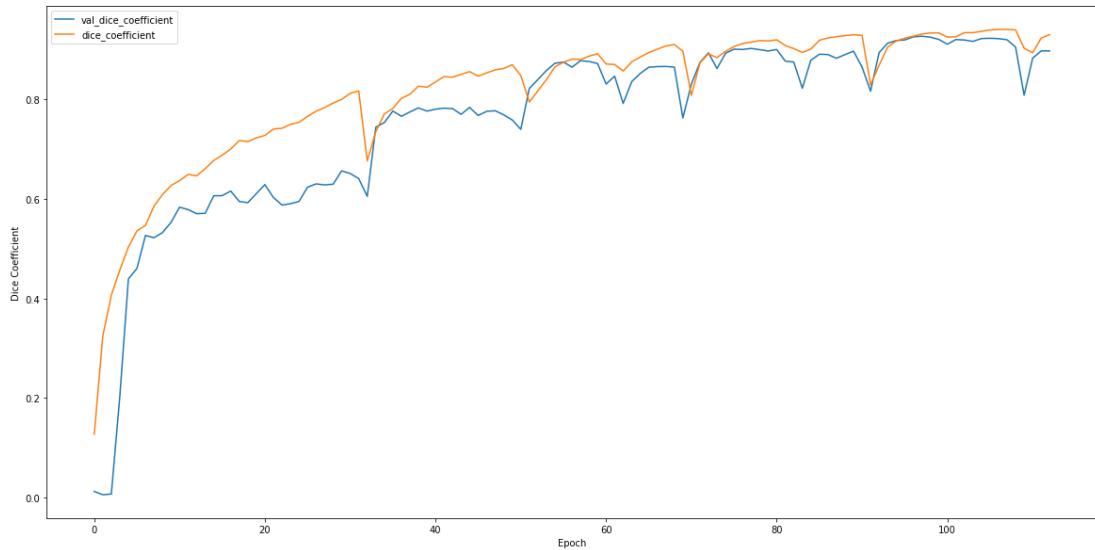
Looking at Figure 57, the model using ResNet-34 as the backbone has the lowest accuracy score for the total number of defect types as well as for each defect type, with results of 0.7 and 0.915 respectively. Meanwhile, the model using EfficientNet-B5 as backbone has the best accuracy score with 0.82 for total defect types, and 0.949 for each defect type. In addition to these 2 options, we have 2 other options with the

same very good results, these are EfficientNet-B3 and EfficientNet-B7 as backbone with 2 pairs of results being 0.804 - 0.943 for accuracy score of total defect types and each defect type.

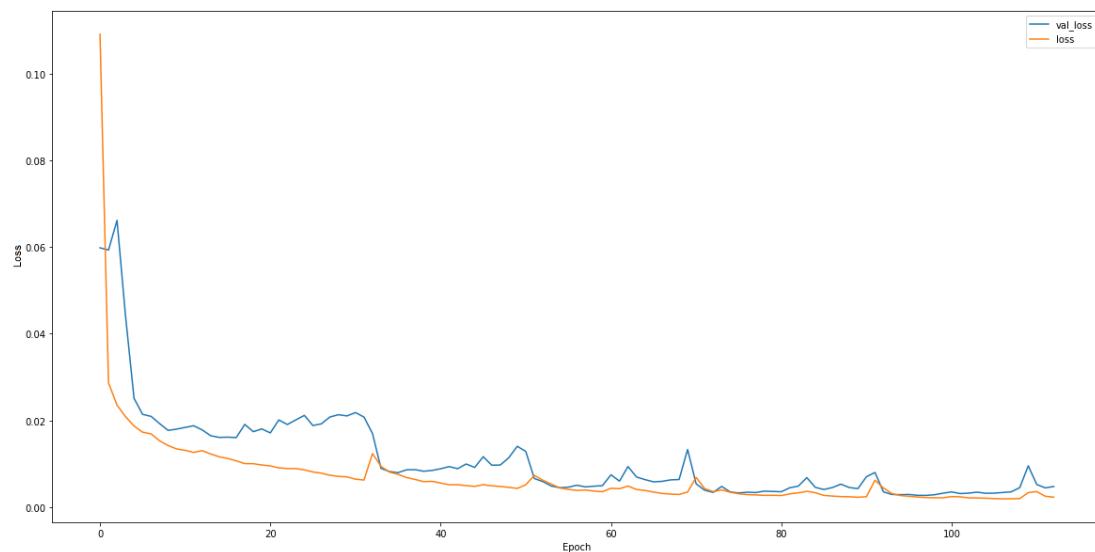
Subsequently, the results of using confusion matrix for model evaluation, overall, all the models have a pretty good result above 0.83. Specifically, there was the lowest result in defect type 3 with the values ranging between 0.832 - 0.908. However, the prediction of defect types 1, 2 and 4 yielded more impressive results, with types 1 and type 2 giving fluctuating results in the range of 0.89 - 0.96, in particular, type 4 gave excellent results at approximately above 0.98.

