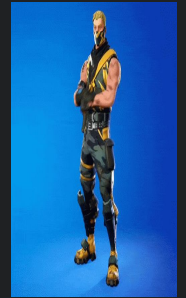




Starting Soon...



# NLP for Beginners:

Build your own summarization model  
for the Talk Tuah Podcast

Benjamin Tsang & Taha Sarfraz

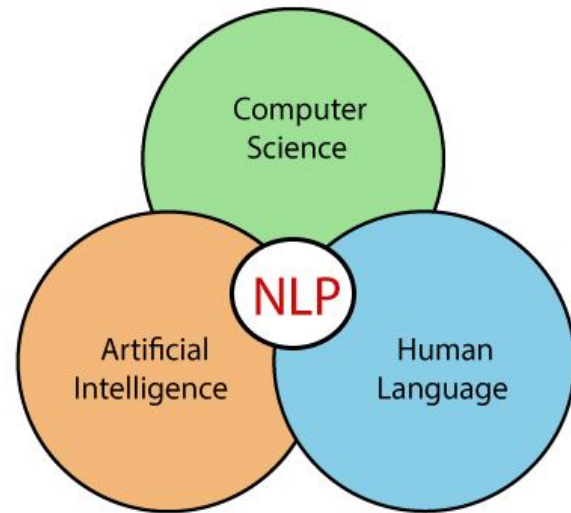


```
...org = filterByOrg ? study.lead_organization == filterByOrg : true  
...status = filterByStatus ? study.status === filterByStatus : true  
...searchStatus) }
```

```
function filterStudies({ studies, filterByOrg  
...lines = studies.filter(study  
...})
```

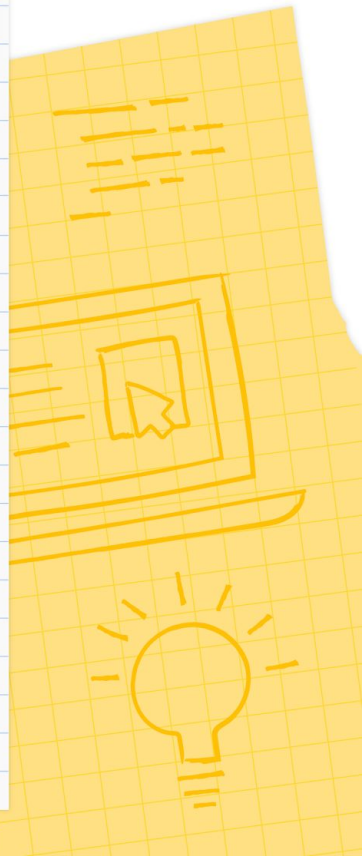
```
children: [  
  icon, color: color  
  container(  
    margin: const EdgeIns  
    child:  
      label  
      style
```

# What is NLP?



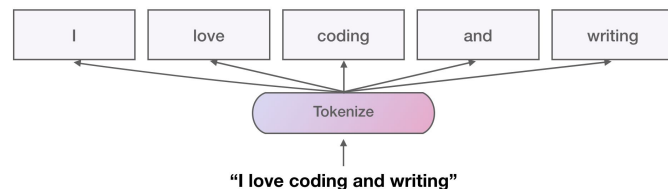
# Building a simple extractive summarization model

- Open [Google Colab Notebook](#)
- Make a copy
- Upload Talk Tuah [transcript](#)
- Follow along, run code for the model
- We will explain concepts at every step!



# Tokenization

- Breaking down text into smaller units - **tokens**
- Sentence, word, subword (coding → “code” + “ing”), character
- Breaking down into smaller pieces to analyze



```
lookup.KeyValue  
f.constant(['em  
=tf.constant([C  
.lookup.StaticV  
_buckets=5)
```



# Feature Extraction

- Word Frequency
- Term Frequency - Inverse Document Frequency (tf-idf)
  - Accounts for the fact that some words are common in every document
- Sentence length
- Positional Importance (starting, ending sentences usually are more important)

$$TF(t, d) = \frac{\text{number of times } t \text{ appears in } d}{\text{total number of terms in } d}$$

