

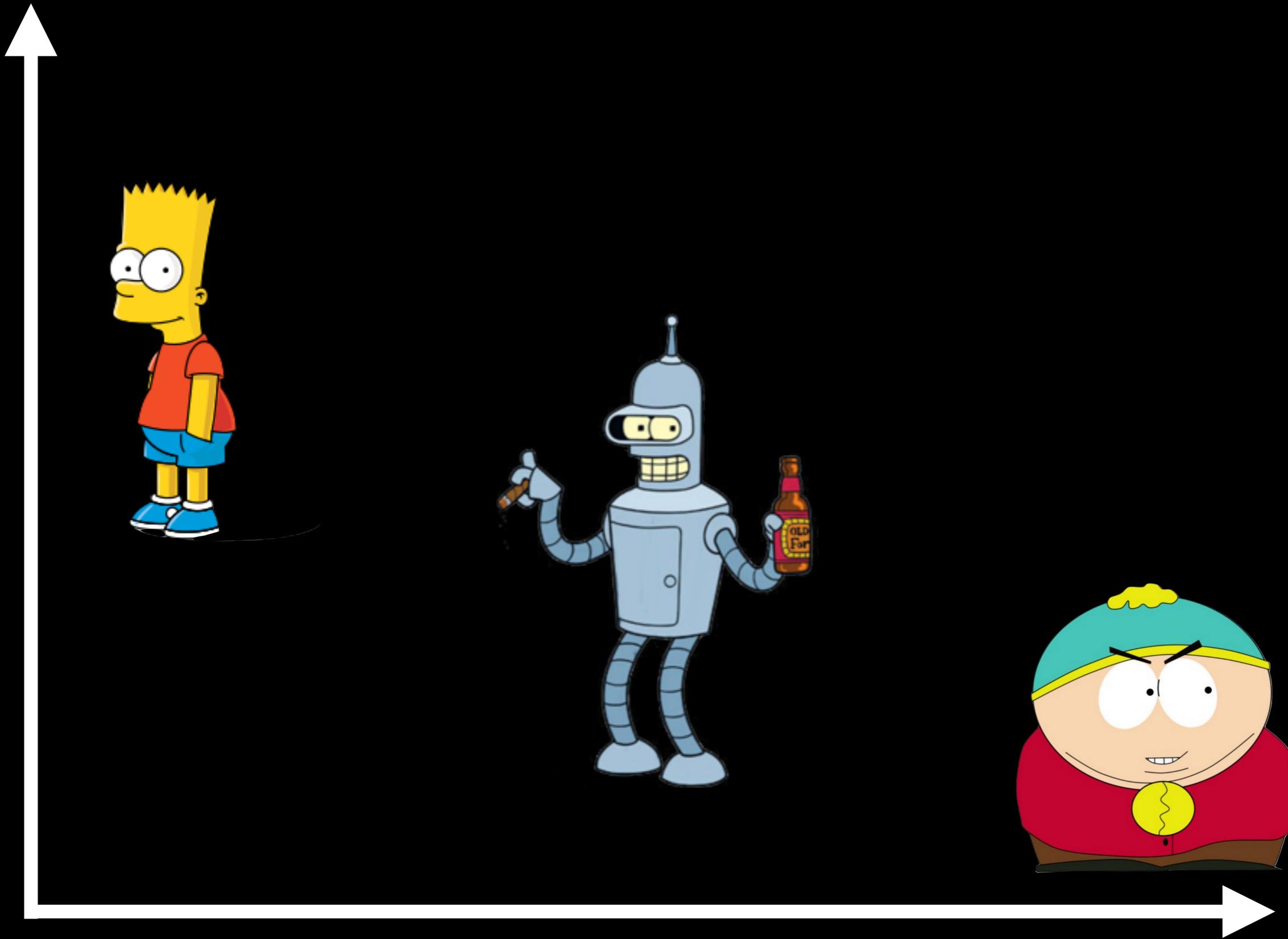
The hero with (several) thousand faces







Agreeableness