Based on the above results, we can conclude that, with the combination of ResNet-34 and U-Net gave a less optimal result, on the contrary, using EfficientNets-B5 and U-Net gave the most optimal results with the highest and closest accuracy values in these versions and a reasonable execution time, it spent 255.26 minutes for training model work. Although the use of EfficientNet-B7 got a slightly better result, we had to trade off the execution time of up to 389.69 minutes for the training model.

The Figure 58, 59 show the results of dice coefficient and loss in both training and validation data set in the combination of ResNet-101 and U-Net model.

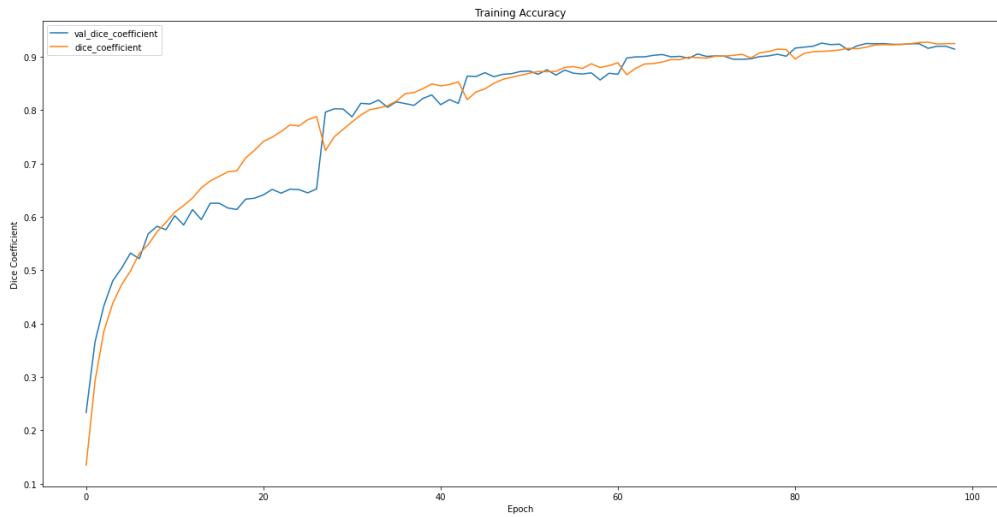


*Figure 58: Dice coefficient of training and validation in ResNet-101 - U-Net model*

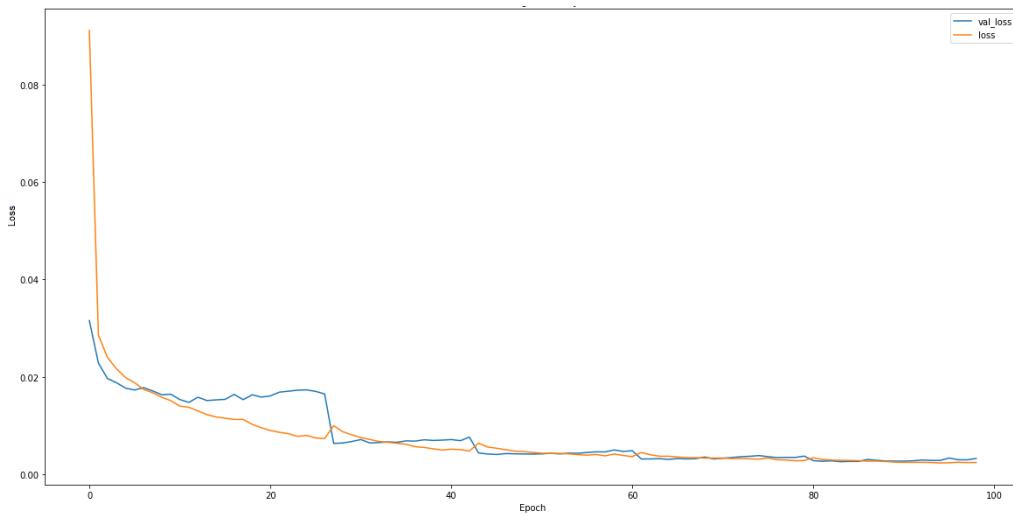


*Figure 59: Loss of training and validation in ResNet-101 - U-Net model*

The Figure 60, 61 show the results of dice coefficient and loss in both training and validation data set in the combination of EfficientNet-B5 and U-Net model.

