



## Reddit post topic modelling

Prepared by:

Abdultawwar Safarji

Najla Bin-Melha

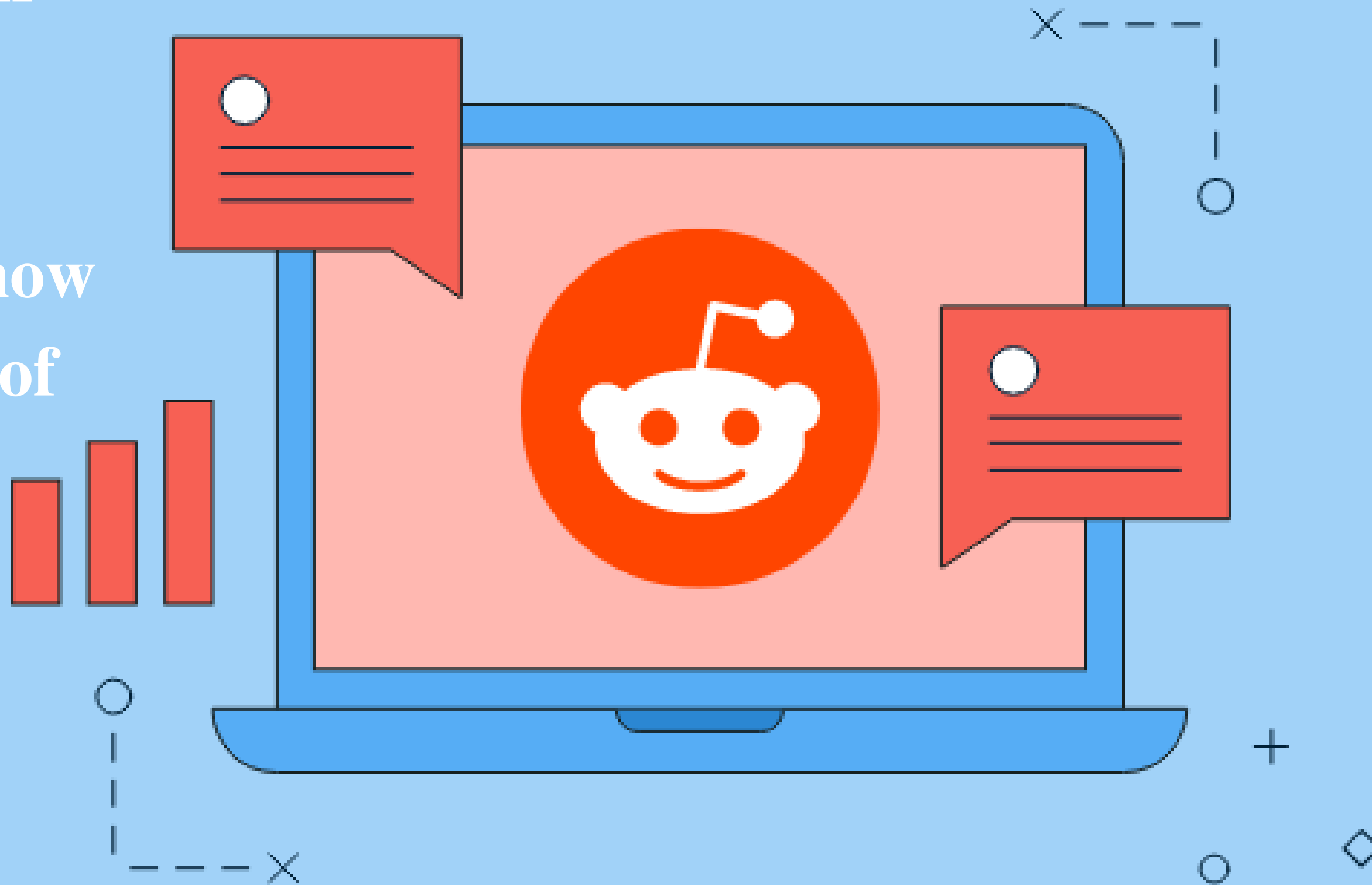
Khalid Alrashed

---

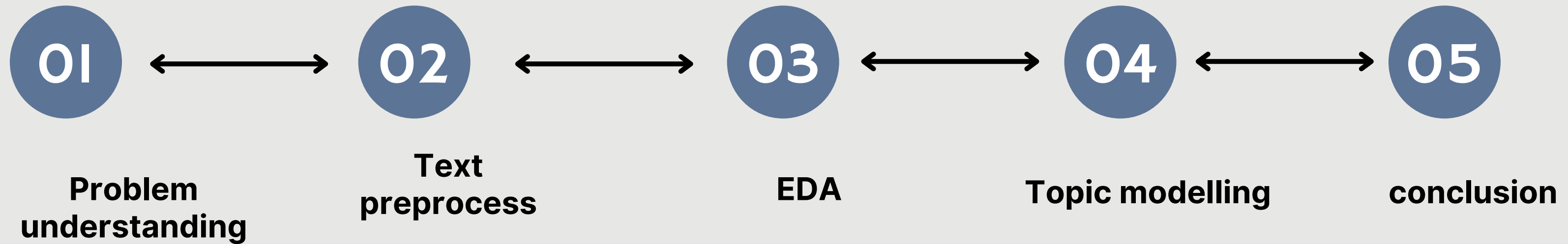


reddit

Have you ever  
asked yourself  
that through  
your use of  
reddit you know  
the behavior of  
users?



# Methodology



# Problem Understanding

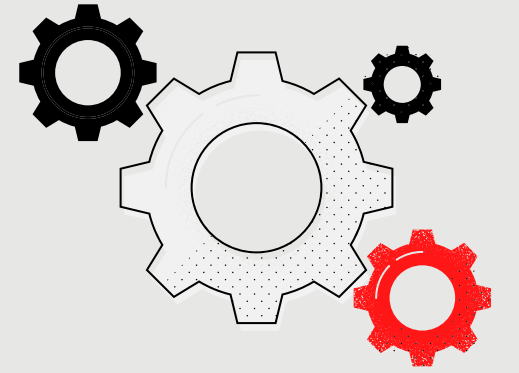
## Why Riddet Toping Modeling ?

- To organize, understand and summarize large collections of textual information.
- Make it easy for the Reddit community to use unsupervised topic modeling and advance in technology as people discuss in text format.

