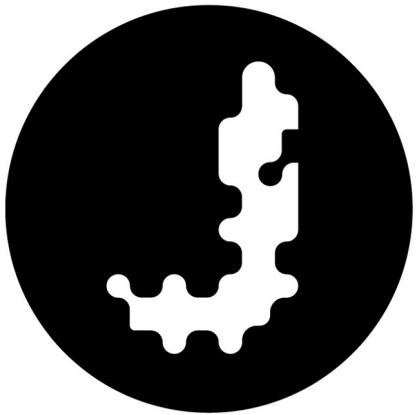


**TOXIC** ■  
■ **CLASSIFIER**

**TEXT** ■  
■ **MINING**

*Data*



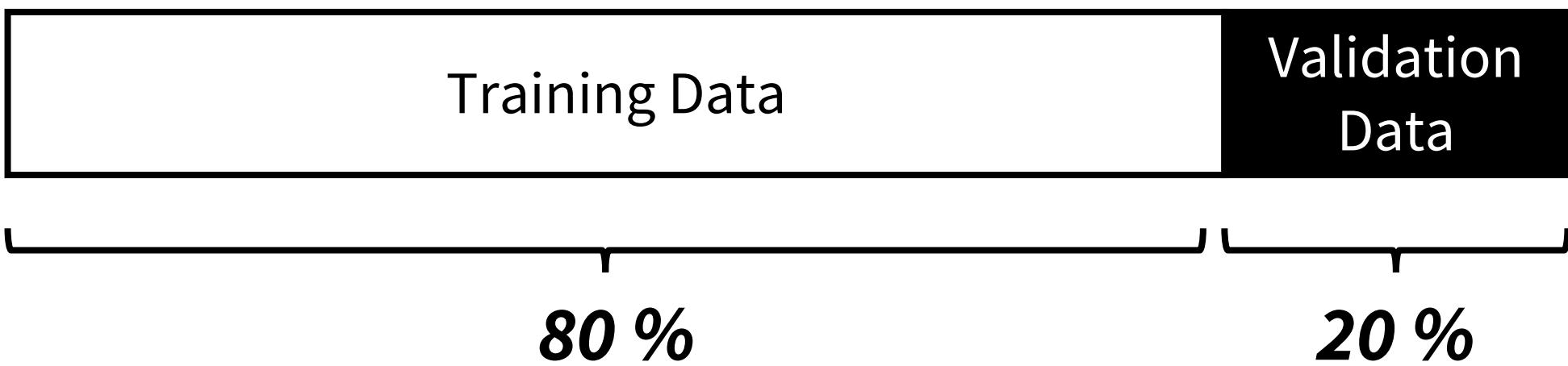
# Jigsaw

# *Data*

***159 571 Observations***

Training Data

# *Data*



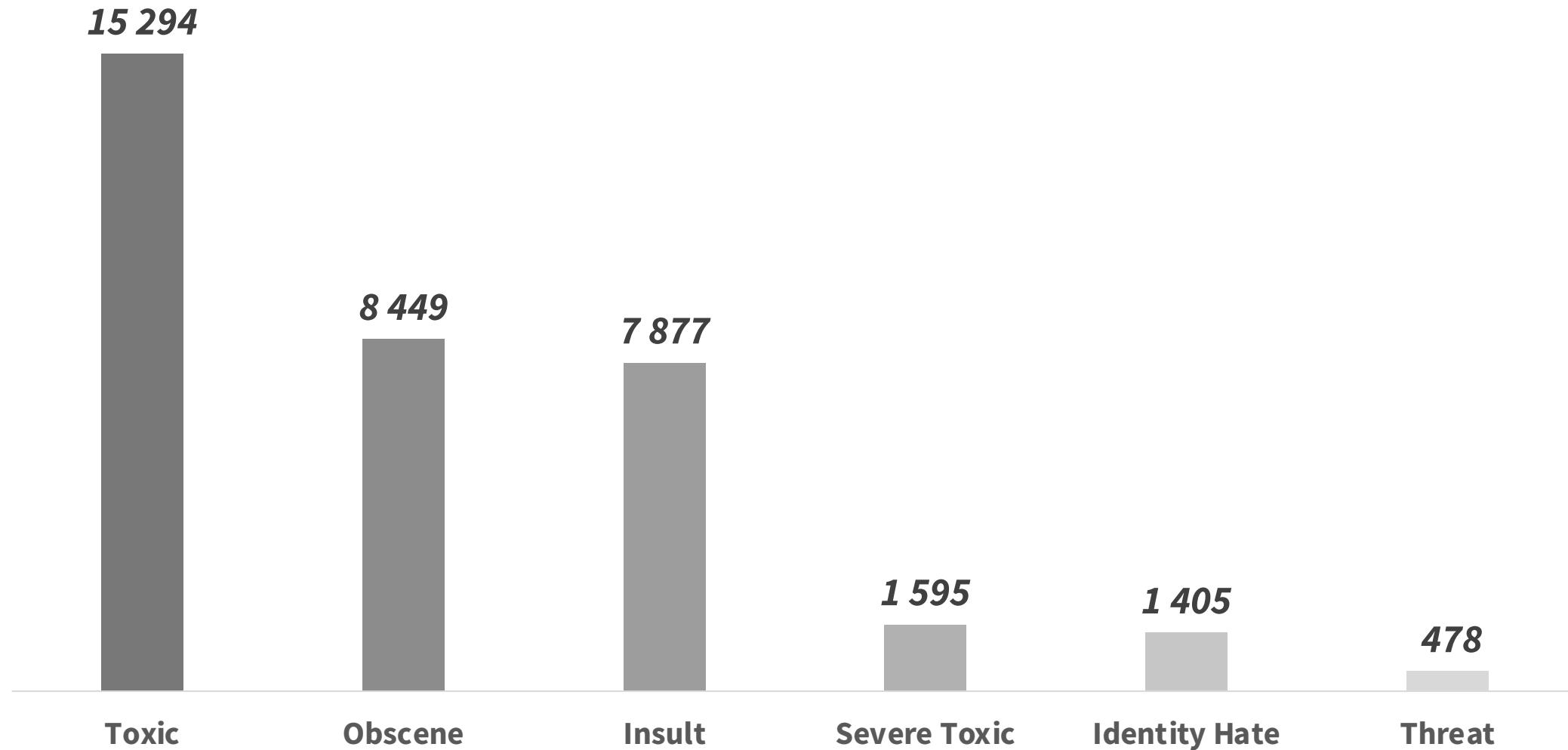
# *Preprocess*

<b>id</b>	<b>Comment Text</b>	<b>Toxic</b>	<b>Severe Toxic</b>	<b>Obscene</b>	<b>Threat</b>	<b>Insult</b>	<b>Identity Hate</b>
1	...	0	0	0	0	0	0
2	...	1	0	0	1	0	0
3	...	0	0	0	0	0	0

