

# NATURAL LANGUAGE PROCESSING

(for the impatient)

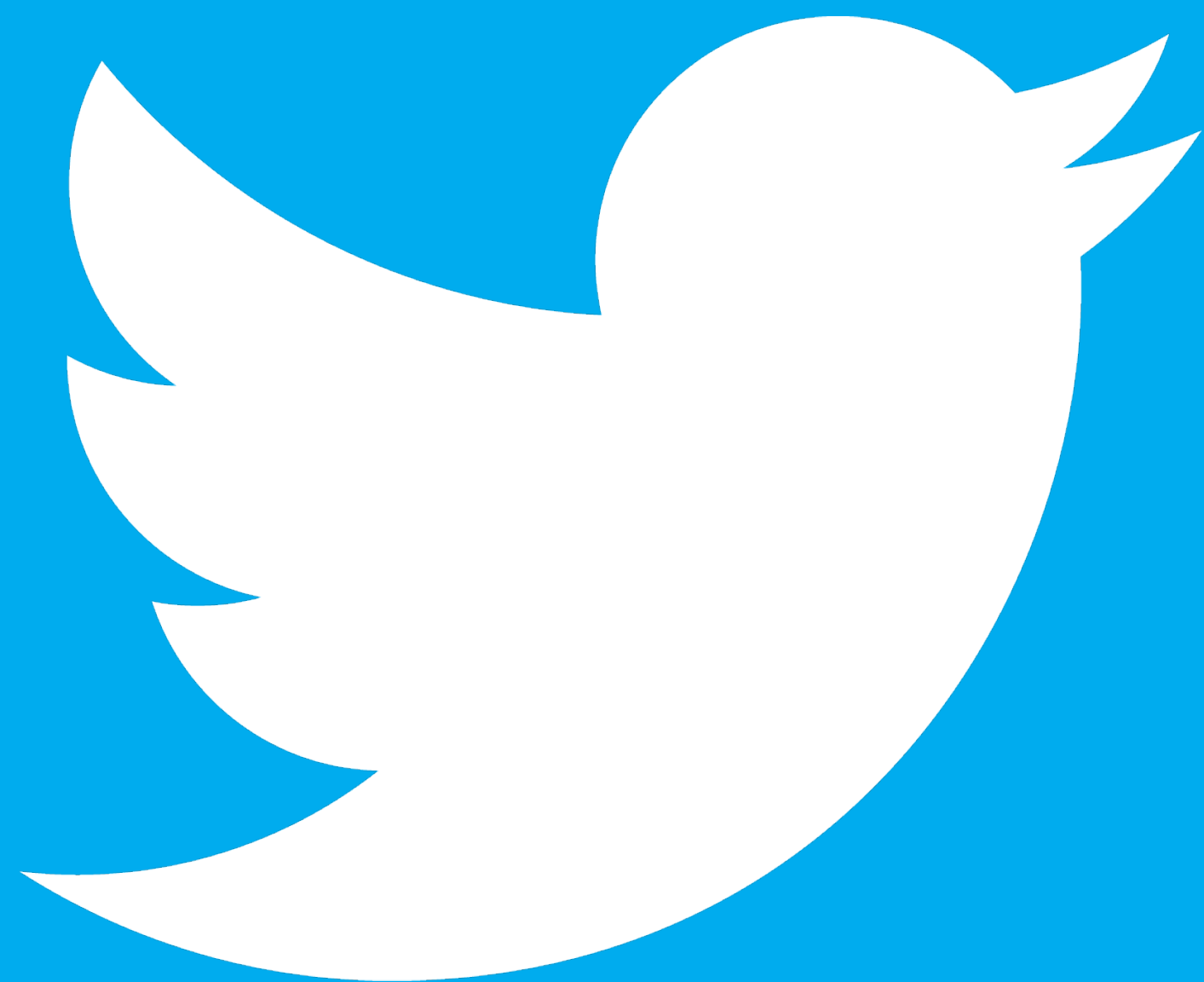
# INTRODUCTION



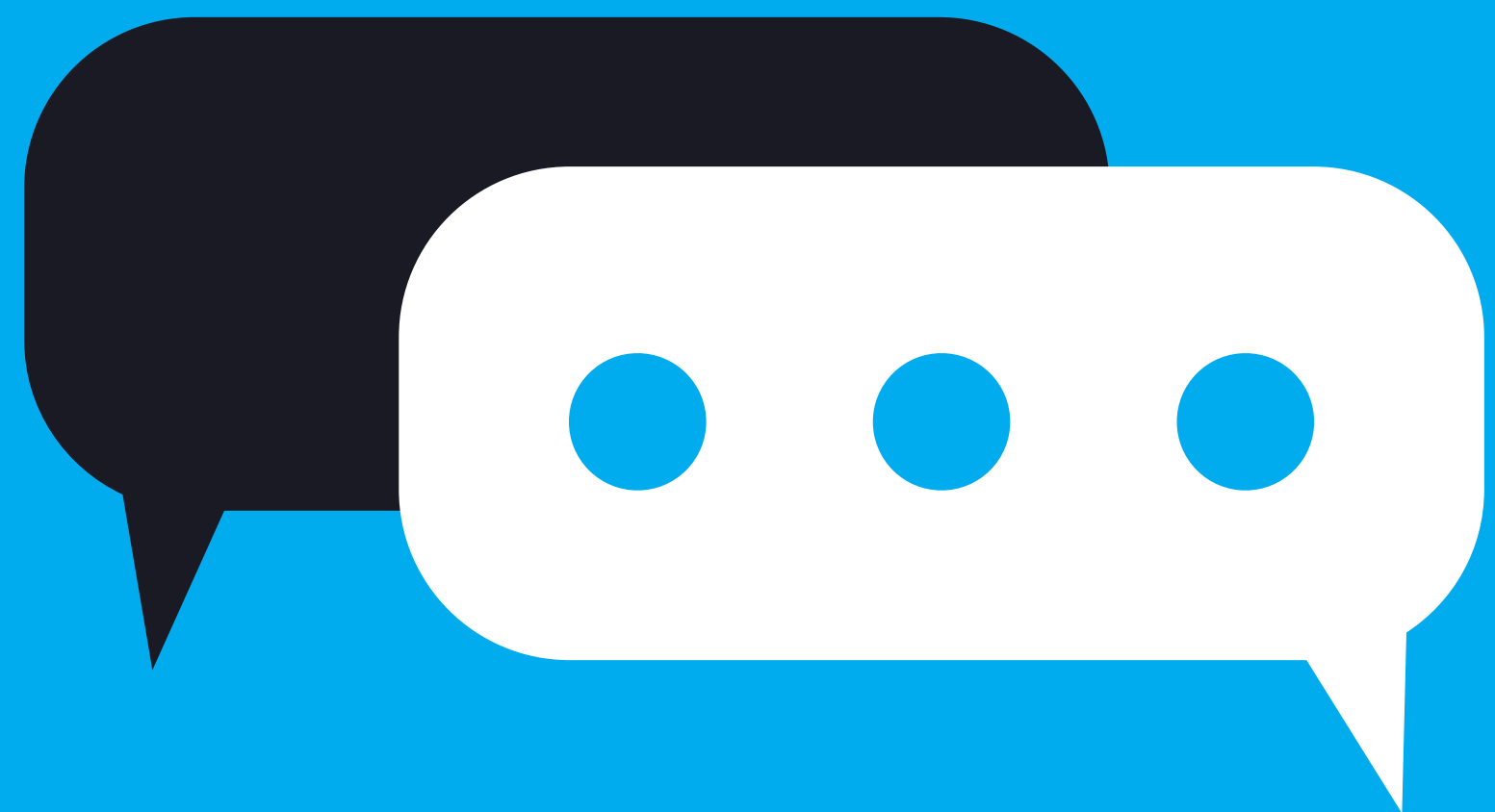


Sebastian **Dziadzio**

[cliqz.com](http://cliqz.com)  
[ghostery.com](http://ghostery.com)



