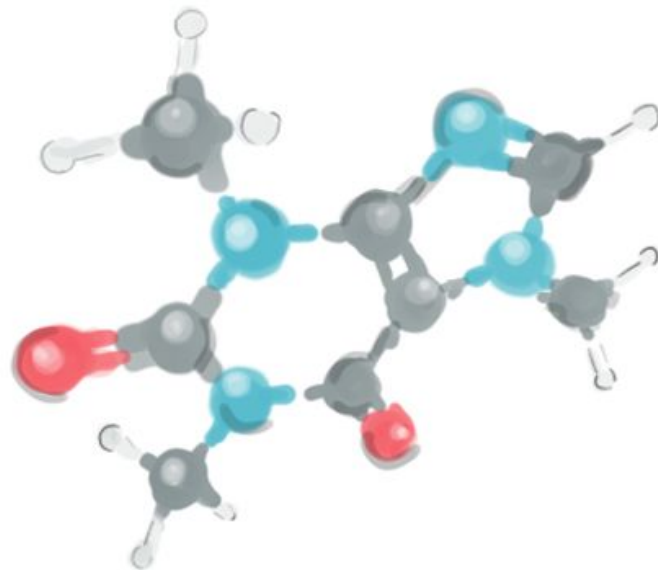


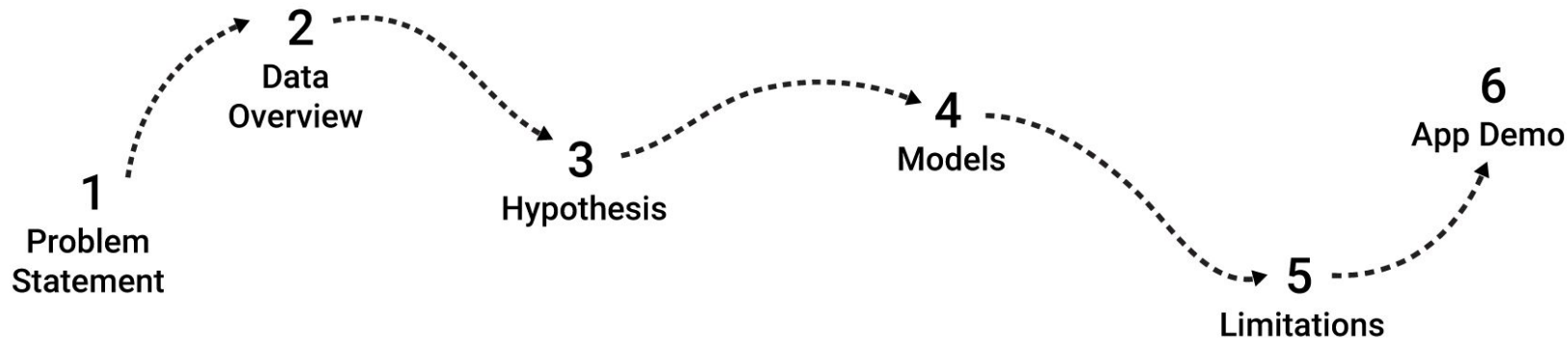
# Multi Class Drug Classification Using Molecular Structures

Vivian Peng

[https://github.com/veeps/molecular\\_classification](https://github.com/veeps/molecular_classification)