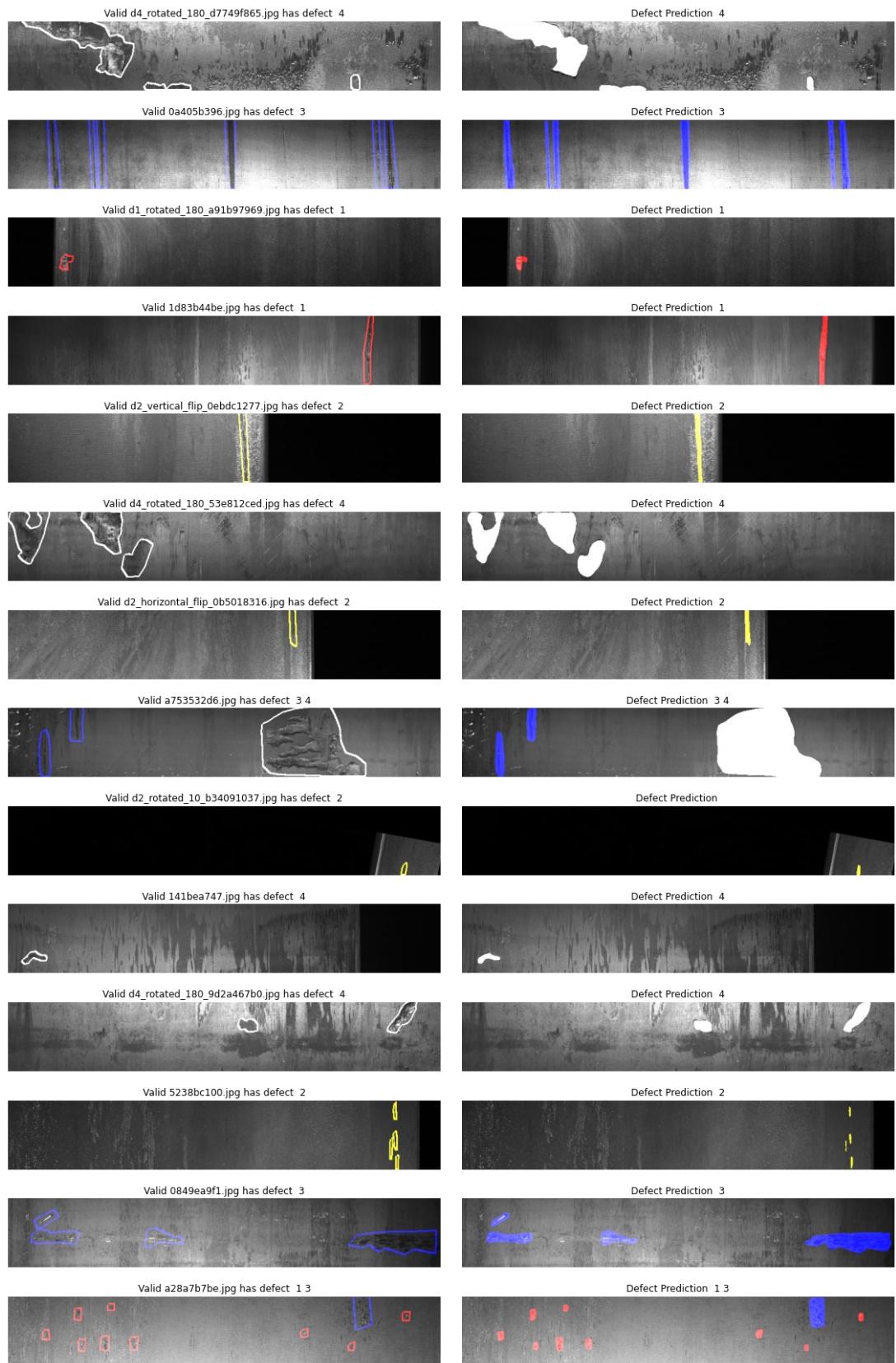


*Figure 60: Dice coefficient of training and validation in EfficientNet-B5 - U-Net model*