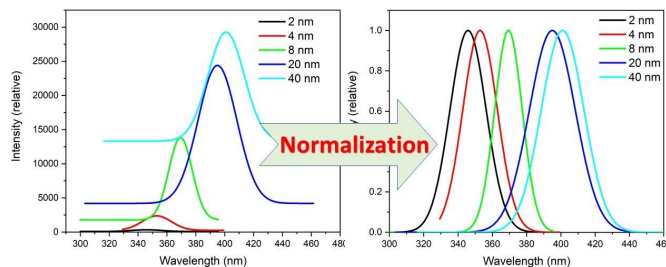
$$IDF(t) = \log \frac{N}{1 + df}$$

$$TF - IDF(t, d) = TF(t, d) * IDF(t)$$



# Sentence Selection

- Normalize, so different scales are comparable
- Combine the different scores
- Select the top N sentences based on combined score.
- This gives us our summary!



```
lookup.KeyValue  
f.constant(['em  
=tf.constant([@  
.lookup.StaticV  
_buckets=5)
```



# Post-processing

- Remove redundancies
- Grammar and punctuation checks
- Simplify or refine overly complex sentences

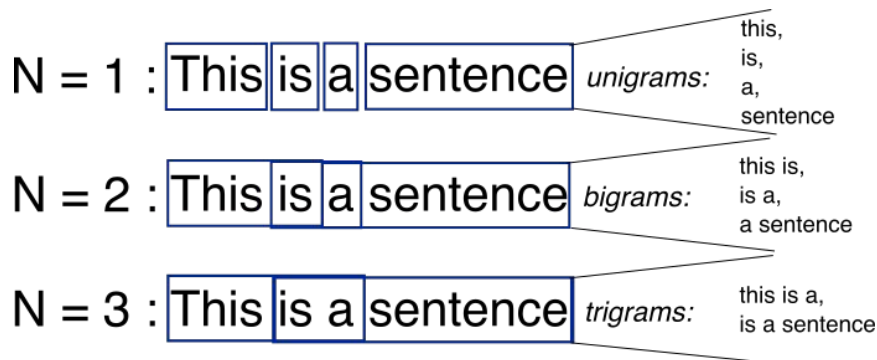
```
lookup.KeyValue  
f.constant(['em  
=tf.constant([G  
.lookup.StaticV  
_buckets=5)
```





# Evaluating the extractive summarization model

- Compression ratio
- Precision
- Recall
- ROUGE-N (overlap of n-grams)
- Intuition



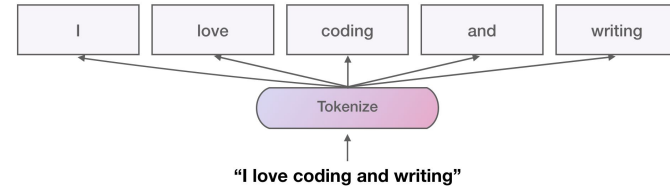
# Limitations of extractive summarization models

- Relies on existing sentences, no rephrasing
- Could be less coherent
- Lack of understanding
- Bias towards sentence position and length
- Potentially meaningless sentences



# Hugging Face

- What is hugging face
  - Pretrained models that you can use.
  - platform for NLP and machine learning (transformers, datasets, spaces, etc.).
- Why is it important
  - Open-source library powering NLP
  - Common models like BART, T5, or Pegasus.



# Use Cases and Applications

- Where Text Summarization is Used:
  - Media: Summarizing news articles or podcasts.
  - Customer Service: Condensing support conversations.
  - Legal: Creating summaries of long legal documents.
  - Education: Summarizing lecture notes or textbooks

```
lookup.KeyValue  
f.constant(['en  
=tf.constant([G  
.lookup.StaticV  
_buckets=5)
```



# Example:



```
from transformers import pipeline

summarizer = pipeline("summarization", model="facebook/bart-large-cnn")
text = " "
summary = summarizer(text, max_length=130, min_length=30, do_sample=False)
print(summary[0]['summary_text'])
```



# Example (2):

```
summarizer = pipeline ("summarization", model="facebook/bart-large-cnn" )
```

- Other things you can do w/ the pipeline class
  - Sentiment-analysis
  - Summary
  - poem

```
lookup.KeyValue  
f.constant(['em  
=tf.constant([G  
.lookup.StaticV  
_buckets=5)
```

## Example (2):

```
summary = summarizer(text, max length=130, min length=30, do sample=False)
```

- do\_sample
  - do\_sample: If False, the model uses greedy decoding (choosing the most probable word at each step).
  - If True, it adds randomness, which can make the output more diverse.

```
lookup.KeyValue  
f.constant(['em  
=tf.constant([G  
.lookup.StaticV  
_buckets=5)
```



# Thank you for listening!



@euphoriakkkkk



```
Lookup.KeyValue  
f.constant(['em  
=tf.constant([@
```



Those who know: