# Application Of Deep Learning For Image Deepfake Detector Using Convolutional Neural Network Algorithm

**Ananda Adhicitta Wangsadidjaja[1]**

[1]*Universitas Buddhi Dharma Imam Bonjol, Tangerang, Indonesia*

[1]nanda143735@gmail.com

**Abstract**

Social media has long been used by the public in general as a means of exchanging information. Behind this commonly exchange of information, hide the malicious intent of those who are not responsible for spreading false information or hoaxes. This false information, which can come in various forms such as images, sounds, or videos, can actually be useful when used as stock photos or simply used as caricatures and satire. Unfortunately, false information often used on famous people instead to make them look like they said or did something that never happened. This certainly needs to be controlled, one of which is by using deepfake detector that aims to recognize false information pattern. Deepfake detector utilizes the computer's ability to self-learn to recognize that invisible patterns in images using one of deep learning algorithms, namely Convolutional Neural Network, which converts images into a collection of arrays containing numbers and then performs mathematical operations repeatedly on each layer. The result of the mathematical operation can then be used as a reference to determine whether an image is real or hoax. Author's deepfake detector application using Convolutional Neural Network, specifically using the Resnet-50 model on hoax images created using AI with the ProGAN model, appears to be able to detect hoax images with the same model, with an accuracy of 85%, precision of 100%, and recall of 65%, but appears to experience decrease in accuracy when used in deepfakes with other models such as StyleGAN and BigGAN.

## I. INTRODUCTION

Social media has become a common means of communication and information exchange used by various levels of society. The ability of social media to eliminate time and regional boundaries in direct communication makes it easier for social media to permeate people's daily lives. According to the "Digital Around the World 2019" report [1], out of 268.2 million Indonesians, 150 million of them or around 56% have used social media. This number also continues to increase as in research by "DataReportal" in January 2022. This research shows an increase in social media users by 21 million users to 191.4 million users when compared to the same period in the previous year [2].

This increase cannot be separated from the development of social media's ability to provide information. Not only in the form of text, social media can also be used as a medium for sharing images, sounds, videos, or documents and other file extensions. The development of internet, especially in terms of upload and download speeds, also allows social media users to share high-quality images, sounds and videos that have large file sizes quickly.

Unfortunately, this ease of sharing information is used by irresponsible parties to spread false information or hoaxes. According to a report by the Masyarakat Anti Fitnah Indonesia in Kompas newspaper [3], the growth of false information in Indonesia has almost doubled by 2020. This false information or hoax is commonly spread either for fun, sensation or fame on the internet, cornering certain parties, or deliberately causing unrest [4].

One of the emerging forms of false information is deepfake. Deepfake is one result of deep learning technology development, a branch of machine learning that uses Neural Network algorithms on large-scale data sets to create false information. This fake information can be in the form of sound, image, or video that will look very similar to the original, thanks to advancement of this machine learning fields. Deepfake usually accepted in the form of stock photos, caricature, or even arts. Unfortunately, deepfake can also be used on famous people to make them look like they said or did something that never happened. One of most well-known deepfake can be seen in figure 1.

Research published in Crime Science [5] show that while still not illegal, deepfake is ranked as most serious Artificial Intelligence (AI) crime threat. This probably happens because people have a strong predisposition to believe their own eyes and ears despite the lengthy history of photographic manipulation, audio and video evidence has typically been accorded a significant lot of credence.



Fig. 1 Deepfake video of former US president Barrack Obama in 2018

In order to prevent this false information from spreading widely and harming various good parties, a method is needed that is able to identify information that is likely to be false information. Problem is that the spreader of false information will certainly do their best to make the recipient of the information believe the information they provide. The recipient of the information will also find it difficult if they have to search for the truth of the information through the person concerned in the information. So, technology is relied upon to weigh the truth of information. Machine learning technology can be used to train machines to learn and recognize patterns in data. By training the machine to learn data on false information, it is expected that the machine can recognize the patterns that false information has so that the machine can identify new information including false information or not.

