

# Logbook

Anie Udofia, Gr. 9, ~~aniedofia8@gmail.com~~ aniedofia8@gmail.com

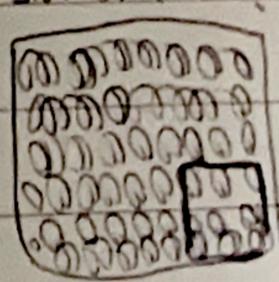
Logbook must include:

- Schedule daily notes / ideas, background research contacts, references,
- Experimental procedure / method, data collection sheets, observations / results, conclusion
- Aspects of the project performed differently
- All change!
- Extreme detail to ensure replication

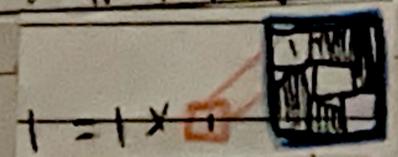
MIT 6.S191/2024 - Source of images

Convolutional Neural Networks (CNNs) are neural network models that process images of data through the extraction of features (important aspects of the input) in order to assign them to a particular class. These features have a set of weights executed on them, known as a filter, which can be of varying size. These filters are implemented on several different features to encourage algorithmical efficiency + ease instead of having individual pixels being processed.

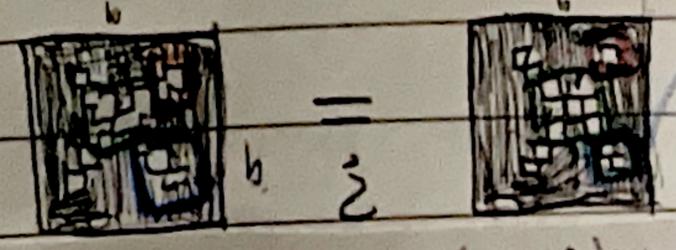
Ex: Filter size is  $3 \times 3$  + 9 different weights



↳ Move the filter by 1 pixels for a subsequent patch + convolution



↳ Tell model that those are both  $X_s$



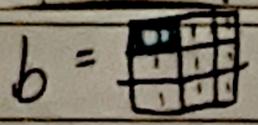
= Features of the images

↳ Then we add a filter ( $3 \times 3$ )

Convolution still allows for the preservation of spatial information because it

in each layer

↳ doesn't join all the values into one large matrix of all the values (they are distributed among the filters)



↳ doesn't multiply them (for each on the grid)

## Abstract:

This scientific article explores the identification of malignant, cancerous tumors via the employment of Convolutional Neural Networks, a deep learning model engineered in Python code. The approach I present develops a machine learning algorithm equipped to attain a threshold of 80% percent within only one week of training, where each day, the CNN is fed

CT scans through <sup>the means of</sup> user input displayed on a personally constructed website. On a daily basis, measures will be implemented to ensure the refinement of the network in order to guarantee optimal performance and results by leveraging the pattern recognition techniques of regularization, performance indicators, optimizers and dataset enhancement, each of which will be meticulously documented in the following chapters. The objective of my novel design\* is to provide a solution to the alarming gap in the availability of rapid and proper resources for individuals affected by cancer. This continues to accelerate the process of diagnosis as well as ~~early~~ offering medical professionals an additional analysis tool among the growing pool of diagnostic frameworks.

\* = change

is to explore + audit the powerful capabilities of CNNs in assisting in the bridging of the alarming gap in the availability of rapid + proper resources for individuals affected by cancer. This continues to accelerate the process of diagnosis as well as offering medical professionals an additional innovative analysis tool among the growing pool of diagnostic frameworks

