# X-ray Classification with Chandra

Shreeya Behera
Adam Broussard
Karthik Prabhu

Project Mentor: Michael Darcy

# Problem statement
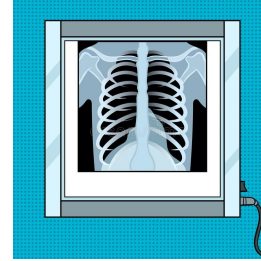
**Data source:**

We use Chest X-ray images from Guangzhou Women and Children's Medical Center, China. The patients are children of one to five years old.
http://www.cell.com/cell/fulltext/S0092-8674(18)30154-5

5860 X-ray scans (1585 - Normal, 4275 - Pneumonia)

**Stakeholders:**

- UNICEF reports that pneumonia kills more children than any other illness.
- It very hard to diagnose X-ray scans, and take decisions deterministically at the right moment.
- Stakeholders are going to be facilities with imaging capabilities who may not have a doctor on hand to interpret X-rays
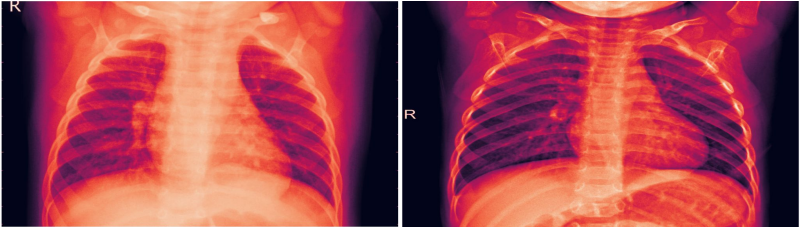
**Chosen Metric**

$$F1 = 2\frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

**Metric balances ensuring we catch all of the pneumonia cases (Recall) with not over-classifying cases as pneumonia (Precision)**
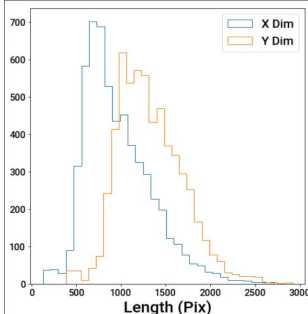
# Exploratory Data Analysis
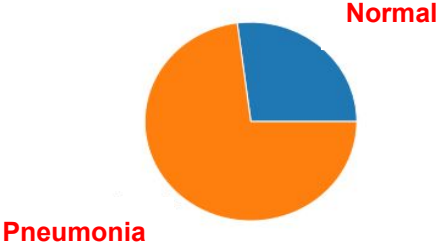


**Pneumonia**       **Normal**

In cases with Pneumonia, images have higher opacity compared to the normal. Therefore we expect the normal cases to have a higher variance

## Challenges with the dataset



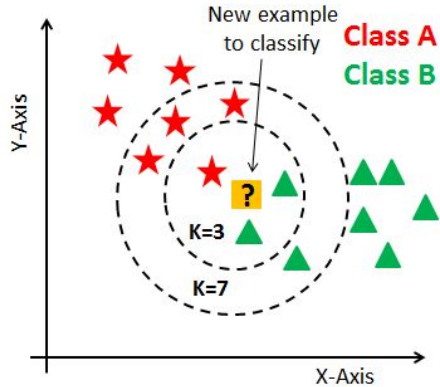**Wide variety of image dimensions and aspect ratios**
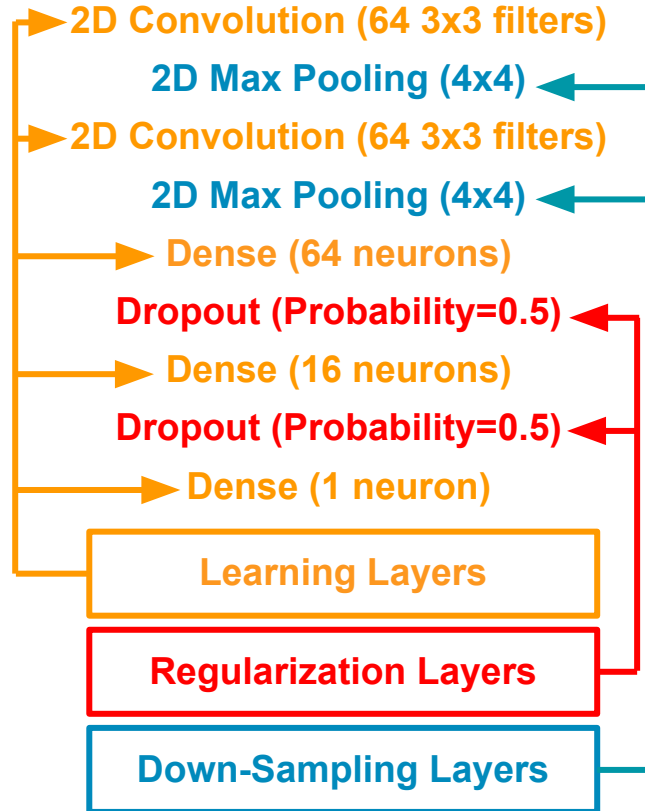