Of the various machine learning methods available, one that has been recognized by many in identifying patterns in data is the Neural Network. This method will be used to generate mathematical equations that can be reused by machines to calculate whether new information is false information or not. Thus, Neural Network is used to develop deepfake detector application which hopefully can help general public identify the truth behind an information.

## II. Related Works/Literature Review

Convulation Neural Network (CNN), moreover ResNet have been used widely to detect pattern in image as seen in Devvi Sarwindaa, et. al. [6], Pulung Adi Nugroho, et. al. [7], or used to detect fake news pattern as seen in Muhammad Umer et. al. [8]. CNN is believed to be widely used because of its fairly good capabilities. Muhammad Umer, et. al. research bring model with stunning metric evaluation of 97.8% accuracy, 97.4 precision, 98.2% recall, and 97.8% F1-score value. Similar result can be seen in Devvi Sarwindaa, et. al. research which produce model with 88% accuracy and 92% sensitivity. Moreover, this research show that the more layer used in ResNet algorithm, the better model produced.

In deepfake detection itself, V. Venkata Reddy, et. al. [9] have tried to detect deepfake using LBP (Local Binary Pattern) method with LBPNet model. Chih-Chung Hsu et. al. research [10] also tried to detect deepfake. Using CFFN (Common Fake Feature Network) which is based on DenseNet enhancements and Siamese network architecture, this method actually produce quite promising result with 90.9% precision and 86.5% recall for BigGAN deepfake or 93% precision and 93.6% recall for SA-GAN (Structure-Aware GAN) deepfake. Sadly, this method is only tested based on same deepfake model as trained model, which mean it has not been tested for its robustness against other deepfakes.

Abdulqader M. Almars [11] states several challenges and problems that need to be considered in developing deepfake detectors, where one of it is the scalability factor where a model must be robust in the sense that it can maintain its ability on other data sets. Siwei Lyu [12] also stated similiar thing, where there are differences in quality between data sets. Some data sets were found to have deficiencies such as color mismatches, low quality images, splicing borders that indicate the editing process is clearly visible, parts of the original image that are overwritten are still clearly visible, and face orientation is inconsistent, so a deepfake detection model needs to pay attention to the robustness of the model.

## III. Methods

### A. *Deepfake*

Deepfake is the result of the development of deep learning technology, a branch of machine learning that uses Neural Network algorithms on large-scale data sets to create false information [13]. This fake information can be in the form of sound, image, or video that will look very similar to the original, thanks to advancement of this machine learning fields.

Famous figures such as politicians and artists are the most frequent targets of deepfake technology, where they are made to say or do something that never happened. As a result, the public's view of the character can change, directly or indirectly, for better or worse.

### B. *Data Mining*

Data mining is the process of using statistical, mathematical, artificial intelligence, or machine learning techniques to extract and identify information and related knowledge from various sources [14]. Often, data mining is also referred to as Knowledge Discovery in Database (KDD), knowledge extraction, or business intelligence data analysis which will be very important for business purposes. These terms actually have different meanings, but are interrelated in explaining the process of interpreting important information in data.

### C. *Neural Network*

Neural Networks (NN), also known as Artificial Neural Network (ANN) or Simulated Neural Network (SNN), are machine learning algorithms inspired by the workings of the human brain [15]. The human nervous system consists of cells called neurons. Neurons are connected using axons and dendrites, where the connections between axons and dendrites are called synapses. In a neural network, each computational unit acts as a neuron that will share information using inputs as synapses and produce outputs for input to other neurons continuously. This process can be illustrated as in Figure 2.

In practice, there is a weight value that will continue to change according to the input in order to increase the similarities of the output with the expected result, which is the label on the data. This process is then done repeatedly (iterative) to produce a learning model that can produce the right output.
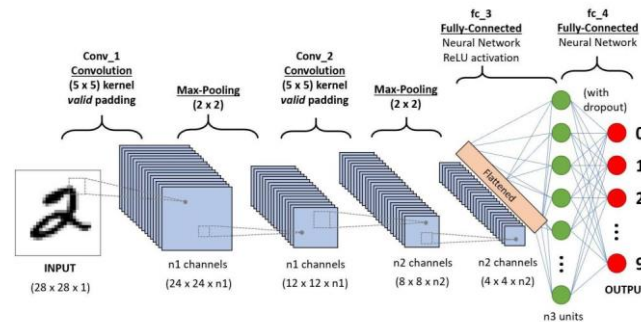


Fig. 2 Illustration of CNN process to detect numbers in images

In addition to the weights, the neural network also recognizes the bias and activation function that also play a role in increasing the output suitability. The initial input will be multiplied by the load value, then summed with the bias, and finally the activation function is used depending on the type of label on the data. The weight and bias values will change throughout the iterations, while the activation function used will remain fixed. This, can then summarized as formula (1), which is basic Neural Network formula.

$$y = g(\sum w_i * x_i - \theta) \tag{1}$$

Information:
  y = output
  g = activation function
  wi = weight for ith input
  xi = ith input
  Θ = bias

### D. *Convolutional Neural Network (CNN)*

Convolutional Neural Network (CNN) is a development of perceptron, which is the very first Neural Network, where CNN architecture has developed many neurons, so a layer is needed that divides these neurons. There are commonly three kind of layers in CNN:

1) *Input Layer*

   Input layer, as the name implies, is the layer containing the initial data that will then be processed in other layers. The processed data must be numeric, so non-numeric data needs to be converted or removed from the data set.

2) *Hidden Layer*

   Hidden layer is a layer where the data manipulation process is carried out using the weight, bias, and activation function values. This layer is called hidden layer because Neural Network users will not be directly related to this layer, but through the other two layers. What distinguishes CNN from other Neural

Networks is that CNN add a convolutional layer in the hidden layer. This convolutional layer allows CNN to be used in image and video processing.

Before, we need to know that an image is a collection of grids (boxes) containing numbers that represent the intensity or color of a pixel. Because it contains numbers, the Neural Network can manipulate these numbers to find patterns and characteristics of the image. In the convolutinal layer, this is done using a filter or kernel.

3) *Output Layer*

Output layer is a layer containing the final result of the data after processing. The result can be several neurons, each containing an object and the probability of that object, or just one neuron.

In addition, it is needed to note that an image is a collection of grids (boxes) containing numbers that represent the intensity or color of a pixel. Because it contains numbers, the neural network can manipulate these numbers to find patterns and characteristics of the image. In the convolutinal layer, this is done using filters. The patterns that the filter can recognize can be edges, dots, or other information so that this collection of filters can eventually recognize a part in the image, such as eyes, nose, mouth, and so on. Afterward, generally the output of the filter will be processed in the pooling layer before being processed again in another filter. This is done to reduce the size of the output by using stride. Filter and pooling layer example can be seen in Figure 3 and Figure 4.
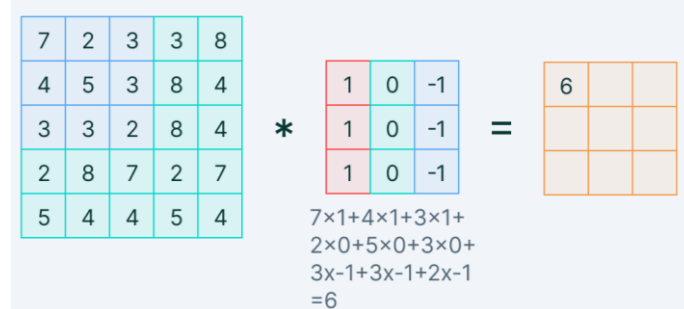


Fig. 3 Illustration of Multiplication Operation of Input with Weight in Filter Produces a Feature Map



Fig. 4 Example of Max Pooling on a 6x6 Dimensional Feature Map using Stride 2

## E. *Algorithm Evaluation*

The selection of metrics in the evaluation of machine learning algorithms is fundamentally determined by the objectives of the algorithm [16]. In classification algorithms, several metrics are known such as Precision-Recall, ROC-AUC, and Accuracy. This application will later use Precision-Recall and Accuracy as base algorithm evaluation.

Before understanding Precision-Recall and Accuracy, it is necessary to know True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN). TP describes the number of true-valued events that the algorithm successfully predicted. Similarly, TN describes the number of false-valued events that the algorithm successfully predicted. In contrast, FP describes the number of true-valued events that the algorithm failed to predict, and FN describes the number of false-valued events that the algorithm failed to predict. These four terms are often used in classification-type algorithm evaluation calculations.

Precision describes the ratio of TP, divided by the total number of positive predictions made by the algorithm (sum of TP and FP). Precision is particularly useful when the algorithm has a high sensitivity to false values such as an email spam filter system. Recall is the ratio of TP, divided by the actual number of positive values (sum of TP and FN). Recall is used when the loss due to FN is greater than FP such as in disease prediction. Lastly, Accuracy is a method to assess the accuracy of algorithm predictions (TP plus TN) compared to the overall algorithm predictions (sum of TP, TN, FP, and FN). A high accuracy means that the algorithm is able to predict something with high accuracy. However, accuracy is not recommended if the data has sensitivity to certain value like the previous example.

$$Precision = \frac{TP}{TP + FP} * 100\% \tag{2}$$

$$Recall = \frac{TP}{TP + FN} * 100\% \tag{3}$$

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} * 100\% \tag{4}$$

Information:
TP = True Positive
TN = True Negative
FP = False Positive
FN = False Negative

## IV. RESULTS

### A. *Program View*

1) *Home Menu View*

Home menu show short explanation, background, configuration used, and developer behind this application. Home menu view can be seen in Figure 5.
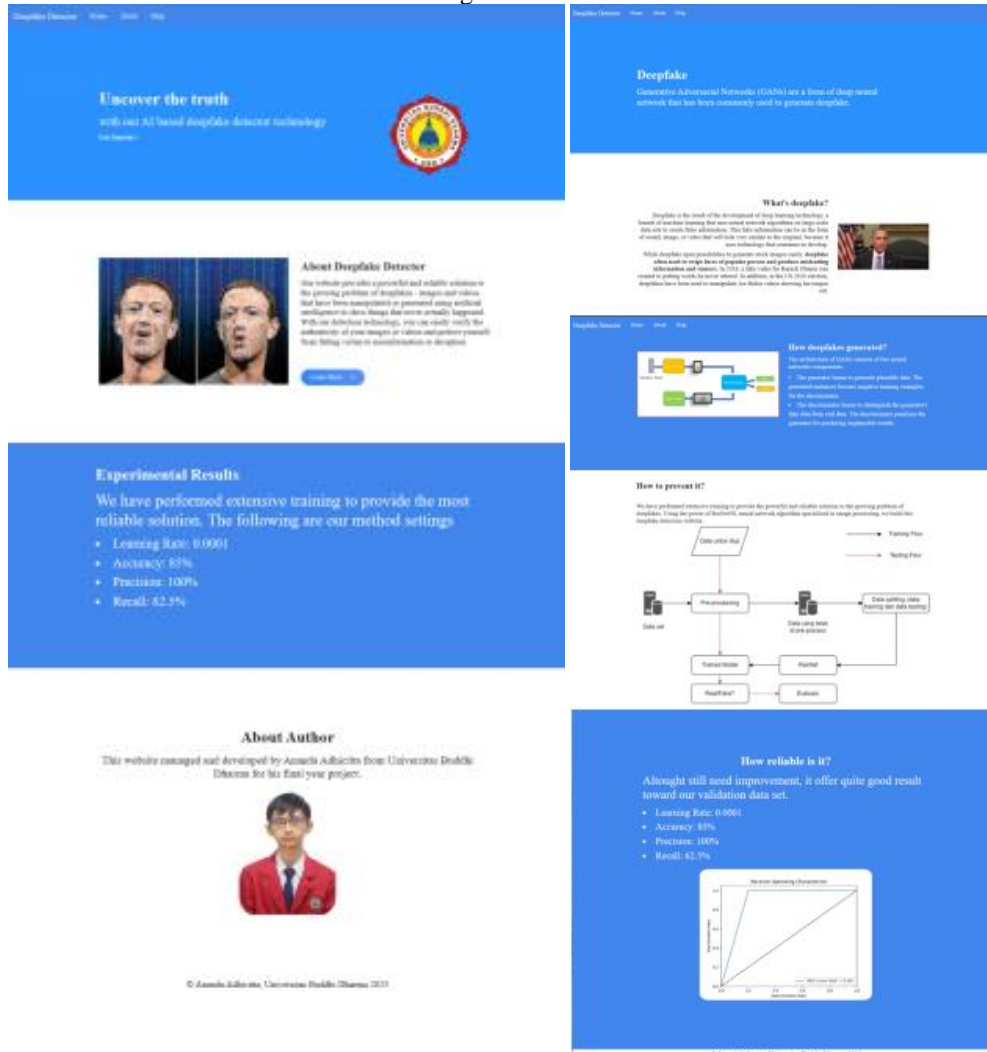


Fig. 5 (left) Home Menu View and Figure 6 (right) About Menu View

2) *About Menu View*

About menu give more detailed explanation about this application, along with flowchart and algorithm evaluation result. About menu view can be seen in Figure 6.

3) *Help Menu View*

Help menu allow user to give suggestion, message, or question regarding this application. Help menu view can be seen in Figure 7.
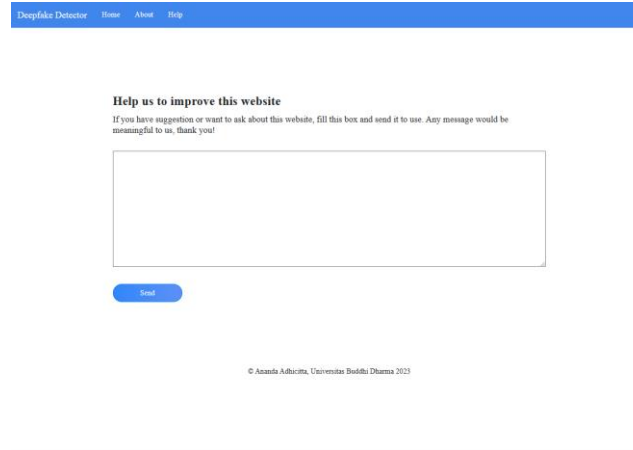


Fig. 7 Help Menu View

4) *Detection Menu View*

Detection menu view is the page where user can input their image to do deepfake detection of their image.Detection menu view can be seen in Figure 8.
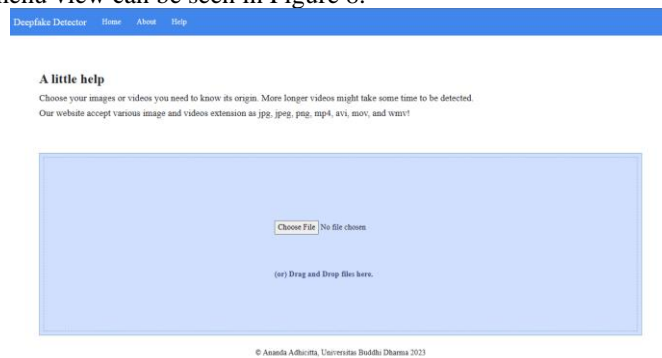


Fig. 8 Detection Menu View

5) *Result Menu View*

Result menu show detection result done by the application. Result menu view can be seen in Figure 9.
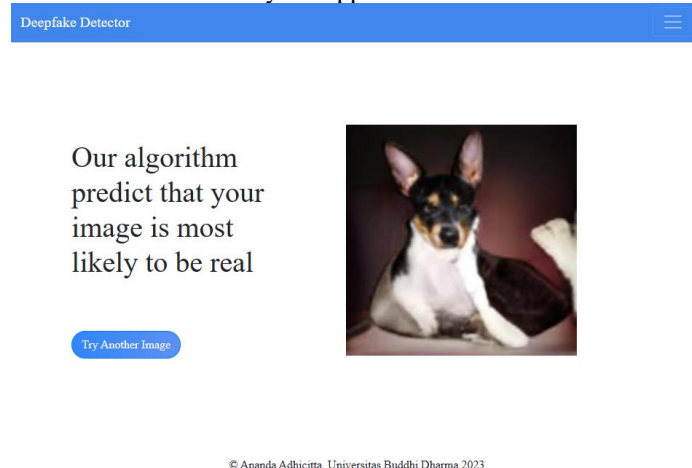