# Overview



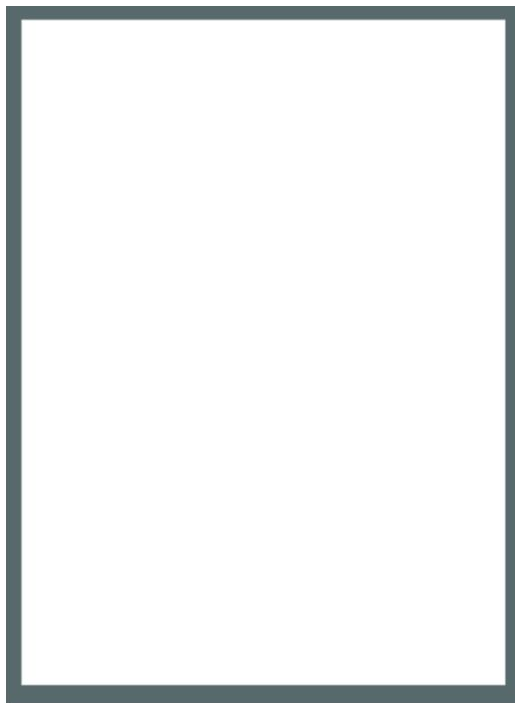


# Problem Statement

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

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

# Neural Network



**Antineoplastic**

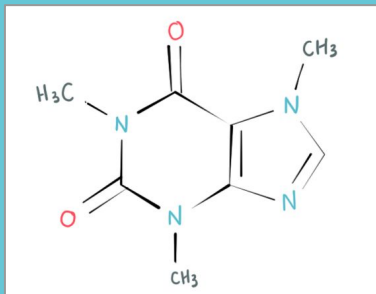
**CNS**

**Cardio**

# Here's what my data looked like:

## CSV file w/ chemical properties


## 2D Images of molecular structures

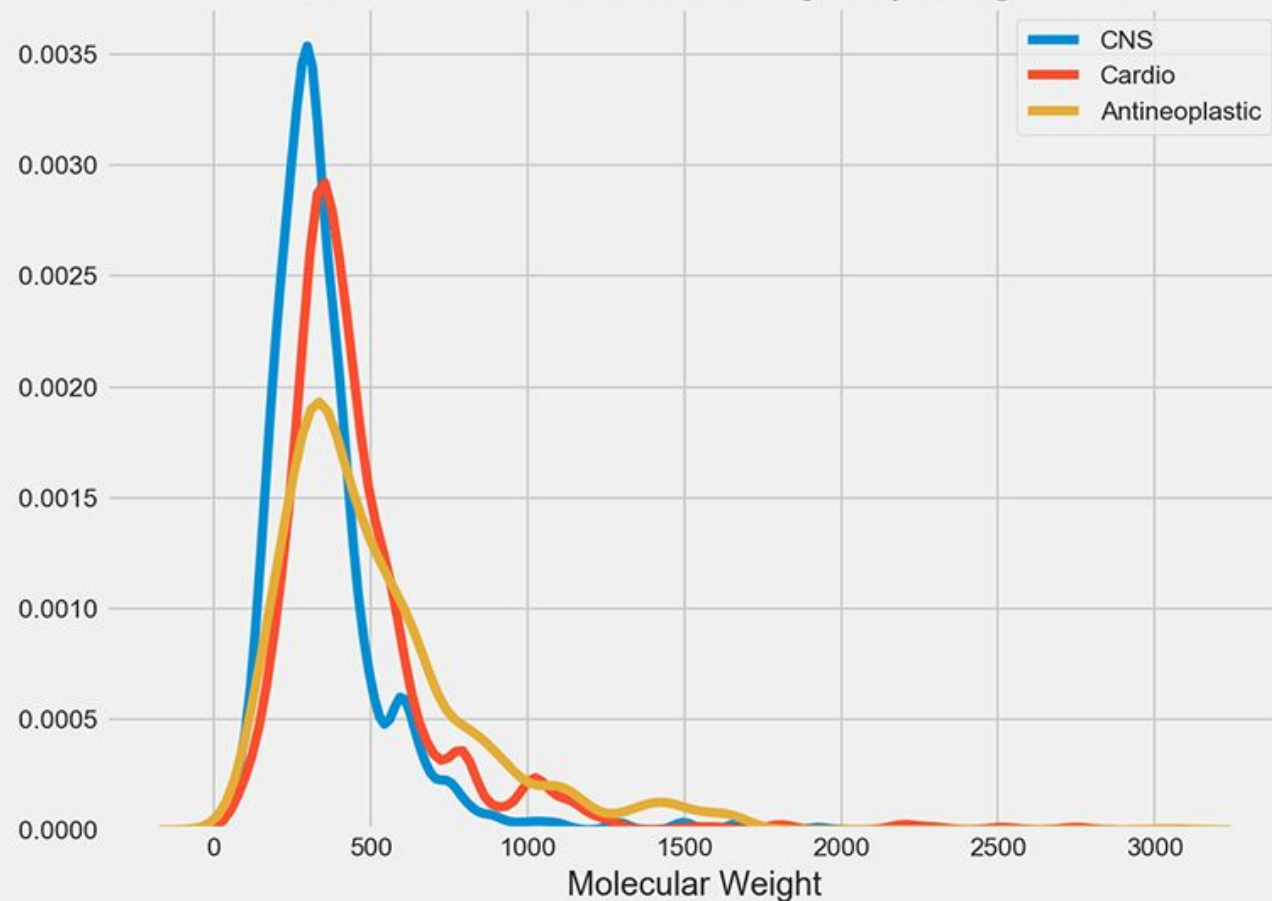


