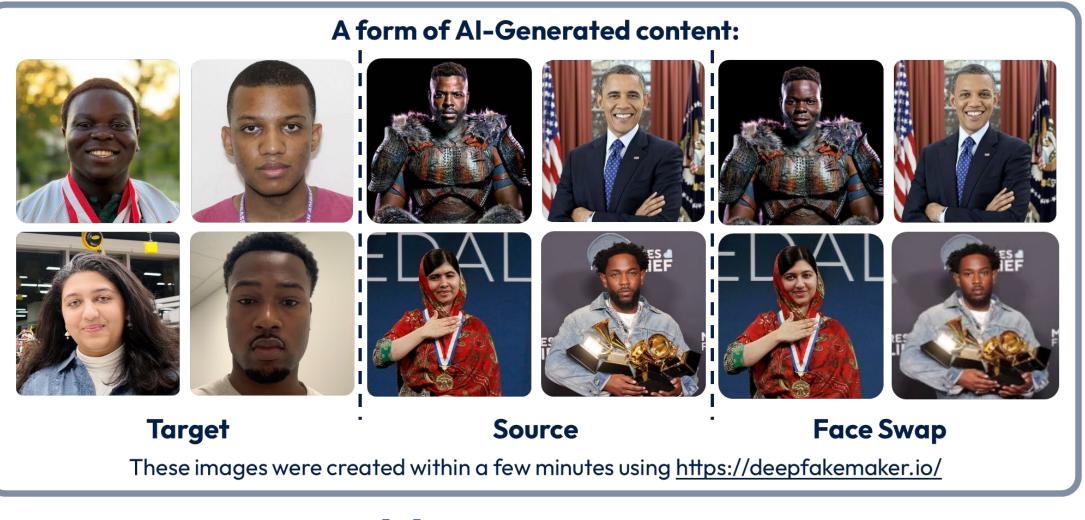


# Team 5: Automated DeepFake Detection

Jouri Ghazi, Ashton Bryant, Jahtega Djukpen, Zacary Louis Instructor: Dr. Maryam Alibeik, Advisor: Dr. Joseph Picone

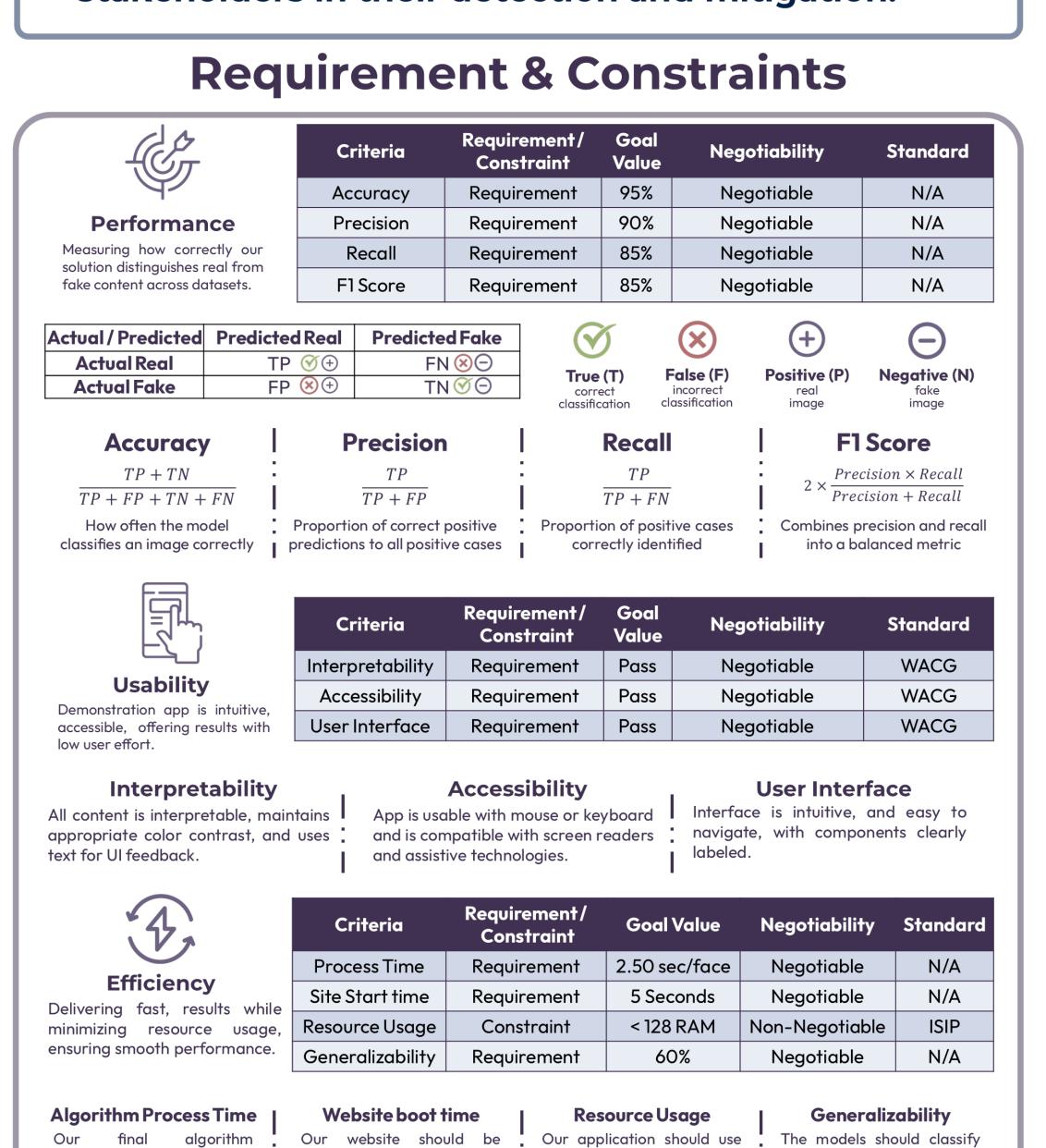
### What is a DeepFake?



#### **Problem Statement**

Artificially generated and manipulated images such as DeepFakes, have become easier to create, and harder to detect. These images pose a serious risk to the credibility of public figures, resulting in great emotional, financial and political turmoil. An accurate and reliable DeepFake detector is needed to protect individuals such as celebrities, politicians, and even corporations who are frequent victims of DeepFake misuse.

DeepFakes pose a risk to everyone, making us all stakeholders in their detection and mitigation.



#### **ISIP** Guideline **ISO 12207 IEEE 7003 ISO 27001** Documenting race or gender naming conventions and bias in the model to confirm the system works as intended. keeping all user data private. which are accessible

the Neuronix Cluster.

#### Our project was developed in 3 portions:

Constraint

Requirement

hosted on the ISIP server.

ISIP Guideline

Algorithmic

Transparency

Website Security



below 2500 ms per face.

**Code Maintainability** 

Ease with which a codebase ca be understood, modified,

tested, and extended over time

**Ethic & Security** 

protecting data, and promoting

fairness and transparency.

commenting style.



Pass

**Web Tool Development** A platform for users to upload images and detect

DeepFake methods.

Non-Negotiable | ISO 12207

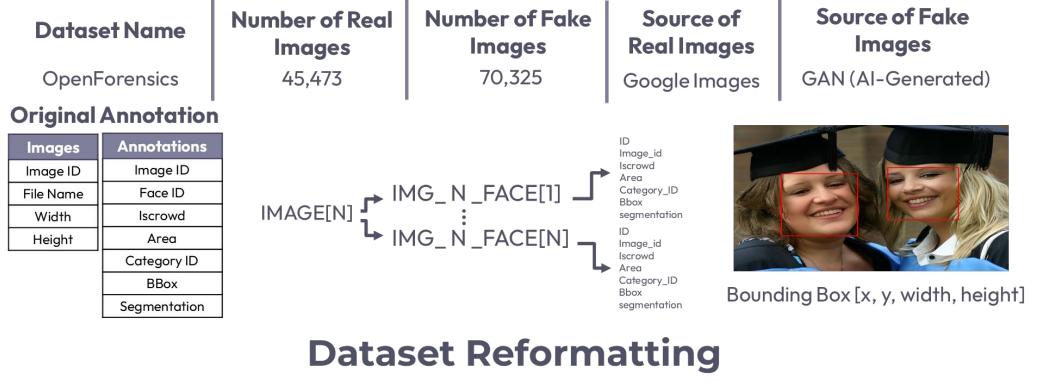
Non-Negotiable | ISO 27001

IEEE7003-

2024

Standard

## **Dataset Management**



A parallel dataset of faces is created to better train and develop the model This dataset will be made up of images that contain one face each.



105.093 **Fake** faces

Filenames ending in \*\_0. jpg are

Source of

#### real, and \*\_1.jpg are fake Used to fine tune model performance metrics

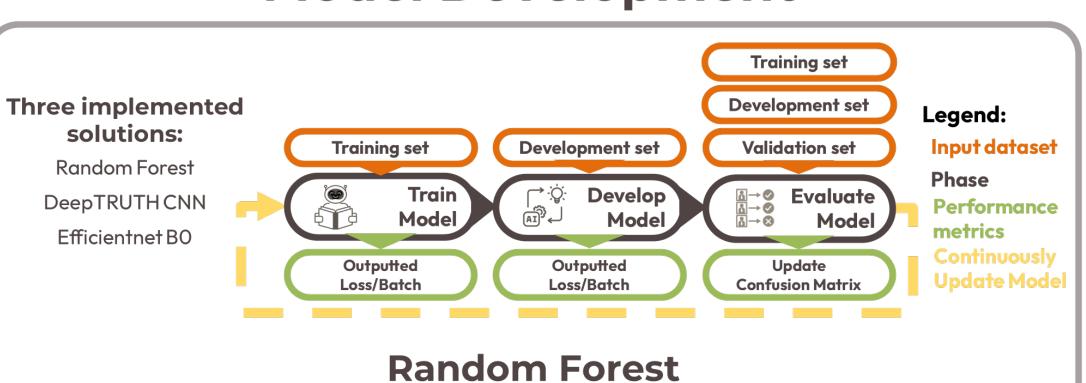
#### **Caltech Dataset**

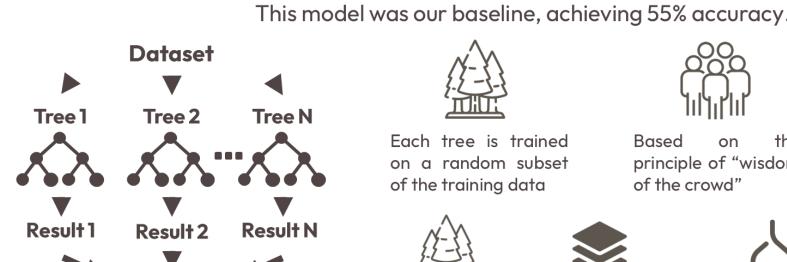
to prevent overfitting

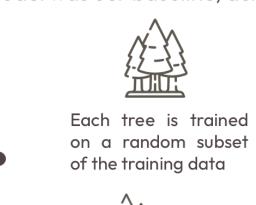
A secondary evaluation dataset is created to assess generalization, testing the model's performance on new and unseen data. This dataset will be 100 Fake images and 100 Real images. The original images are sourced from Caltech. DeepFakes made from Photoshop and DeepfakeMaker.

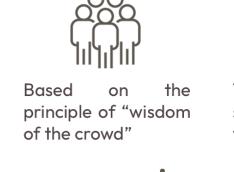


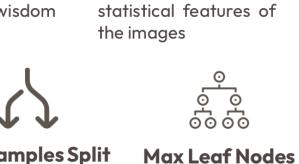
## **Model Development**



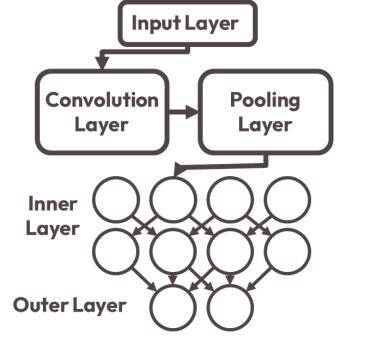








## **Convolutional Neural Network**



MBConv1, 3x3

MBConv6, 3x3

MBConv6, 3x3

MBConv6, 5x5

MBConv6, 5x5

MBConv6, 3x3 MBConv6, 3x3

MBConv6, 3x3

MBConv6, 5x5

MBConv6, 5x5

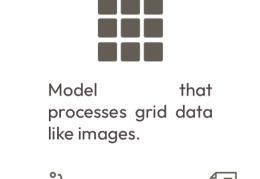
MBConv6, 5x5

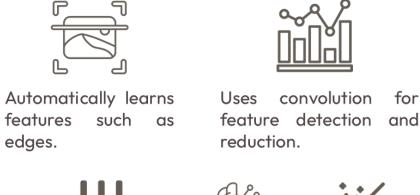
MBConv6, 3x3

Feature Map

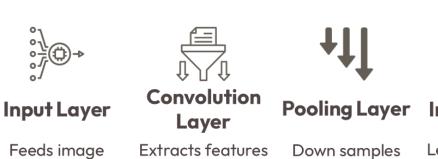
Averaging

Classification

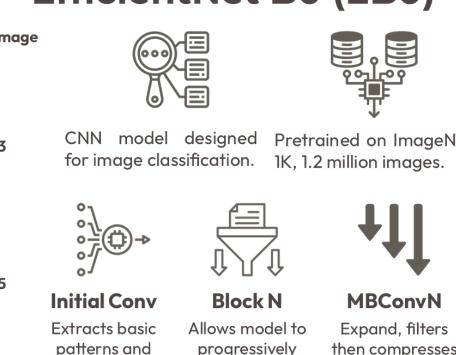


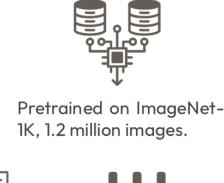


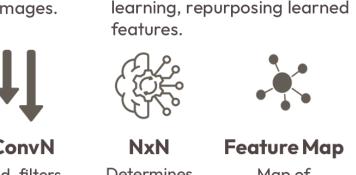
tree is allowed to split

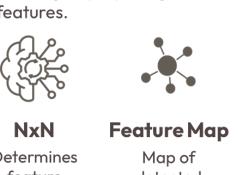


## EfficientNet B0 (EB0)



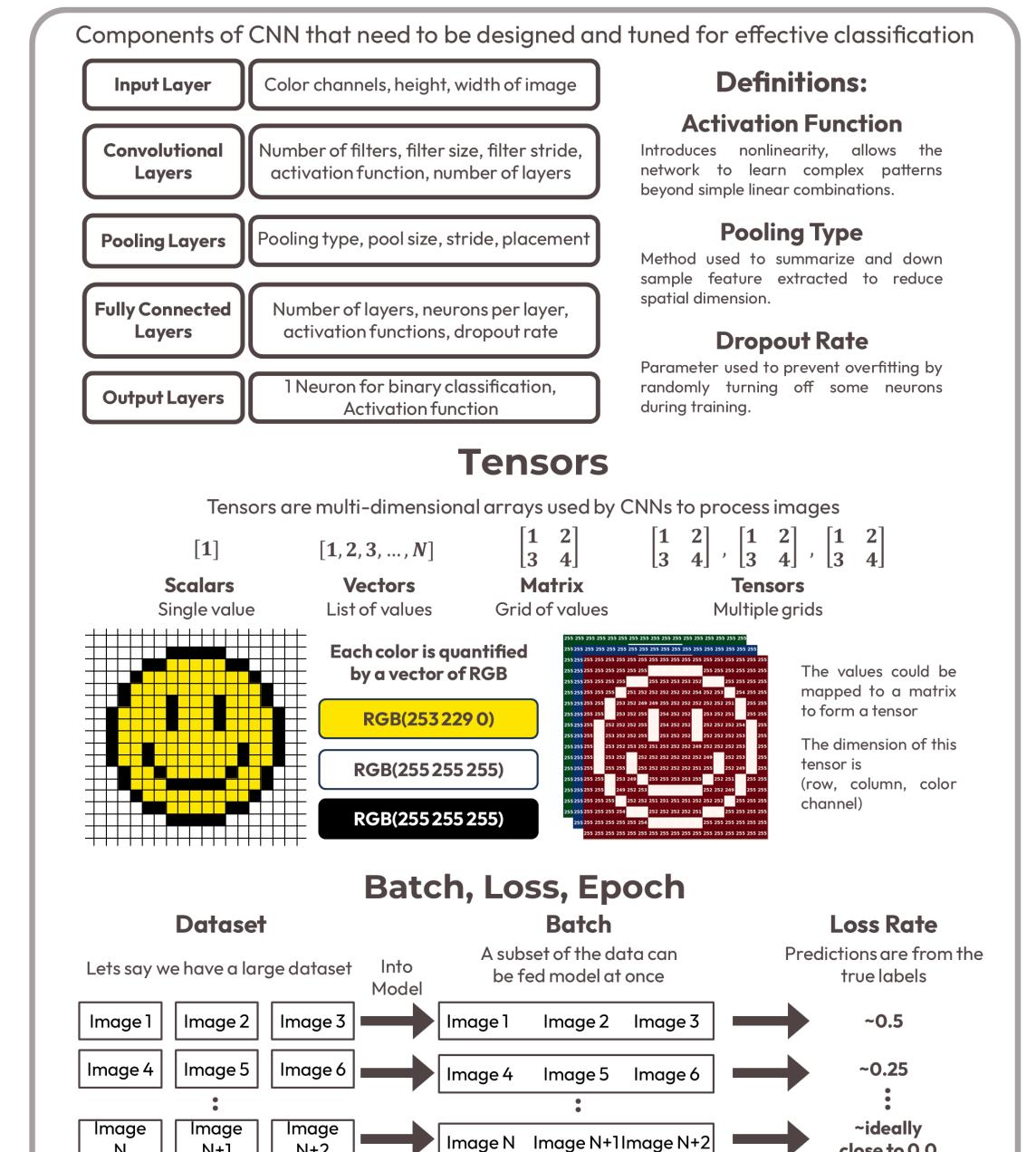






 $\widehat{\Theta} = P_2 - P_1$ 

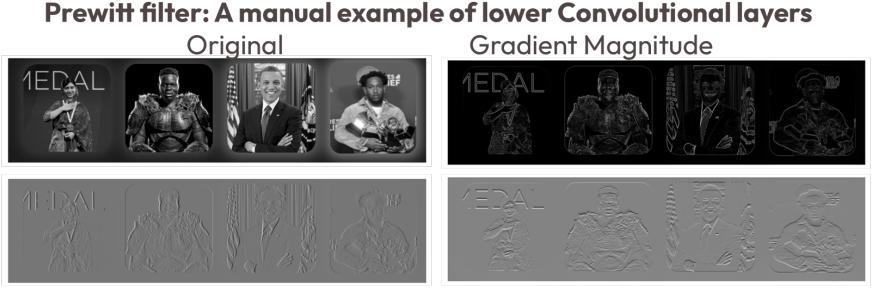
## **Convolutional Neural Network**



#### Epochs are complete pass throughs in the entire training dataset within the model. The models performance should increase with more epochs **CNN Edge Detection**

Lower Layers: edges, lines, colors Middle Layers: patterns, shapes, textures convolutional layer: Higher Layers: faces objects signs of manipulation

Prewitt filter: A manual example of lower Convolutional layers



Y-Direction Derivative X-Direction Derivative

## **Model Performance**

The dataset is split into 3 parts, when model is trained then each subsection is decoded.

Model	Random Forest				DeepTRUTH				Goal				
Model	Train	Train Development		<b>Validation</b>	Train	ain Development Validat		Train	Development	<b>Validation</b>	Value		
Accuracy	98.90%	54.60	O%	54.70%	98.54%	96.41%	94.44%	98.97%	99.36%	98.90%	95%		
Precision	N/A				99.46%	95.01%	89.84%	99.89%	99.33%	98.10%	90%		
Recall	N/A				97.97%	96.46%	92.13%	98.31%	99.16%	98.36%	85%		
Fl	N/A				98.71%	95.73%	90.97%	99.10%	99.25%	98.23%	85%		
	Deen'	TRUTH	EB	O Go	al								
Model		TECH	CalTi										
Accuracy	60.	.64%	60.6	4% 60	<mark>%</mark> Т	The CalTECH dataset is used to determine the models generalizability							
Precision	70.	.07%	96.6	0% N/	<b>~</b>								
Recall	65.	.95%	60.4	.3% N/	Δ m								
Fl	67.	.76%	74.3	5% N//	Δ								

#### **Statistical Significance**

The measure that the difference in accuracy between the two results is meaningful improvement, unlikely due to random chance.

	N	$P_1$	$P_2$	$\widehat{m{\Theta}}$	SE	Z	$CI_{low}$	$CI_{up}$	Sig
Dataset	Total	Accuracy		Difference	Standard		CI	CI	Significant?
Dalasei	Faces	DeepTRUTH	EBO	Dillerence	Error	Statistic	Lower	Upper	Significant:
Train	150866	98.54%	98.97%	0.43%	4.0E-04	10.65	0.35%	0.51%	Yes
<b>Development</b>	49718	96.41%	99.36%	2.95%	9.1E-04	32.50	2.77%	3.13%	Yes
Evaluation	15345	94.44%	98.90%	4.46%	2.0E-03	21.94	4.06%	4.86%	Yes
Caltech	223	60.64%	60.64%	0.00%	4.6E-02	0.00	-9.07%	9.07%	No
						4.06	0 = 0 /	C	

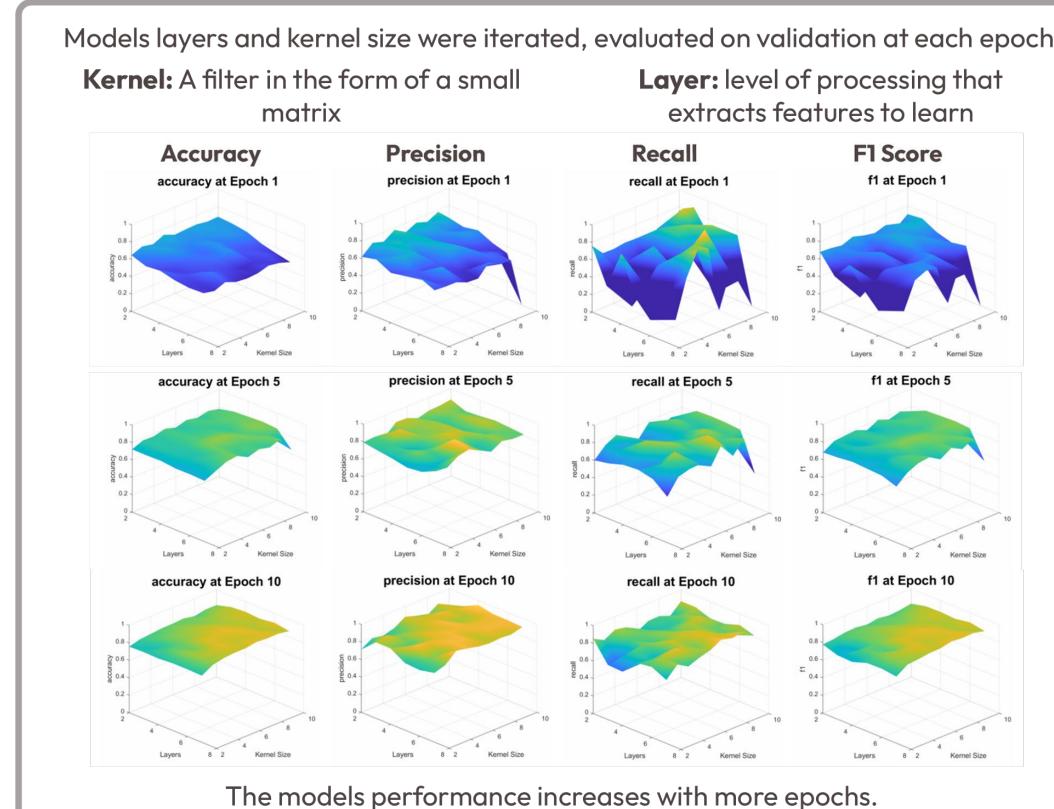
## $Z_{\alpha/2} = 1.96,95\%$ confidence $CI = Confidence\ Interval$