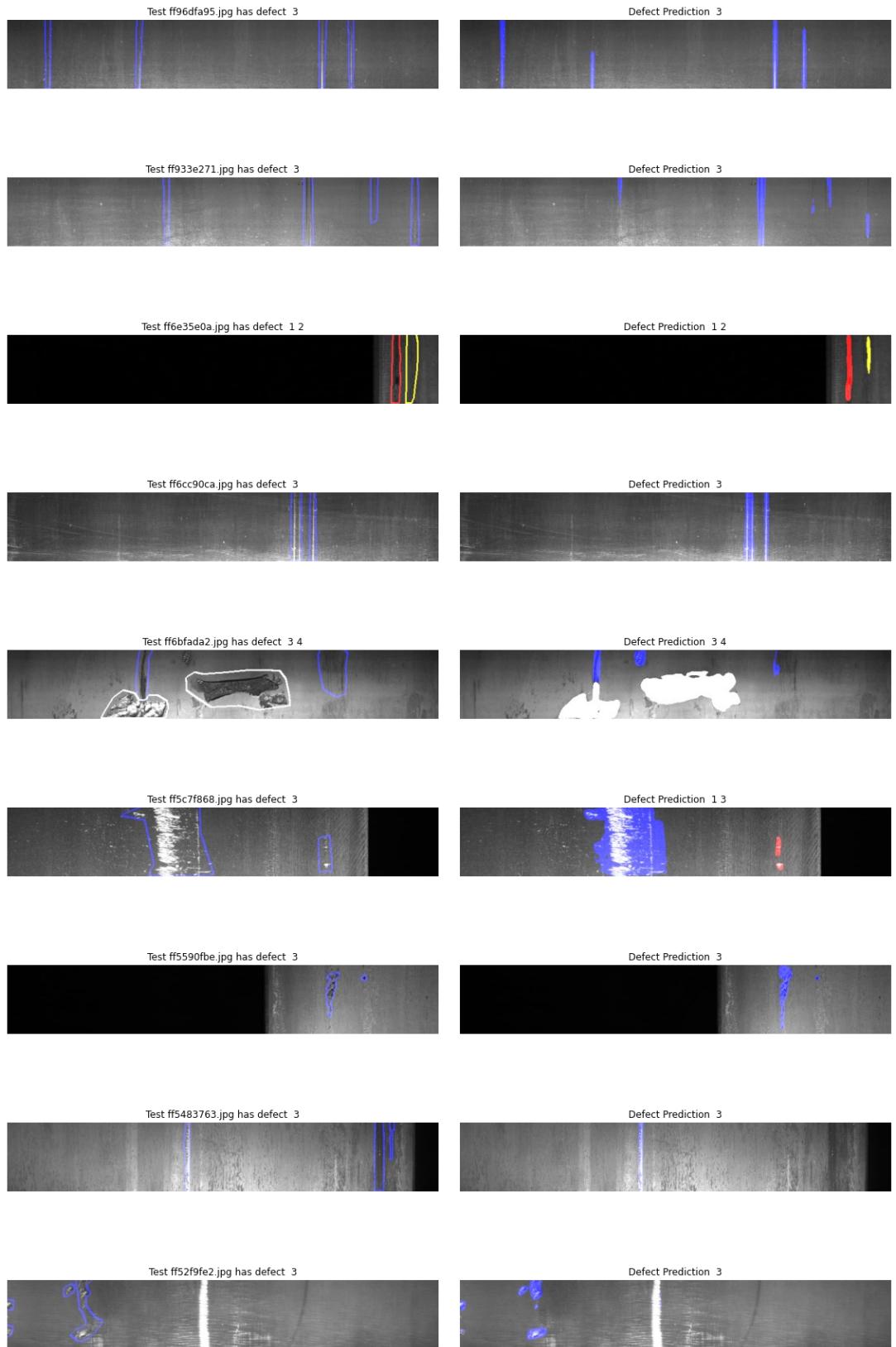


*Figure 61: loss of training and validation in EfficientNet-B5 - U-Net model*

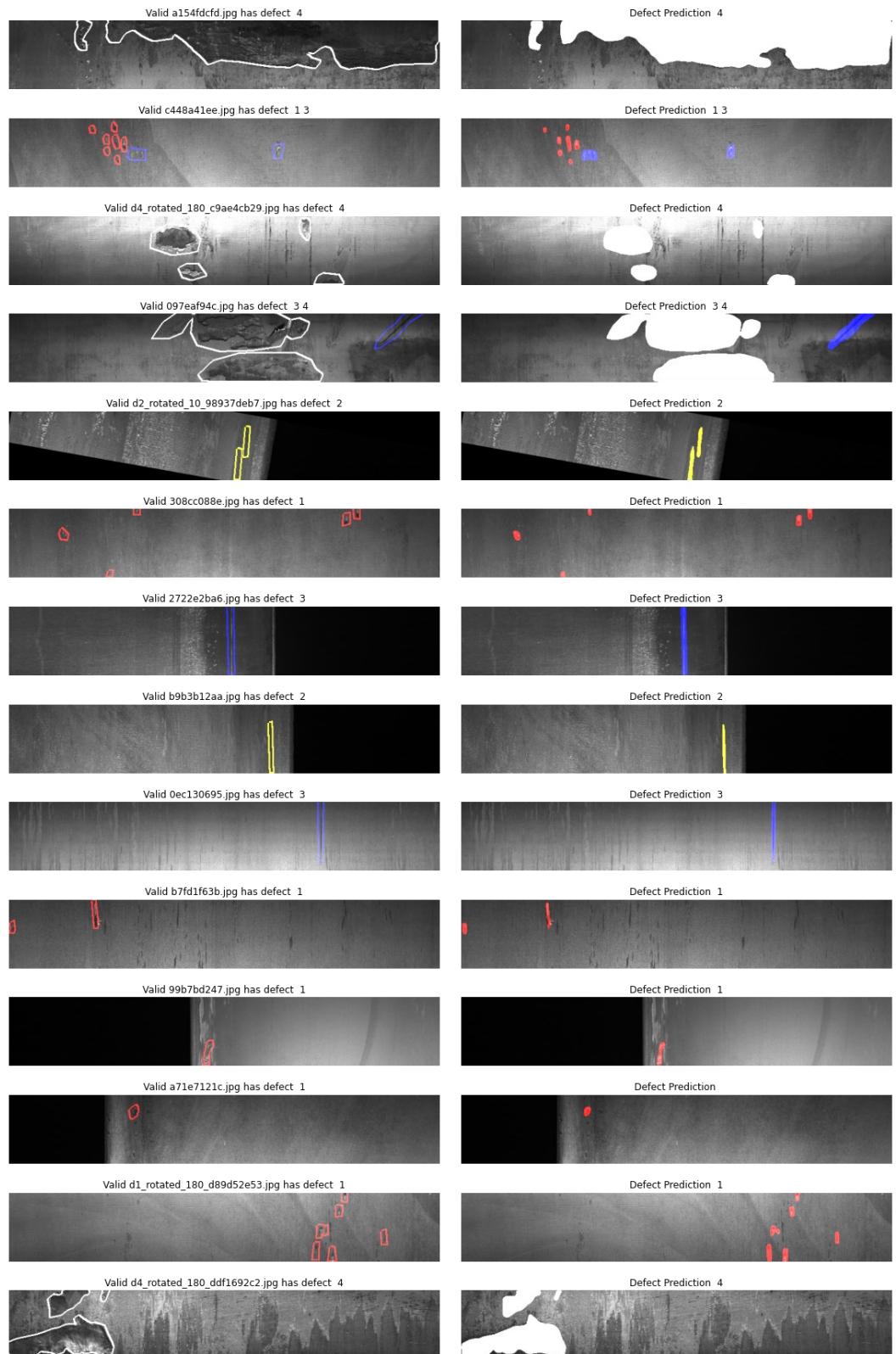
Figure 62 and 63 show the results of image prediction in validation and test data set in the combination of ResNet-101 and U-Net model and figure 64 and 65 show the results of image prediction in validation and test data set in the combination of EfficientNet-B5 and U-Net model. The predicted label based on the number of pixel that model marked as the defect, the defect type will be supposed a defect if there are more than 200 pixels of being predicted defect.