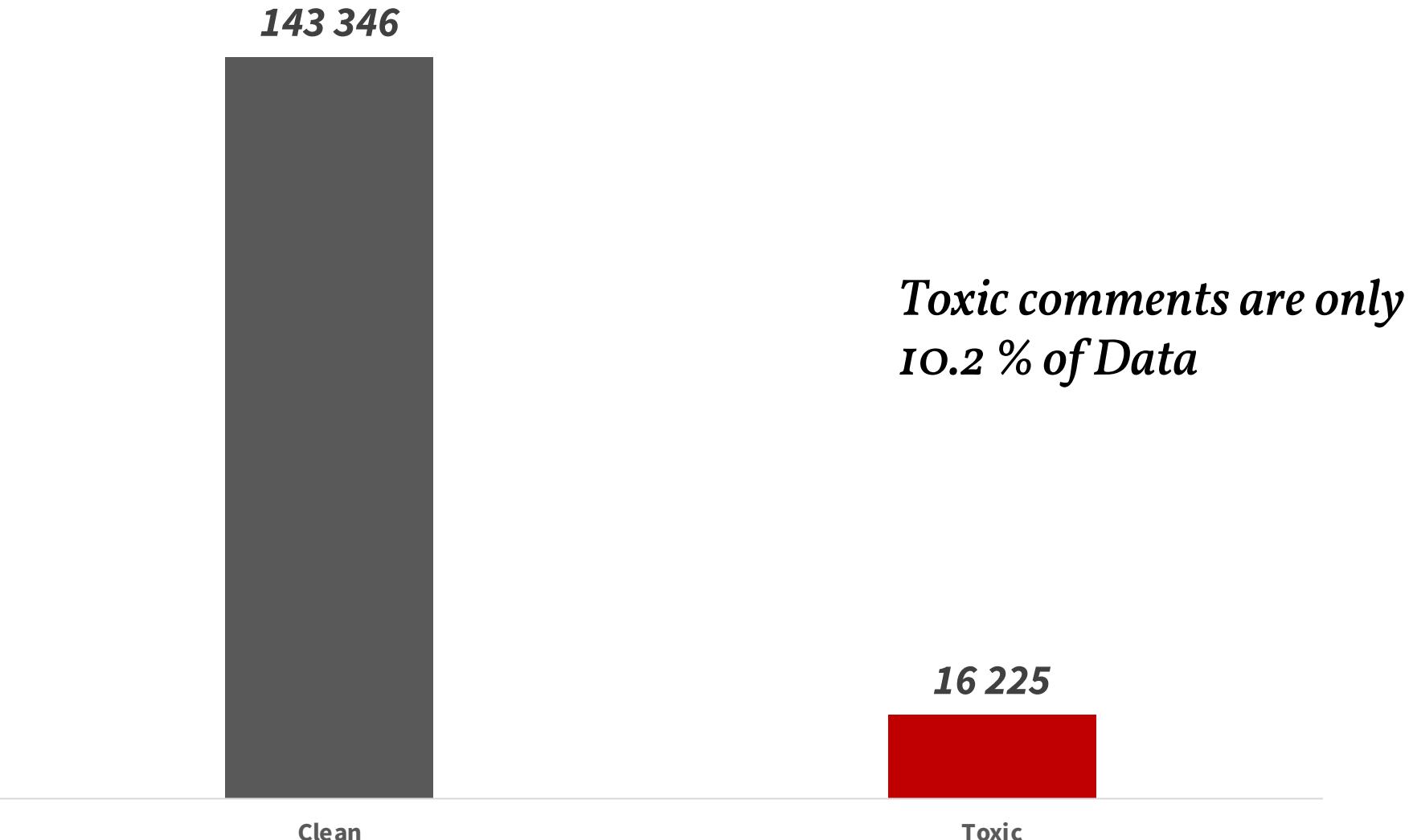
# *Preprocess*

<b>id</b>	<b>Comment Text</b>	<b>Toxic</b>	<b>Severe Toxic</b>	<b>Obscene</b>	<b>Threat</b>	<b>Insult</b>	<b>Identity Hate</b>	<b>Target</b>
1	...	0	0	0	0	0	0	0
2	...	<b>1</b>	<b>0</b>	<b>0</b>	<b>1</b>	<b>0</b>	<b>0</b>	<b>1</b>
3	...	0	0	0	0	0	0	0

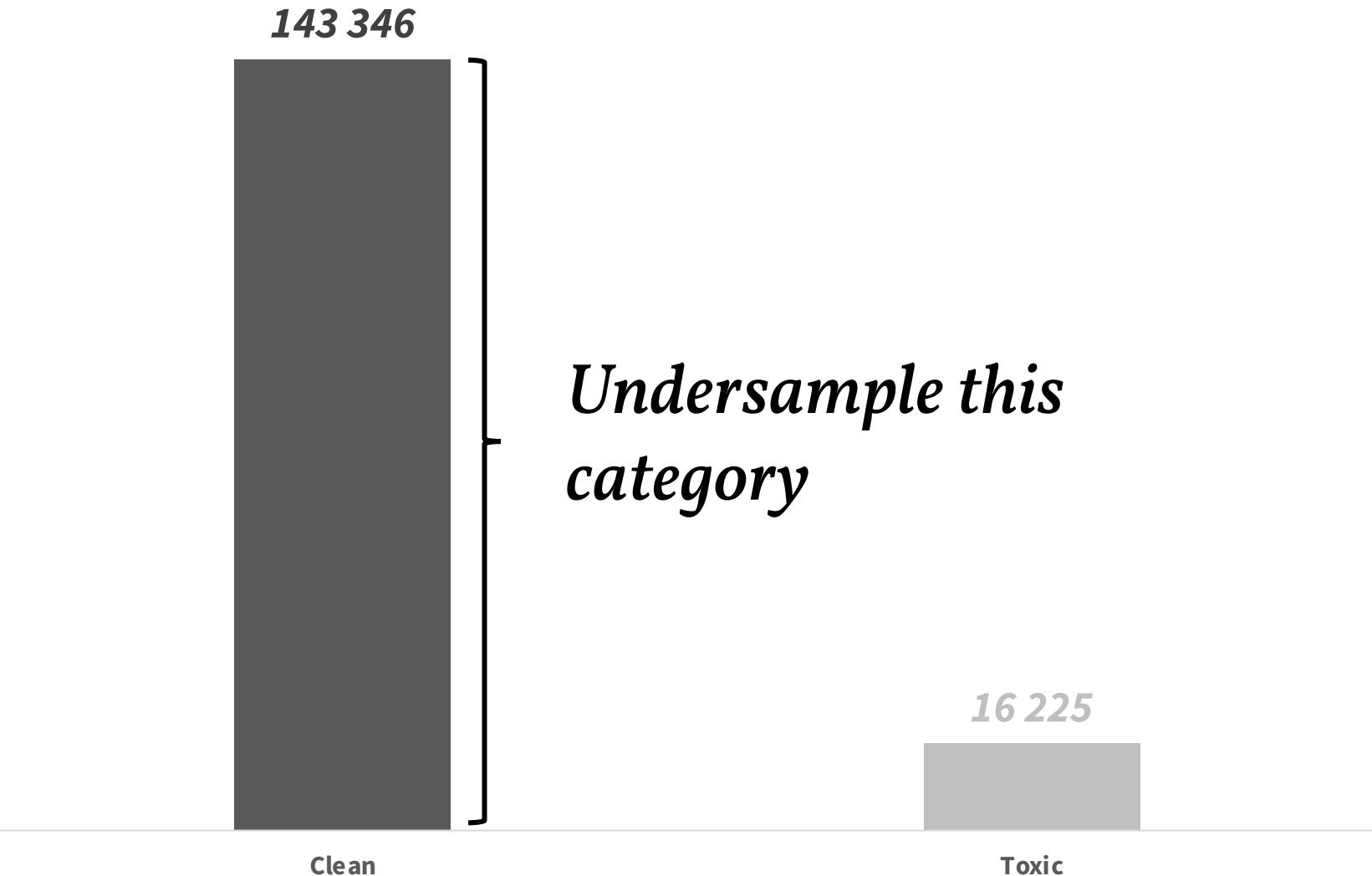
# *Preprocess*



# *Preprocess*

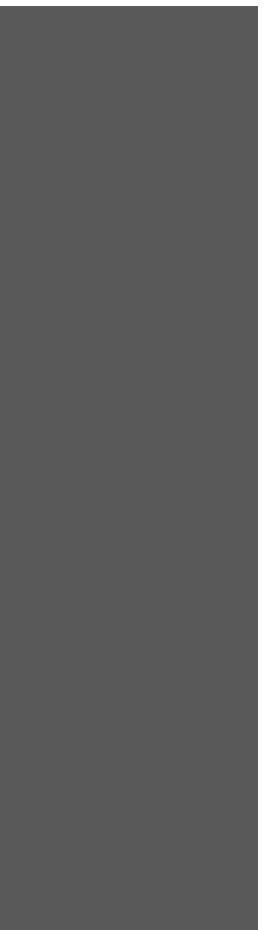


# Preprocess



# *Preprocess*

**16 225**



**Clean**

**16 225**



**Toxic**

**32 450**  
*Observations*

# *Preprocess*

## ***1. Balance Dataset***

# *Preprocess*

## *1. Balance Dataset*

## **2. New Features**

- Character count
- Word count
- Word density
- Punctuation count
- Title word count
- Upper case word count