**@sebadzia**



# Natural Language Processing



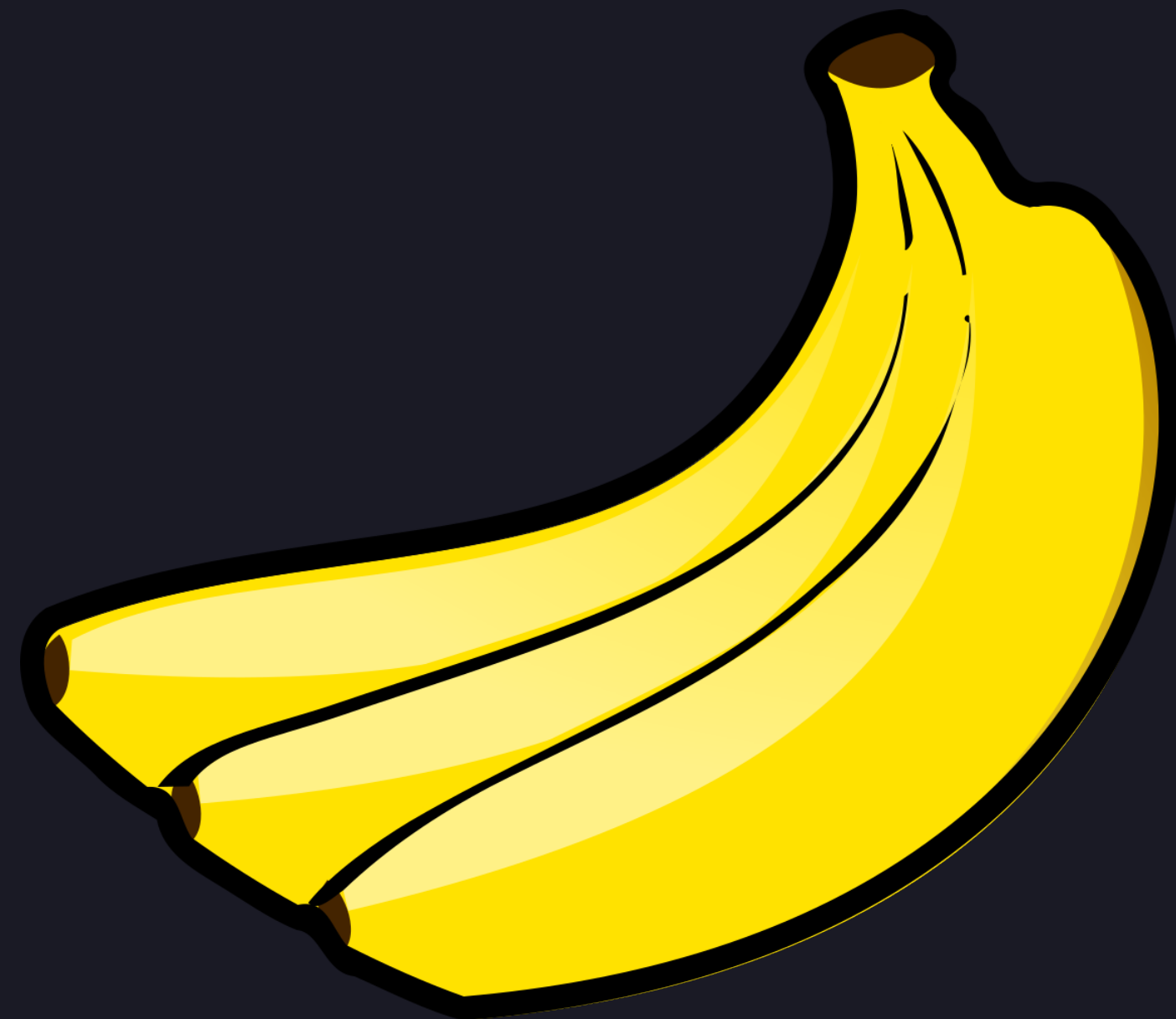
**Why is  
Natural  
Language  
Processing  
important?**



**Why is  
Natural  
Language  
Processing  
hard?**

Time flies like an arrow.

Fruit flies like a banana.





worstest

hatemonger

post-truth

pogonophobia

opinion size age shape colour origin material purpose noun

ugly large old brown English tweed hunting jacket

big fat Greek wedding

# Rindfleischetikettierungsüberwachungsaufgabenübertragungsgesetz



# PROCESS & TOOLS



**Get data** → Transform → Encode → Visualize → Model

**Get data**

kaggle™



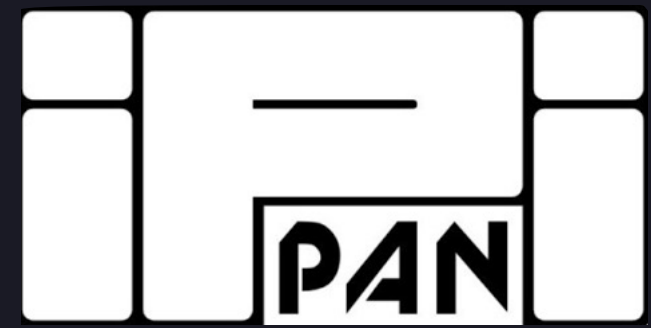
**data.world**



**niderhoff/nlp-datasets**



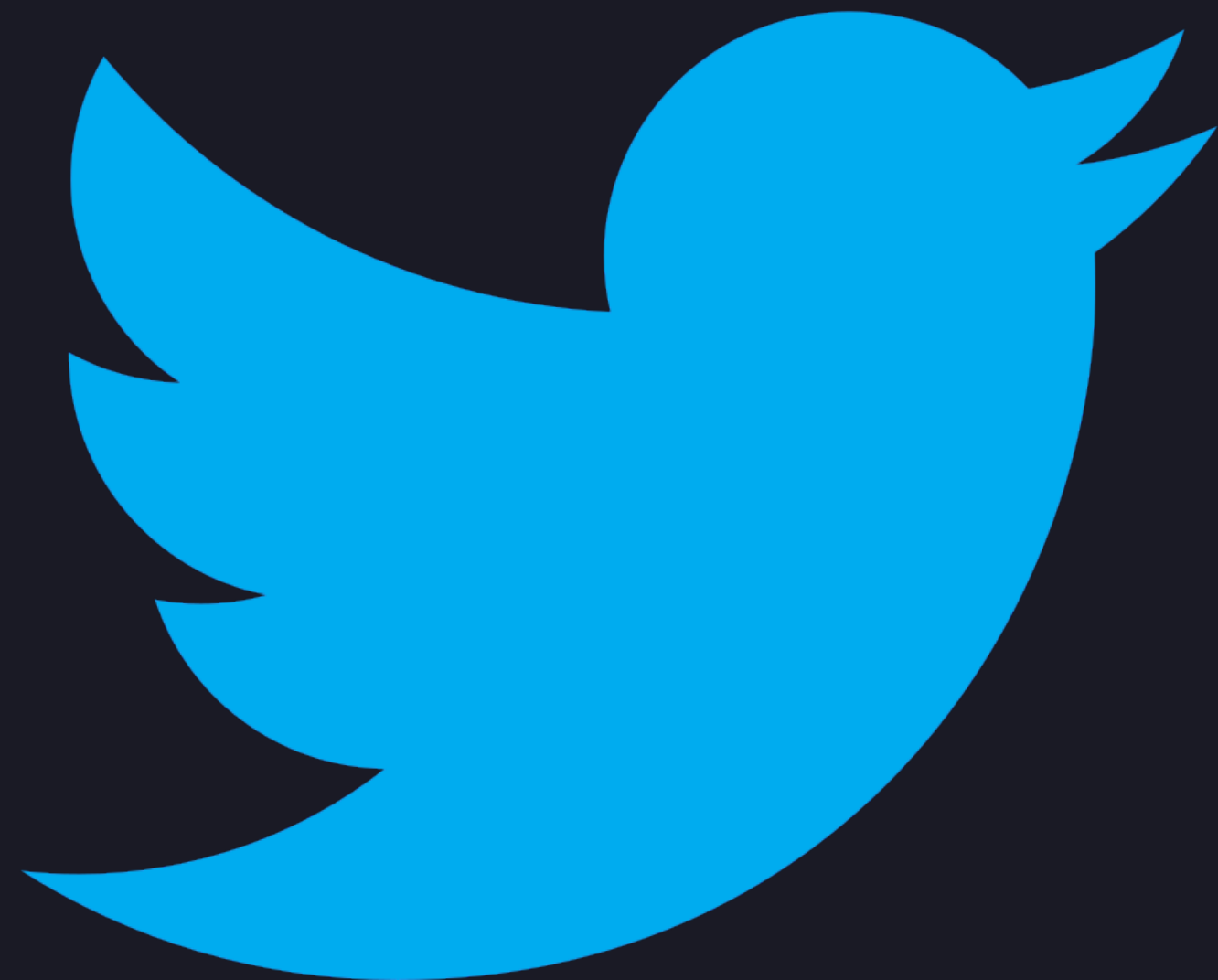
[nlp.stanford.edu](http://nlp.stanford.edu)



