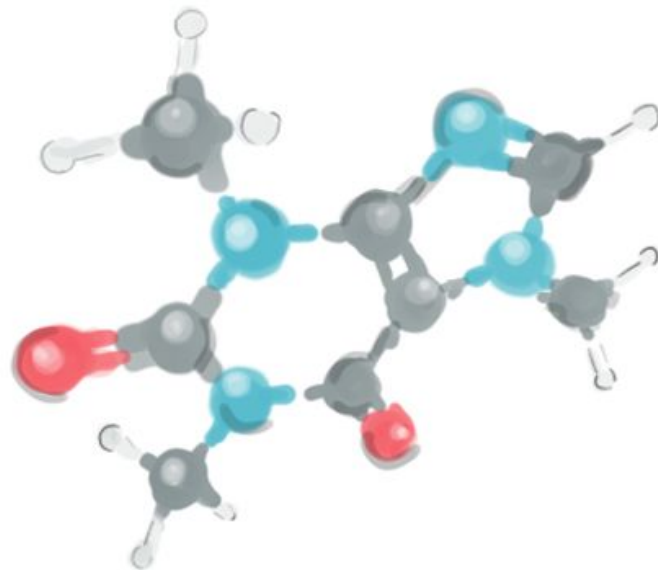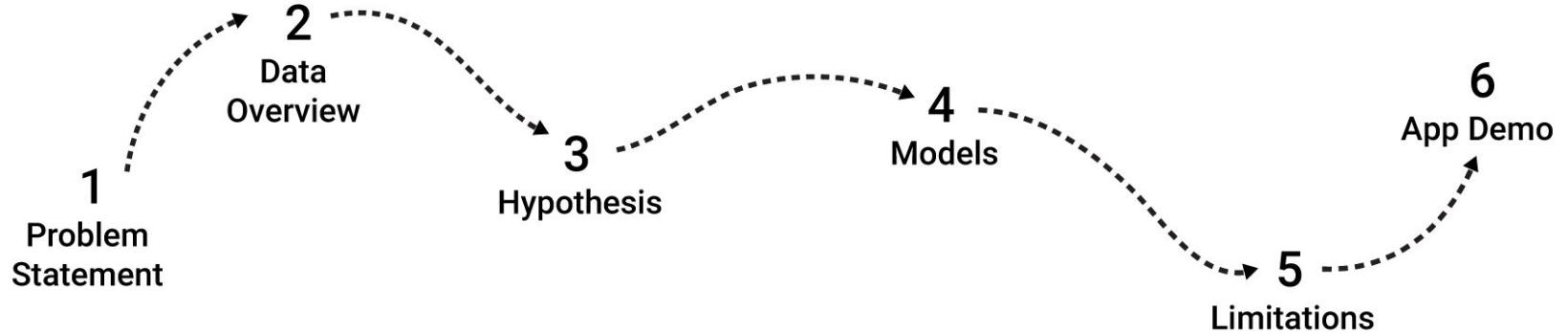# Multi Class Drug Classification Using Molecular Structures

**Vivian Peng**

https://github.com/veeps/molecular_classification

# Overview



1 Problem Statement

2 Data Overview
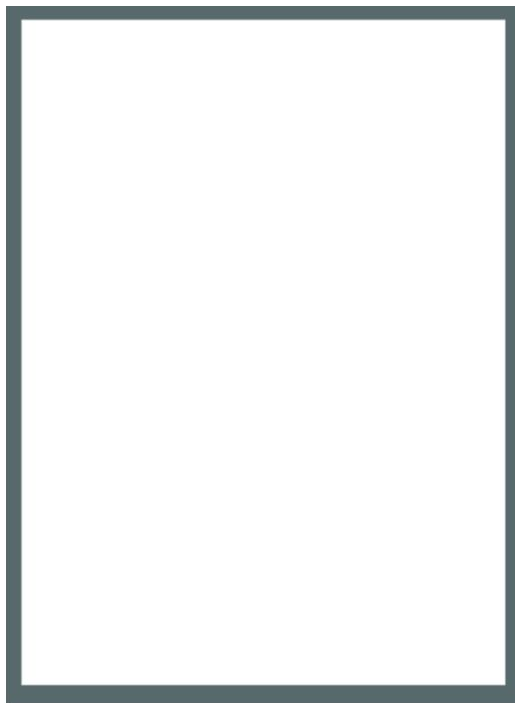
3 Hypothesis

4 Models

5 Limitations

6 App Demo

# Problem Statement

Current drug R&D is long and costly...
And it doesn't incentivize for research into treatments for neglected diseases.