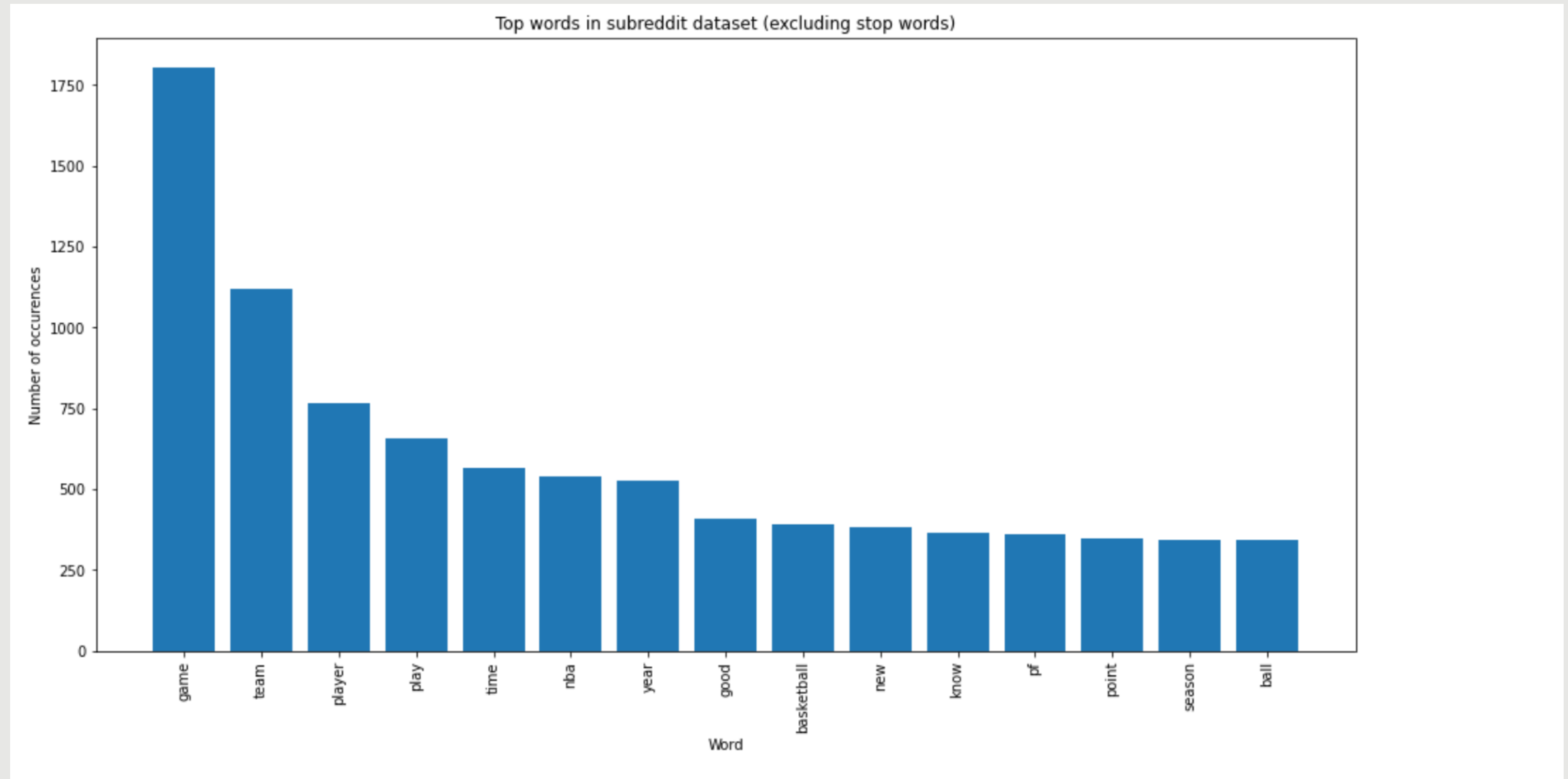


# Text preprocess

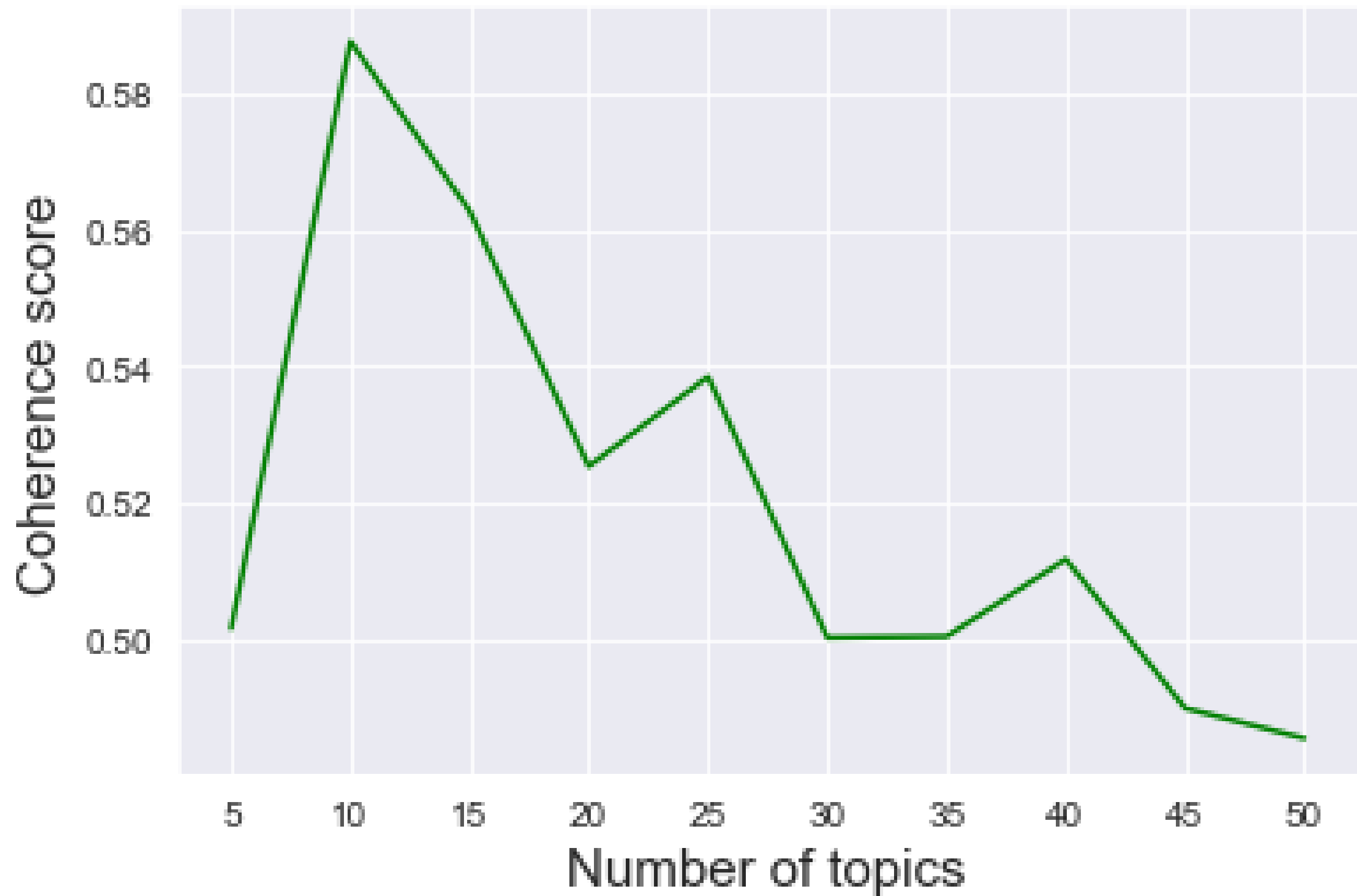
- URL links
- Numbers  
&punctuation(etc..)
- Tokenization
- Stop words
- lemmatization



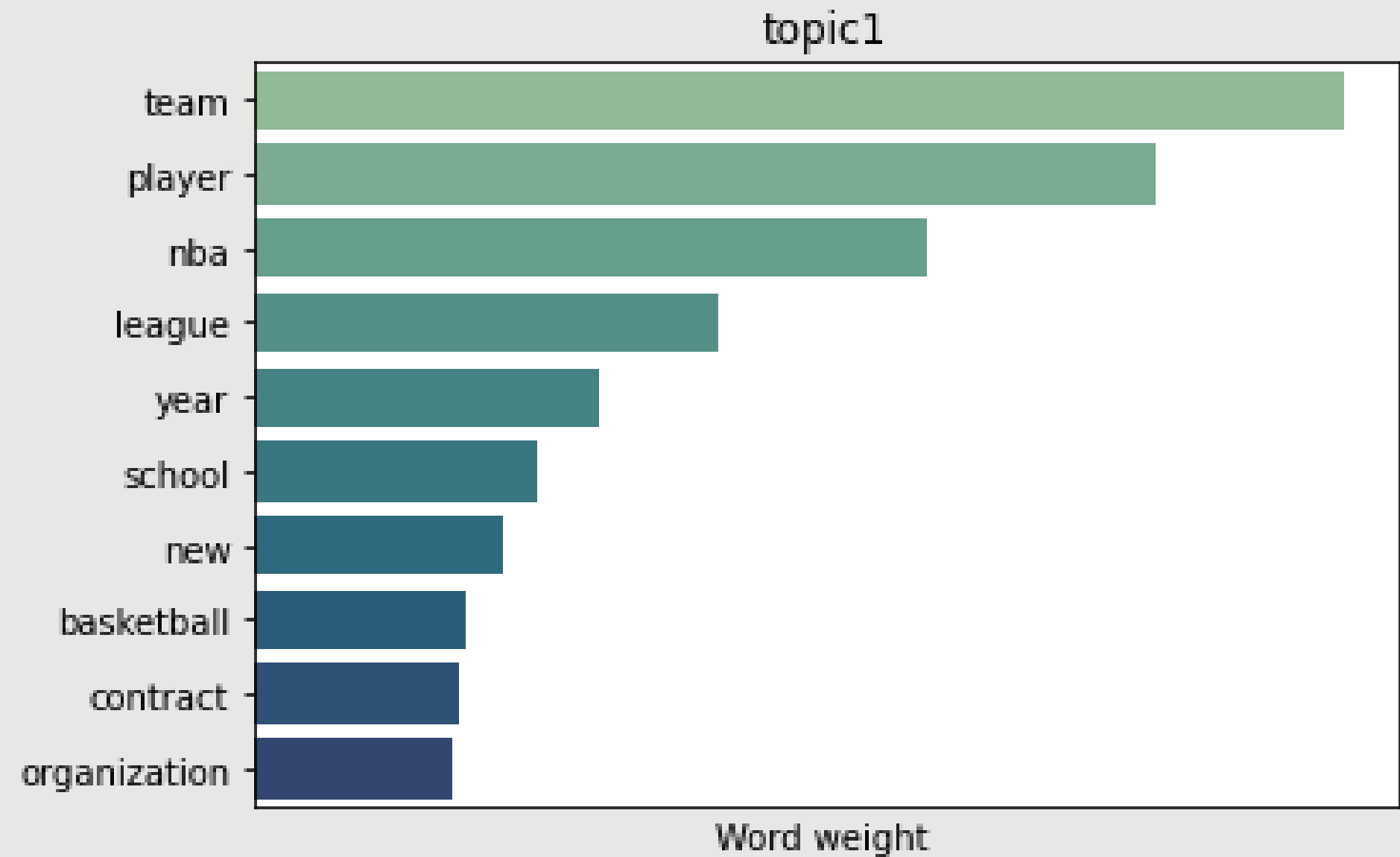
# Exploratory Data Analysis



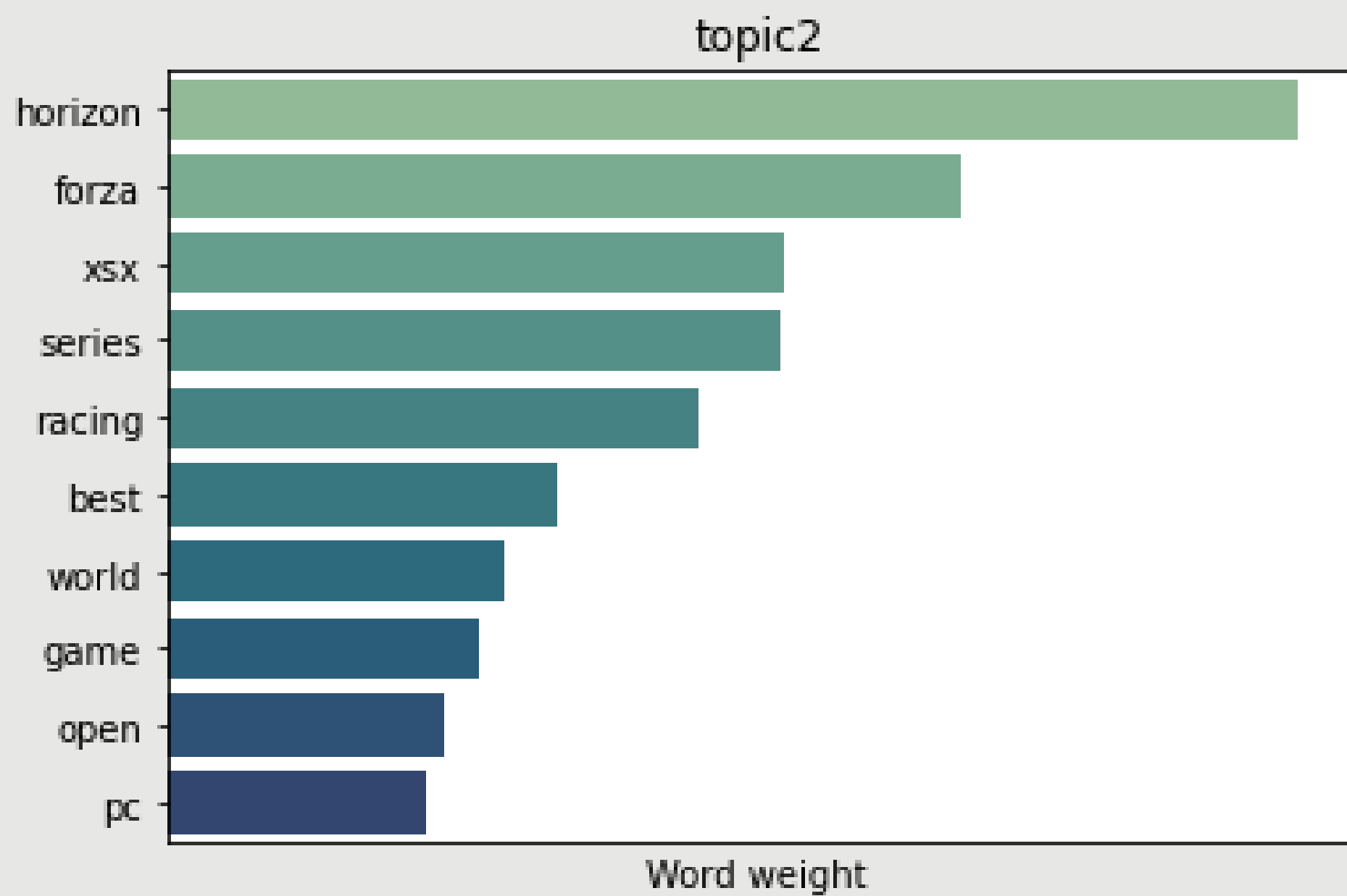