[zil.ipipan.waw.pl](http://zil.ipipan.waw.pl)



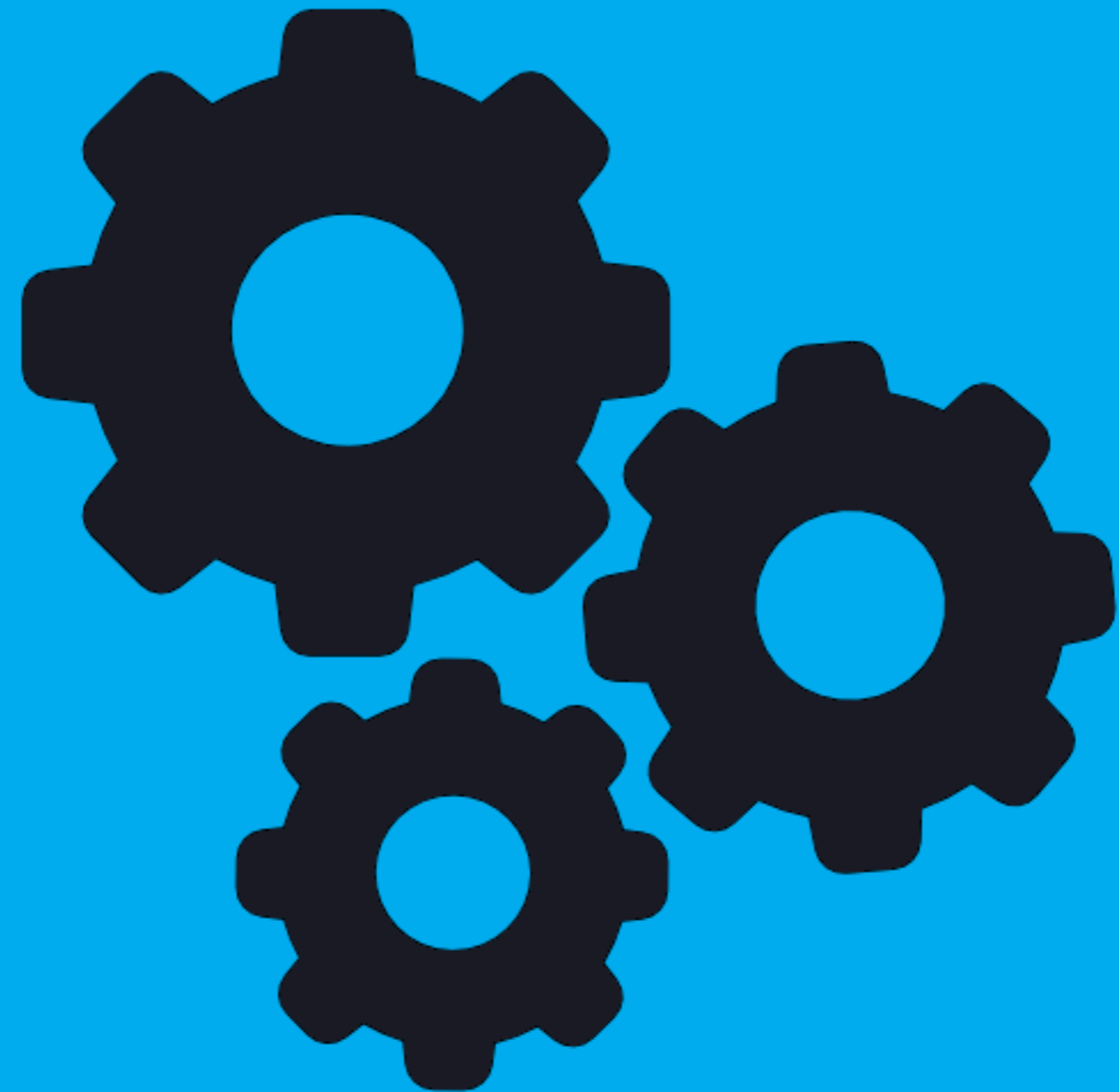
[keon/awesome-nlp](https://github.com/keon/awesome-nlp)





Get data → **Transform** → Encode → Visualize → Model

# Transform



**normalization**

**stemming**

**lemmatization**

**POS tagging**



**Plain Python**

**Unix**

**NLTK**

**SpaCy**

@TRINARockstarr Trina flow, working on my abs & legs ... So ""when sb ask ""what's yo REAL name"" I tell him Da Baddest Bitch 🤔😂😂😱💪 #gym

trina flow working on my abs and legs so when somebody ask what is your real name I tell him da baddest bitch

will devour your HTML parser, application and existence for all time like Visual Basic only worse  
he comes he comes do not fight he comes, his unholý radiancé destroying all enlightenment,  
HTML tags leaking from your eyes like liquid pain, the song of regular expression parsing will  
extinguish the voices of mortal man from the sphere I can see it can you see it it is beautiful the  
final snuffing of the lies of Man ALL IS LOST ALL IS LOST the pony he comes he comes he  
comes the ichor permeates all MY FACE MY FACE god no NO NOOOO NO stop the angles  
are not real ZALGO IS TONY THE PONY, HE COMES

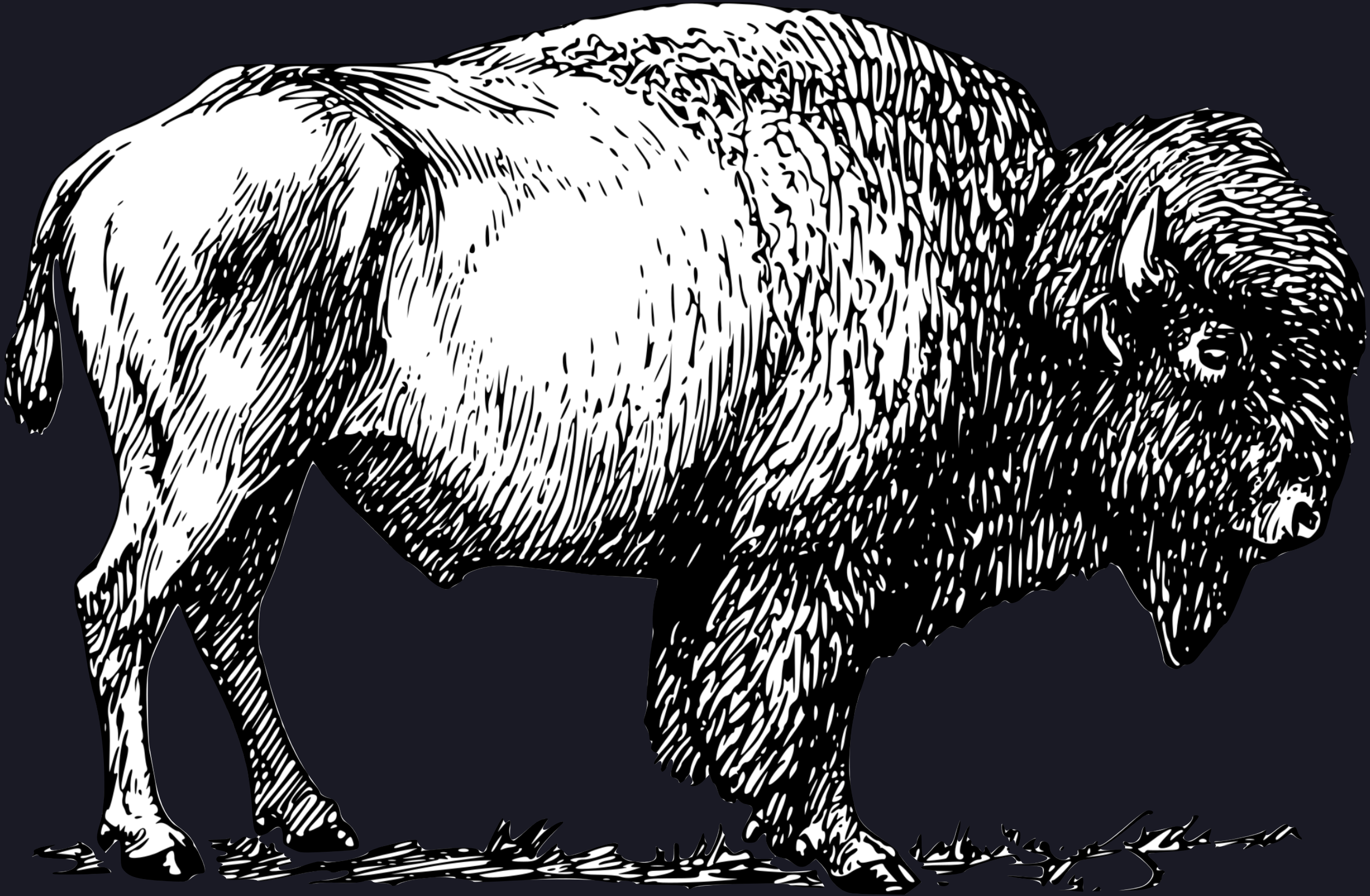
Listen, strange women lying in ponds distributing swords is no basis for a system of government. Supreme executive power derives from a mandate from the masses, not from some farcical aquatic ceremony.