## SMILES of molecular structures

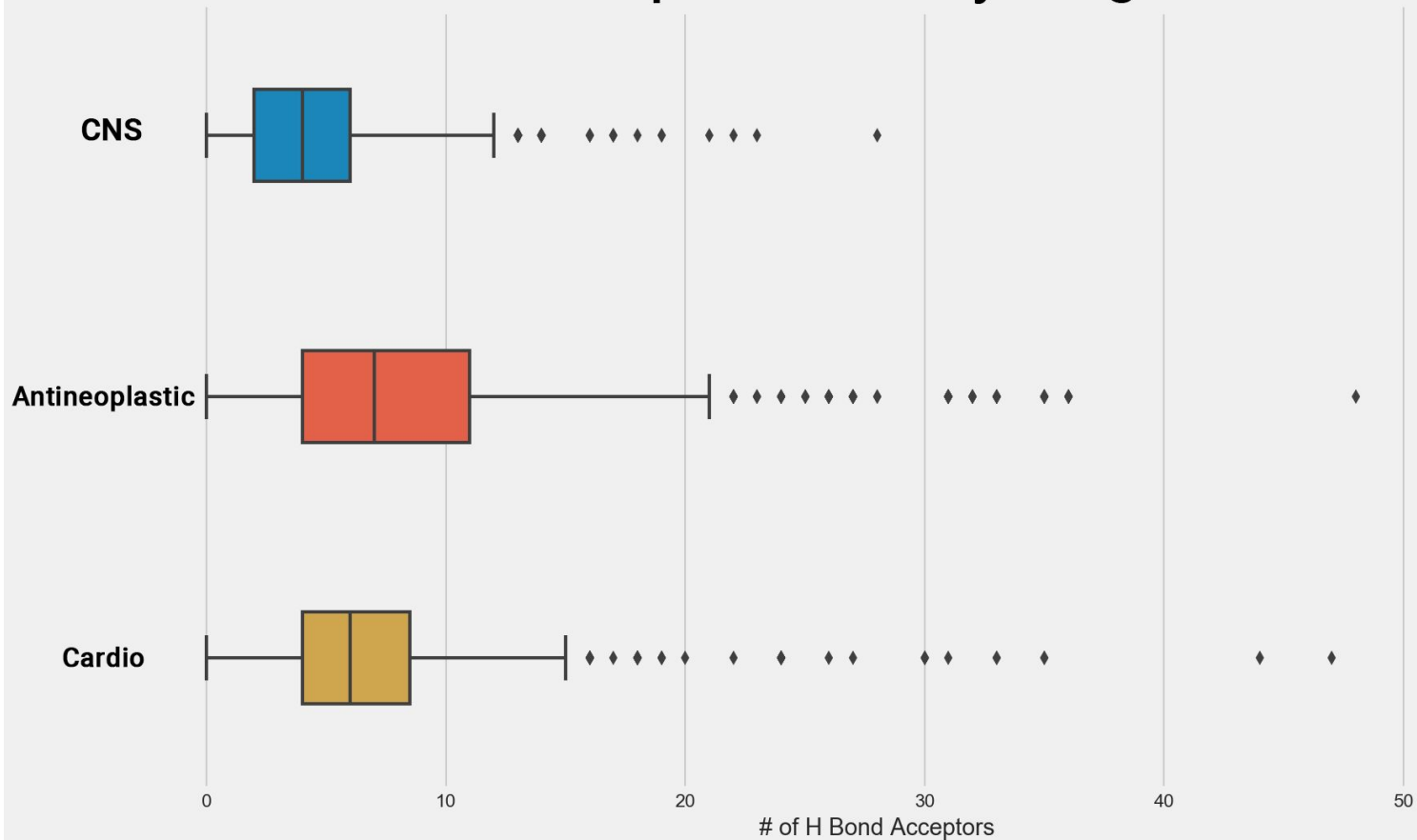
SMILE  
CN1C(=O)N(C)C2NCN(C)C2C1=O

SMILE  
CN1C(=O)N(C)C2NCN(C)C2C1=O

Distriubtion of Molecular Weight by Drug Class



# H Bond Acceptor Count by Drug Class





# Ran 3 types of models:

CSV file w/ chemical  
properties

**SVC**

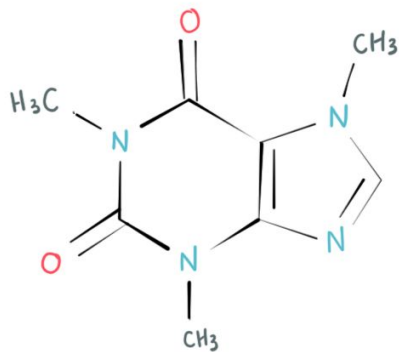
2D Images of  
molecular structures

**CNN**

SMILES of  