 $CI_{up} = \widehat{\Theta} + Z_{\alpha/2} \cdot SE$ 

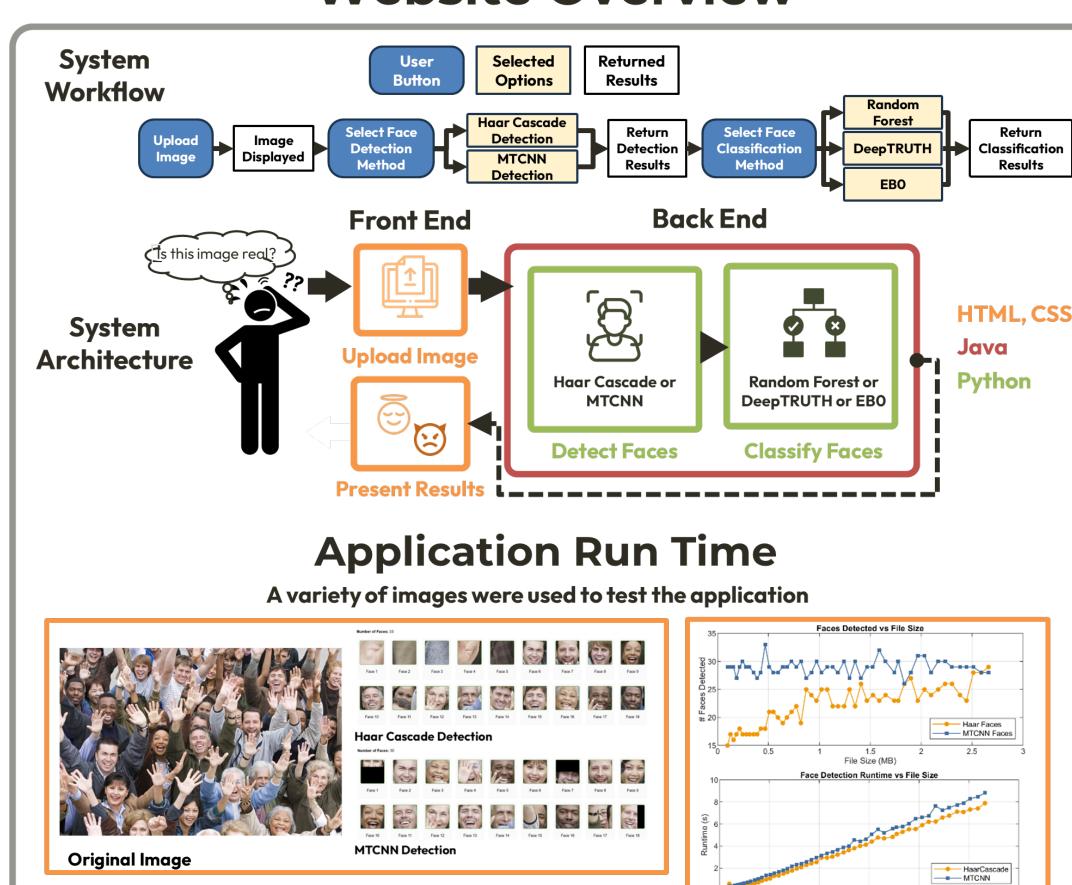
application lies in its ability to classify unseen data.

The performance of our

## CNN DeepTRUTH Model



#### **Website Overview**



## **Project Results**

 
 2.38
 24
 307.79
 312.29
 309.79
 29
 276.00
 279.72
 302.48
Average 156.10 161.57 156.77 Average 139.33 143.33 161.60

All possible combinations perform better than goal amount

	Criteria	Goal Value	<b>Achieved Value</b>	Justification			
Performance	Accuracy	95%	98.90%				
	Precision	90%	98.10%	The FDO Medales are former addless hand			
	Recall	85%	98.36%	The EBO Model performed the best			
	F1 Score	85%	98.23%				
	Algorithm Process time	1500 ms/Face	< 161.60 ms/Face	MTCNN+EBO took the longest			
F.66: -:	Website boot time	5 Seconds	< 1 Second	Webtool Available Online			
Efficiency	Resource Usage	< 128 GB RAM	<16 GB RAM	Runs on laptop			
	Dataset Generalizability	60%	60.64%	Performance on Caltech Dataset			
Code	ISIP Guide Compliance	Pass	Pass	Working on producing the proper			
Maintainability	Documentation	Pass	Pass	documentation for software			
	Interpretability	Pass	Pass	Feedback console & confidence information			
Usability	Accessibility	Pass	Pass	Keyboard tested			
	User Interface	Pass	Pass	High contrast colors + simple UI			
	Algorithmic	Pass	Danding	Actively working on finding notantial his			
Ethics & Security	Transparency	Fuss	Pending	Actively working on finding potential bias			
	Website Security	Pass	Pass	All dependencies on CSV's removed, photo data cleared			

This project was entirely software based, supported by the NeuroNix Cluster and required no additional cost.

#### **Future works**



Image was scaled up by 1.1x per iteration