Fig. 9 Result Menu View

## B. *Algorithm Test*

Neural Network algorithm training and usage developed using CRISP-DM (Cross-Industry Standard Process for Data Mining) standard which is a universal concept that can be used in various types of data mining algorithms (Peter Butka, et al., 2020: 6). This structure can be described by Figure 10. Furthermore, the trained algorithm model has a total of 6 main layers that can be subdivided into 50 layers, briefly can be seen in Figure 11. Website development is carried out on 2 amazon ec2 server instances, where the frontend uses an instance of type t2.micro, while the backend uses an instance of type t2.large. In general, the frontend uses react.js and the backend uses python as its programming language. The trained algorithm model then placed in backend, which secured by CORS so that algorithm model can only be accessed from frontend.
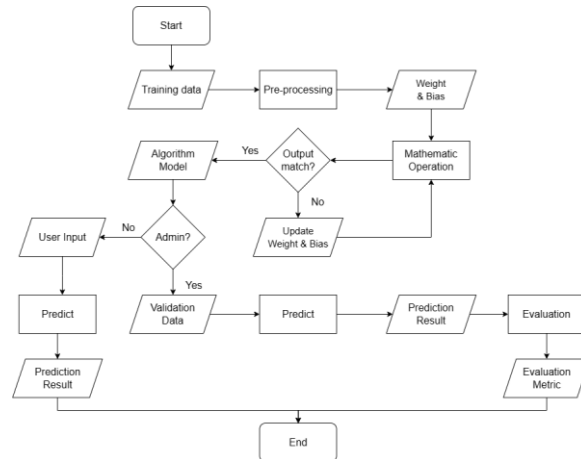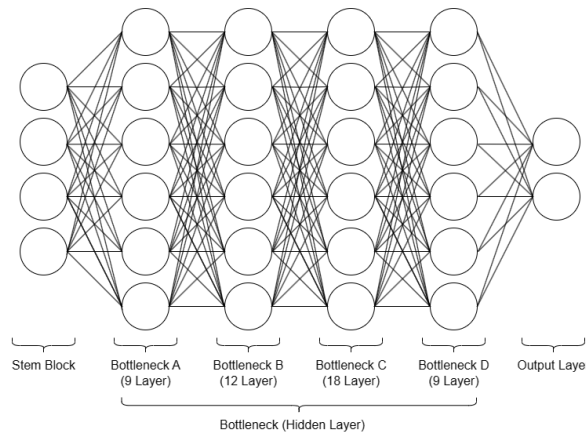


Fig. 10 Algorithm Training and Usage Flowchart



Figure 11Trained Algorithm Model

Data set used is the data set also used previously by Sheng-Yu Wang, et al (2020) in their paper for CVPR 2020, a virtual annual conference for computer vision. The data set consists of 720,000 training data and 4,000 validation data, which are divided into 360,000 original labeled training data and 2,000 original labeled validation data, as well as 360,000 false labeled training data and 2,000 false labeled validation data. The false labeled data itself is an image generated using ProGAN. ProGAN was chosen because it can produce high-quality images even though it uses CNN algorithms that tend not to be complex.

Data preprocessing is done before training by equalizing the pixels of each image to 224x224 pixels. Next, the data set was divided into 20 categories based on what the images depicted. Each category have 20 image with false and true, result in 800 images data used for training. This equalization is important so the model will not have imbalance which result in defective model.