number of topics vs coherence score



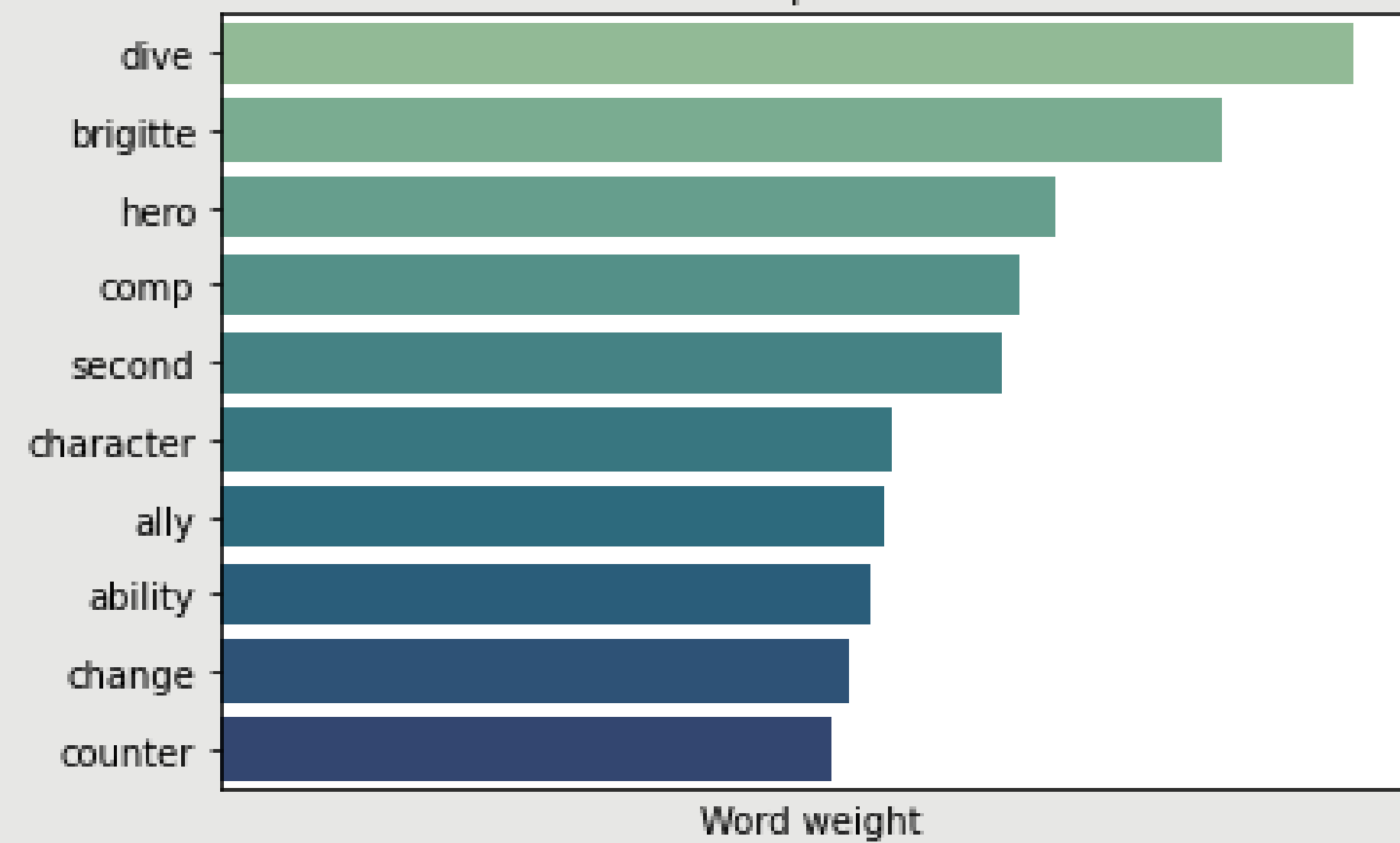
# Results

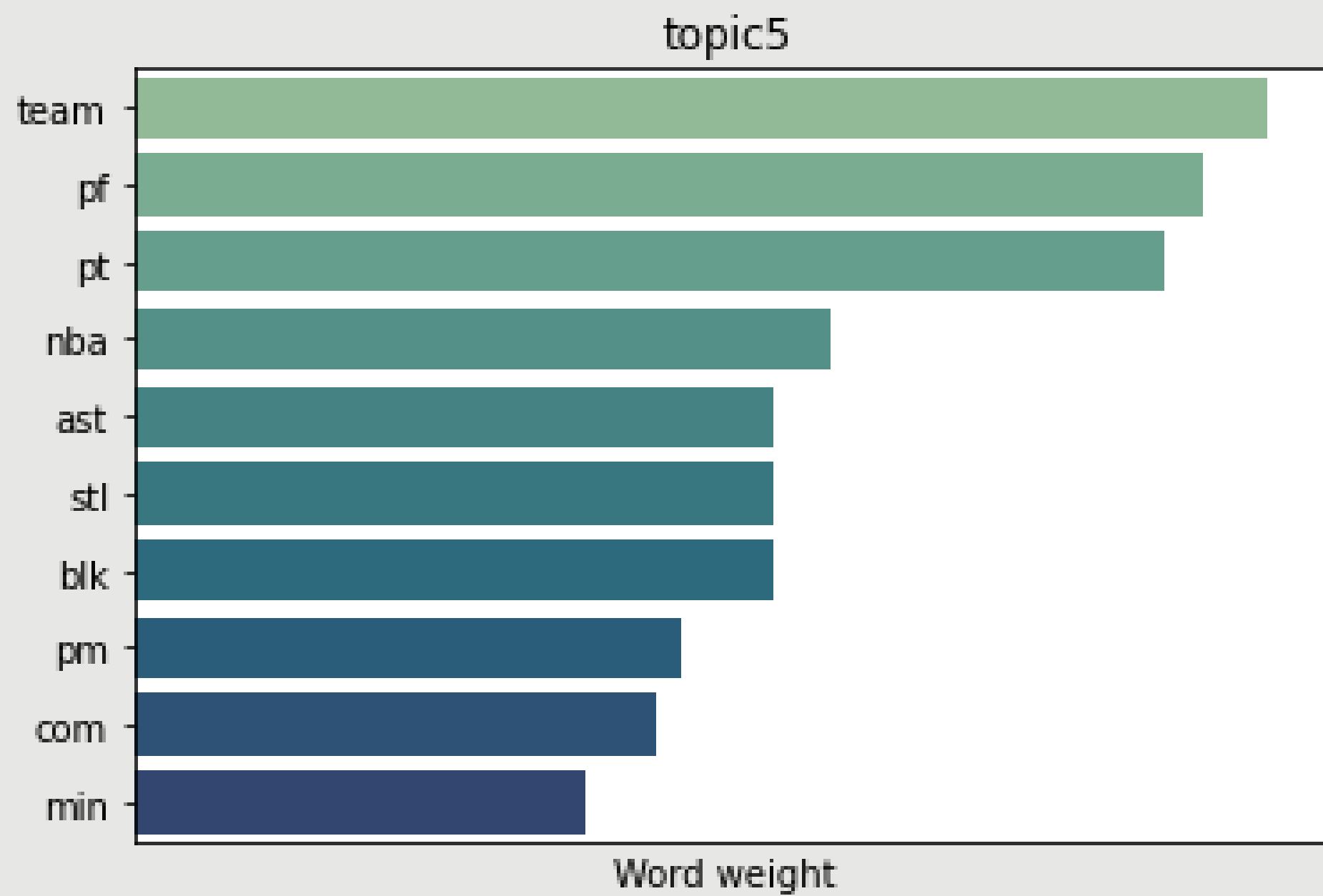