Sociopathy

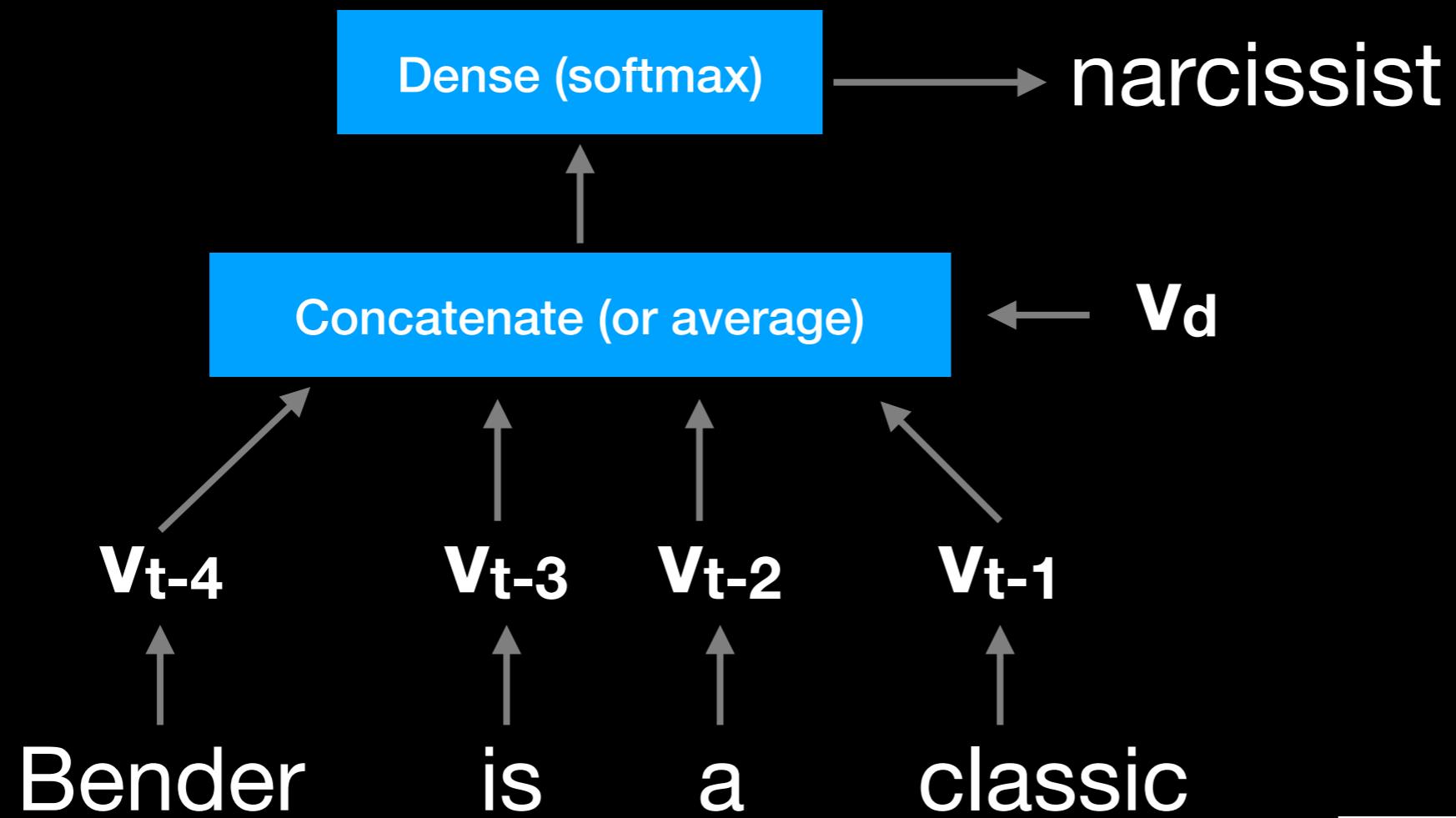


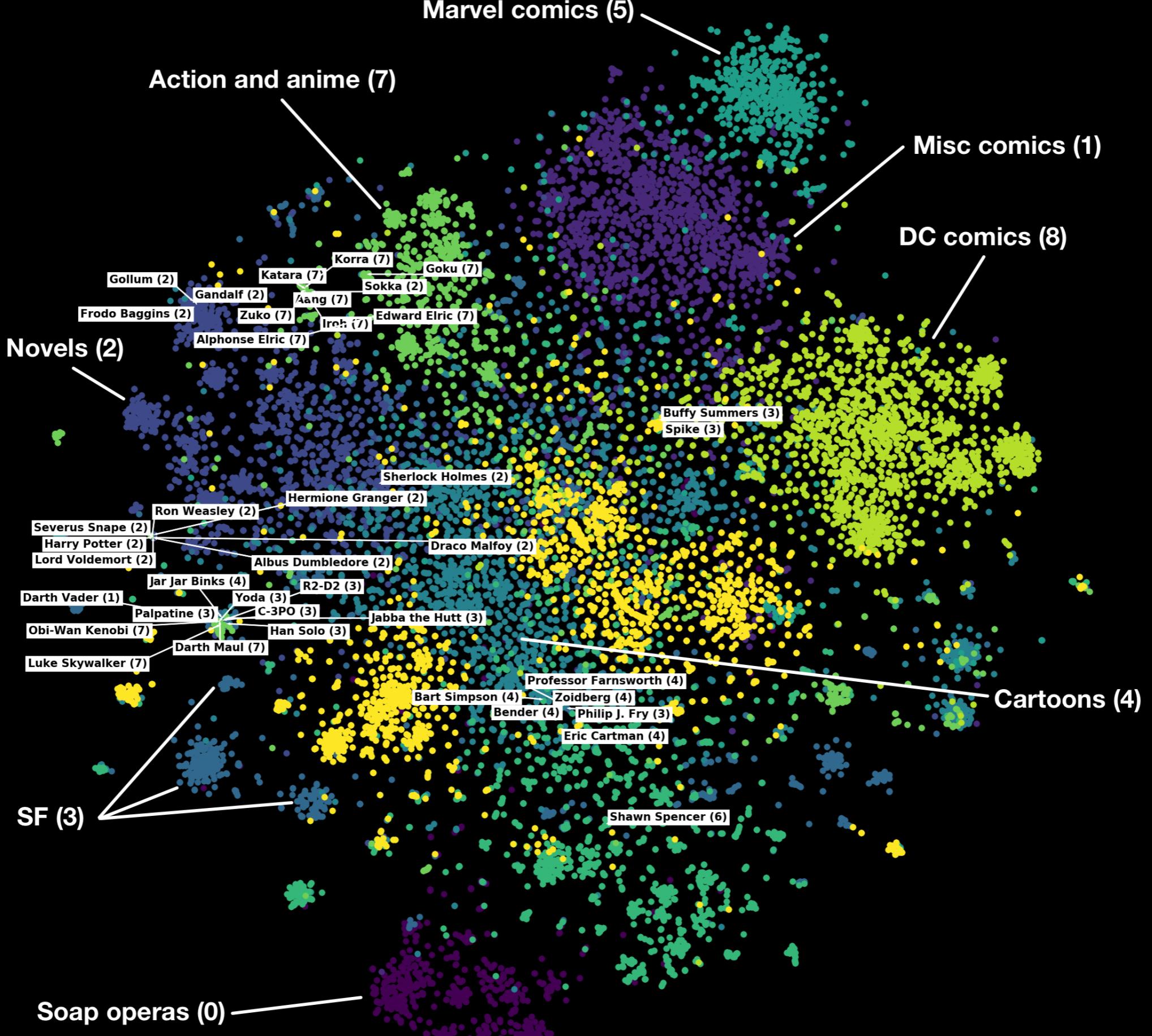
# Data set

- ~14,000 character articles from Wikipedia
- 23,616,309 raw words
- 1 to 5 million relevant ones