# *Preprocess*

## *1. Balance Dataset*

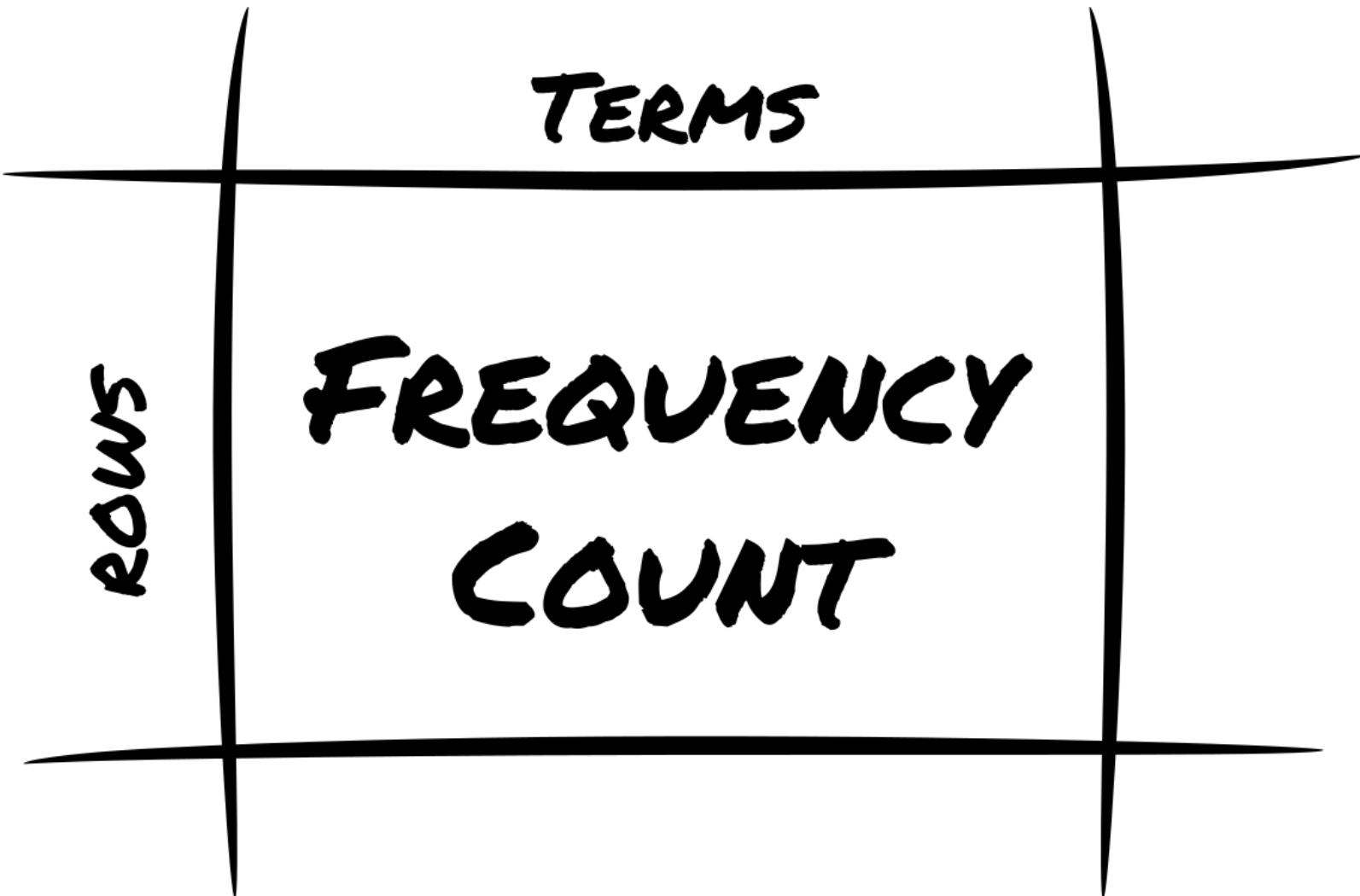
## *2. New Features*

- Character count
- Word count
- Word density
- Punctuation count
- Title word count
- Upper case word count