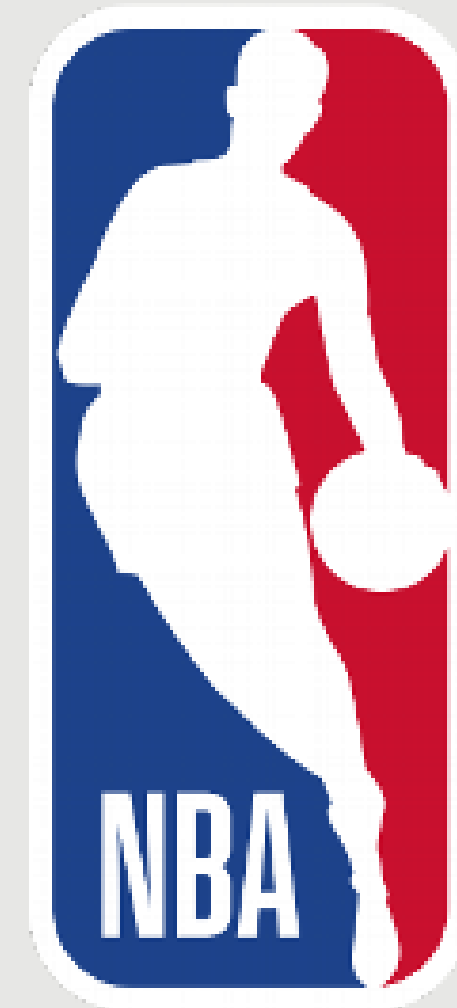


topic7

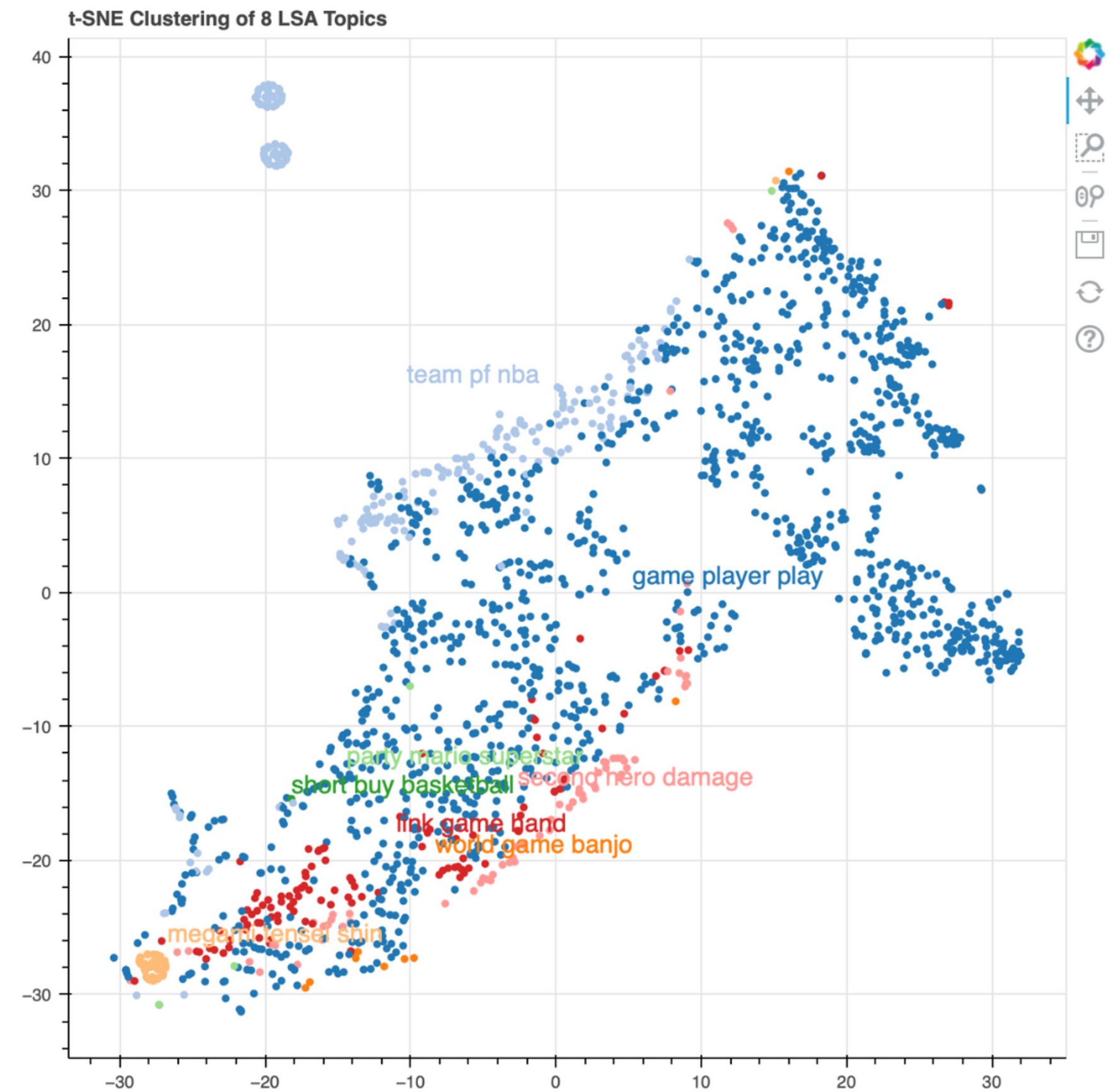
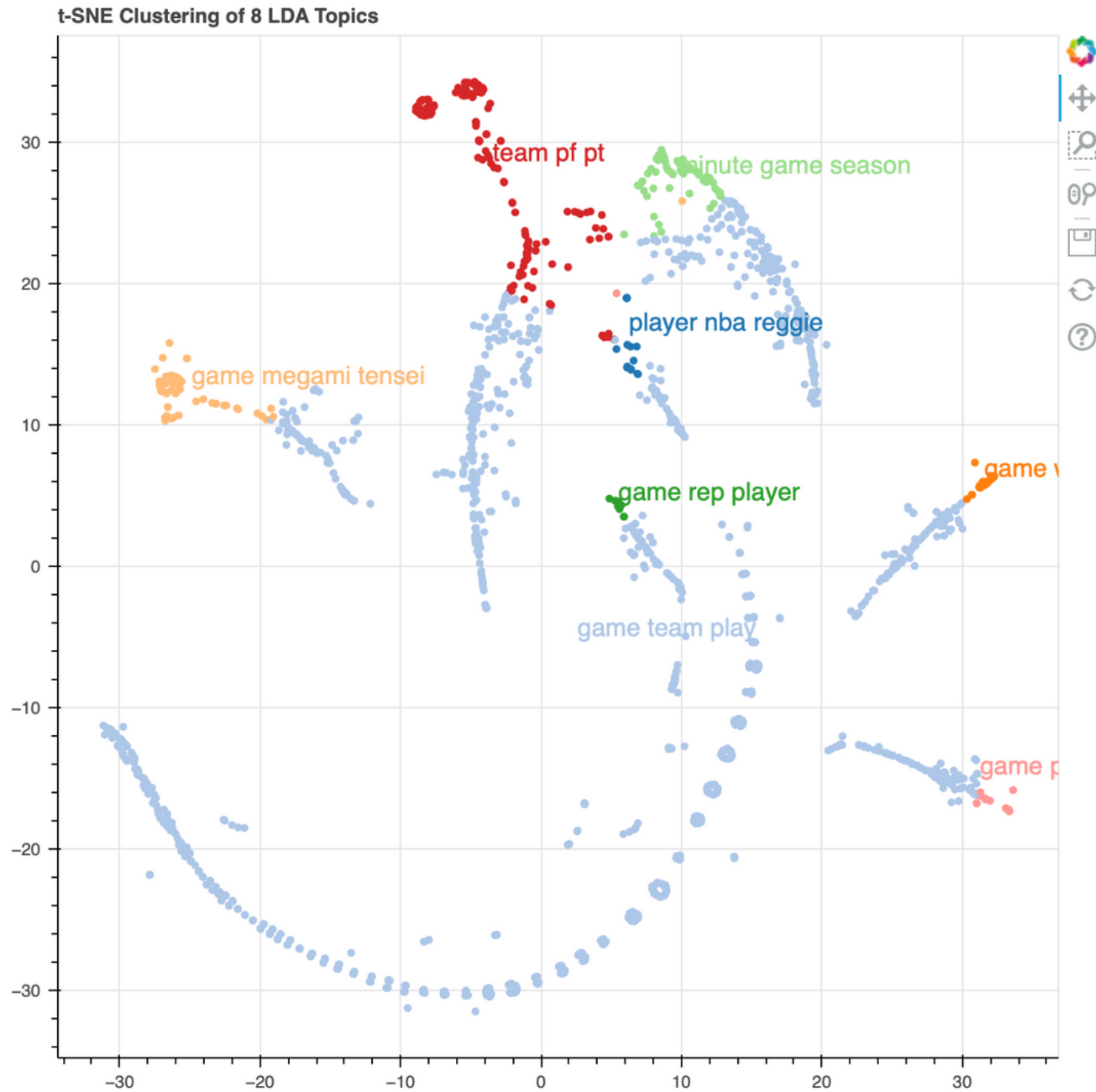