molecular structures

**RNN**

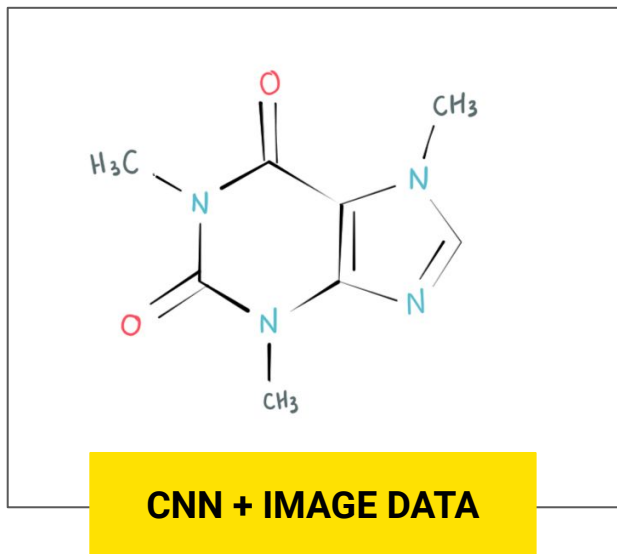
# How do these neural networks compare in classifying drug classes using molecular structure data?



**CNN + IMAGE DATA**

SMILE  
CN1C(=O)NC2=C1N(C)C(=O)N2C

**RNN + TEXT 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)

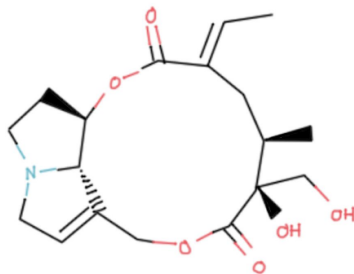
## HYPERPARAMETERS:

- Ridge penalty
- Linear kernel
- Gamma = “scale”

# CNN

Convolutional  
Neural Network

## DATA PROCESSING:



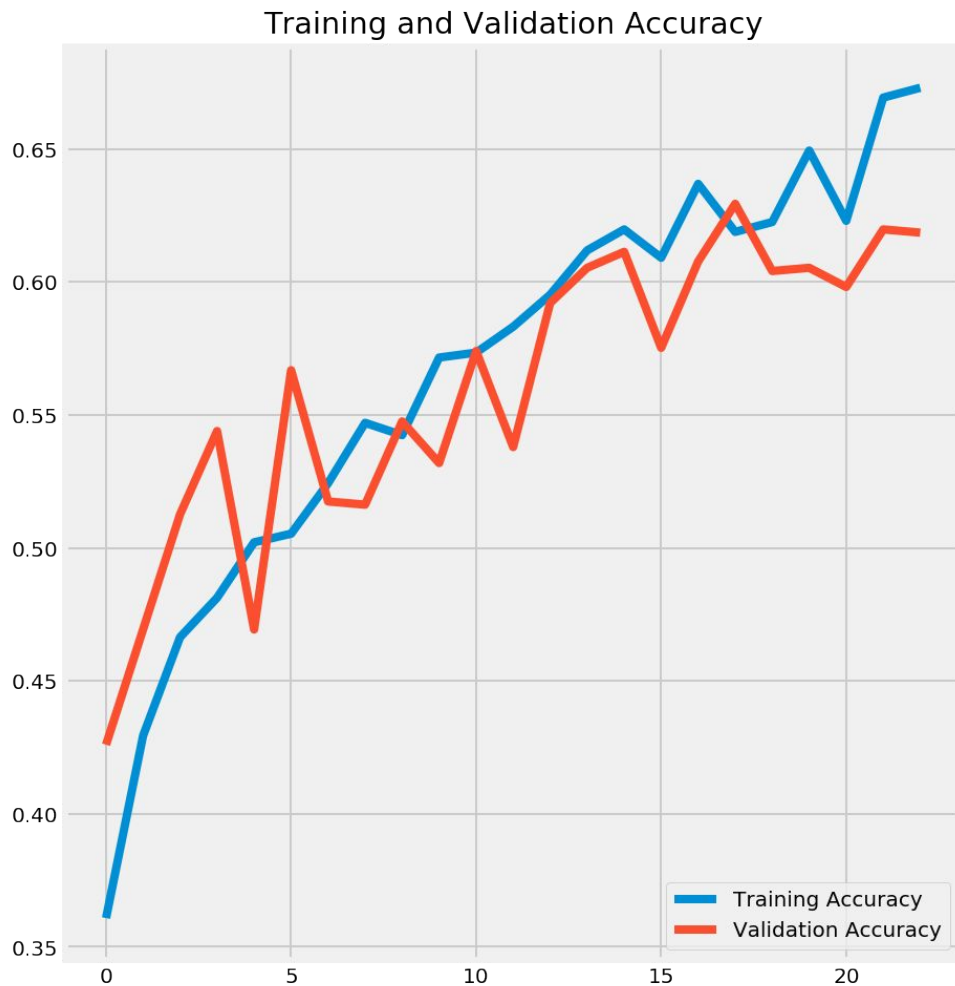
```
array([[1., 1., 1.],  
       [1., 1., 1.],  
       [1., 1., 1.],  
       ...,  
       [1., 1., 1.],  
       [1., 1., 1.],  
       [1., 1., 1.]])
```