1

What if we could improve the R&D process by reducing redundancy, and screen for multiple therapeutic uses in parallel?
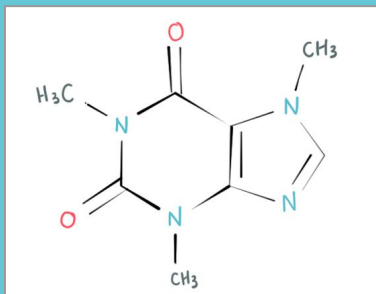
## Neural Network

Antineoplastic

CNS

Cardio

# Here's what my data looked like:

Distribution of Molecular Weight by Drug Class

# Ran 3 types of models:

**CSV file w/ chemical properties**

SVC

**2D Images of molecular structures**

CNN

**SMILES of molecular structures**
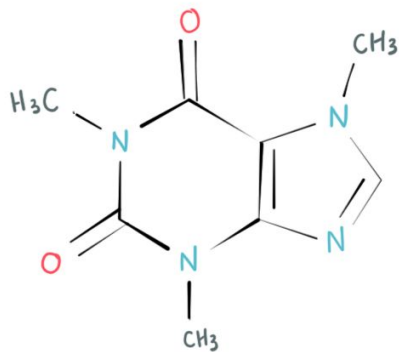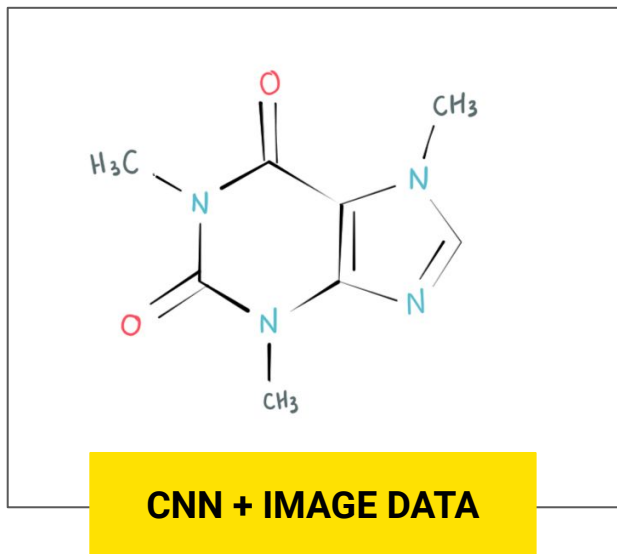
RNN

# How do these neural networks compare in classifying drug classes using molecular structure data?



**CNN + IMAGE DATA**



SMILE

CNlC(=O)N(C)c2NCN(C)c2Cl=O

**RNN + TEXT DATA**

CNN + IMAGE DATA

My hypothesis was that a **CNN model** using image data would be better at classifying than an RNN model using SMILES.

# Models

# SVC

Support Vector Classification

## FEATURES:

- **Molecular Weight**
- **Hydrogen Bond Acceptors**
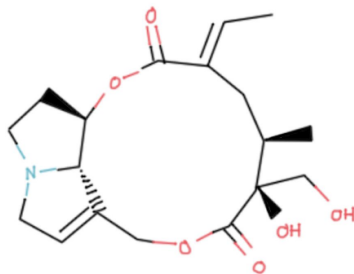- **Hydrogen Bond Donors**
- **XlogP (Solubility)**

## HYPERPARAMETERS:

- **Ridge penalty**
- **Linear kernel**
- **Gamma = "scale"**

# CNN

Convolutional Neural Network

## DATA PROCESSING:
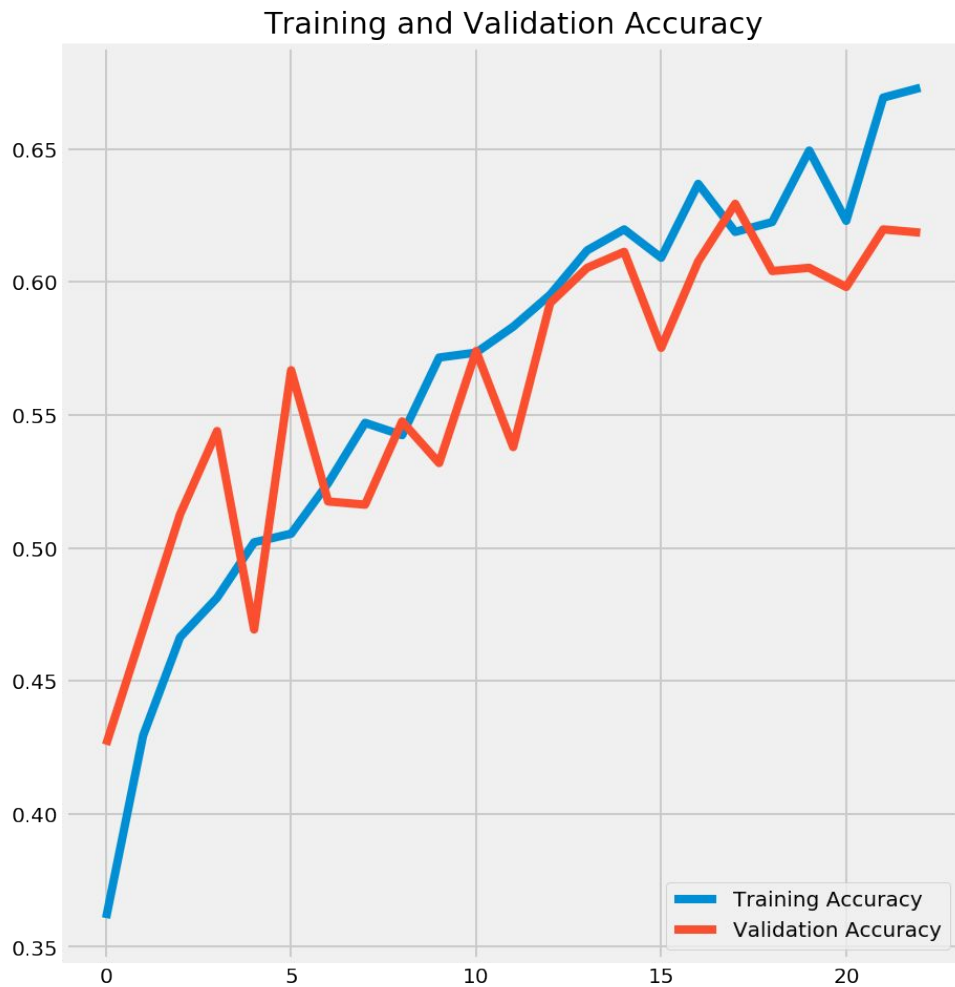


```
array([[[1., 1., 1.],
        [1., 1., 1.],
        [1., 1., 1.],
        ...,
        [1., 1., 1.],
        [1., 1., 1.],
        [1., 1., 1.]],
```

## MODELS:

- **Custom convolutional neural net**
- **Pre-trained VGG16 model**
  - 2 hidden layers (256 and 128 neurons)
  - Adam optimizer