ORB	DRB	REB	AST	STL	BLK	TO	PF
0	2	2	3	0	0	1	4



# Models Testing



# LDA Tunning

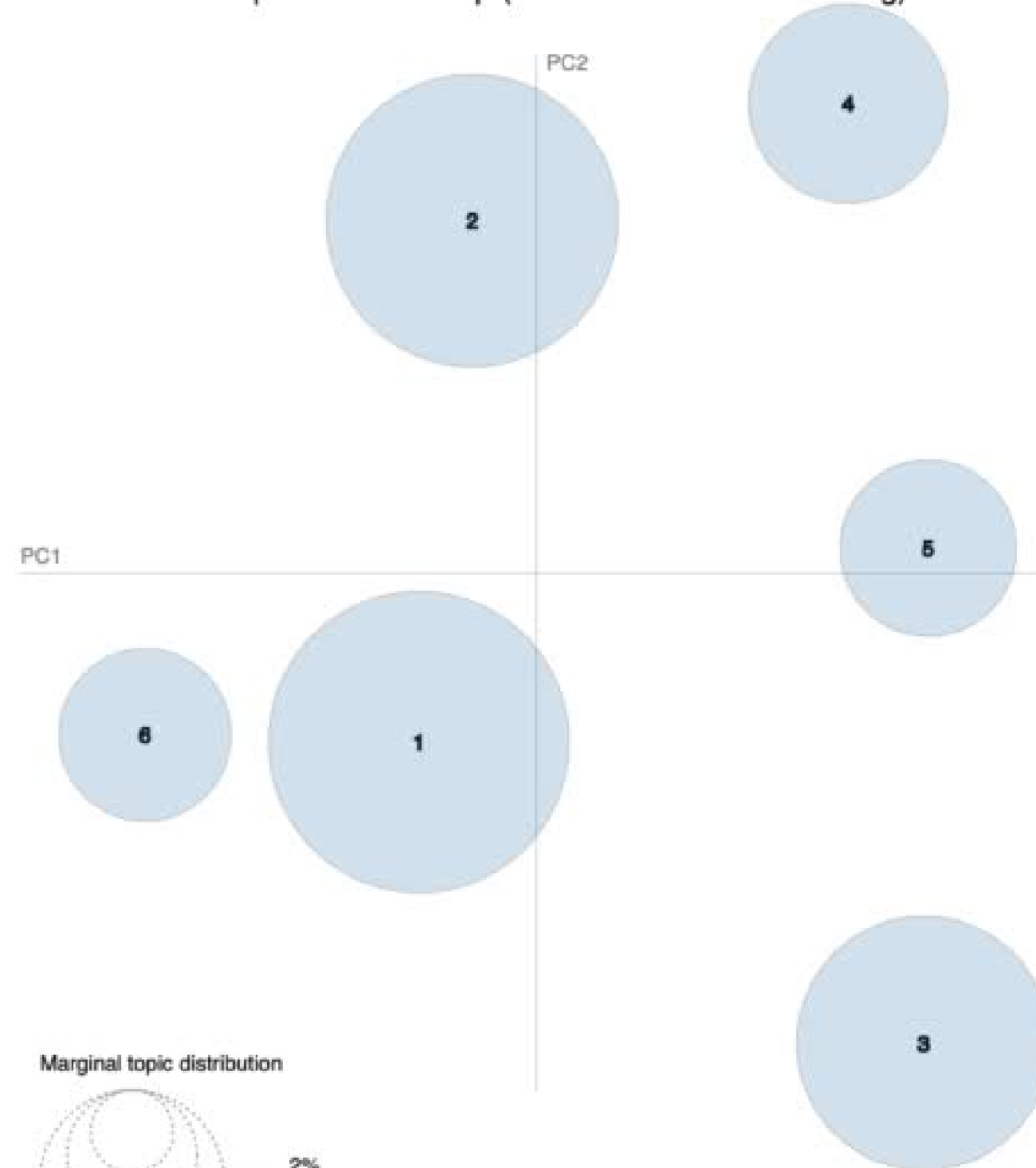
Selected Topic:  Previous Topic Next Topic Clear Topic