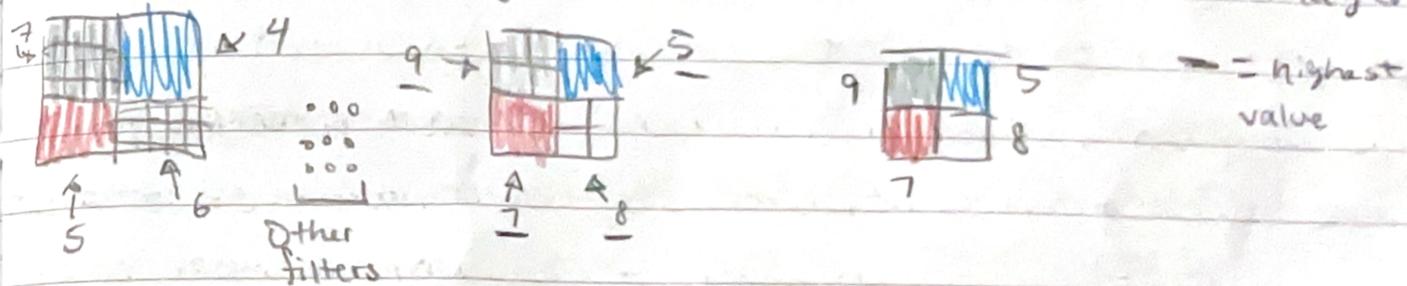
The similarity of the two initial X's is determined by the element wise multiplication that sums up every location's results (which becomes the output).

Equation of convolution (4x4 filter):

$$\sum_{i=1}^4 \sum_{j=1}^4 w_{ij} X_{i+p, j+q} + b$$

Summing over indices i and j from 1-4  
input from prev layers (weighted sum of filter's weights and the corresponding input values)  
for neuron (p,q) within the hidden layer

Maxpooling - The process of extracting the highest, or maximum value of every patch created by the filters to the next convolutional layer



Pooling helps to minimize the dimensionality of the features, allowing for the lack of abundance of hyperparameters, and the preservation of spatial invariance

## Non-Linearity

- CNNs often use ReLU activation functions due to their maximum operation feature that maintains a continuous positive gradient for inputs that return positive and a reaction of zero for negative values

(Not strictly 0 and 1 but  $> 1$  and 0)

- Dismantles (solves) vanishing gradient\*

- Leaves room for nuance + mirrors perception

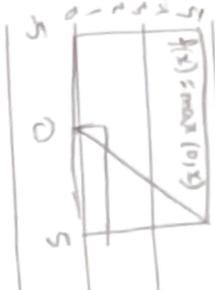
is easily combatted by ReLU because ReLU employs a system of values higher than 0. Vanishing gradient causes some weights to not increase at all because the weight updates gradually become unnoticeable, halting learning and bringing about a lot of frustration and unhappiness.

## NOTE TO SELF

\* = Problem during backprop when the weights become very small as they are fed backwards from the gradient of the loss function (how does a small change in weight affect the final loss). Sigmoid, the activation function, causes this because it squishes values between 0 and 1 so they have even smaller derivatives for input values. This problem

(but better)

\*ReLU is applied to every pixel, replacing all negative values (less than 0) to zero itself



↳ My assumption is that this is to reduce the gravity of unnecessary weights to remain at equilibrium, and to not overfit to overfitting over weights that don't seem more important when they are not

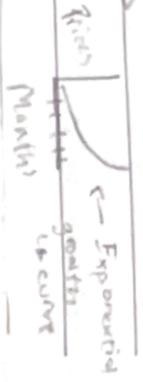
- ReLU is applied after every convolution

Layer to each individual unit, only offering non-negative values

- Nonlinearity allows networks to be accustomed to

the data of the real world where values don't always have to be straight lines

Classification



Once the advanced level features of the previous layers have been outputted, the output is flattened to be interpreted by the final fully connected layer (Dense), which used said features for image classification.

- The classification of the input image is expressed in the form of a probability (how likely this image belongs to — class)

↳ How likely this image depicts the presence of cancer

- In order to convert the output into a probability, we use the sigmoid function, converting the set of values into probabilities between 0 and 1

softmax ( $y_i$ ) =  $\frac{e^{y_i}}{\sum_j e^{y_j}}$

softmax can help inform us of the distance between two probabilities distributions (between the predicted

output and the or accurate output), so it is given context to the mind process of the machine (at

negative values are the same as 0 because there is no activation of weights to carry into further layers, so replacing negatives with zero makes more sense

also because it eliminates unnecessary exceptions

(No activation = 0 and  $\neq < 0$ , etc.). Simplicity

is more sought after to be more efficient,

especially with the creation of a model within

ably the span of a week.

The model currently includes

Fixed CAM in the code

because it visualizes the

productions of the model

while employing the gradient

of the output with respect to

the activations of the CNN

layer. Research says it is

easier to use compared to the

boxes because it is more precise

(think of dots).

equally used in loss quantification to our overall loss

activation = 'softmax'

\* For model ver .h5 file loc → useful for scientific + numerical data (weights, hyper-params) (Gives performance) > Load the files in Python

↳ Told me to use .keras instead because '.h5 format is legal. Please use the .keras format instead'

### Methods:

#### 1 Image-Preprocessing:

- Hardest part is actually uploading images to use on my computer to be classified by the CNN model. The process is very tedious and requires a lot of troubleshooting. Currently, I am trying to

change my directories and file environments to accommodate to the demanded paths of the program (VS code). I already had the files and the folders set, but the paths I tried to use did not work due to not

being set to work in my intended environment. I navigated to my directory → cd "C:\Users\... \code\CNN\imageDownloader"

and adjusted my files to include this path (I think I made the path recognize the files by specifying through Powershell actually).

Every time I tried to run my script: python download-images.py, it gave me one error code after another (realised it was the way

I formatted my file [folder + text; Python file = .py, not an additional

py file which turned to a .py.py file], so I fixed that and went to my

Image Downloader Folder, created a Text Document, opened it with notepad and created a set of code to open the file of the images I

downloaded from The Cancer Archive → IDC-IDR1. I was able to see the names and amount of files I had under my directory by using the dir function in Powershell, after activating the %CD%\Image Downloader

path. When I tried to run my script, it showed another error message regarding my URLs, saying, "500 Server error: for url: ..., Error down-loading: ...", and I believe this might be due to invalid URLs or limited access <sup>to the server</sup>. This is only the image pre-processing <sup>part</sup>.

Thankfully, I have finally found a dataset that does not require me to go through the unnecessary hassle of transcribing the images to a path to download that I do not understand what to do, and this catalogue is: [IQ-OTH/NCCD - Lung Cancer Dataset](#). The dataset includes labeled CT scans of lung cancer, so my assumption is that I must find another dataset where the images are unlabeled. My goal right now is to feed my CNN % of the files tomorrow (there are 1100 images) because it is at its early stage of classification and I want it to be easier for me to understand what is going on. There are a total of 110 cases divided into 3 major categories: benign, malignant and normal (40m, 15b, 55n → cases). Description says they were collected in DICOM format, but on my files, showed up as JPGs

- 561 malignant images
  - 416 normal images
  - 120 benign images
  - 600 - 700 images per day
- Last day (day 7) is the determinant of the model's accuracy. For my CNN, I am using Streamlit as the website for my model, acting as the backend of my interface, while the more decorated Reddymag page introduces the user to my CNN: [ONCO SCANS](#)

The Reddymag houses the network address that the user can then go to (there are 2 websites joint together: front-end, back-end). Streamlit doesn't

allow you to directly edit your application on the website. you have to write a script of code describing the website (appearance), option for image upload, and the model (I used <sup>vscode</sup> Notepad). Then I run the script with Windows Powershell `cas admin`, then I select the directory and environment I used.

What I will do for code:

```
import streamlit (st), pandas (pd), tensorflow (tf), sequential, load_model, Dense, Flatten (keras.models), layers, numpy (np), cv2, os, matplotlib? → all using cv2 and st, so no need (for visualization)
def create_model (input_shape):
    model = Sequential(), model.add (Flatten (input_shape = input_shape))
    model.add (Dense (128, activation = 'relu'), model.add (Dense (1, activation = 'sigmoid')) # Cancer or no cancer [or] activation = 'relu'
```

```
if I wanna also (malignant, benign, normal).
model.compile (optimizer = 'adam', loss = 'binary_crossentropy', metrics = 'accuracy'), return model
```

```
from keras
IMAGE_HEIGHT, IMAGE_WIDTH = 150, 150
EPOCH = 10 (but it can be from min=1 to max=100)
BATCH_SIZE = 32
```

```
MODEL_FILE = 'lung_cancer_detection_model.keras'
base_data_dir = 'os.path.join (os.getcwd(), 'data')
train_data_dir = 'os.path.join (base_data_dir, 'train')
val_data_dir = 'os.path.join (base_data_dir, 'val')
test_data_dir = 'os.path.join (base_data_dir, 'test')

# I have both a python script + a streamlit app
script, st = st.app
more into detail, but not sure why have the python script, I will just keep it, think it will still be useful

# I have both a python script + a streamlit app
script, st = st.app
more into detail, but not sure why have the python script, I will just keep it, think it will still be useful
```