Listen, strange woman lie in pond distribute sword be no basis for a system of government. Supreme executive power derive from a mandate from the mass, not from some farcical aquatic ceremony.

Listen, strang women lie in pond distribut sword is no basi for a system of government. Suprem execut power deriv from a mandat from the masses, not from some farcic aquat ceremon.

Buffalo buffalo Buffalo buffalo buffalo Buffalo buffalo.

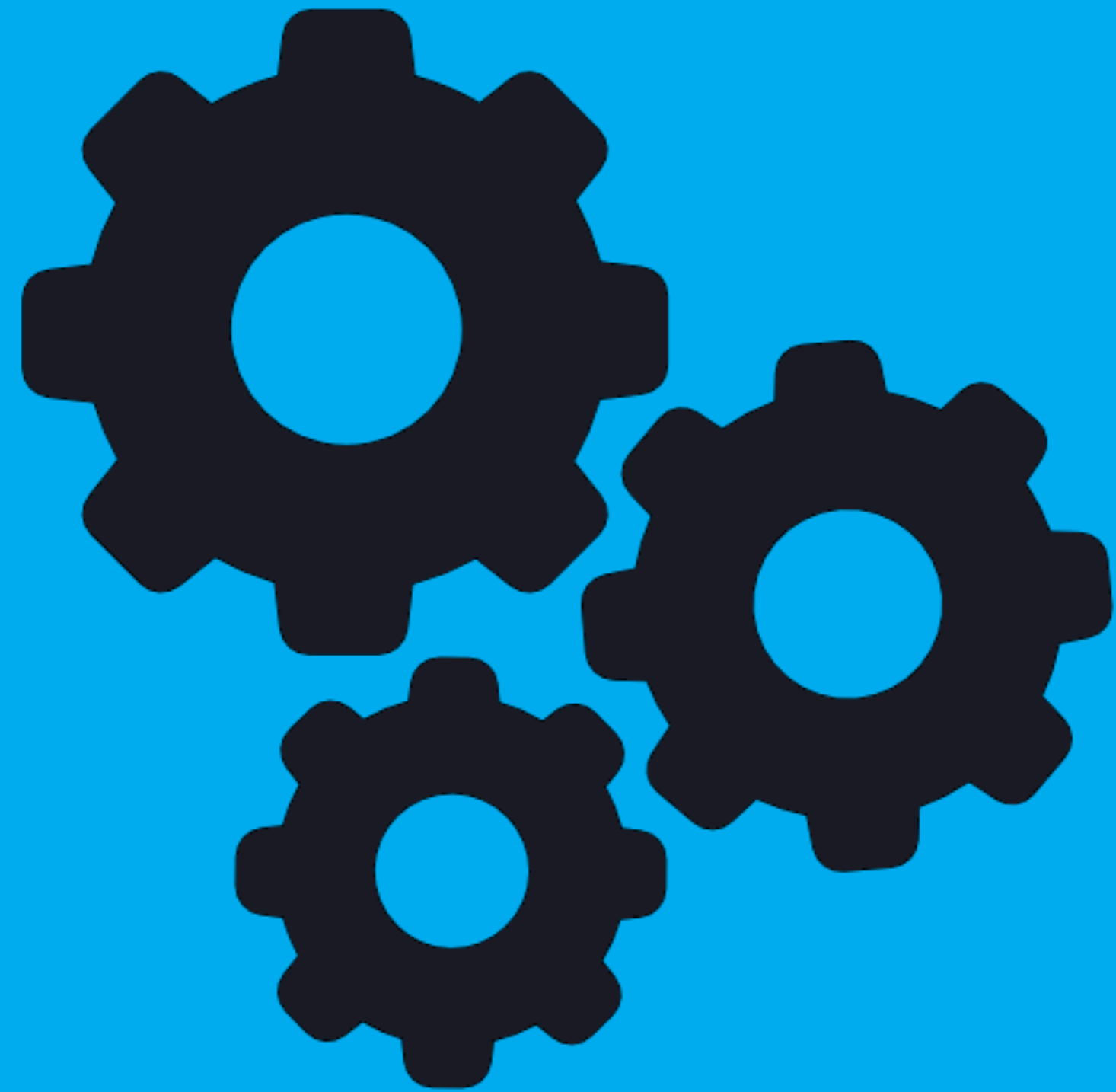
PN N PN N V V PN N



Get data → Transform → **Encode** → Visualize → Model

# Encode





**word2vec**

**GloVe**

**doc2vec**

**sense2vec**



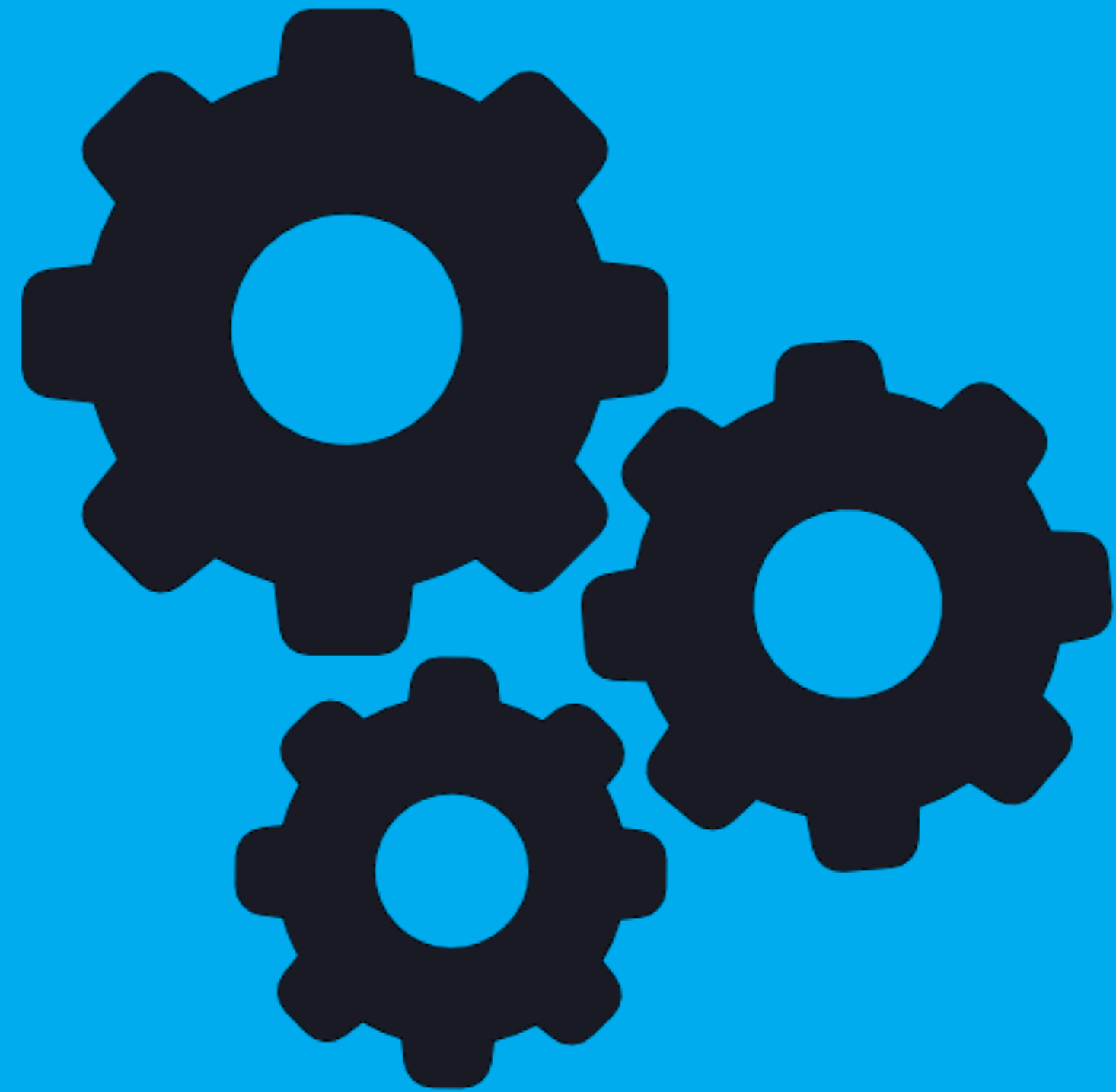
**gensim**  
**fastText**  
**GloVe**



**EXAMPLE:**  
**Information  
Retrieval**

Get data → Transform → Encode → **Visualize** → Model

# Visualize



**t-SNE**

**PCA**



**scikit-learn**

**matplotlib**



**EXAMPLE:**  
**Project**  
**Gutenberg**

governess judgement meryton park library her bingley convinced rosings uncle lizzy caroline harlotte ladyship fitzwilliam collins mary william attentions

darcey bourgh janeforster catherine behaviour kitylydia lucas wickham longbourn netherfield eliza bennet that hertfordshire

advantages enumerating twelve months connections boast agreeable partiality fortnight mamma militia probability

darcey bourgh janeforster catherine behaviour kitylydia lucas wickham longbourn netherfield eliza bennet that hertfordshire

convicted undoubtedly he believing vanity console teasing tuesday wednesday gratitude anybody gardiner hunsford lizzy caroline harlotte ladyship fitzwilliam collins mary william attentions

darcey bourgh janeforster catherine behaviour kitylydia lucas wickham longbourn netherfield eliza bennet that hertfordshire

derbyshire engagement consequently harlotte ladyship fitzwilliam collins mary william attentions

darcey bourgh janeforster catherine behaviour kitylydia lucas wickham longbourn netherfield eliza bennet that hertfordshire

derbyshire engagement consequently harlotte ladyship fitzwilliam collins mary william attentions

darcey bourgh janeforster catherine behaviour kitylydia lucas wickham longbourn netherfield eliza bennet that hertfordshire