Slide to adjust relevance metric:<sup>(2)</sup>

$\lambda = 1$

0.0 0.2 0.4 0.6 0.8 1.0

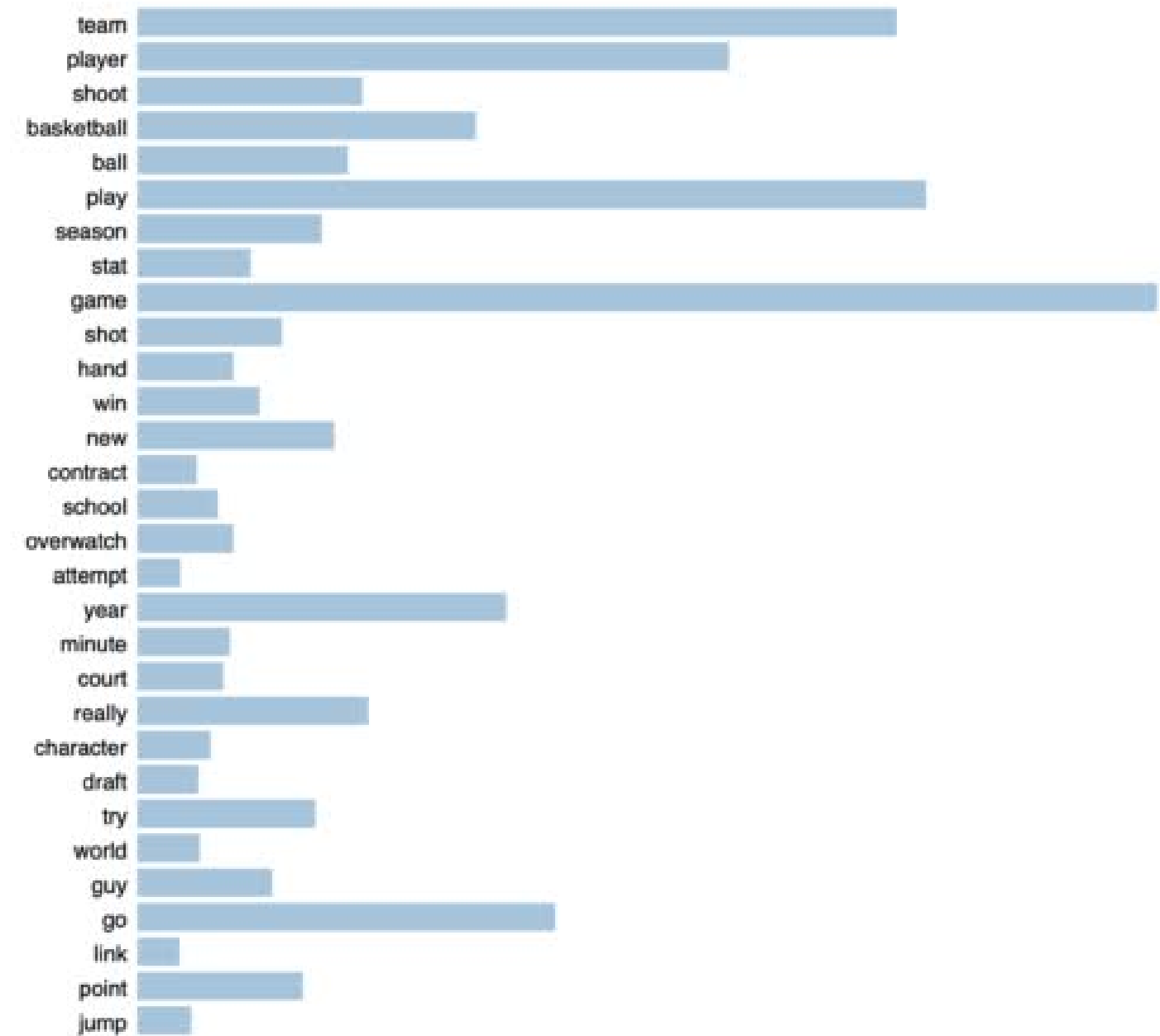
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Salient Terms<sup>1</sup>



Overall term frequency

Estimated term frequency within the selected topic

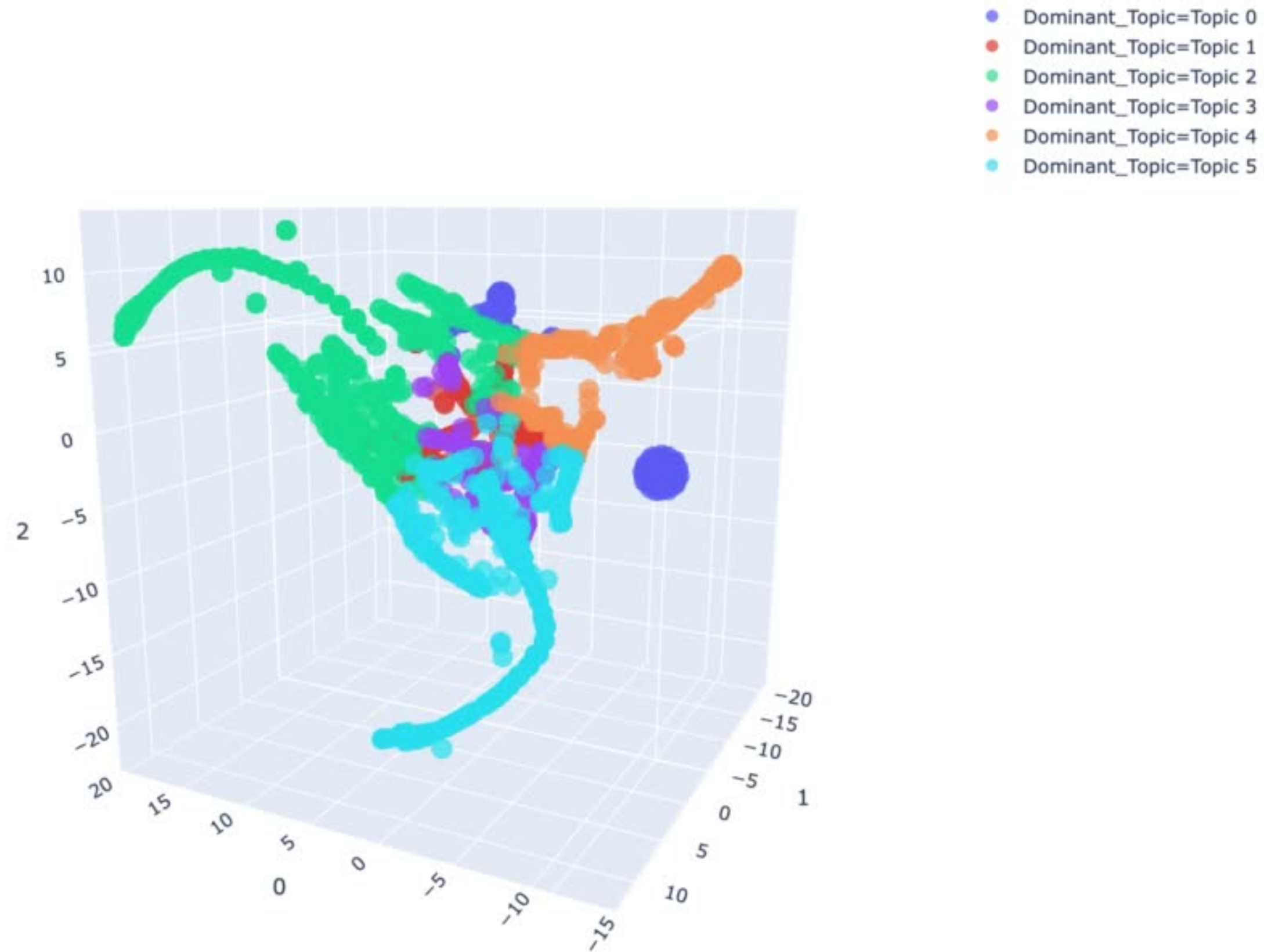
1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w)/p(t))]$  for topics  $t$ ; see Chuang et. al (2012)