**Dataset skewed towards Pneumonia cases**
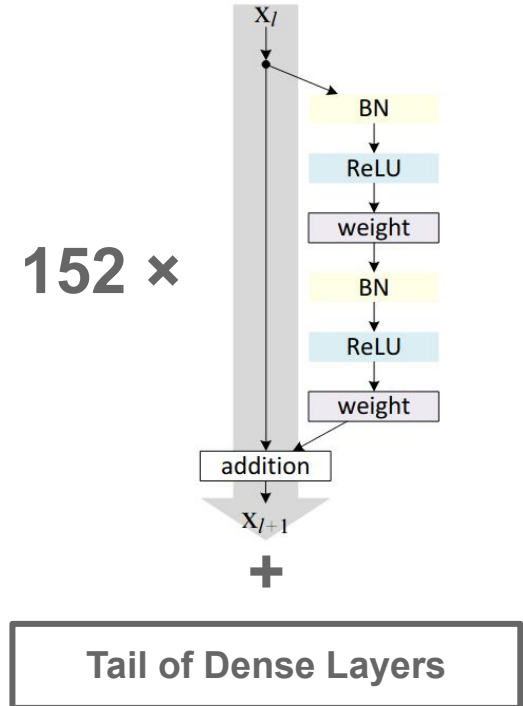


**Some images were in color**

# Model Selection



**Model 1:**
**K- nearest neighbors**
**(KNN Model)**

**Model 2:**
**Convolutional Neural Network**
**(CNN Model)**

**Model 3:**
**Transfer Learning**
**(TL Model)**

**2D Convolution (64 3x3 filters)**
**2D Max Pooling (4x4)**
**2D Convolution (64 3x3 filters)**
**2D Max Pooling (4x4)**
**Dense (64 neurons)**
**Dropout (Probability=0.5)**
**Dense (16 neurons)**
**Dropout (Probability=0.5)**
**Dense (1 neuron)**

**Learning Layers**

**Regularization Layers**

**Down-Sampling Layers**

$x_l$
BN
ReLU
weight
BN
ReLU
weight
addition
$x_{l+1}$

**152 ×**

**+**

**Tail of Dense Layers**

# Data Preprocessing

| Redistributing Train, Test, and Validation Datasets | Homogenizing Image Dimensions | Transforming Categorical Variable into Binary |
|---|---|---|