## MODELS:

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

# CNN

Convolutional  
Neural Network



# 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.]])
```

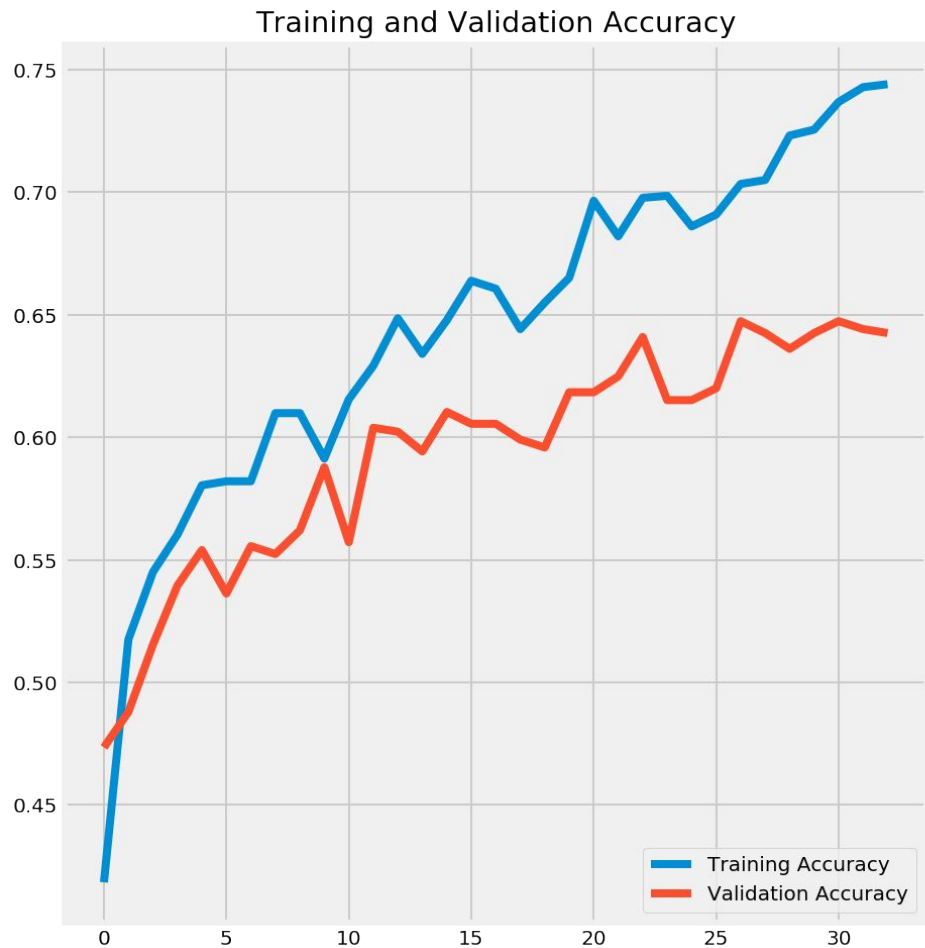
## MODEL:

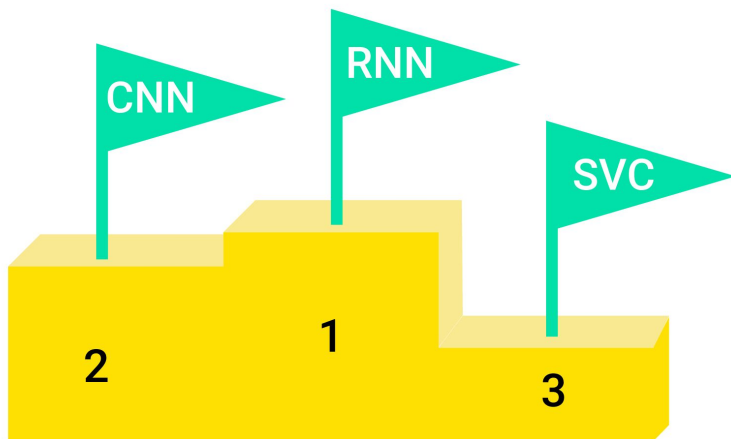
- Simple RNN
  - 5 hidden layers
  - Adam optimizer
- LSTM



