

# **ESL TEXT CLASSIFIER**

**HOW TO FIND THE READING THAT'S RIGHT FOR YOU**



---

**“TO LEARN A LANGUAGE,  
YOU NEED A  
DREAM”**

**BY ME**

- **DATA**
- **REPETITION**
- **EXPERIENCE**
- **ACTION**
- **MOTIVATION**

**DATA**

**REPETITION**

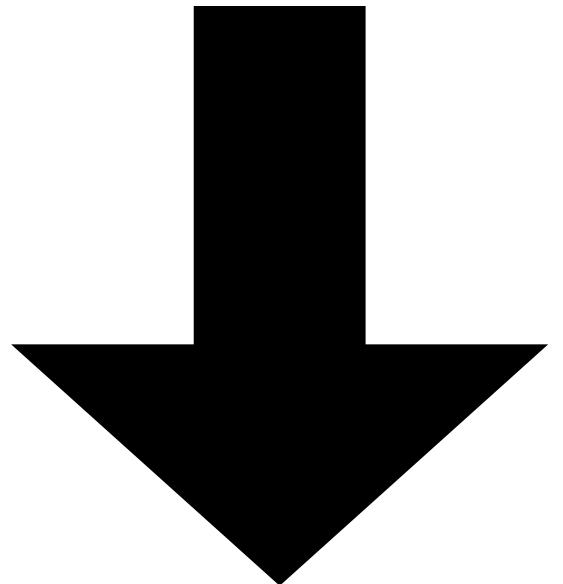
**EXPERIENCE**

**ACTION**

**MOTIVATION**

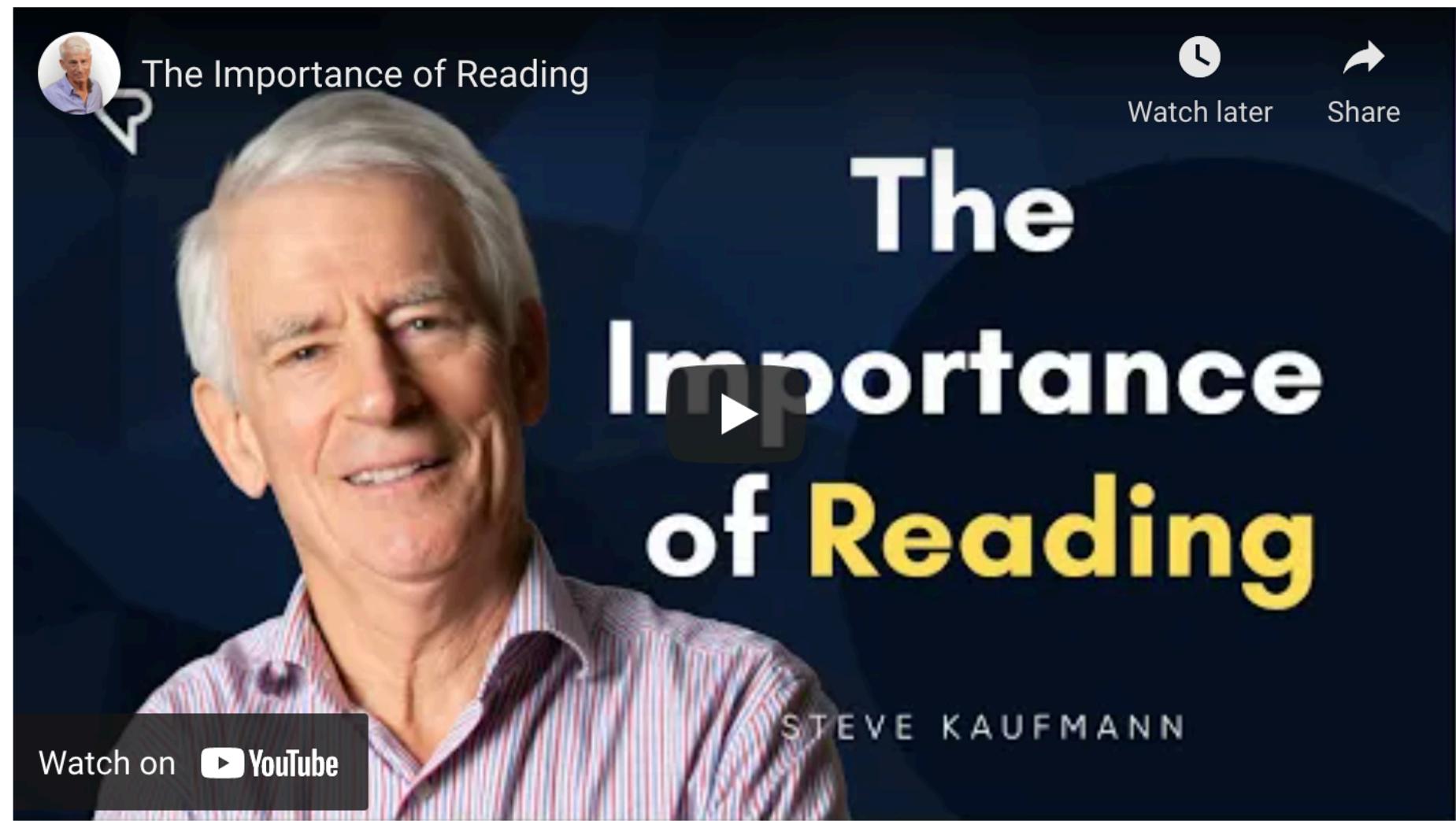
**Data**

**Input**



**Output**

## The Importance of Reading



# Data

## Reading

## Listening

# LEVEL APPROPRIATE READING

# LEVEL APPROPRIATE READING

Storytelling

Context

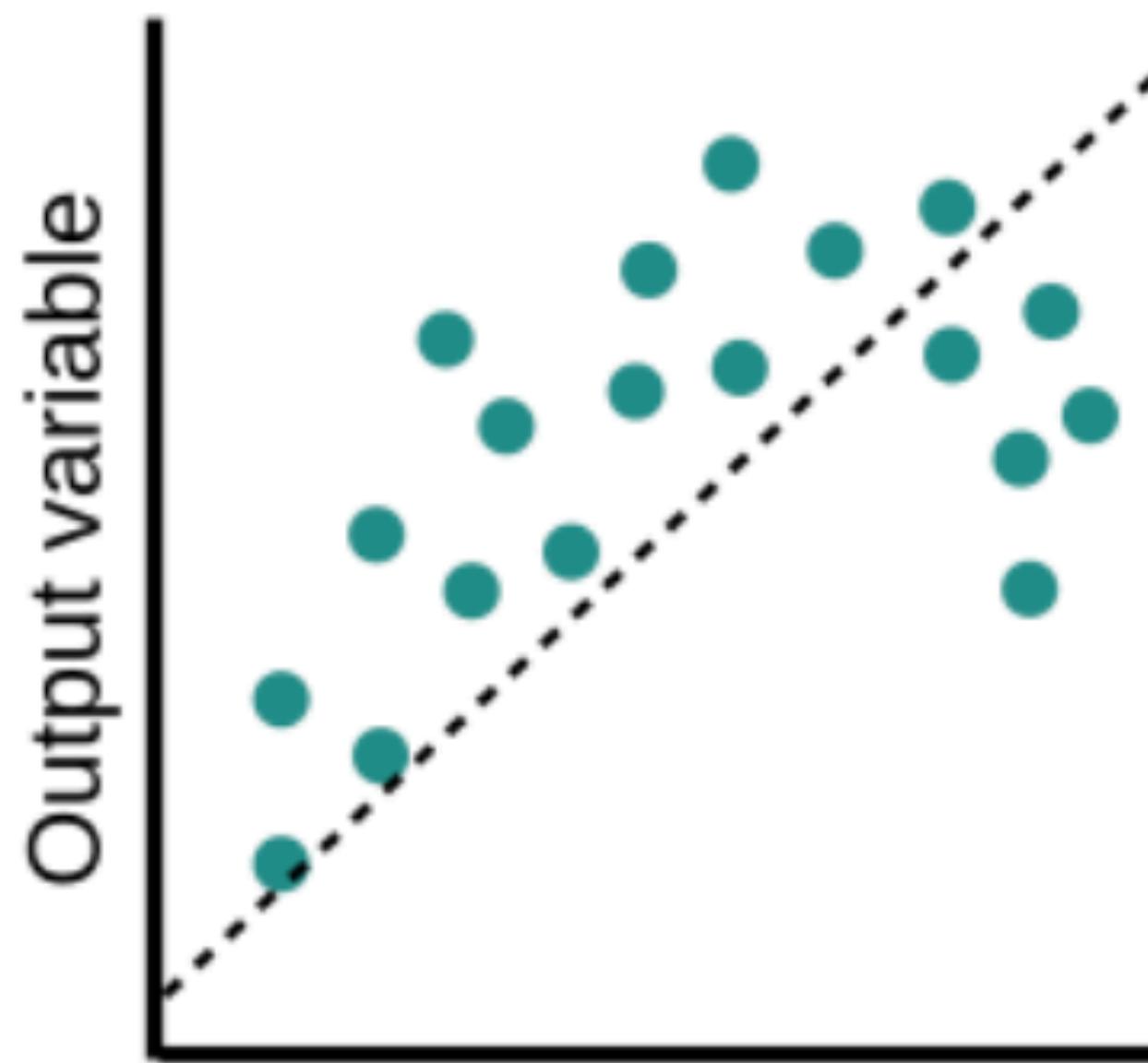
Data

Grammar

Vocabulary

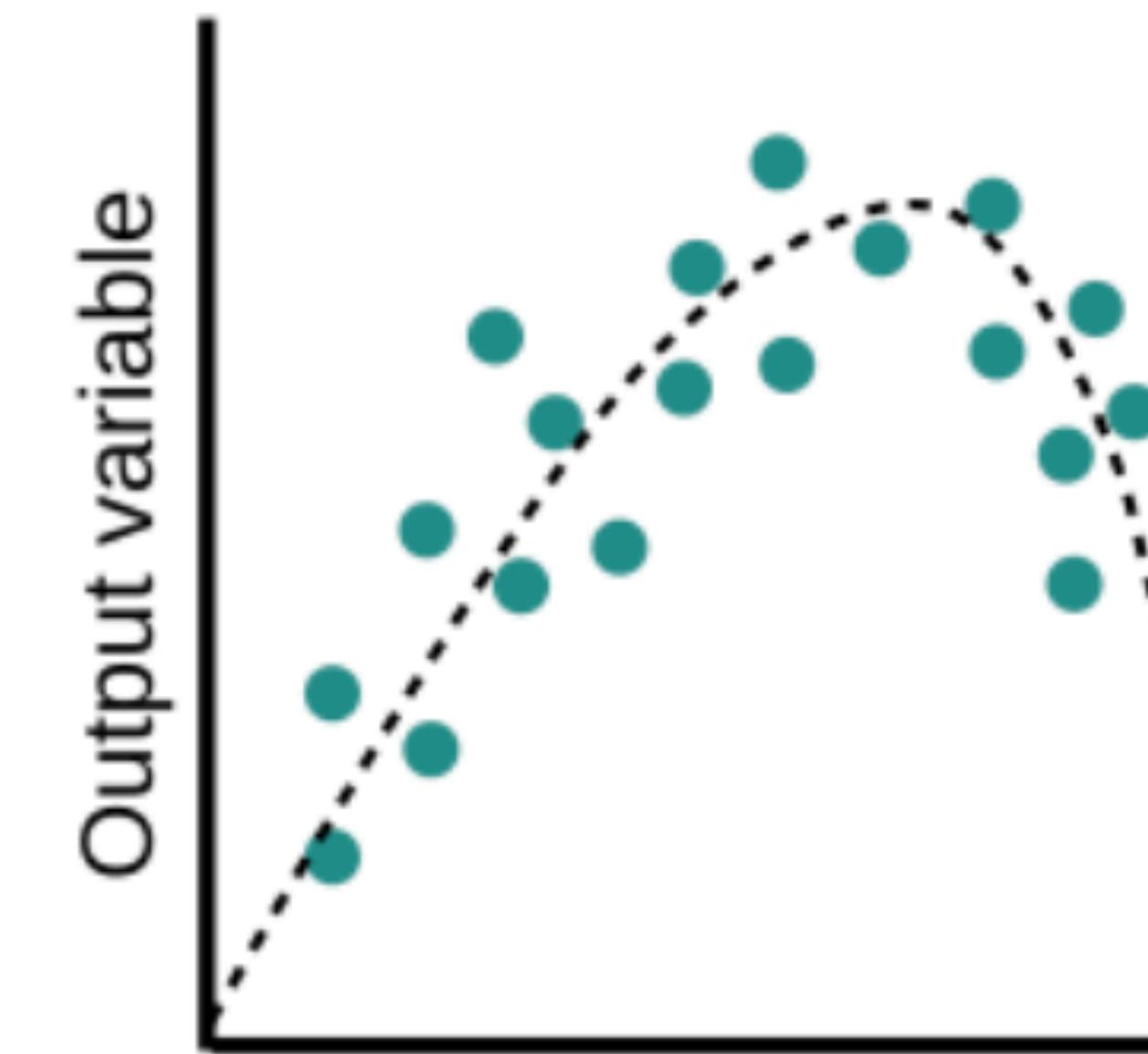
# ROBUST READING TEXT CLASSIFIER

Underfit