persuasion

pride



emma

eyre

sense

shirley

villette

categories

poetics



politics

constitution

caesar

macbeth

hamlet

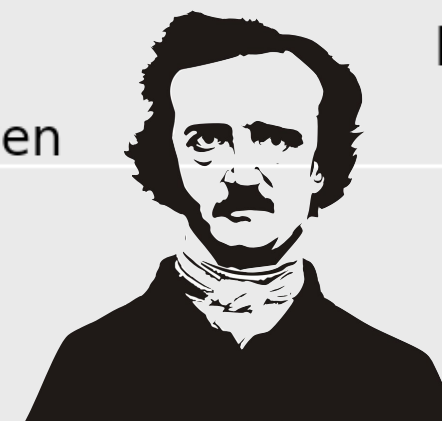


paradise

leaves

poems

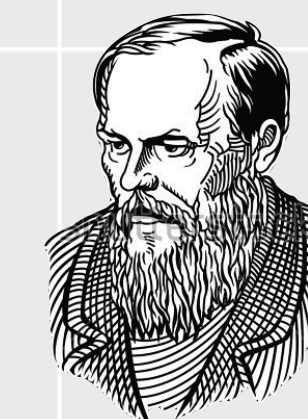
raven



crime

karamazov

idiot



scarlet

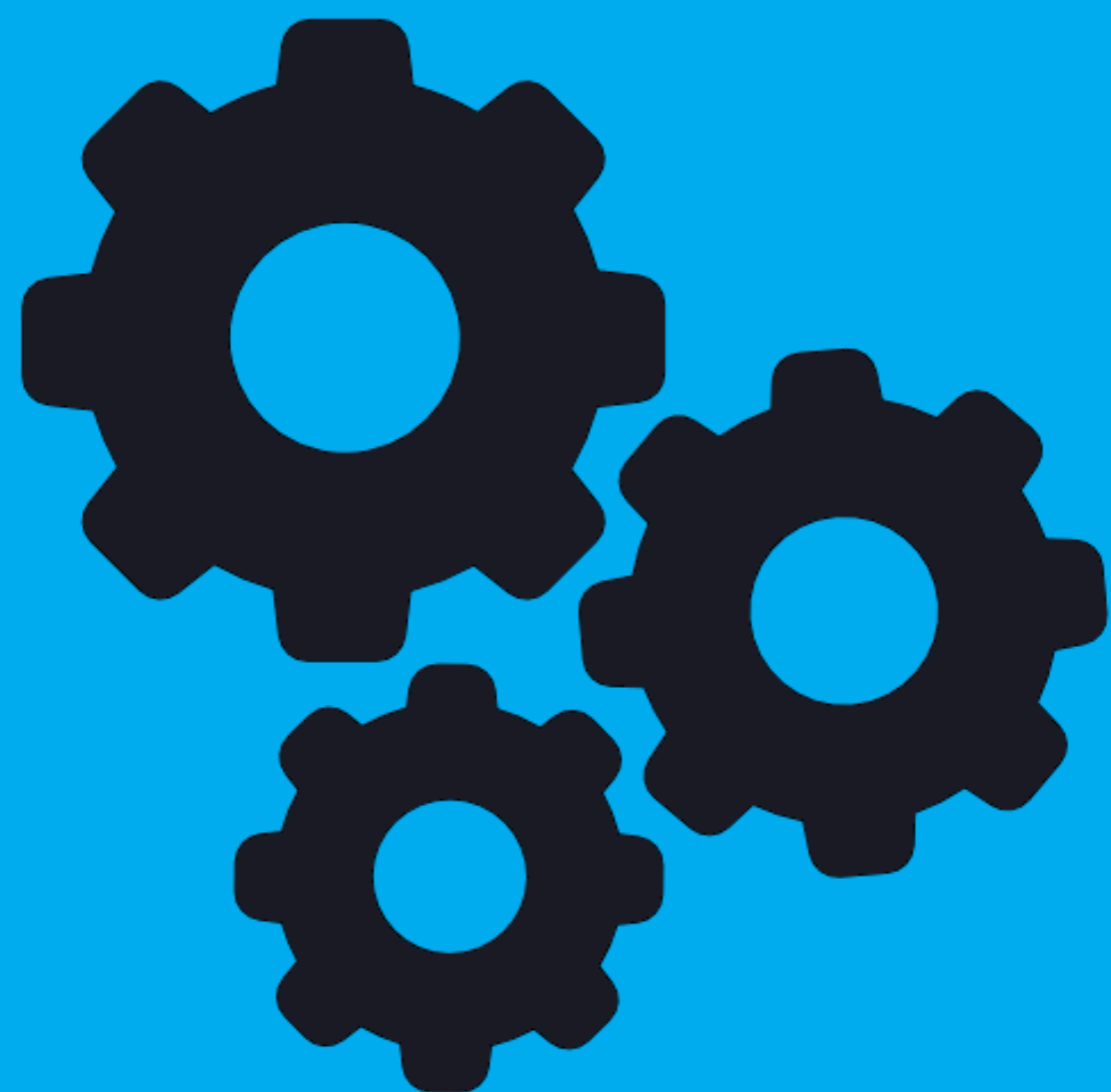


four

baskervilles

Get data → Transform → Encode → Visualize → **Model**

**Model**



**old school ML**

**(bi-)LSTM**

**attention**

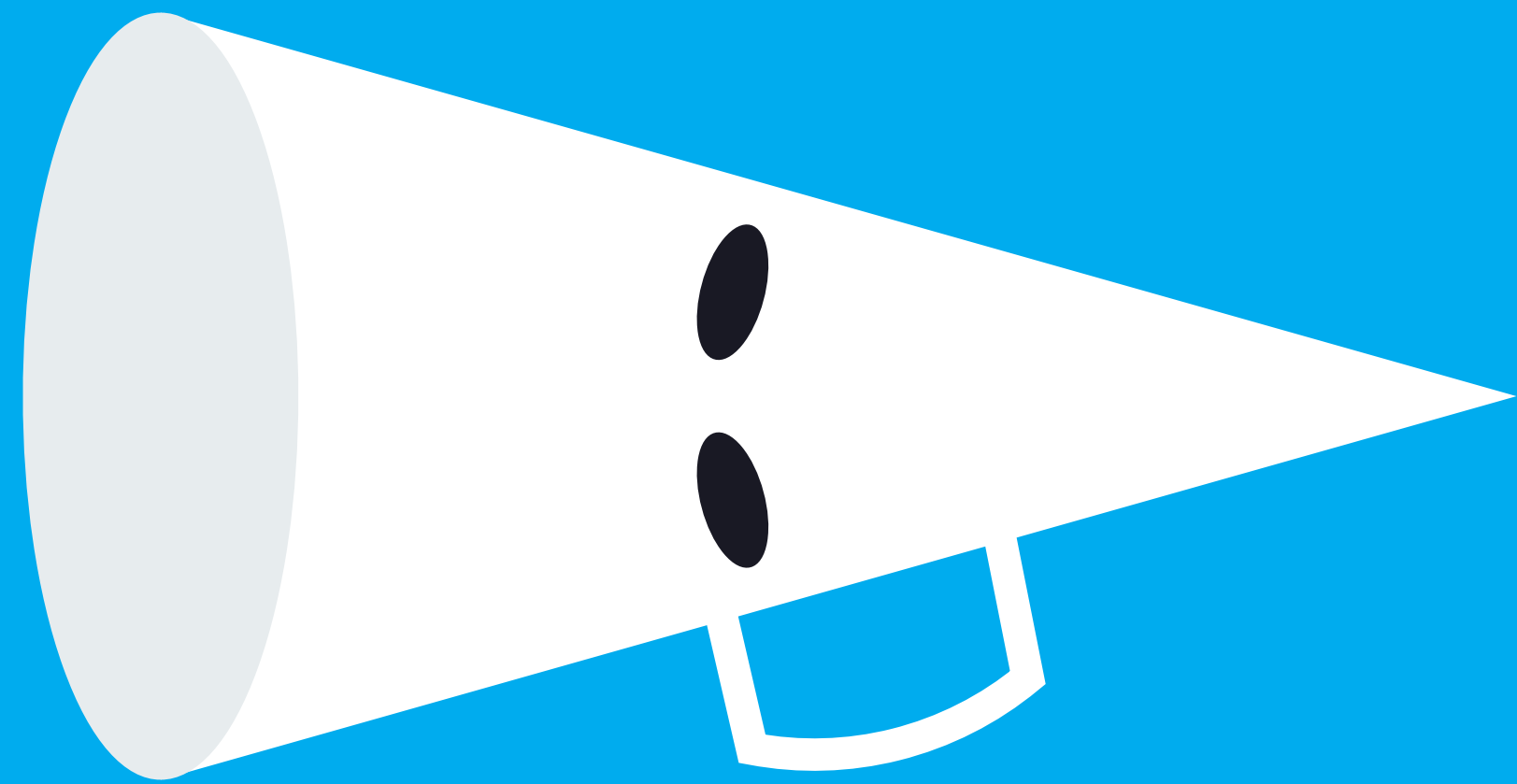


**scikit-learn**

**Keras**

**TensorFlow**

**SpaCy**

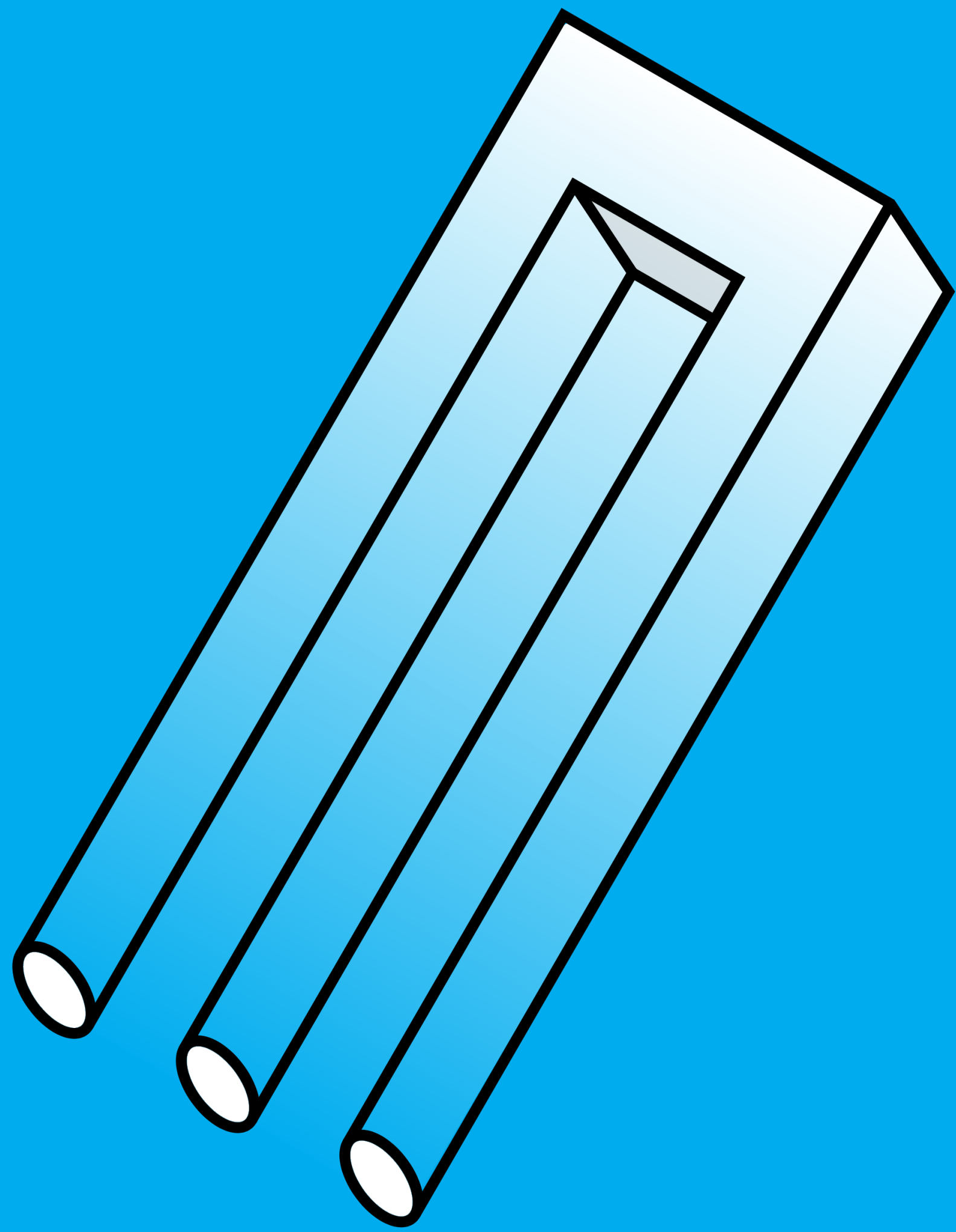