# RNN

Recurrent Neural  
Network

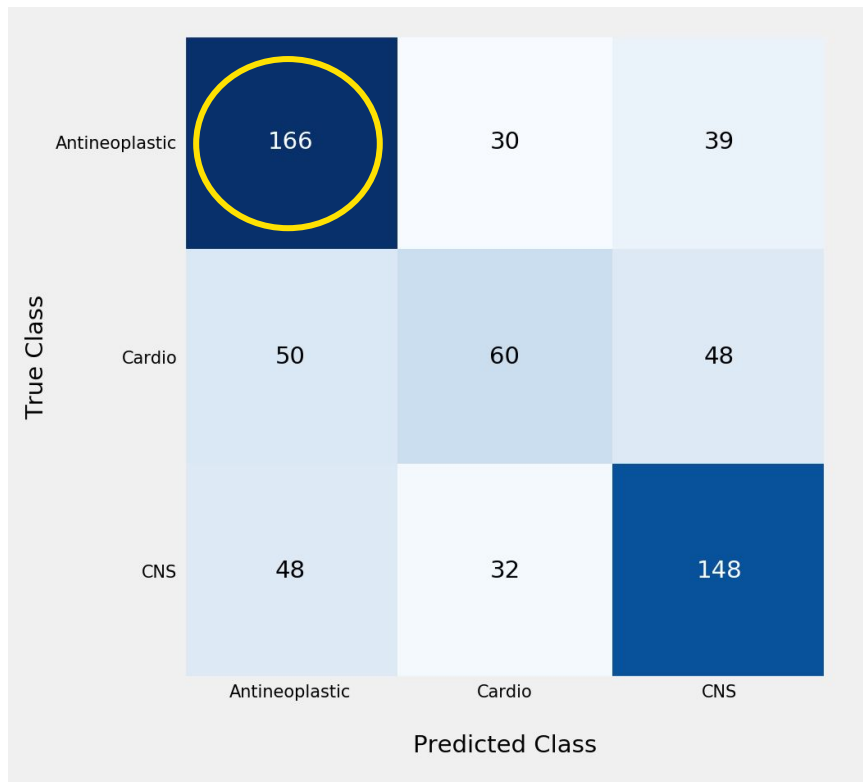




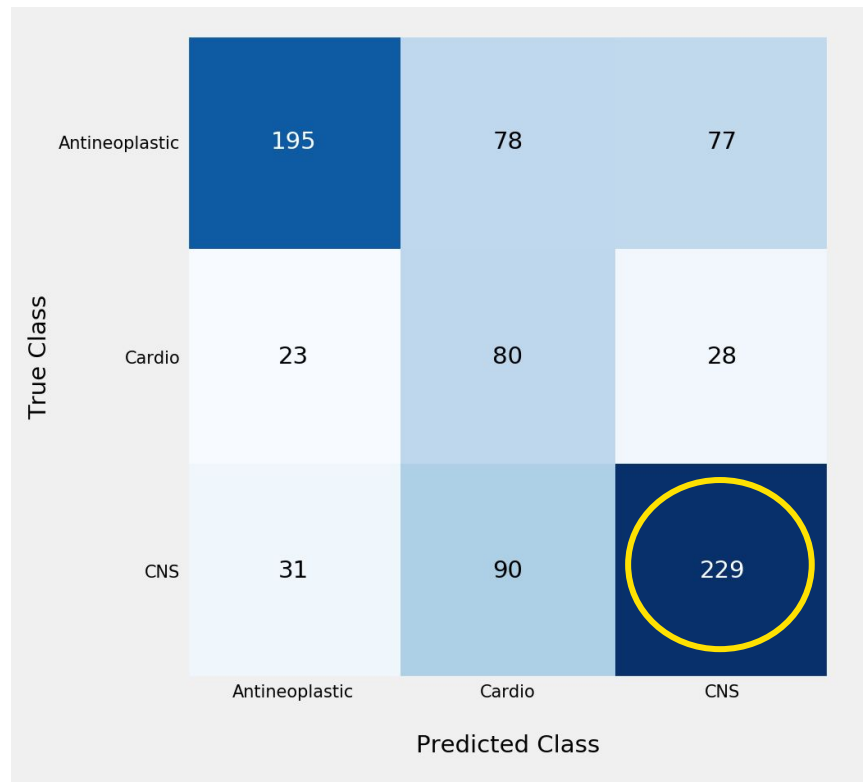
Model	Score
RNN	0.64
CNN	0.62
SVC	0.53
Baseline	0.37

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

What's a likely  
drug class for  
your chemical?

Upload an image!

Choose File No file chosen

Upload

**Thank you!**