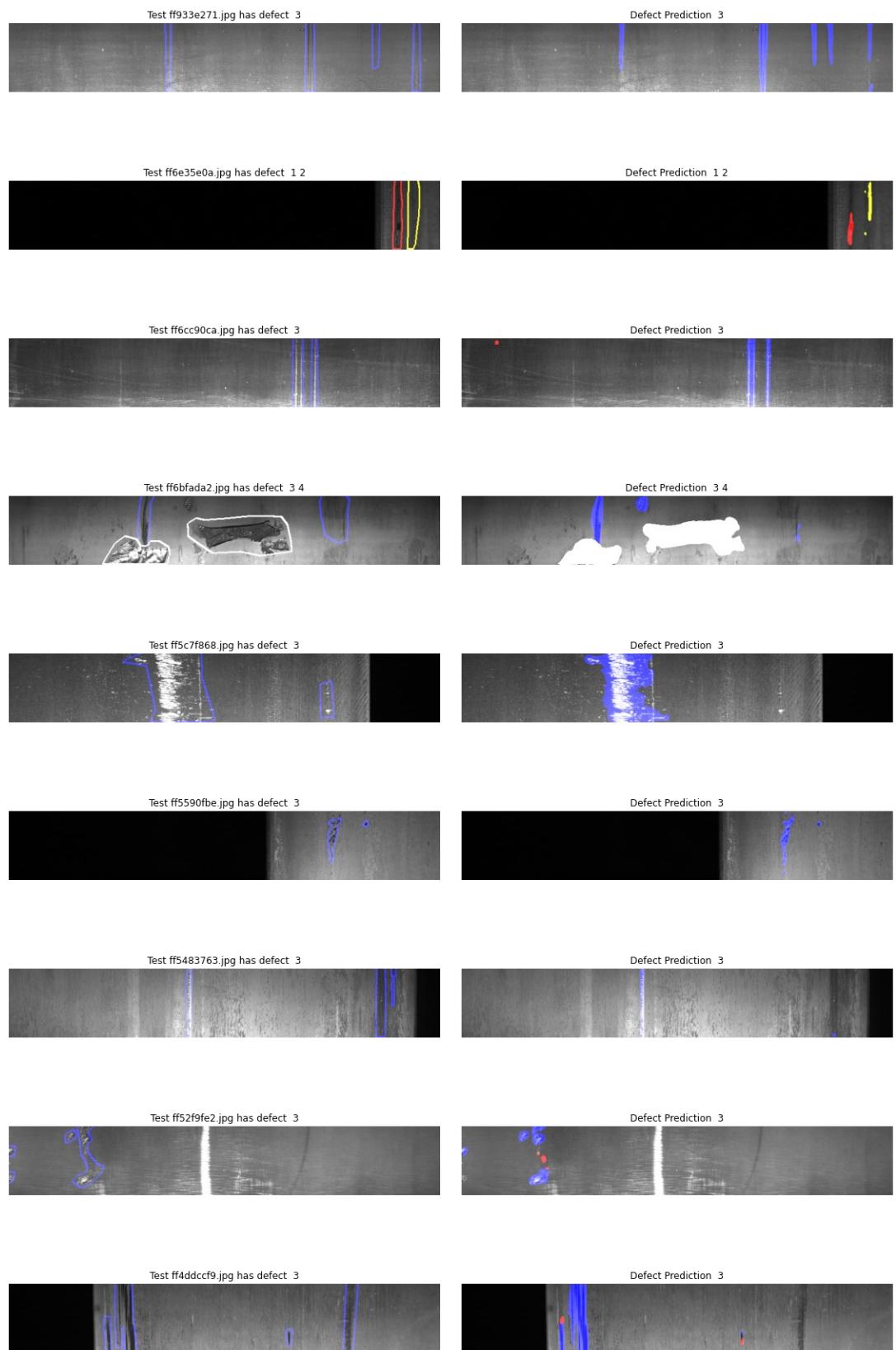
*Figure 62: Prediction result on validation data in ResNet-101 – U-Net model with ground-truth compare to predicted masks.*



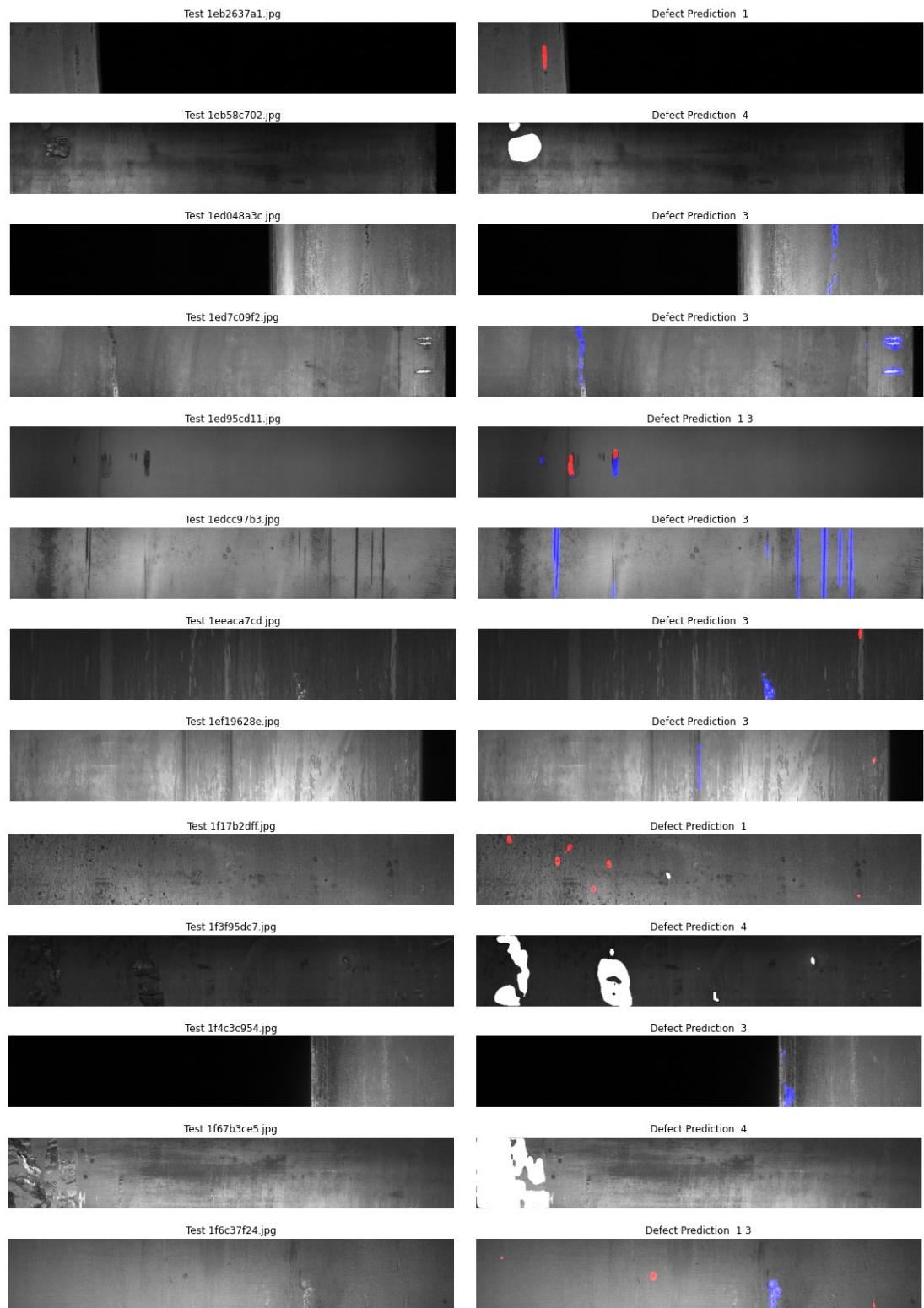
*Figure 63: Prediction result on test set in ResNet-101 – U-Net model with ground-truth compare to predicted masks.*