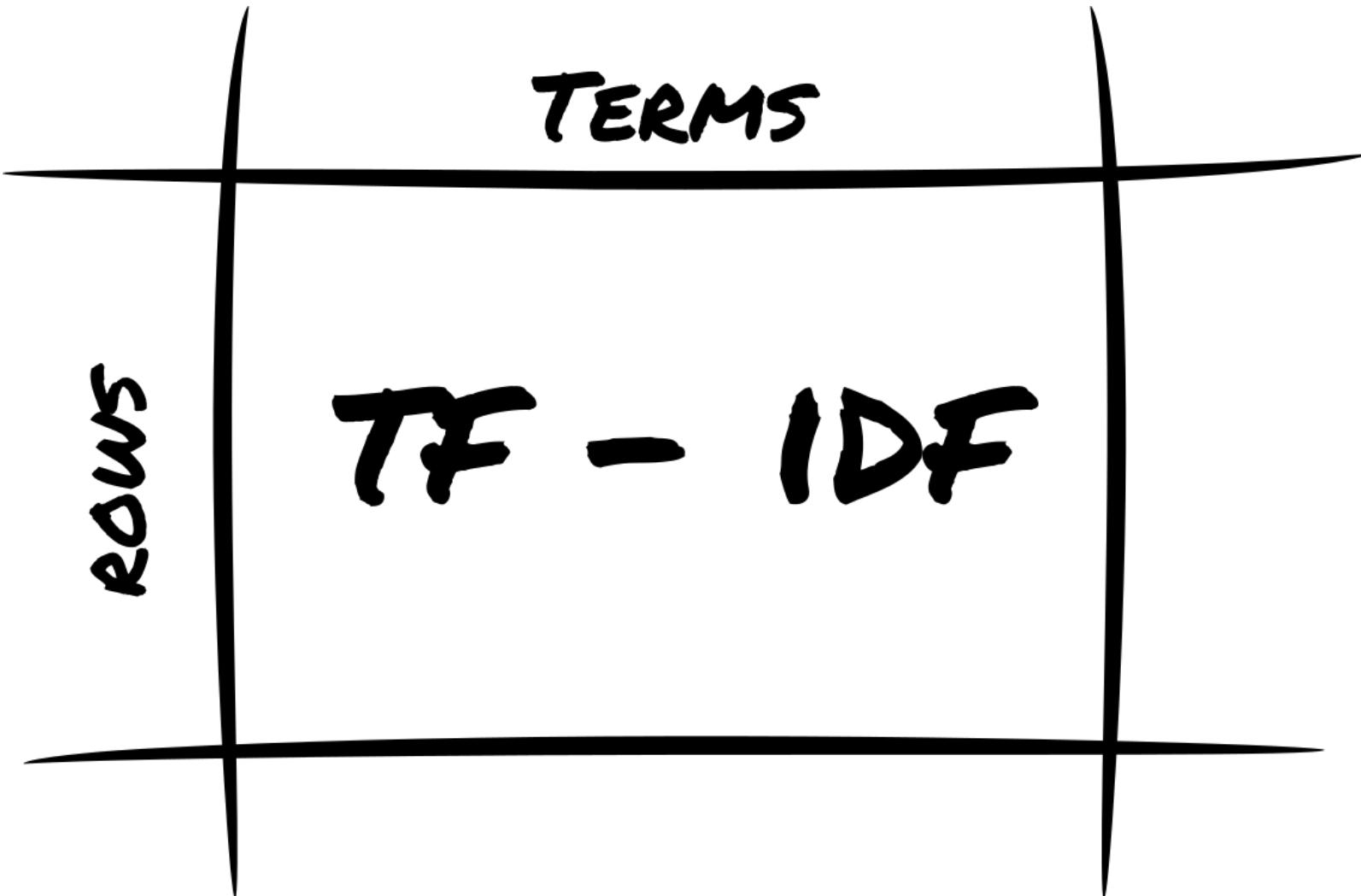
## *3. Clean Text*

- Lower Case
- Remove punctuation
- Tokenize
- Remove stopwords
- Stem words

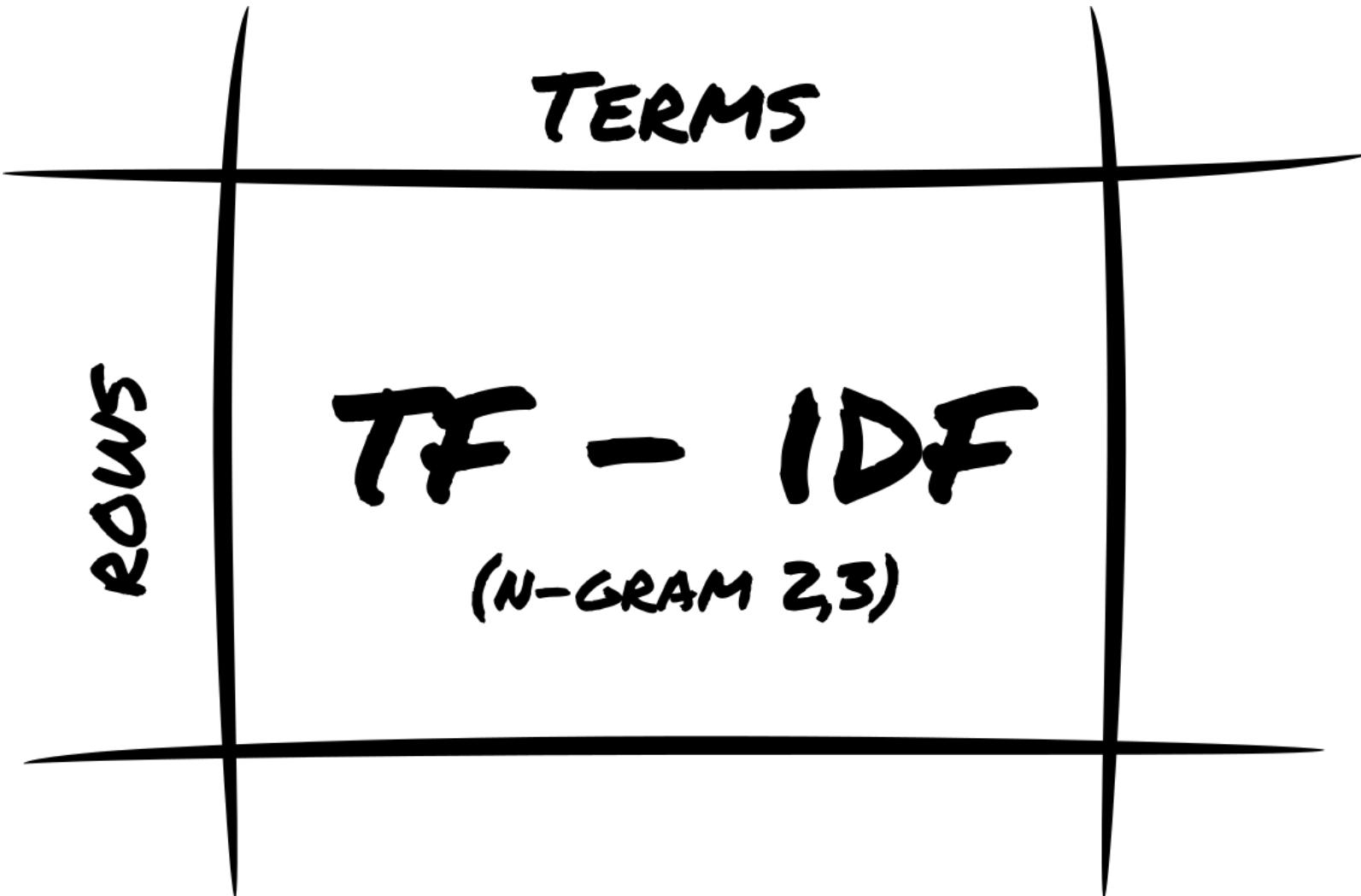
# *Feature Engineering - Vectorizers*



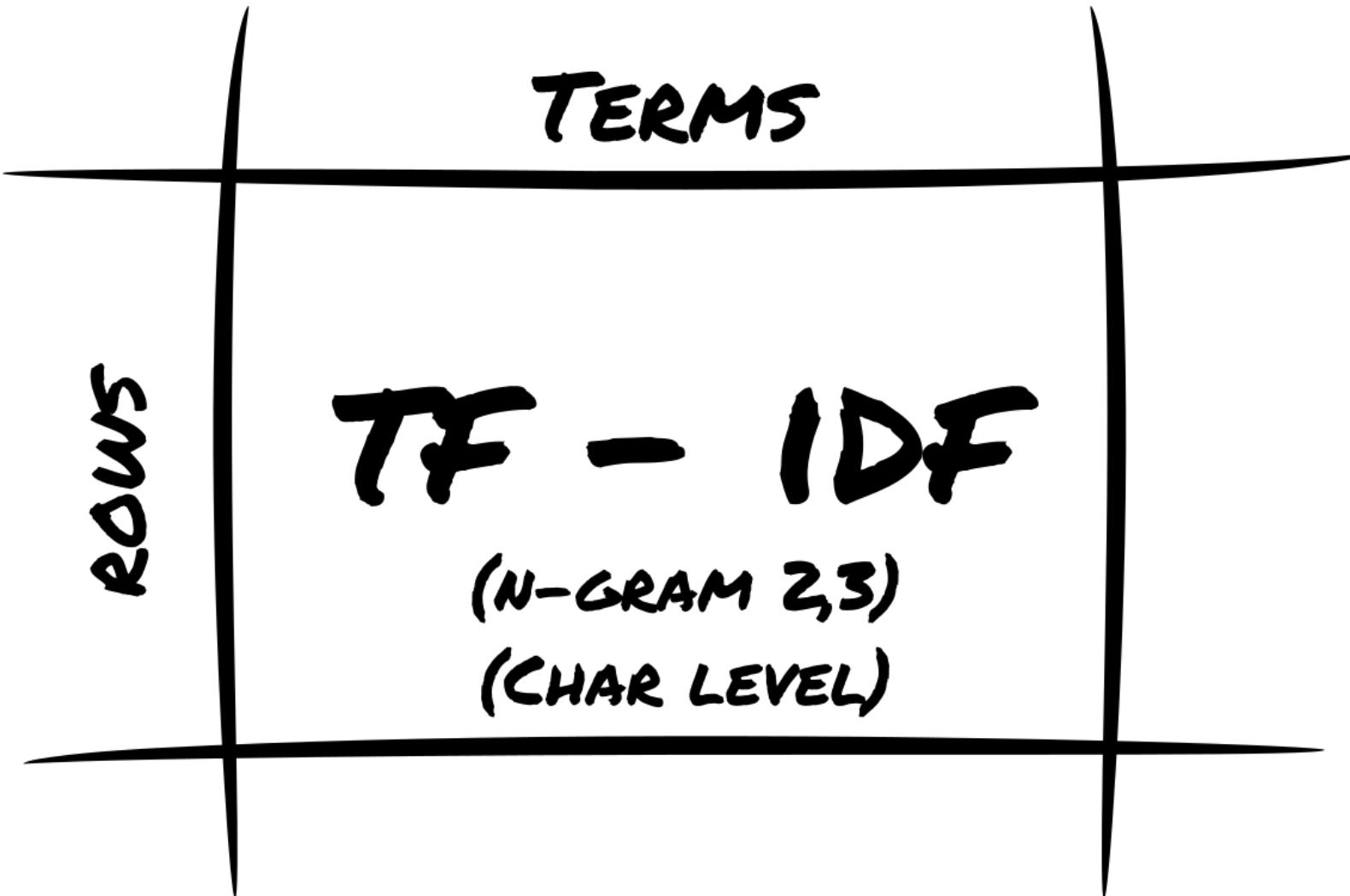
# Feature Engineering - Vectorizers



# Feature Engineering - Vectorizers



# Feature Engineering - Vectorizers



# *Models*

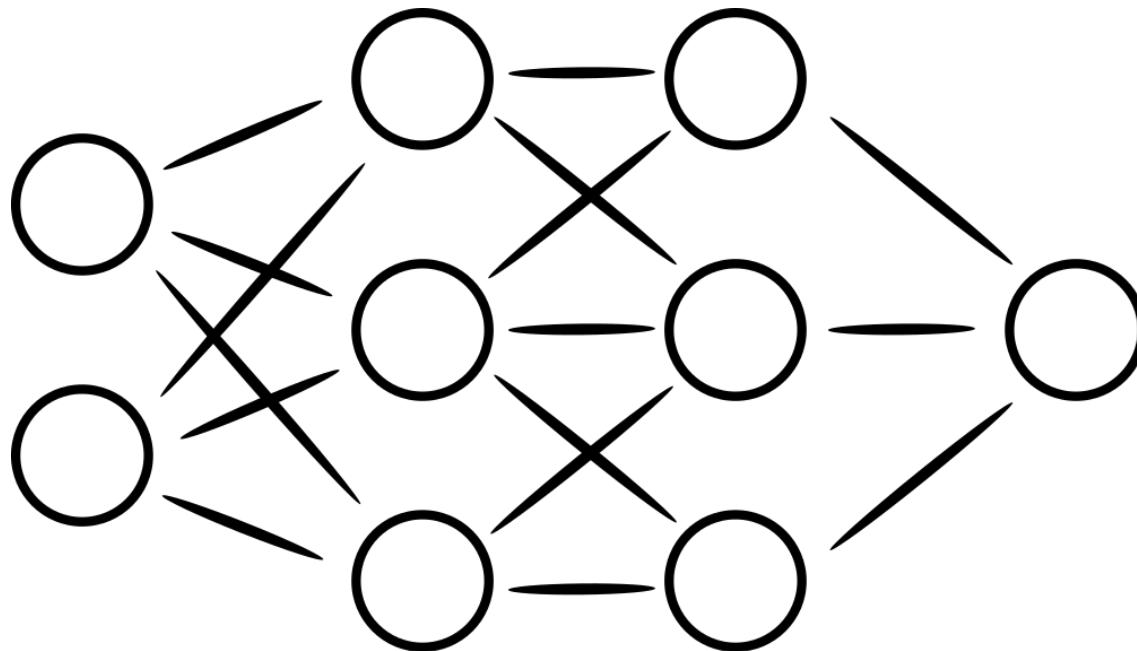
*We tested some simpler models*

- Logistic Regression
- Multinomial NB
- Random Forest Classifier
- Stochastic Gradient Descent Classifier
- Decision Tree Classifier
- XGBoost Classifier

