

IS Datathon

Team 5

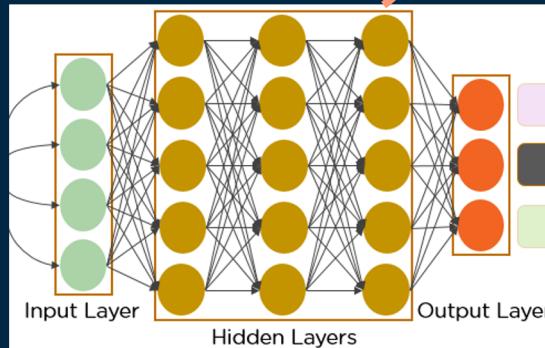
Wang-Han Li
Chung-Hao Lee
Fabienne Yang

H E A R S T

Goal



Input



CNN Model

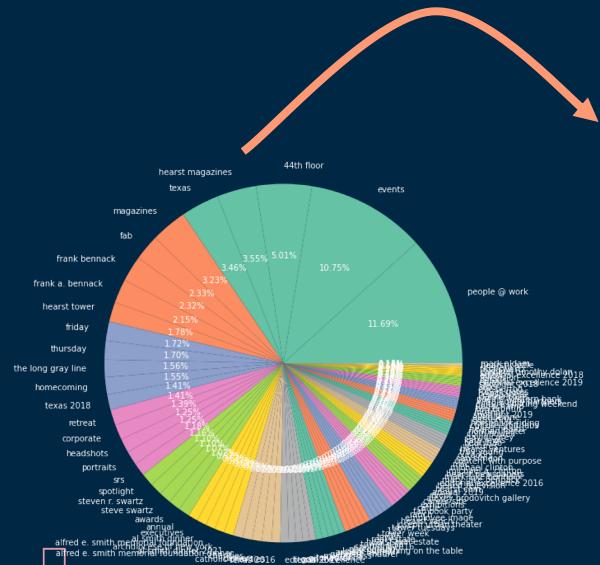
Output

Keywords:

- headshots
- portraits
- srs
- steve swartz
- steven r. swartz

Keyword Exploration

Filter out four digit Year and keyword counts < 93 (1% data points)



1388 keywords

A horizontal bar chart titled "Top 112 Keywords by Count". The x-axis is labeled "Count" and ranges from 0 to 7,000. The y-axis lists keywords with their corresponding ranks. The bars are blue, and the exact count value is displayed at the end of each bar.

Keyword	Ranking	Count
people @ work	1	7,059
events	2	6,493
44th floor	3	3,025
hearst magazines	4	2,143
texas	5	2,090
magazines	6	1,950
fab	7	1,408
frank bennack	8	1,401
frank a. bennack	9	1,301
hearst tower	10	1,076
...		
editorial excelle..	103	113
details	104	112
television	105	108
editorial excelle..	106	102
scouts	107	101
cardinal timothy..	108	95
tina swartz	108	95
hearst castle	110	94
scenery	110	94
mark aldam	112	93

112 keywords

Data Processing

Step1: Split keywords

Portraits,Headshots,2017,Steven R. Swartz,revised 2022



Portraits	Headshots	2017	Steven R. Swartz	revised 2022
-----------	-----------	------	---------------------	--------------

Step3: Filter out year and low frequency keywords

portraits headshots ~~2017~~ steven r. swartz revised 2022



portraits headshots steven r. swartz

Step2: Convert keywords to lowercase

Portraits	Headshots	2017	Steven R. Swartz	revised 2022
-----------	-----------	------	---------------------	--------------



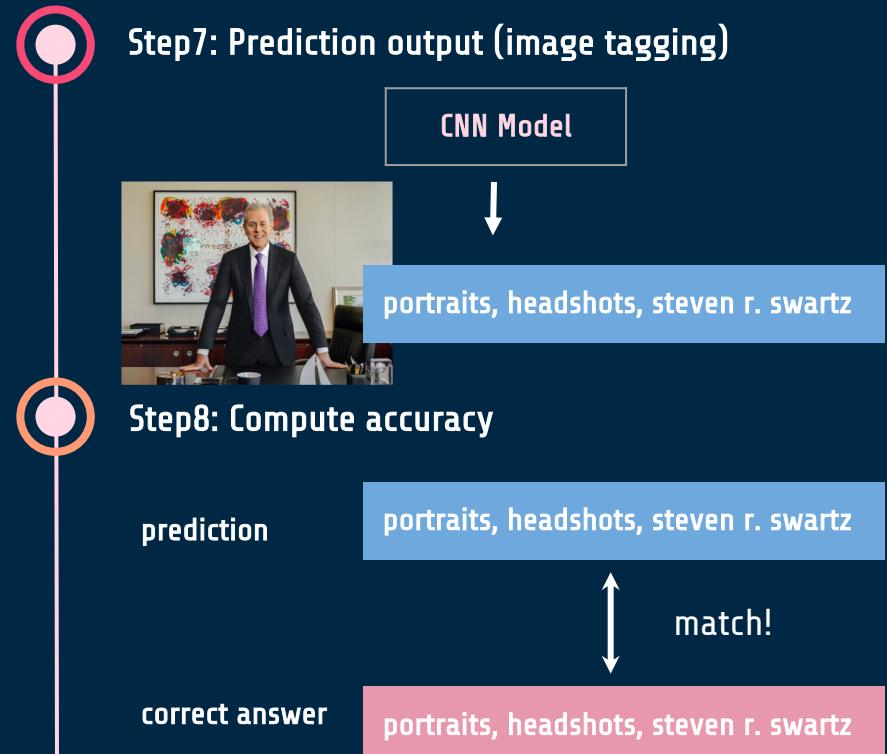
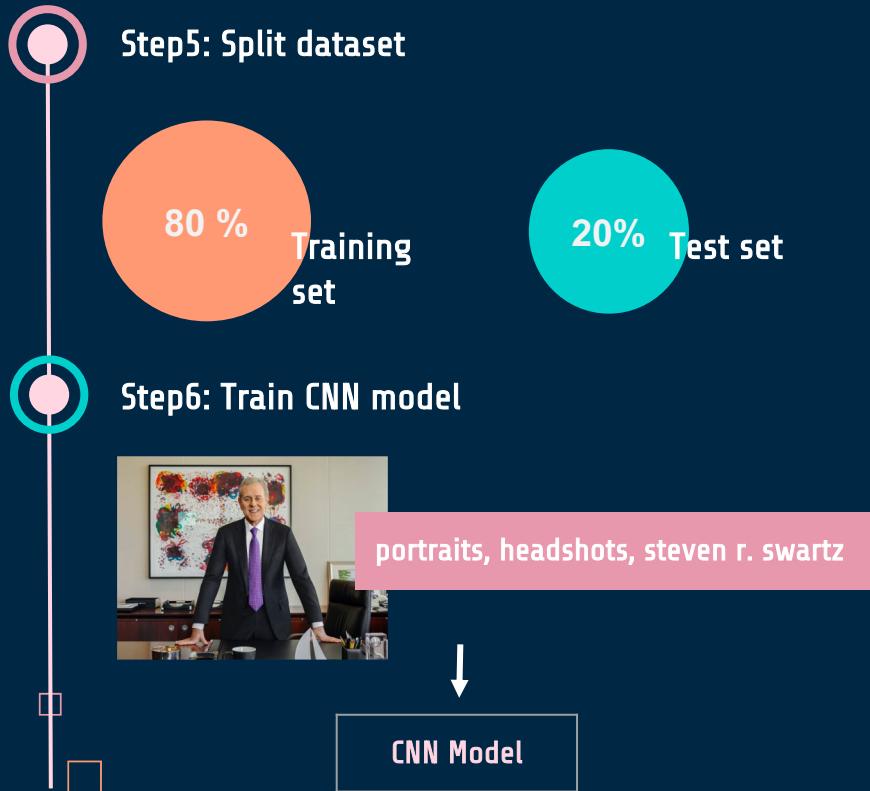
portraits headshots 2017 steven r. swartz revised 2022

Step4: Connect image with keywords



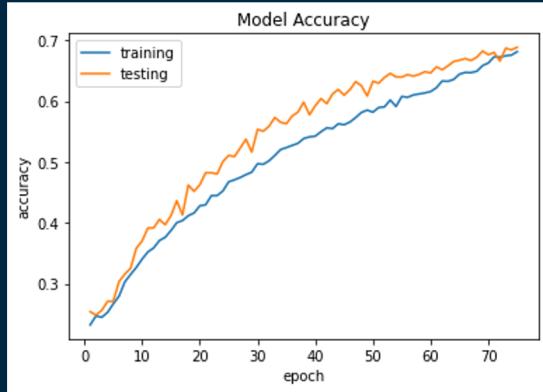
portraits, headshots, steven r. swartz

Build Image Classification Model



Model Performances

Best model has highest testing accuracy: 68.86%



Output examples:



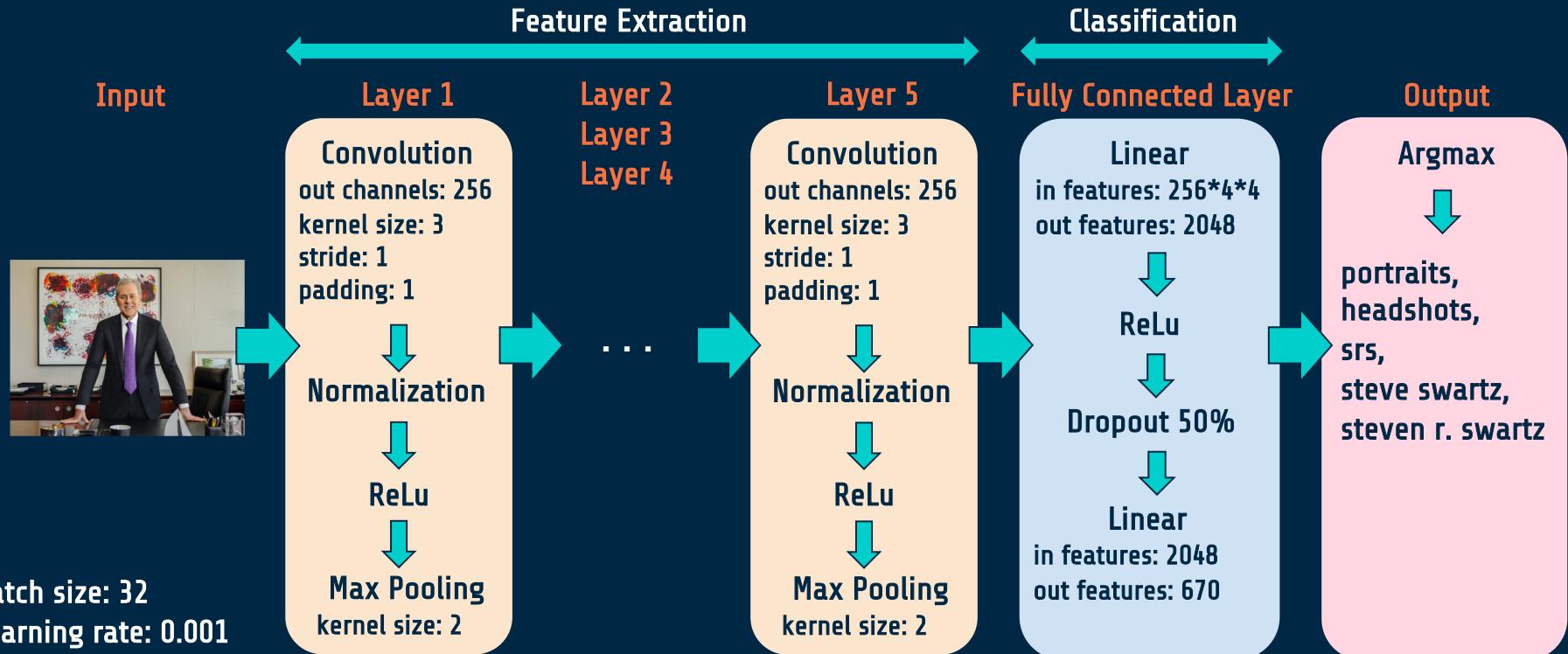
portraits, headshots, steven r.
swartz



people @ work, the long gray line,
events, homecoming, 44th floor

	Dropout	Augmentation	Keywords filtering	Accuracy
Best model	●	●	●	68.86%
Model A		●	●	65.83%
Model B	●		●	64.69%
Model C	●	●		51.31%

Structure of Our Best CNN Model

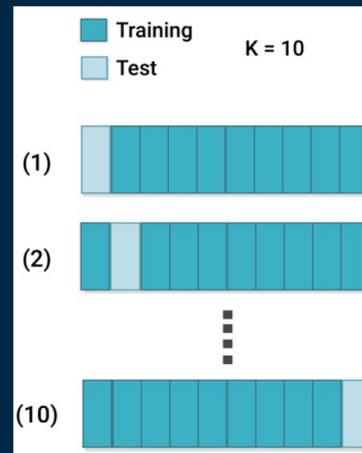


Next Step: Improve Model Accuracy

- Add more convolutional layers

Number of Layers	Testing Accuracy
3	49.38 %
4	63.28 %
5	68.86 %

- Use k-Fold cross-validation



- Transfer learning

Summary

1. Successfully created an image classifier help editors better search and describe media at scale.
1. Change hyperparameters can enhance the model performance.
 - a. **Dropout** increased accuracy by **3.03 %**
 - b. **Augmentation** increased the accuracy by **4.17 %**
 - c. **Keyword filtering** increased the accuracy by **17.55 %**
 - d. **5 layers model** increased the accuracy by **19.48 %** (compared to 3 layers model)
1. Trial and error method and some practices to find the best model.

Thank You



<https://github.com/lch99310>
<https://www.linkedin.com/in/lch99310>
chunghao.lee@marylandsmit.umd.edu



<https://github.com/whl0217>
<https://www.linkedin.com/in/wang-han-li-5b789913b>
wang-han.li@marylandsmit.umd.edu



<https://github.com/FabienneYang>
<https://www.linkedin.com/in/fabienne-yang/>
fabienne.yang@marylandsmit.umd.edu