*Figure 64: Prediction result on validation data in EfficientNet-B5 – U-Net model with ground-truth compare to predicted masks.*



*Figure 65: Prediction result on test data in EfficientNet-B5 – U-Net model with ground-truth compare to predicted masks.*



*Figure 66: Prediction result on some original test data without ground-truth masks in EfficientNet-B5 – U-Net model.*

## V. Future Work Scope

With a pretty good result when the dice coefficient of validation data reaches 0.94 and 0.71 on the test set with ground-truth masks taken from the training set, this model should be further improved in the future. I plan to find a way to fine-tuning the model for a better result, or another approach based on other state-of-the-art works. Actually, in this project with this same data set, many people or teams have come up with better solutions with the result of dice coefficient of validation reaching 0.98 even though there is no verifiable result on test data set, because the supplier did not give out a test data set with ground-truth masks to be able to re-evaluate the model in the most objective way. Also, based on this result, maybe in the future I will learn about and develop them to a backend server to provide a direct API for mobile devices or IOT devices, specifically here devices including camera, which aims to predict steel images containing defects through camera scanning. And another scope is learning and researching how to integrate inference model function on Android platform based on TensorFlow Lite that google provides, it helps to integrate offline model on android devices which is also an interesting idea.

## **VI. Conclusion**

In conclusion, with the problem posed to detect defects on images of steel surfaces, along with basing on the background knowledge of ML as well as DL, and applying state-of-the-art works in DL, specifically, the combination of ResNets and U-Net architectures, and the combination of EfficientNets and U-Net architectures, the thesis gives a relatively satisfactory result for this problem, which is using semantic segmentation based on these complex architectures with binary cross entropy of loss and dice coefficient for evaluation metrics, fine-tuning with pre-trained weights of "imagenet" corresponding to backbone type of model and Adam as an optimizer. And cannot be lacking of making training data balanced. This strategy of implementation gave out the best validation of loss about 0.0018 and best validation of dice coefficient about 0.94, although it is not the best result, it is suitable in research and academic scope. Besides, because my ability is within the scope of study and research, the thesis cannot avoid shortcomings and limitations in achieving an excellent result. Hopefully, this thesis can contribute something in the practical application of the problem of identifying defects on the surface of steel and materials in general, and a better solution will be found to reach a higher improvement in the future.