# *Models*

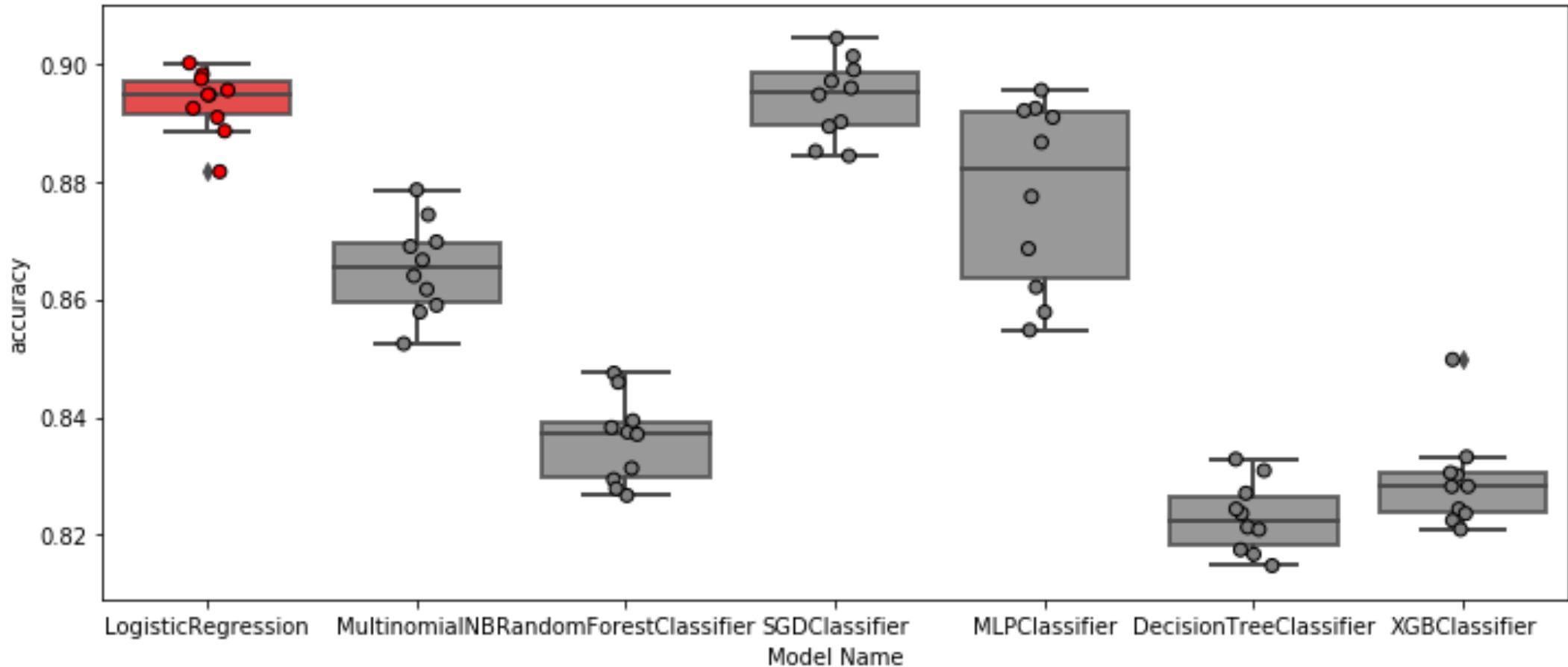
*And some more complicated - NN*

- MLP
- Deep Learning applications using Keras



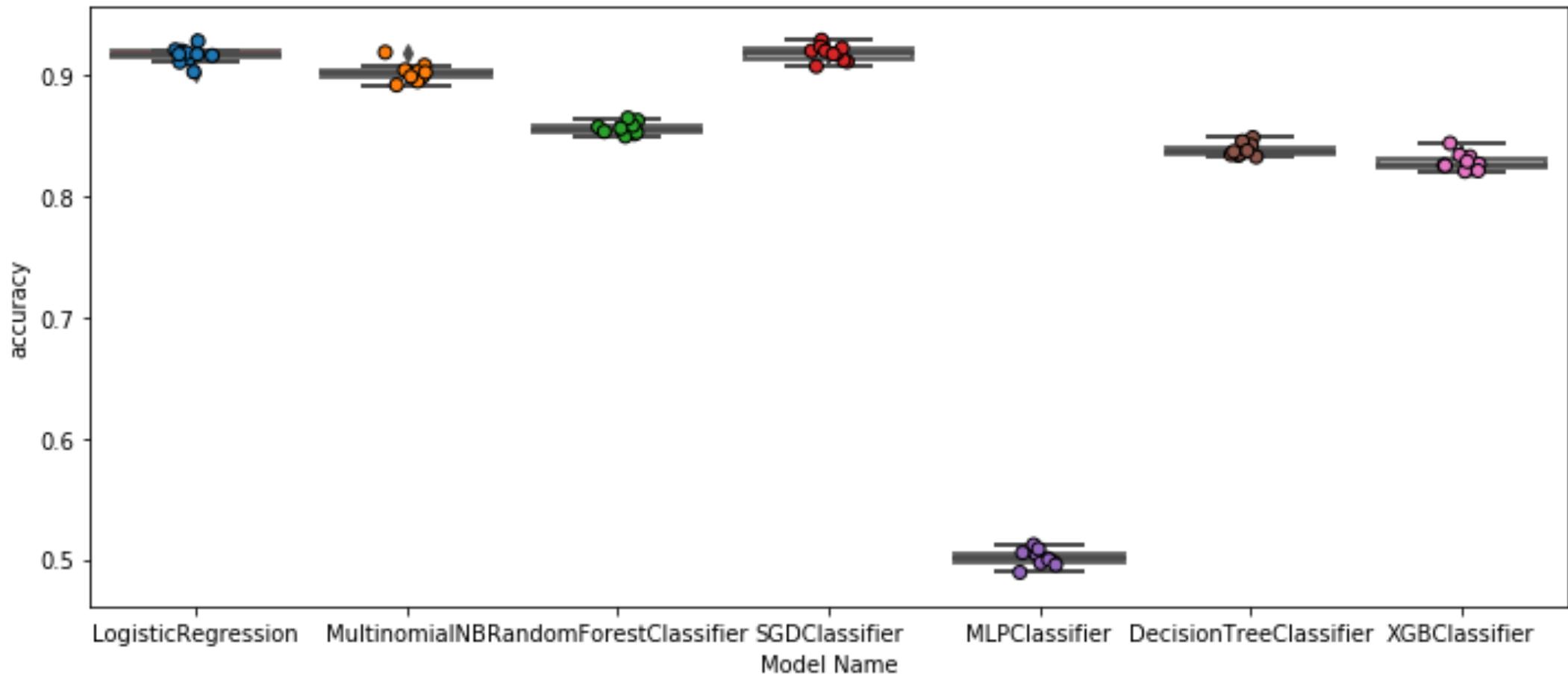
# *Models*

## *Without feature selection*



# *Models*

## *With feature selection*



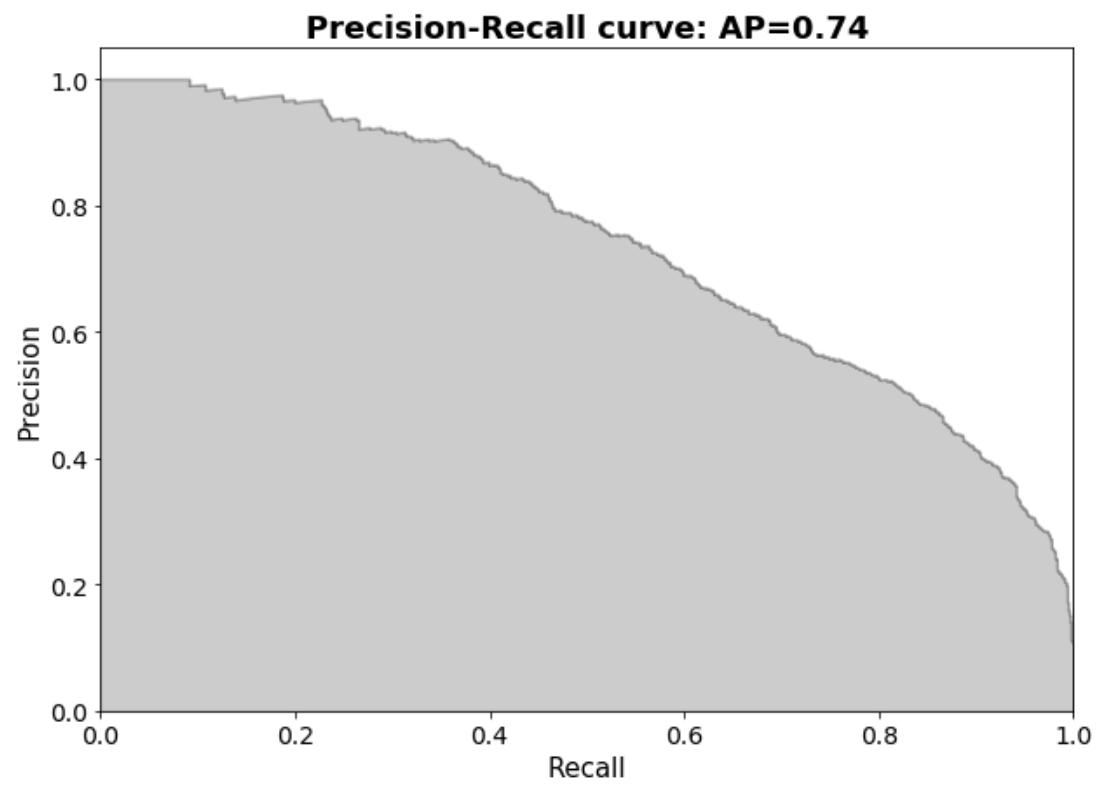
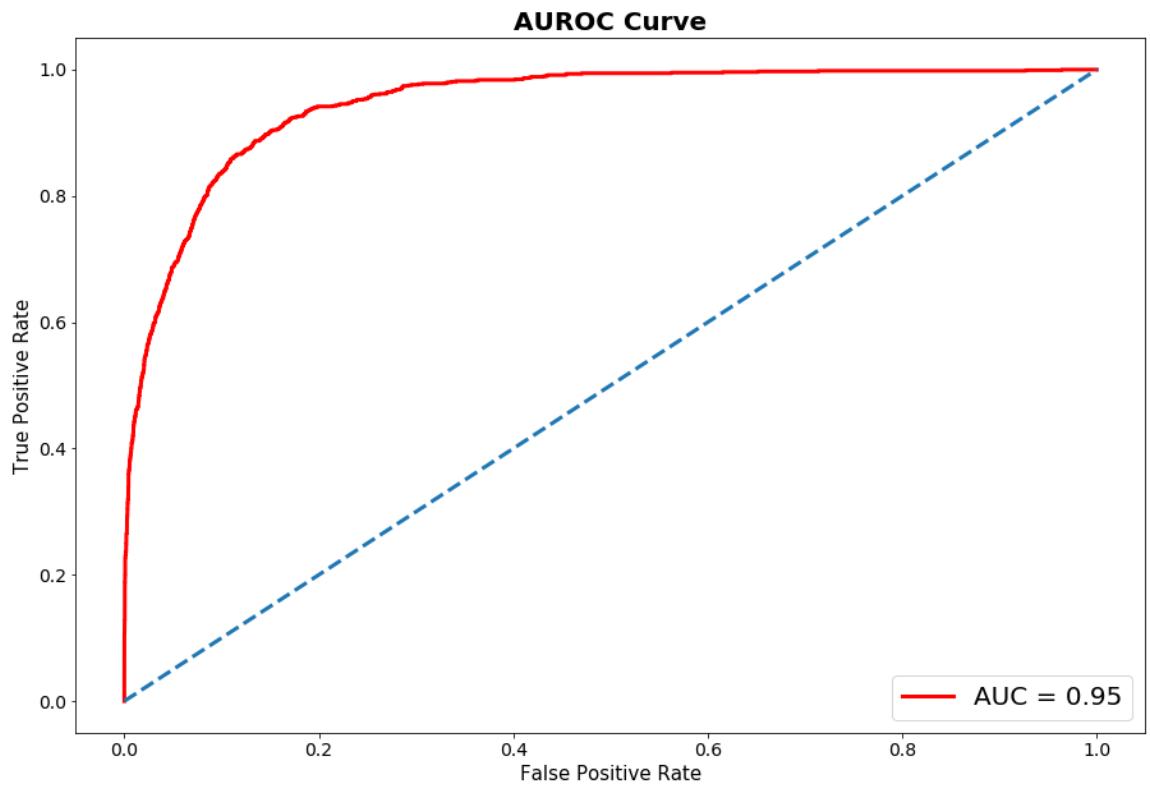
# *Models*