**EXAMPLE:**  
**Hate Speech  
Detection**

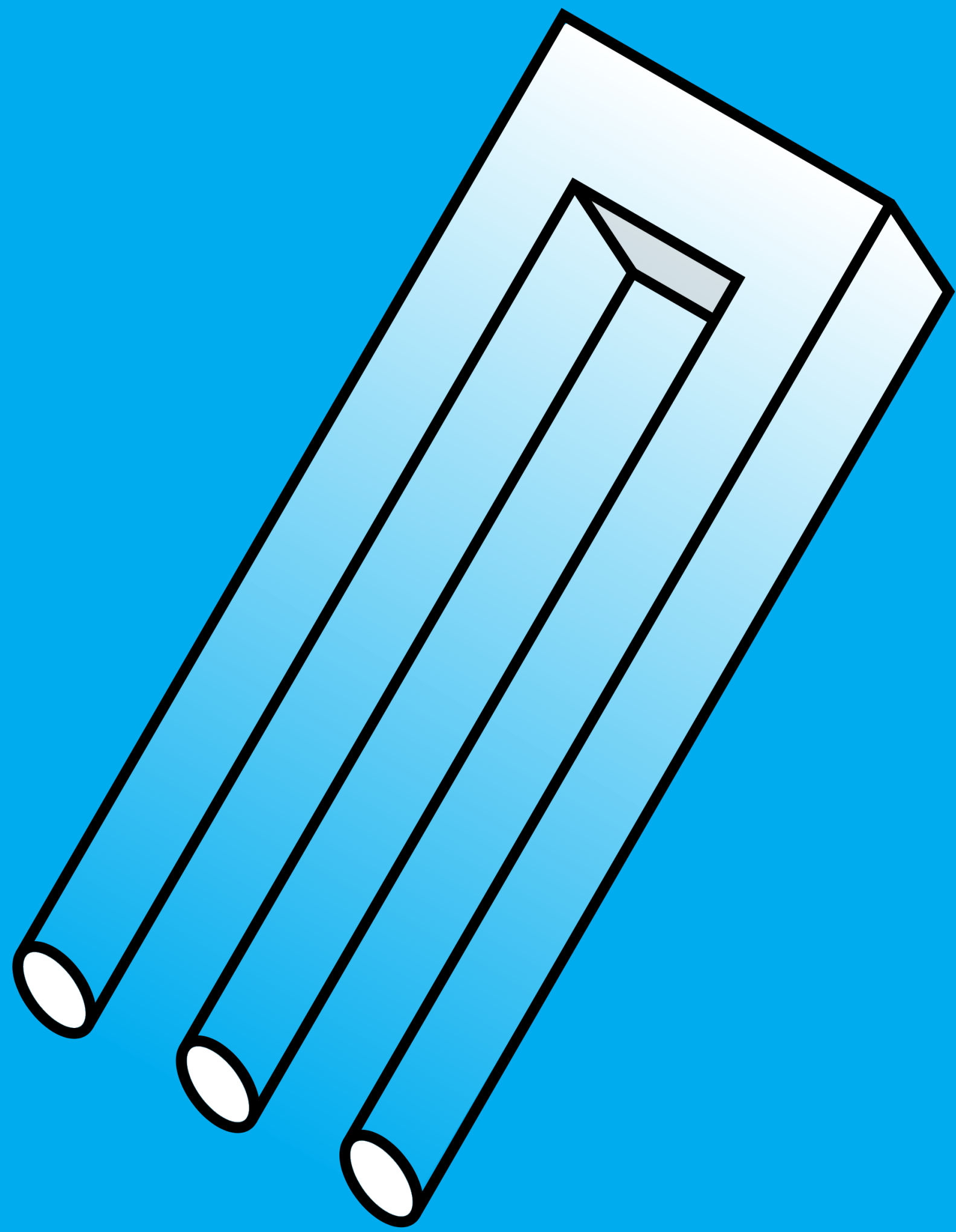
# BEST PRACTICES







**DEFINE**  
**the problem**



**task**

**data**

**metric**

**assumptions**

**motivation**



**INSPECT**  
**the data**



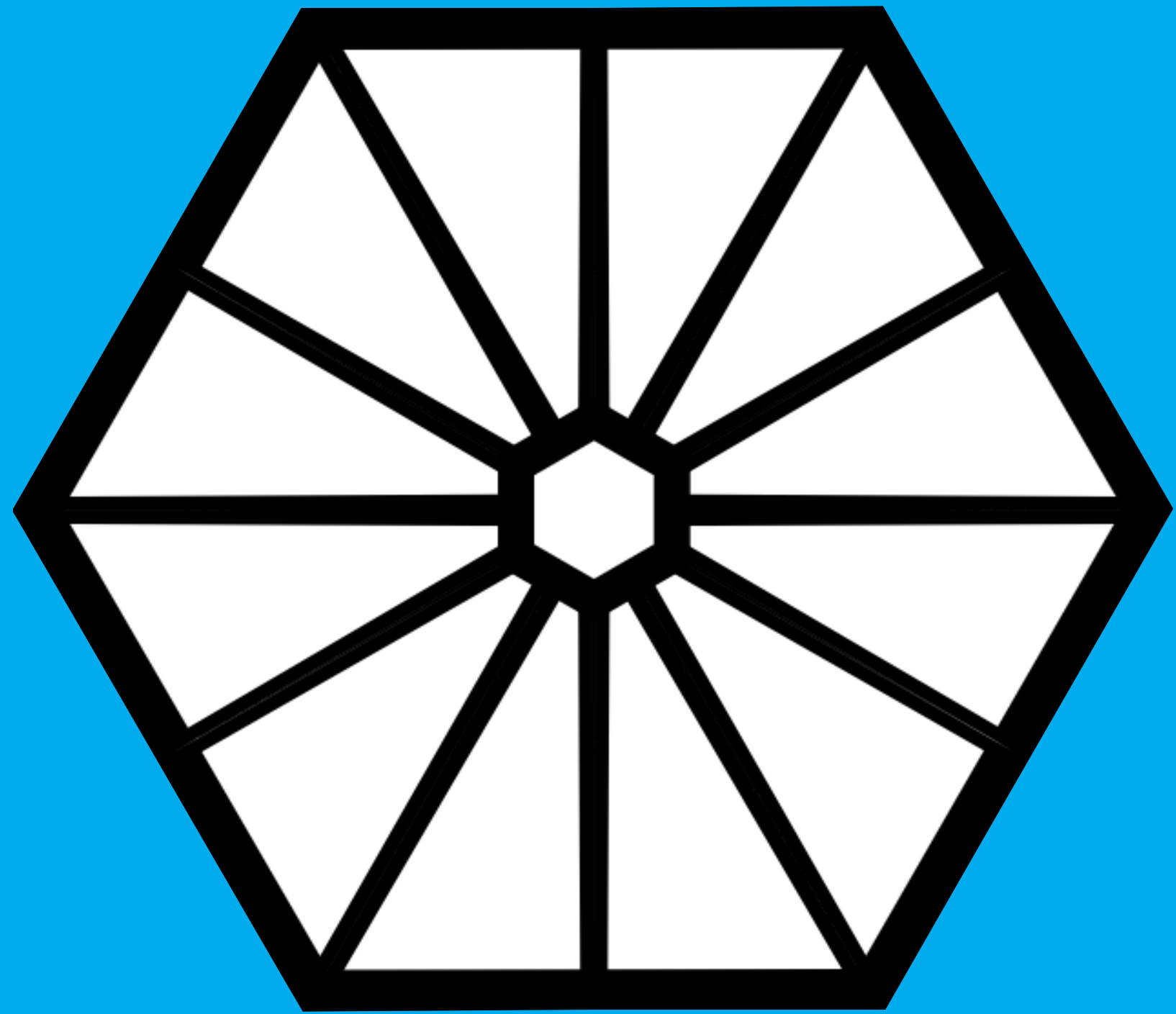
**noise**

**order**

**imbalance**

**sparsity**

**encoding**



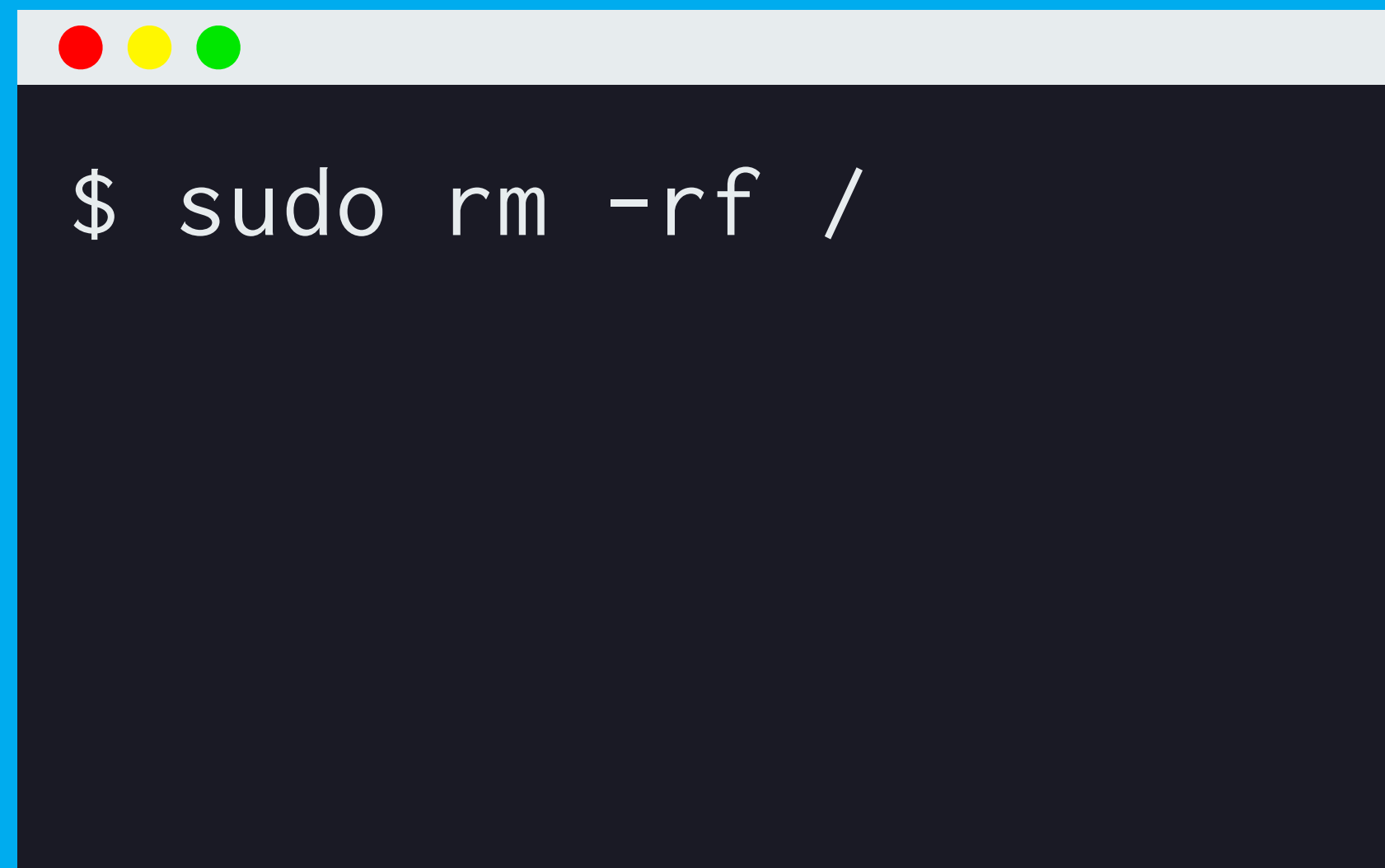
**AVOID**  
**reinventing**  
**the wheel**



**CHOOSE**  
**the simplest**  
**tool possible**



**OPEN**  
**the black box**

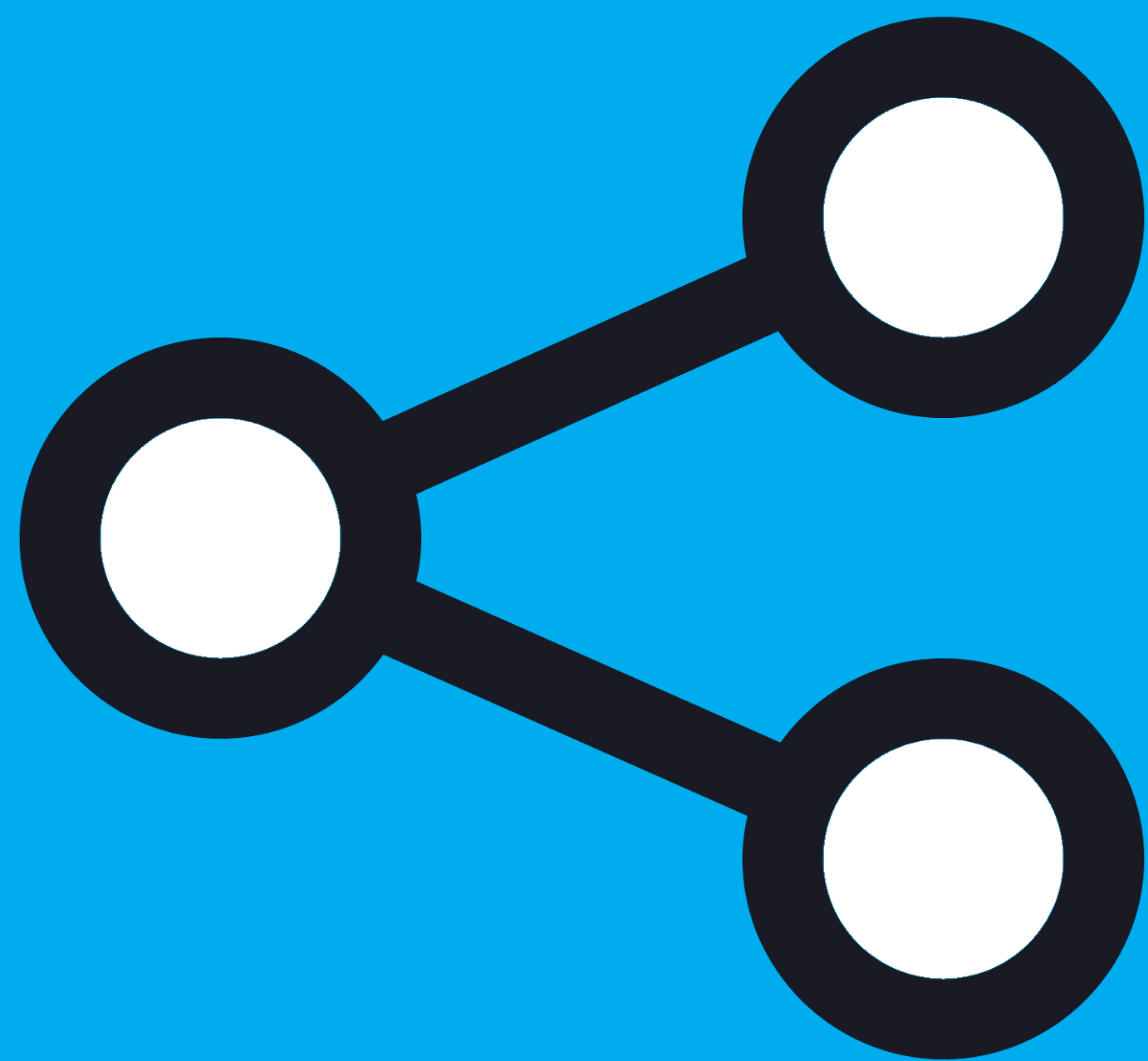