ST = same thing

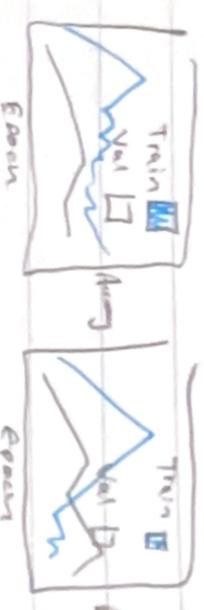
Function + flow - from directory (train-dir, target-size = (224, 224))

batch\_size = BATCH\_SIZE (32), class\_mode = binary

preprocess\_img = load\_img (img-path, target\_size = ST)

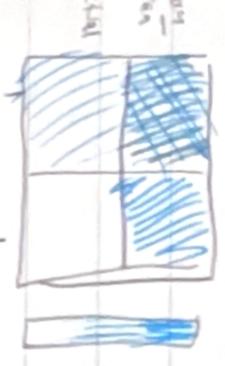
img\_array = img\_to\_array (img), image\_array = np.expand\_dims (img\_array, axis=0)

model.compile for training = Train Accuracy, Val Accuracy, Train Loss, Val Loss



Maybe I should add emojis to make interface more interesting and more fun, less boring and more appealing... (25)

sketch - Learn + Plot its for overall = Precision, Recall, F1, Accuracy



F1 =  $\frac{2 * Precision * Recall}{Precision + Recall}$   
Accuracy =  $\frac{TP + TN}{TP + TN + FP + FN}$



Prediction = model.predict (processed\_image)

result = "Cancerous" if prediction [0][0] > 0.5 else "Non-Cancerous" & cause like

st. subheader ("Prediction Result @")

store the model predicts that the image is {result}

no clue what last conv-layer name so maybe I should do a backward math thing like, last-conv-layer-name = dense\_layer

[-2] or synth similar

cache! if st. button ("Clear Cache"):  
st. cache\_data = clear!  
st. success ("Cache cleared successfully")

btn changed my mind about datasets being only one. Found 3 other datasets on Kaggle to use as well, and I have them downloaded. I have a folder of data my device so I can update my data after every 3 days.

## Day 4: Observations

Windows Powershell model shows table =

Layer (type)	Output shape	Params
Flatten (Flatten)	(None, 15328)	0
Dense (Dense)	(None, 128)	197248
Dense-1 (Dense)	(None, 1)	129

Total Params = 57 603 525 (220.5 MB)  
Trainable Params = 19 267 841 (73.5 MB)  
Non-trainable Params = 0  
Optimizer Params = 38 335 684 (147 MB)

So far right now: 33 steps/epoch

Accuracy: 1.0000, 0.9876, 0.9968, 0.9983, 0.9984 ~ biggest to smallest

Loss: 0.0000e+10, 0.1982, 11.0870, 0.1000, 11.0470, 0.0033, 7.0330e-19, 5.3226e-21, 3.29820e-21, 2.4744e-21, 3.7319e-21, 1.7108e-21, 6.5832e-21

Val-Accuracy: 0.9962, 0.6921

Val-Loss: 1.0174, 1.0199, 1.0200, 0.4470, 0.1944, 0.5785, 0.5155, 0.5860

Notes: 100% accuracy might be an indicator of overfitting

Research says that CNNs are prone to overfitting "due to their high complexity and large number of parameters."

Especially since the loss value is so low, Training loss became as low as 0%, demonstrating the model's ability to learn the training data, but could also be another case of overfitting with all the hyperparameters. The validation accuracy is a good start, showing the model identified patterns that were similar in the training data; can make the strong connections we are looking for. Validation loss could be better though, but the model's relation to unseen data is promising, and I am very excited to see how this will go.

Today, I fed my model 20 480 images (20 epochs, 32 steps per epoch or 20 x 32 x 32 = 20480) -> But the images are repeated so a 2000 images

finish all my data of your streamer app + computer.





# Day 6: Observations

TP = 4 371 764

Tr-P = ~~40000~~ ~~Stitcher~~ 528 193 Num 1 found it

N-TP = 4 049 571

input-target	None, 211231	0
efficient-nr-1	None, 17, 2849	10413
global-nr-1	None, 1280	0
dense	None, 152	317 911
dropout	None, 152	0
None-1	None, 1	152

Accuracy: 0.1291, 0.7295, 0.7500, 0.7574, 0.6562, 0.7630, 0.8438, 0.7592, 0.8438, 0.7984, 0.782, 0.805, 0.8280, 1.0000, 0.4834, 0.5319, 0.5772

Loss: 0.8382, 0.1462, 0.7258, 0.7230, 0.6338, 0.4850, 0.1044, 0.2487, 0.7826, 0.8289, 0.9132, 0.2286, 0.1009, 0.8897, 0.7595, 0.3112, 0.2003

Val-Accuracy: 0.7246, 0.7293, 0.7650, 0.8238, 0.7545, 0.8668, 0.8349, 0.8495, 0.8654, 0.8766, 0.8845, 0.4975, 0.7625, 0.8038, 0.6753, 0.5765, 0.7675, 0.818, 0.627

Val-Loss: 0.2905, 0.3154, 0.3075, 0.3739, 0.3335, 0.3755, 0.3673, 0.3147, 0.7120, 0.1275, 0.7805, 0.1023

Notes: Turns out the repository didn't update all the images. Now that I have updated every image in the dataset from the original to contain images from Dataset 3, the website is taking time to load on my computer and its accuracy is lower than what is was on Day 4.

Plan is to finally do the test section on Day 7 because my own test is when I just try out the app and when my friends do it. Or maybe it is going slower because I'm listening to YT on my browser + its memory usage is 403 MB... probably not right lol. I also added the early stopping instead of manually stopping the training myself. It is slower because of my limited space on Git LFS too. Never always been able to see the actual image, so changing it from 150x150 to 224x224, mirroring those of predevel type.

Chaos from tensorflow. Later in day = Purchased 50 GB of Git LFS and I am using EfficientNetB0 as a form of transfer learning to reach my goal. I added a st-button to show the model layer names because errors keep occurring about input shape and attribute error (button acts as debugging tool) *grad* *Parameter for*

also noticed that 425 to 15/5step -> 570, 144, 280

# Day 7: Observations

\* LAST DAY + Testing \*

TP = 4 571 764

Tr-P = 528 193

N-TP = 4 649 571

OP = 0 *if I still want to want data to work more*

input-target-1	None, 211231	0
efficient-nr-1	None, 17, 2849	10413
global-nr-1	None, 1280	0
dense	None, 152	317 911
dropout	None, 152	0
None-1	None, 1	152

Accuracy: 0.8159, 0.9859, 0.9488, 0.9065, 0.9026, 0.9067, 1.0000, 0.9128, 0.8954, 0.9132, 0.9426, 0.9625, 0.8995, 0.8885, 0.9435, 0.9972

Loss: 0.6620, 0.2170, 0.2200, 0.1097, 0.3889, 0.3631, 0.2400, 0.7041, 0.6959, 0.6938, 0.6994, 0.7041, 0.7052, 2.4095

Val-Accuracy: 0.8962, 0.8108, 0.8583, 0.9072, 0.9945, 0.9029, 0.8954, 0.9554, 0.9785, 0.9575, 0.9125

Val-Loss: 0.3028, 0.2479, 0.4165, 0.9549, 0.1275, 1.00259, 0.1279, 1.1027, 0.2784, 0.5075, 0.2264

Test-Accuracy: 0.9034, 0.9839, 0.9436, 0.9825, 0.9275, 0.9340, 1.0000

Test-Loss: 0.2350, 0.6936, 0.2954, 0.9823, 0.1026, 0.2843, 1.7289, 0.1299, 0.5642, 0.7726

Overall Result: Hypothesis was correct!

GPU should promising results, it did good

Out of 10

Precision: 81.72%, 80.39%, 74.77%, 15.73%, 79.41%

82.86%, 80.73%, 79.25%

Accuracy: 89.23%, 85.44%, 86.14%, 84.92%, 74.21%, 82.05%, 79.05%

Recall: 10000, 8095, 7821, 78.26%

81.19%, 78.26%, 76.55%, 76.4%

80.38%, 81.69%, 82.42%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

\* Still takes some time to load

Fl Score: 61.27%, 84.33%

81.19%, 78.26%, 76.55%, 76.4%

80.38%, 81.69%, 82.42%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

Purchasing?