- Dataset is pre-sorted into train, test, validation sets
- Combine all of these into a single dataset and re-distributing them, giving full control over data splitting and stratification
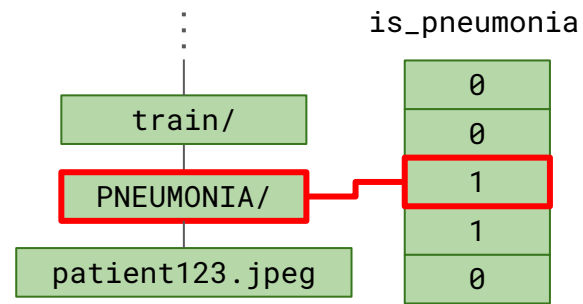
- Resize the images, maintaining axis ratio to prevent distortion and loss of information
- Crop the long axis of the image from the center to protect the region with the most information

- Read the folder structure to generate a binary indicator of pneumonia
- Save resized images along with references to their original files

Initial     Final

Test
Val
Train

968x592    370x224    224x224

is_pneumonia

train/

PNEUMONIA/

patient123.jpeg

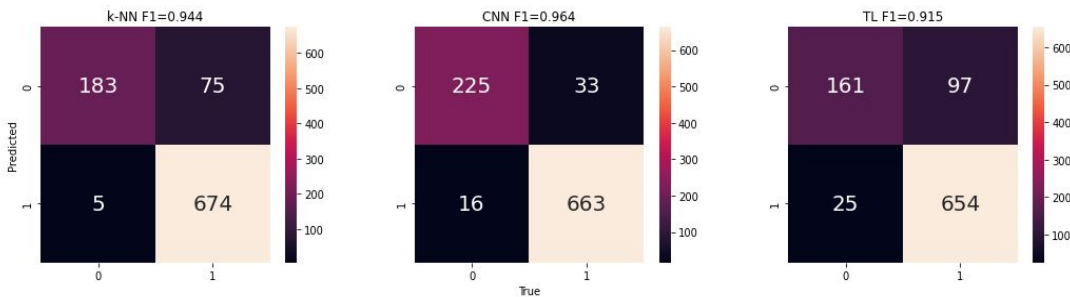| |
|---|
| 0 |
| 0 |
| 1 |
| 1 |
| 0 |

# Model Selection and Validation

**k-NN Model:** Struggles to correctly identify a sizeable portion of the pneumonia cases

**CNN Model:** Greatly improves false negatives with only a small increase in false positives

**TL Model:** More false negatives and false positives than the other two models
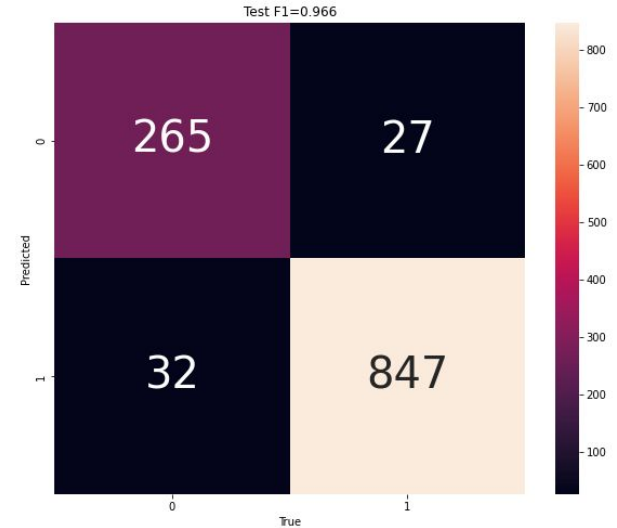
**F1-Score of the three models on validation sample**

| *k*-NN | CNN | TL |
|--------|-----|-----|
| 0.944 | 0.964 | 0.915 |



We choose the custom-trained CNN as our final model as this achieves the highest **F1 score** of **0.964**.

# Final Conclusions

- F1 Score balances finding all pneumonia cases against correctly classifying pneumonia
- Trained three models of varying complexity: a *k*-Nearest Neighbors model, a convolutional neural network (CNN) model, and a CNN with transfer learning
- Recommend CNN model to healthcare providers
  - Low incidence of false negatives means patients get the care they need
  - Low incidence of false positives means radiologists don't waste time treating healthy patients



CNN with a test-sample F1 score of 0.966