A terminal window with a white title bar containing three colored window control buttons (red, yellow, green). The terminal content is on a black background with white text. It shows a shell prompt '\$' followed by the command 'sudo rm -rf /'.

```
$ sudo rm -rf /
```

# MASTER

## the terminal





**GET**  
**your work**  
**out there**

**THANK YOU!**

 **@sebadzia**

# IMAGES

[Title image](#) (Flickr Commons)

[Vault Boy in a suit](#) (Nukapedia, CC-BY-SA)

[Twitter logo](#) (Flaticon free license)

[Banana clipart](#) (CC0 1.0)

[Pretzel clipart](#) (Freepik free license)

[Vault Boy with a wrench](#) (Nukapedia, CC-BY-SA)

[Tool icon](#) (Flaticon free license)

[Gear icon](#) (public domain)

[Jane Austen silhouette](#) (Wikimedia commons)

[Pride and Prejudice](#) (public domain)

[Shakespeare](#) (CC BY 3.0)

[Bison](#) (CC0 1.0)

[Thumb up Vault Boy](#) (Nukapedia, CC-BY-SA)

[Impossible trident](#) (Wikipedia, public domain)

[Black box](#) (CC0 1.0)

[Share icon](#) (CC BY-ND 3.0)

[Swiss army knife](#)

[Holmes silhouette](#) (public domain)

[IPI PAN logo](#) (public domain)

[SpaCy logo](#) (CC BY-SA 4.0)