*Final Model: SGDClassifier with grid search*

	Precision	Recall	F1	Support
0	0.99	0.85	0.91	8954
1	0.41	0.90	0.56	1046
<b>Weighted AVG</b>	0.93	0.85	0.87	10000

Accuracy: 0.852

# Models



# *Applications*

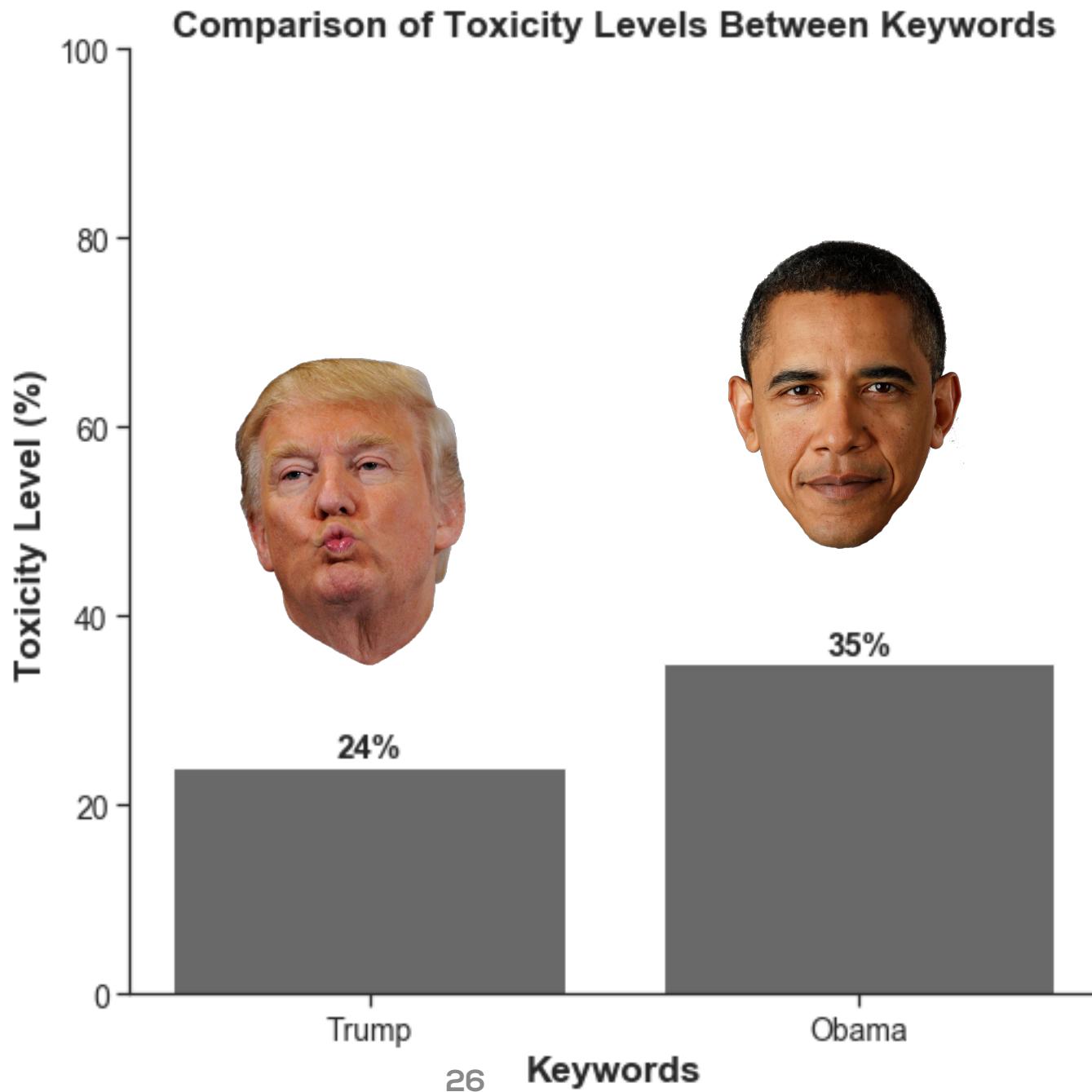
## SOCIAL NETWORKS

# *Applications*

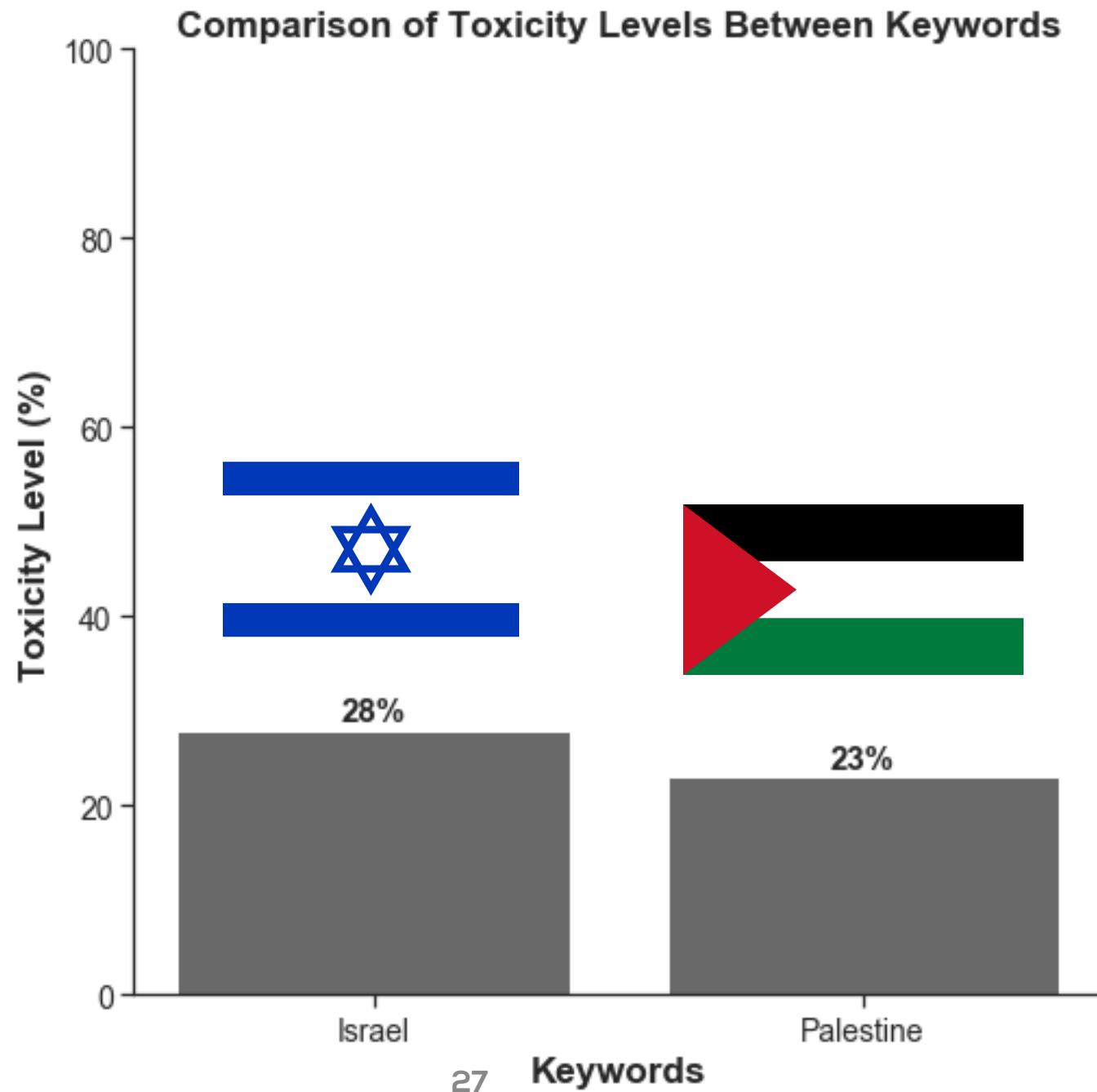
SOCIAL   
NETWORKS

BATTLE ROYALE

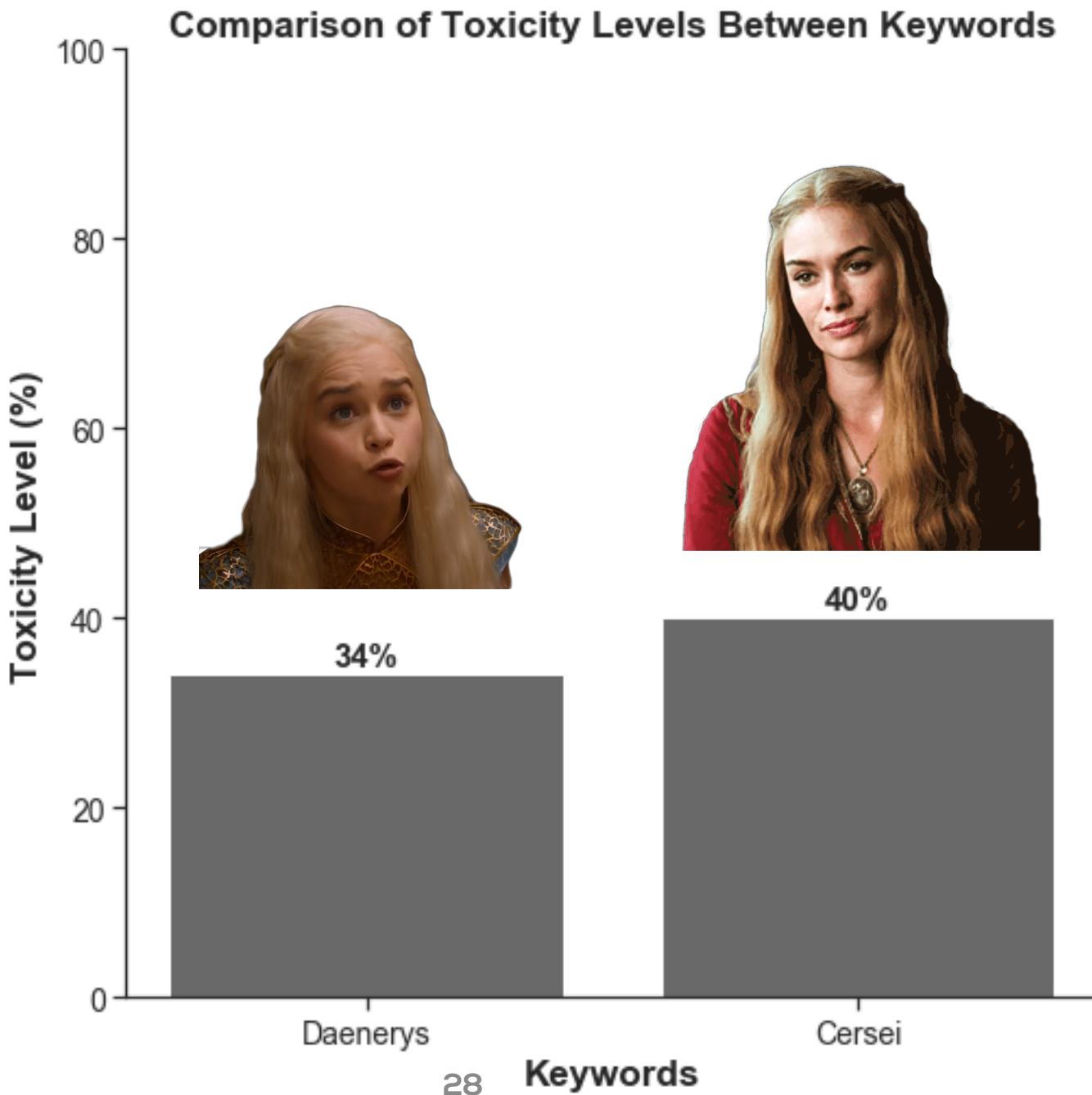
# *Applications*



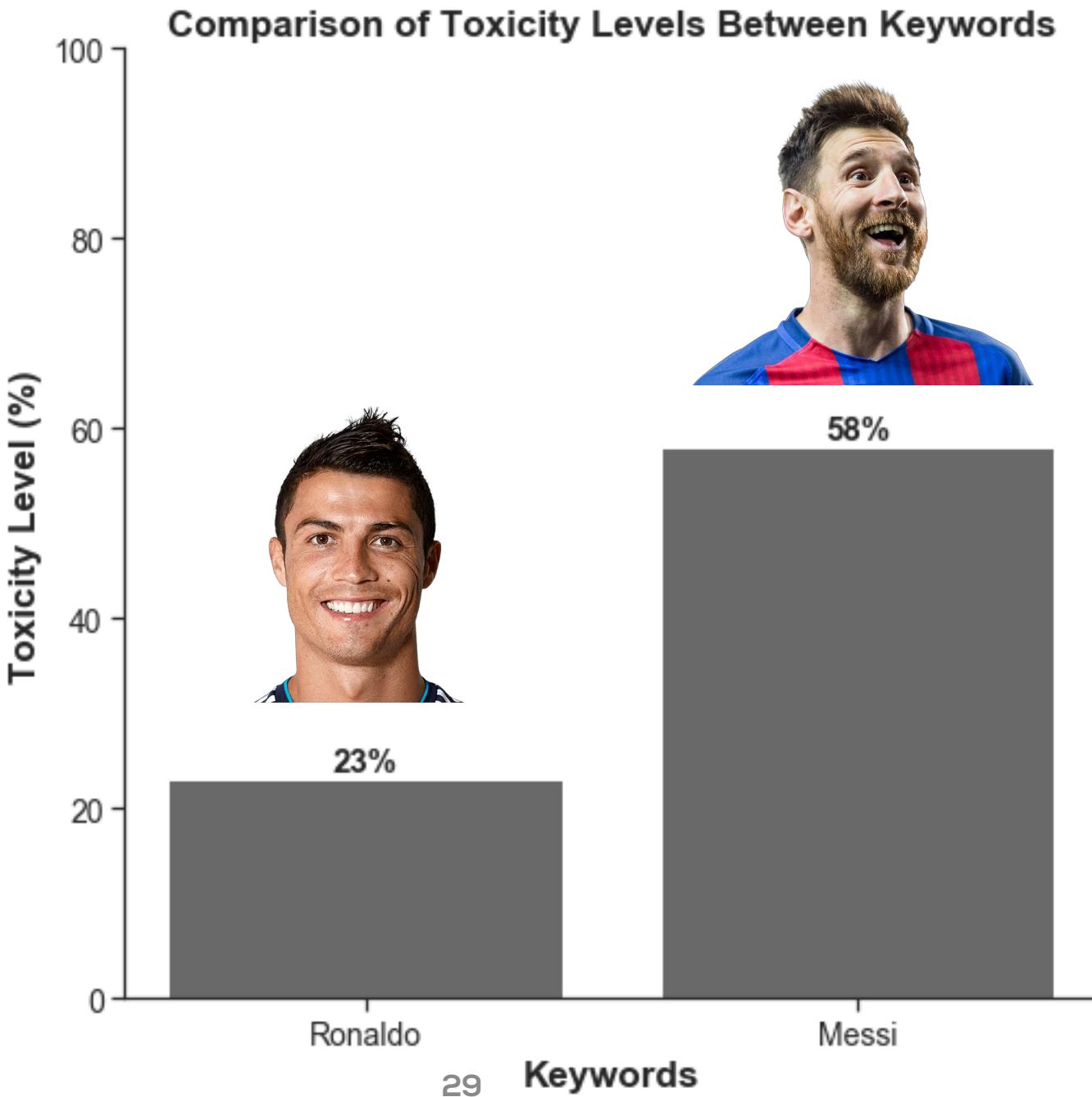
# *Applications*



# *Applications*



# *Applications*



# *Applications*

**MOVIES**   
**OFFENSIVENESS**

**The Internet Movie Script Database (IMDb)**

**IMDb**  
The web's largest movie script resource!

**Search IMDb**

## All Movie Scripts on IMDb (A-Z)

**Alphabetical**

#	A	B	C	D	E	F	G	H
I	J	K	L	M	N	O	P	Q
R	S	T	U	V	W	X	Y	Z

**Genre**

Action	Adventure	Animation
Comedy	Crime	Drama
Family	Fantasy	Film-Noir
Horror	Musical	Mystery
Romance	Sci-Fi	Short
Thriller	War	Western