2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

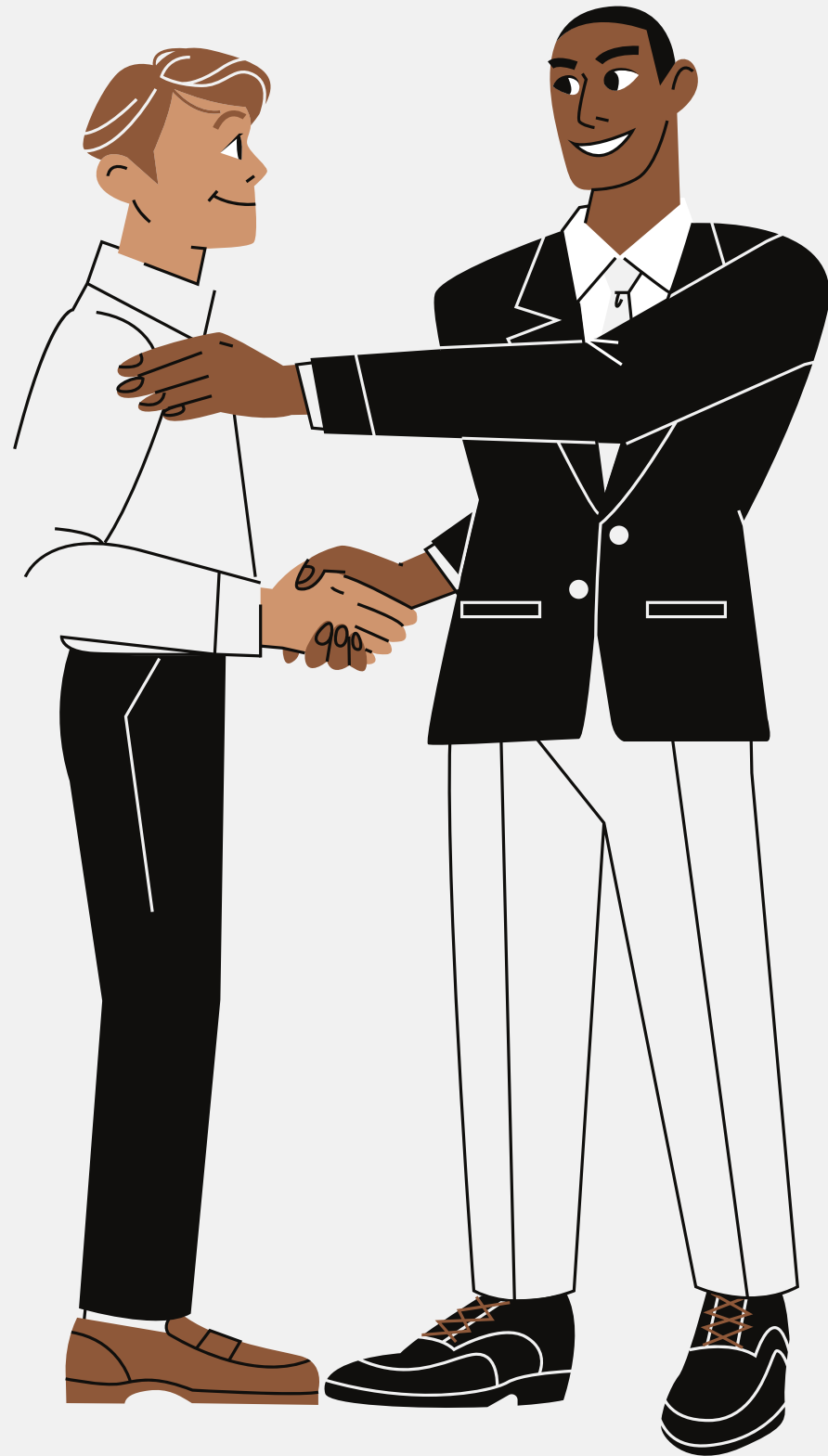


# LDA 3D

3d TSNE Plot for Topic Model

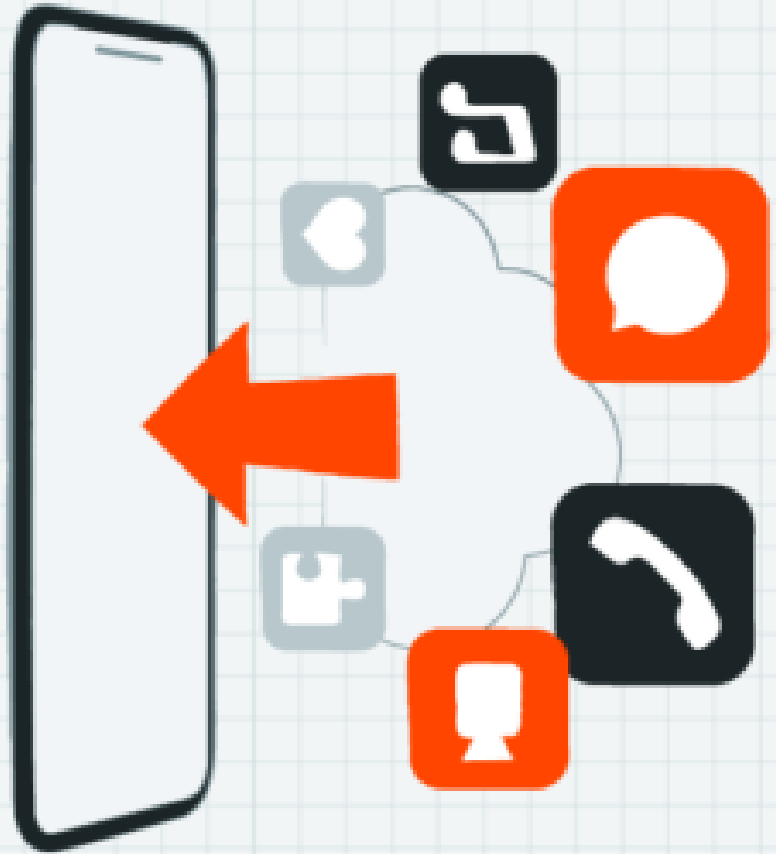


# Conclusion



- **Reddit topics and posts need continuous cleaning and gathering.**
- **Need to explore more in more algorithms in topic modeling for Reddit.**
- **Ensure text quality modeling.**

# Thank You For Listening.



## Do you have any questions?