# CNN

## Convolutional Neural Network

Training and Validation Accuracy

Training Accuracy
Validation Accuracy

# RNN

Recurrent Neural Network

## DATA PROCESSING:

Brc1c(NC2=NCCN2)ccc2nccnc12
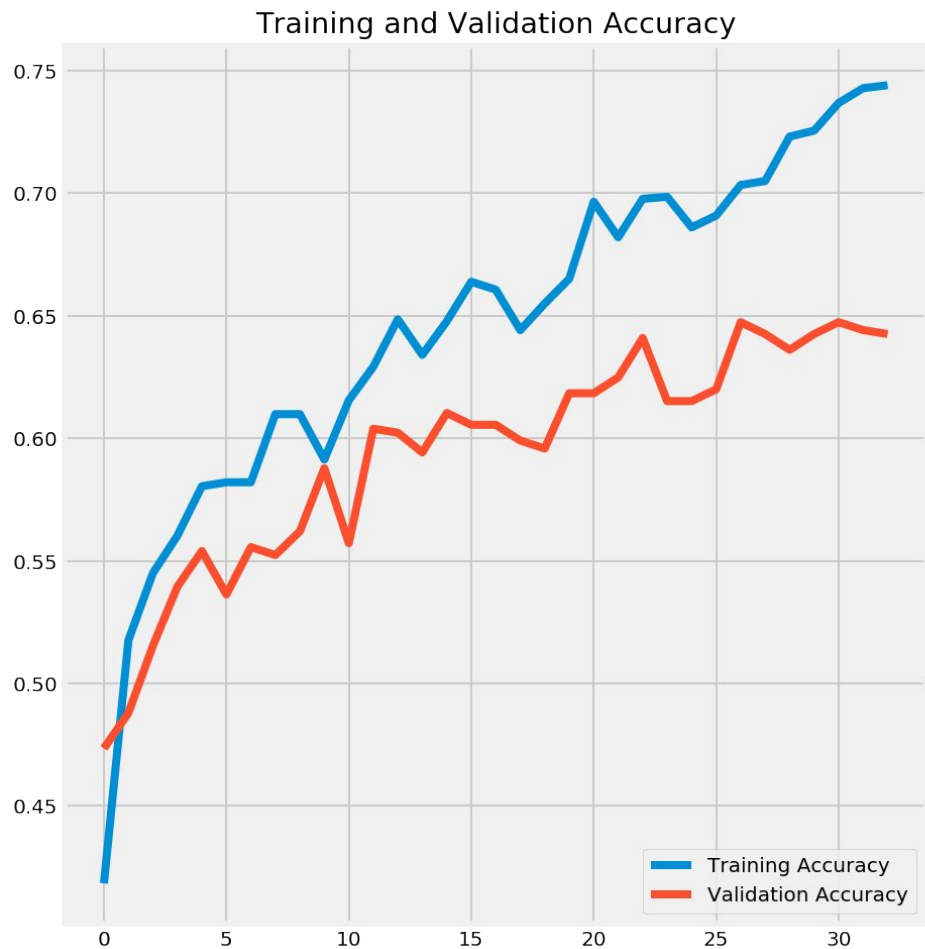
```
array([ 0.,  0., 10.,  6.,  2.,  2.,  1.,  0.,  2.,  0.,  0.,  1.,  0.,
        2.,  2.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,
        0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,
        0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.])
```
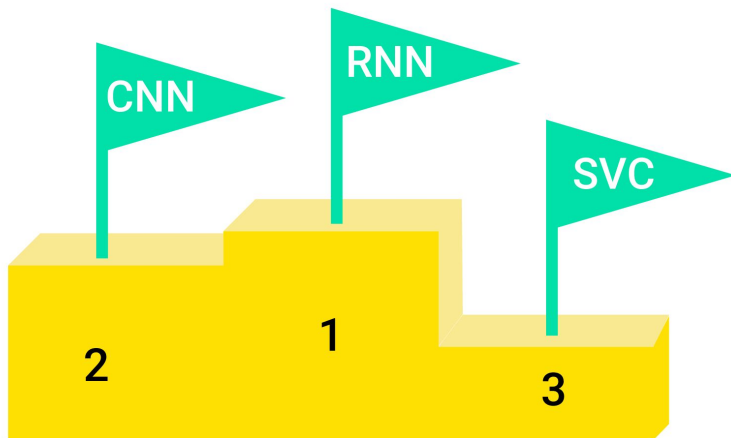
## MODEL:

- **Simple RNN**
  - 5 hidden layers
  - Adam optimizer
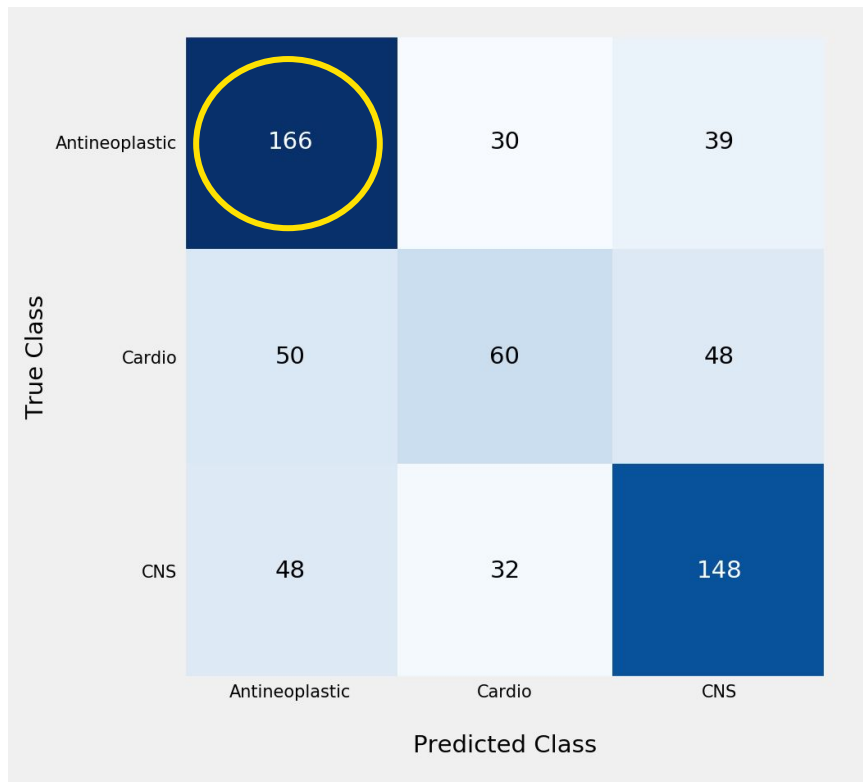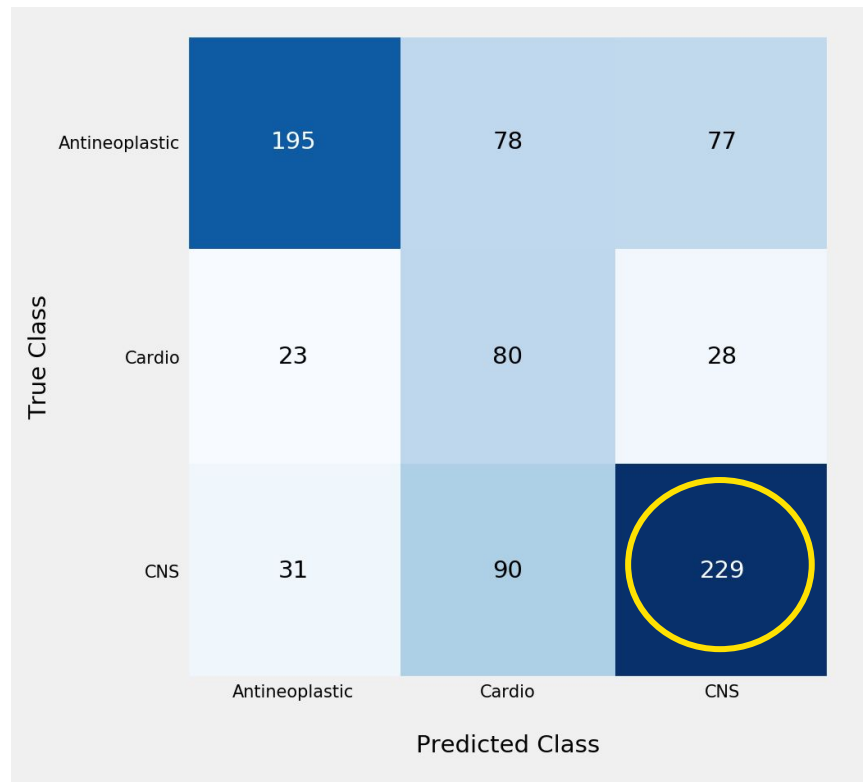- **LSTM**

**RNN**

Recurrent Neural
Network

Training and Validation Accuracy

Training Accuracy
Validation Accuracy

| Model | Score |
|---|---|
| RNN | 0.64 |
| CNN | 0.62 |
| SVC | 0.53 |
| **Baseline** | **0.37** |

Both types of neural networks performed **relatively the same** in predicting drug classes. **This is really interesting** because running an RNN with text data is **computationally much lighter** than running a CNN with image data.

**RNN**

**CNN**

Models are good at predicting different types of drug classes.

# Limitations

- Unable to tokenize SMILES to keep two-letter elements as one unit
- Lack of chemical expertise
- Unable to deploy app to Heroku

# App Demo

# Thank you!