**Sponsor**

[1492: Conquest of Paradise](#) (1991-09 Draft)  
Written by Roslyn Bosch

**TV Transcripts**

[Futurama](#)  
[Seinfeld](#)  
[South Park](#)  
[Stargate SG-1](#)  
[Lost](#)  
[The 4400](#)

**International**

[French scripts](#)

**Movie Software**



WinX DVD Ripper  
+ HD Converter (free)

[Rip from DVD](#)  
[Rip Blu-Ray](#)

[10 Things I Hate About You](#) (1997-11 Draft)  
Written by Karen McCullah Lutz, Kirsten Smith, William Shakespeare

[12](#) (Undated Draft)  
Written by Lawrence Bridges

[12 and Holding](#) (2004-04 Draft)  
Written by Anthony Cipriano

[12 Monkeys](#) (1994-06 Draft)  
Written by David Peoples, Janet Peoples

[12 Years a Slave](#) (Undated Draft)  
Written by John Ridley

[127 Hours](#) (Undated Draft)  
Written by Simon Beaufoy, Danny Boyle

[1492: Conquest of Paradise](#) (1991-09 Draft)  
Written by Roslyn Bosch

[15 Minutes](#) (Undated Draft)  
Written by John Hertzfield

[17 Again](#) (2007-10 Draft)  
Written by Jason Filardi

[187](#) (1996-11 Draft)  
Written by Scott Yagemann

[2001: A Space Odyssey](#) (1989-02 Draft)  
Written by Stanley Kubrick, Arthur C. Clarke

[2012](#) (2008-02 Second draft)  
Written by Roland Emmerich, Harald Kloser

[25th Hour](#) (2001-04 Draft)  
Written by David Benioff

[30 Minutes or Less](#) (2009-12 Draft)  
Written by Michael Diliberti, Matthew Sullivan

[42](#) (2012-07 Revised draft)  
Written by Brian Helgeland

[44 Inch Chest](#) (Undated Draft)  
Written by Louis Mellis, David Scinto

[48 Hrs.](#) (Undated Draft)  
Written by Steven E. De Souza, Walter Hill, Roger Spottiswoode, Larry Gross, Jeb Stuart

**IMDb**  
The web's largest movie script resource!

**Search IMDb**

**Alphabetical**

#	A	B	C	D	E	F	G	H
I	J	K	L	M	N	O	P	Q
R	S	T	U	V	W	X	Y	Z

**Genre**

Action	Adventure	Animation
Comedy	Crime	Drama
Family	Fantasy	Film-Noir
Horror	Musical	Mystery
Romance	Sci-Fi	Short
Thriller	War	Western

**Sponsor**

**TV Transcripts**

[Futurama](#)  
[Seinfeld](#)  
[South Park](#)  
[Stargate SG-1](#)  
[Lost](#)  
[The 4400](#)

**International**

[French scripts](#)

**Movie Software**

 WinX DVD Ripper + HD Converter (free)

[Rip from DVD](#)  
[Rip Blu-Ray](#)

**Latest Comments**

"PULP FICTION"

By Quentin Tarantino & Roger Avary

PULP [pulp] n.

1. A soft, moist, shapeless mass or matter.
2. A magazine or book containing lurid subject matter and being characteristically printed on rough, unfinished paper.

American Heritage Dictionary: New College Edition

INT. COFFEE SHOP - MORNING

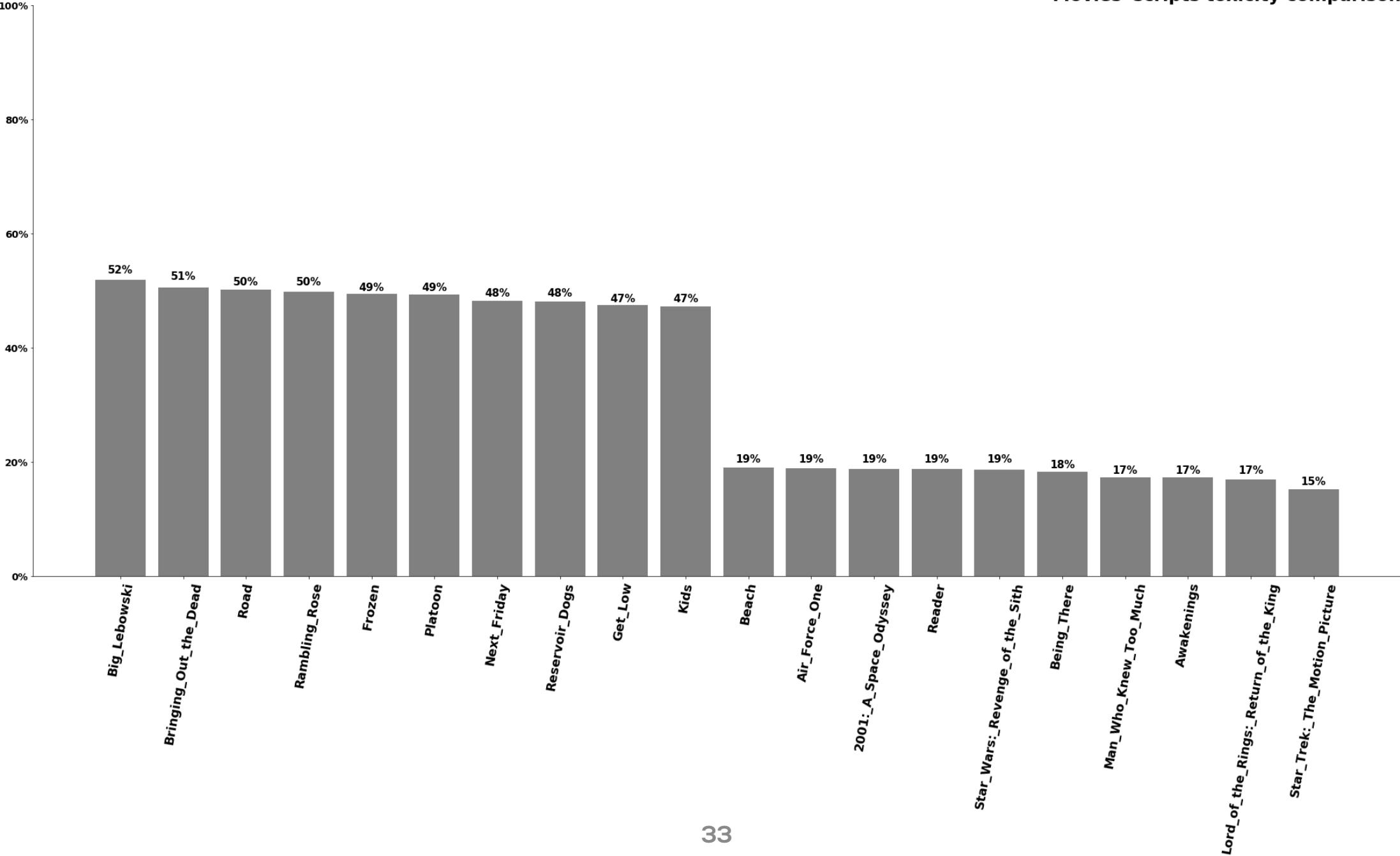
A normal Denny's, Spires-like coffee shop in Los Angeles. It's about 9:00 in the morning. While the place isn't jammed, there's a healthy number of people drinking coffee, munching on bacon and eating eggs.

Two of these people are a YOUNG MAN and a YOUNG WOMAN. The Young Man has a slight working-class English accent and, like his fellow countryman, smokes cigarettes like they're going out of style.

It is impossible to tell where the Young Woman is from or how old she is; everything she does contradicts something she did. The boy and girl sit in a booth. Their dialogue is to be said in a rapid pace "HIS GIRL FRIDAY" fashion.

YOUNG MAN  
No, forget it, it's too risky. I'm through doin' that shit.

## Movies' scripts toxicity comparison



# *Thank you*

*Rita Franco*      20180081  
*Rodrigo Umbelino*    20180060  
*Vitor Manita*       20180054