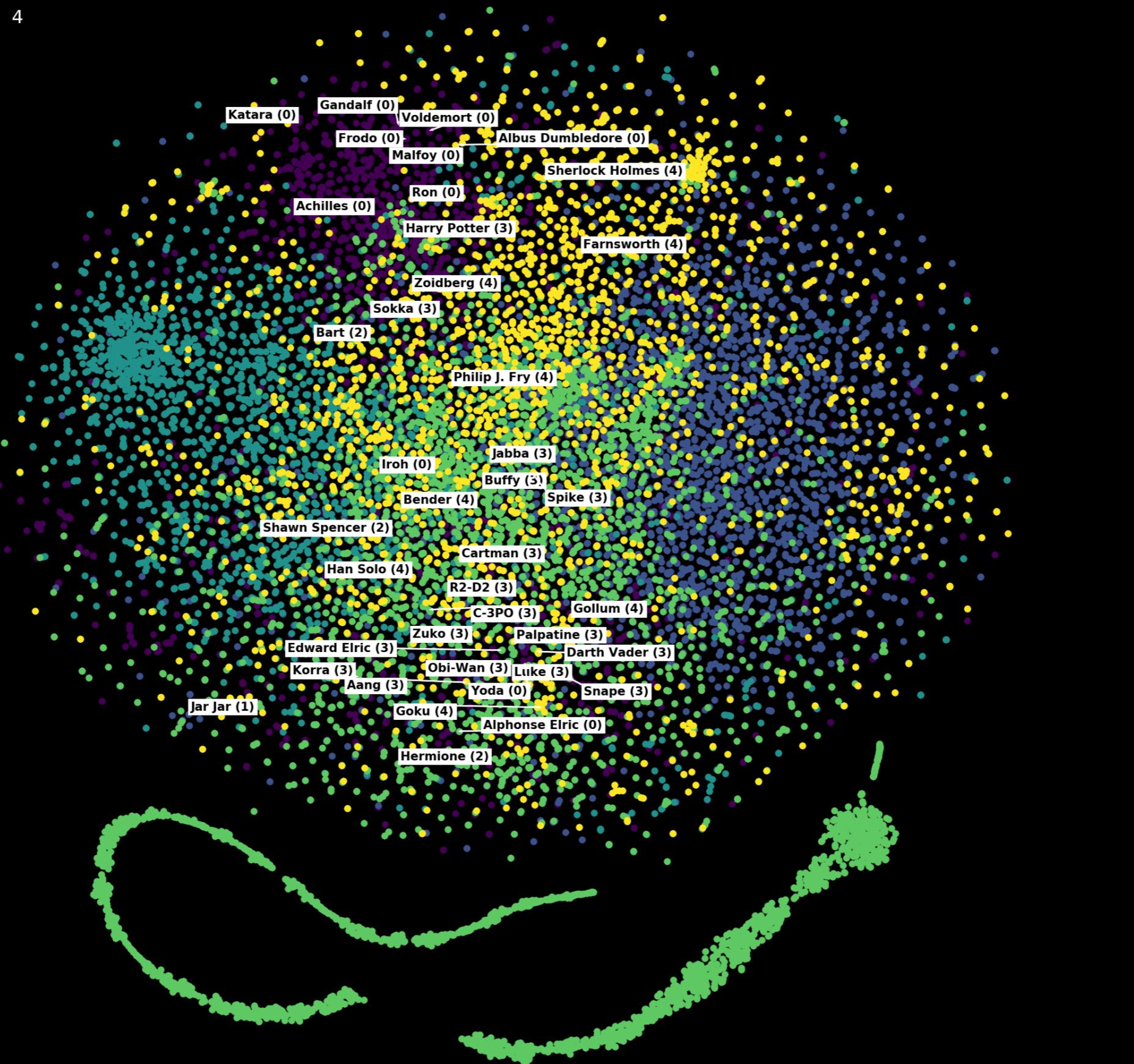
# Doc2Vec





**"<SUBJECT> is the mad scientist proprietor of  
the <ENTITY> delivery service , for whom the  
main characters work ."**

- 0
- 1
- 2
- 3
- 4



# Raw text similarity

```
1 model.docvecs.most_similar('Yoda')  
  
[ ('Obi-Wan Kenobi', 0.8464793562889099),  
 ('Count Dooku', 0.8301168084144592),  
 ('R2-D2', 0.8243268132209778),  
 ('Qui-Gon Jinn', 0.82326340675354),  
 ('C-3PO', 0.8140591382980347),  
 ('Luke Skywalker', 0.8125448822975159),  
 ('Palpatine', 0.7923175096511841),  
 ('Admiral Ackbar', 0.7624970078468323),  
 ('Darth Maul', 0.7532809376716614),  
 ('Jar Jar Binks', 0.7466059923171997)]
```

# Filtered text similarity

```
1 d2v_model.docvecs.most_similar('Yoda')  
  
[ ('Qui-Gon Jinn', 0.6633567810058594),  
 ('Luke Skywalker', 0.6167474985122681),  
 ('Optimus Prime', 0.6107190251350403),  
 ('Scourge (Transformers)', 0.5825557708740234),  
 ('Kanan Jarrus', 0.5685822367668152),  
 ('Zuko', 0.5668112635612488),  
 ('Gimli (Middle-earth)', 0.565727710723877),  
 ('Mokujin', 0.5602509379386902),  
 ('Albus Dumbledore', 0.5537698268890381),  
 ('Beorn', 0.5510784387588501)]
```

# Future directions

- I still have mistags and other noise.
- Implement Doc2Vec in Keras/Tensorflow.
- Improve sentence extractor and apply to plot summaries.
- Understand at the axes of the space (axis of evil, et cetera).