To fulfill the need to be robust or able to be used not only on one data set, the application was tested using three GAN models, namely ProGAN, StyleGAN, and BigGAN. Data with real labels uses flickr as its data source. Data sets with fake labels and the ProGAN model are validation data obtained from the same source as the training data, namely Sheng-Yu Wang, et al (2020) in their paper for CVPR 2020. This data set has a dimension of 224x224 pixels. Data sets with fake labels and StyleGAN models are divided into two sources, data sets obtained through the Kaggle site with the website page https://www.kaggle.com/datasets/awsaf49/artifact-dataset with 200x200 pixels dimensions and author's data set which obtained through the AI Image Generator on the pages https://replicate.com/orpatashnik/styleclip and https://replicate.com/yuval-alaluf/sam, both have 1024x1024

pixels dimensions. Data sets with fake labels and StyleGAN models are also divided into two sources, namely data sets obtained through the Kaggle site with the website https://www.kaggle.com/datasets/awsaf49/artifact-dataset, with dimensions of 200x200 pixels, and author's data set which obtained through Tensorflow's AI Image Generator on the page https://github.com/tensorflow/hub/blob/master/examples/colab/biggan_generation_with_tf_hub.ipynb which has 256x256 pixels dimensions.

The results of the test can be seen in the following confusion matrix.

TABLE 1
ALGORITHM CONFUSION MATRIX

| | | Prediction | |
|---|---|---|---|
| | | Real | Fake |
| Actual | Real | 20 | 0 |
| | Fake | 12 | 48 |

Using formula shown before, the precision, recall, and accuracy of this application are obtained as follows.

$$\text{Precision} = \frac{TP}{TP + FP} * 100\%$$

$$\text{Precision} = \frac{20}{20 + 0} * 100\%$$

$$\text{Precision} = \frac{20}{20} * 100\%$$

$$\text{Precision} = 100\%$$

$$\text{Recall} = \frac{TP}{TP + FN} * 100\%$$

$$\text{Recall} = \frac{20}{20 + 12} * 100\%$$

$$\text{Recall} = \frac{20}{32} * 100\%$$

$$\text{Recall} = 62.5\%$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} * 100\%$$

$$\text{Accuracy} = \frac{20 + 48}{20 + 48 + 0 + 12} * 100\%$$

$$\text{Accuracy} = \frac{68}{80} * 100\%$$

$$\text{Accuracy} = 85\%$$

## V. DISCUSSION

Image Deepfake Detector using Convolutional Neural Network (CNN) with ResNet-50 model shown that deepfake detector can be build as website, not only as runnable code in code editor as shown in most research, while maintaining its accuracy, precision, and recall. With this, hopefully deepfake detector can be more widely spread and accessible by common people, thus contributing in fight against hoax news, specially in form of image.

Based on algorithm evaluation performed, Image Deepfake Detector using Convolutional Neural Network (CNN) with ResNet-50 model able to maintain performance when used in deepfake data set, compared with Chih-Chung Hsu et. al. Research, but overall having degradation compared to CNN based model like in Devvi Sarwindaa, et. al. or Muhammad Umer et. al..

This research also meant to prove robustness of CNN algorithm against different data set, which in this case using data set of ProGAN deepfake image against StyleGAN and BigGAN deepfake image. This research still using small variety of deepfake image due to limition of public deepfake data set, but future research might use even more diverse data set to further test model robustness. The world of deepfake is also still vast and still under development, which mean deepfake detector as prevention tool of hoax news, specially in form of image, need to be developed too.

Another noteworthy topic for another research is, there still many interesting factor worth developed, such as image compression which often applied in social media, integration with commonly used social media as REST API, and so on. Hopefully, next research can also detect hoax news in other form, such as article, sound, or video while still maintaining ease of access and use aspect.

## VI. CONCLUSIONS

The Image Deepfake Detector application using the Neural Network algorithm with the ResNet-50 model has good accuracy when applied to deepfakes that use similar models in their development with the training data model. This research proved image deepfake detector can be build as website, not only as runnable code in code editor as shown in most research. To be precise, this application use 2 Amazon EC2 server to serve frontend and backend request. This research also manage to test CNN algorithm, specifically ResNet-50, robustness against another data set, where the application can detect deepfakes with ProGAN models that are similar to the training data with 100% accuracy and shows a decrease in accuracy when used in deepfakes with StyleGAN and BigGAN models. Overall, the app achieved 85% accuracy, 100% precision, and 65% recall.

However, this application can still be further developed. Artificial Intelligence is still rather new things for most people and currently is in rapid development, which mean there is still room for deepfake detector to improve along with deepfake itself. There is also many interesting factor worth developed, such as image compression which often applied in social media, integration with commonly used social media as REST API, and so on.

## REFERENCES

[1] W. K. Pertiwi, "Separuh Penduduk Indonesia Sudah 'Melek' Media Sosial," Kompas.

[2] L. Jemadu, "Jumlah Pengguna Media Sosial Indonesia Capai 191,4 Juta per 2022," Suara.

[3] W. K. Pertiwi, "Jumlah Hoaks di Indonesia Meningkat, Terbanyak Menyebar lewat Facebook," Kompas.

[4] A. R. Hamsah, *Kejahatan Berbahasa (Language Crime)*, 1st ed. Tasikmalaya: Langgam Pustaka, 2022.

[5] M. Caldwell, J. T. A. Andrews, T. Tanay, and L. D. Griffin, "AI-enabled future crime," *Crime Sci*, vol. 9, no. 1, p. 14, Dec. 2020, doi: 10.1186/s40163-020-00123-8.

[6] D. Sarwinda, R. H. Paradisa, A. Bustamam, and P. Anggia, "Deep Learning in Image Classification using Residual Network (ResNet) Variants for Detection of Colorectal Cancer," *Procedia Comput Sci*, vol. 179, pp. 423–431, 2021, doi: 10.1016/j.procs.2021.01.025.

[7] P. A. Nugroho, I. Fenriana, and R. Arijanto, "Implementasi Deep Learning Menggunakan Convolutional Neural Network (CNN) pada Ekspresi Manusia," *ALGOR*, vol. 2, no. 1, pp. 12–21, Sep. 2020.

[8] M. Umer, Z. Imtiaz, S. Ullah, A. Mehmood, G. S. Choi, and B.-W. On, "Fake News Stance Detection Using Deep Learning Architecture (CNN-LSTM)," *IEEE Access*, vol. 8, pp. 156695–156706, 2020, doi: 10.1109/ACCESS.2020.3019735.

[9] V. V. Reddy, P. Priyanka, D. K. Priyanka, P. R. Vishnu, A. D. Kumar, and S. B. Gole, "Fake Image Detection Using Machine Learning," in *International Journal of Advanced Research in Computer and Communication Engineering*, Guntur: IJARCCE, Jan. 2022, pp. 138–144.

[10] C.-C. Hsu, Y.-X. Zhuang, and C.-Y. Lee, "Deep Fake Image Detection Based on Pairwise Learning," *Applied Sciences*, vol. 10, no. 1, pp. 370–384, Jan. 2020, doi: 10.3390/app10010370.

[11] A. M. Almars, "Deepfakes Detection Techniques Using Deep Learning: A Survey," *Journal of Computer and Communications*, vol. 09, no. 05, pp. 20–35, 2021, doi: 10.4236/jcc.2021.95003.

[12] S. Lyu, "Deepfake Detection: Current Challenges and Next Steps," in *2020 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, IEEE, Jul. 2020, pp. 1–6. doi: 10.1109/ICMEW46912.2020.9105991.

[13] C. Ategeka, *The Unintended Consequences of Technology: Solutions, Breakthroughs, and the Restart We Need*, 1st ed. New Jersey: John Wiley & Sons, 2021.

[14] Mulaab, *Data Mining : Konsep dan Aplikasi*, 1st ed. Madura: Media Nusa Creative (MNC Publishing), 2017.

[15] J. Stone, *Artificial Intelligence Engines: A Tutorial Introduction to the Mathematics of Deep Learning*, 1st ed. Sheffield: Sebtel Press, 2019.

[16] D. Mukunthu, Shah Parashar, and W. H. Tok, *Practical Automated Machine Learning on Azure*. Sebastopol: O'Reilly Media, 2019.