Recall: 10000, 8095, 7821, 78.26%

81.19%, 78.26%, 76.55%, 76.4%

80.38%, 81.69%, 82.42%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

84.54%

Email: Williams.Udofia.Ehaviates.Ca

Phone: 2K18@hava20214good

Service account details (email) = my-app-account@modelanddata.lan.generative

account.com

1 Ways to improve score (Google search):

1. Higher value of weights to cancerous images (cancerous images are the

exception) + class-labels = [0, 1] # cancer has stronger activation (1), than

no-cancer (0) - " should I do other way around?

2. # of layers: Maybe a random generator hyperparameter tuner (but I think it will take too long maybe) -> could also end up lowering performance

3. Cross validation: further divide my data to x and y (symbolizing features and labels from HOG) then split dataset into k subsets (no live what that is). Eval.

4. Model is trained on this model followed by an additional validation metric

4. HOG (Histogram of Oriented Gradients), which uses edges and shapes to capture

distinguishing features of C-scans and NC-scans, helpful to improve performance

of CNN when performance is low

5. Transfer Learning, where I find a CNN structure to implement so as to

train the CNN on how to see lung cancer images in particular, as it

already has had practice on other types of images -> Ex = DenseNet

(121, 169, 201, etc.), VGG16, ResNet, etc.

or I'll do DenseNet; saw grad ATM implementation already + works

good on datasets with not too many images; faster and intuitive

or V-DenseNet is too heavy and took too long to load, spouted many errors,

and I placed my app with problems. Now I switched to EfficientNetB0 and it

is significantly more lightweight; allows for a smoother run time on Databricks

and still has a high accuracy (F1 score) when implemented into code, already 61.27% from the get-go

Trying out test feature

Currently, when testing:

(Day 4)

Confusion matrix

	Actual C	Actual N-C
Predicted C	197	185
Predicted N-C	149	149

144 True C (P)  
219 True N-C (P)  
155 Wrong N(N)  
197 Wrong C (N)

	Actual C	Actual N-C
Predicted C	197	185
Predicted N-C	149	149

Possibly points towards a bias the model has towards C-C score

Accuracy =  $\frac{144 + 219}{144 + 219 + 177 + 185} = \frac{363}{725} \approx 0.49924$

Precision =  $\frac{144}{144 + 197} = 0.421$

Precision =  $\frac{149}{149 + 149} = 0.499$

149 True C  
224 True N-C  
192 Wrong C  
180 Wrong N-C

Precision =  $\frac{149}{149 + 192} = 0.435$

$\approx 0.44$  or  $44\%$



	Model	CNN	DenseNet	MobileNetV2	EfficientNet
Makes					
F1		1	1	1	1
Accuracy		1	1	1	1
Less		1	1	1	1
heavy		1	1	1	1
slow		1	1	1	1
to load		1	1	1	1
to train		1	1	1	1

CNU Names

(Order of names in)

ONCOSCAN (air exists)

Endings

First

PULMOSCAN (air exists)

1. SCAN X

1. ONCO

PULMONUM

2. NET X

1. PULMO

ONCOVISION

3. AI

1. PULMO  
2. LUNGA

Pulmo AI

4. VISION X

4. RESPIRA  
5. BRONCHI

ONCOCAN

5. CAN X

1. BREATH  
2. LARYNG-

PULMO AI (air exists!)

4. PNEUM-

Pulmo Net X air exists.

(It looks like latin on the site help

Pulmo CNU?

RESPIRASCAN

ONCOlung NET

Respira Scan

ThoraxSense

ONCO AI

WIMSCAN?

## Analysis

- Custom CNN model, once AI, attains a score of 80% accuracy, displaying the amazing abilities of CNNs and their ability to learn at an incredible rate, proving my proposed hypothesis to be correct. 8/10 tries, model reached and surpassed the 80% threshold
- I employed (1) Normal, 3 layer, then 4 layer CNN structure, then (2) the VGG-16, but then I quickly replaced it with the DenseNet 121 because I searched up online which structure is more compatible with CT scans, or better at accessing CT scans, and all pointed in favor to DenseNet.