## Bibliography

- [1] G. Cybenko, “Approximation by superposition of a sigmoidal function,” *Math. Control Signals Systems*, 1989.
- [2] F. Chollet. Deep Learning with Python. Manning, Shelter Island, NY, USA, 1<sup>st</sup> edition, 2017
- [3] I. Goodfellow, Y. Bengio, and A. Courville. Deep Learning (Adaptive Computation and Machine Learning Series. MIT Press, Cambridge, MA, USA, 2016.
- [4] <https://www.ibm.com/cloud/learn/gradient-descent>
- [5] Sebastian Ruder. An overview of gradient descent optimization algorithms
- [6] <https://builtin.com/data-science/gradient-descent>
- [7] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. Cottrell. Understanding convolution for semantic segmentation. In IEEE Winter Conference on Applications of Computer Vision (WACV), pages 1451 - 1460, March 2018.
- [8] Kunihiko Fukushima, “Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position,” *Biological cybernetics*, vol. 36, no. 4, pp. 193–202, 1980.
- [9] <https://cs231n.github.io/convolutional-networks/>
- [10] Jonathan Long, Evan Shelhamer, and Trevor Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of*

the IEEE conference on computer vision and pattern recognition, 2015, pp. 3431–3440.

[11] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi, “You only look once: Unified, real-time object detection,” in Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 779–788.

[12] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” 2012.

[13] J. Han and C. Moraga, “The influence of the sigmoid function parameters on the speed of backpropagation learning,” in Natural to Artificial Neural Computation. IWANN 1995. Lecture Notes in Computer Science. Berlin, Heidelberg: Springer, 1995, pp. 195–201.

[14] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[15] B. Karlik and A. Vehbi, “Performance Analysis of Various Activation Functions in Generalized MLP Architectures of Neural Networks,” *International Journal of Artificial Intelligence and Expert Systems (IJAE)*, vol. 1, no. 4, pp. 111–122, 2011.

[16] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier, “Language Modeling with Gated Convolutional Networks,” *arXiv*, 2017.

[17] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines,” Haifa, 2010, pp. 807–814.

[18] M. D. Zeiler, M. Ranzato, R. Monga, M. Mao, K. Yang, Q. V. Le, and G. E. Hinton, “On rectified linear units for speech processing,” in

International Conference on Acoustics, Speech and Signal Processing.  
IEEE, 2013, pp. 3517–3521

[19] X. Glorot, A. Bordes, and Y. Bengio, “Deep Sparse Rectifier Neural Networks,” in International Conference on Machine Learning, 2011

[20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. “Deep residual learning for image recognition”

[21] <https://viso.ai/deep-learning/resnet-residual-neural-network/>

[22] <https://towardsdatascience.com/understanding-and-visualizing-resnets-442284831be8>

[23] <https://medium.com/@nainaakash012/efficientnet-rethinking-model-scaling-for-convolutional-neural-networks-92941c5bfb95>

[24] <https://ai.googleblog.com/2019/05/efficientnet-improving-accuracy-and.html>

[25] Mingxing Tan, Quoc V. Le, “EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks”

[26] <https://www.mygreatlearning.com/blog/fcn-fully-convolutional-network-semantic-segmentation/>

[27] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. “U-net: Convolutional networks for biomedical image segmentation”.

[28] <https://towardsdatascience.com/optimizers-for-training-neural-network-59450d71caf6>

[29] <https://ruder.io/optimizing-gradient-descent/index.html#momentum>

[30] [https://www.researchgate.net/publication/324701535\\_Cross-Validation](https://www.researchgate.net/publication/324701535_Cross-Validation)

- [31] Bhakti Baheti, Shubham Innani, Suhas Gajre, Sanjay Talbar, “EffUNet: A Novel Architecture for Semantic Segmentation in Unstructured Environment”
- [32] <https://medium.com/@wilburdes/semantic-segmentation-using-fully-convolutional-neural-networks-86e45336f99b>
- [33] <https://medium.com/@danqing/a-practical-guide-to-relu-b83ca804f1f7>
- [34] <https://medium.com/@nishanksingla/unet-with-resblock-for-semantic-segmentation-dd1766b4ff66>
- [35] <https://towardsdatascience.com/understanding-binary-cross-entropy-log-loss-a-visual-explanation-a3ac6025181a>
- [36] <https://www.kaggle.com/c/severstal-steel-defect-detection/>
- [37] <https://www.expert.ai/blog/machine-learning-definition/>
- [38] <https://www.ibm.com/topics/computer-vision>
- [39] <https://xd.adobe.com/ideas/principles/emerging-technology/what-is-computer-vision-how-does-it-work/>
- [40] <https://medium.com/analytics-vidhya/image-classification-vs-object-detection-vs-image-segmentation-f36db85fe81>
- [41] <https://viso.ai/deep-learning/image-segmentation-using-deep-learning/>
- [42] <https://developer.nvidia.com/deep-learning>
- [43] <https://www.geeksforgeeks.org/object-detection-vs-object-recognition-vs-image-segmentation/>