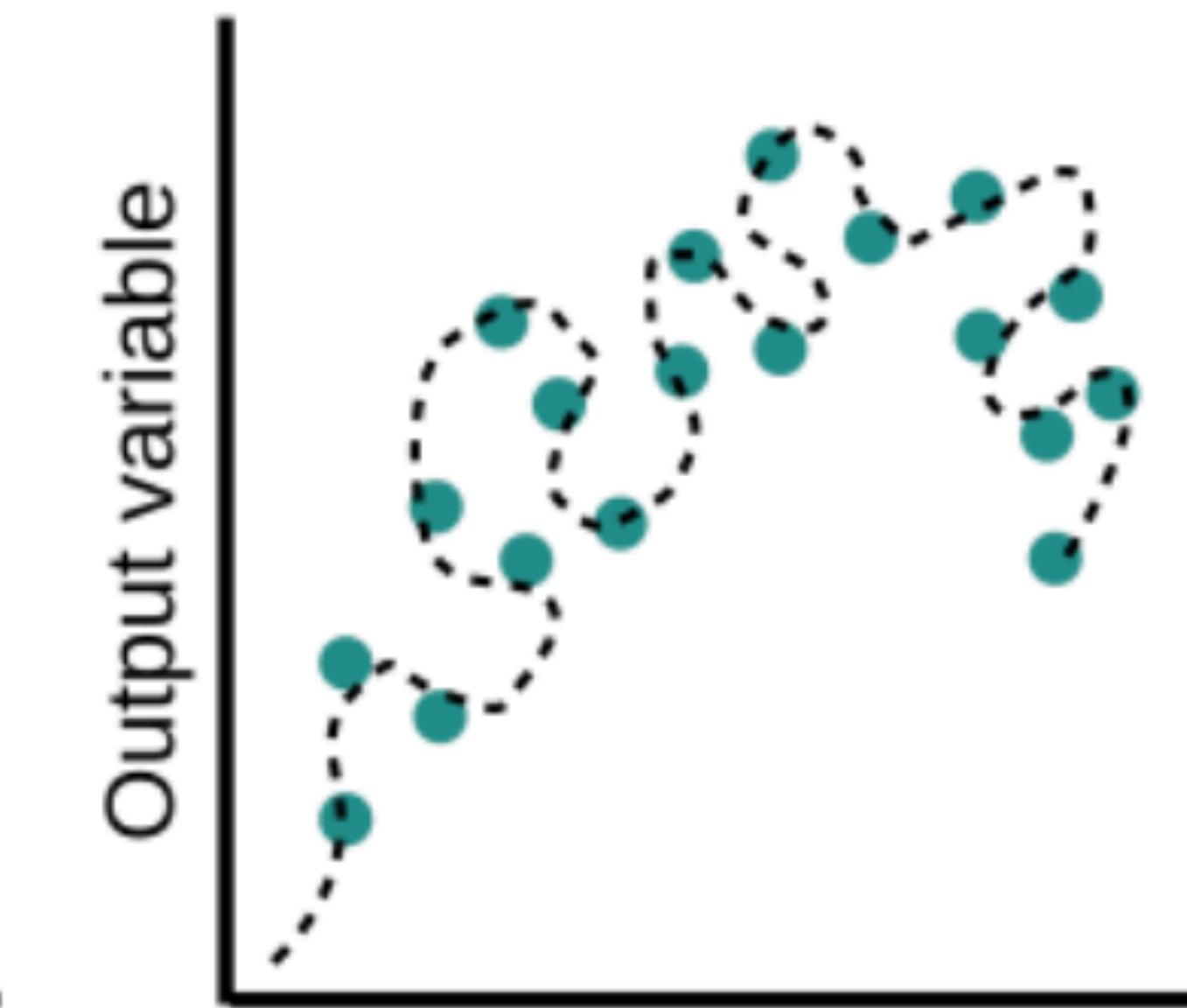
I hardly understand anything!

Optimal



It's just right

Overfit



I understand everything!

Pandas



# Data Acquisition

Natural Language  
Analyses with NLTK



# Data Cleaning

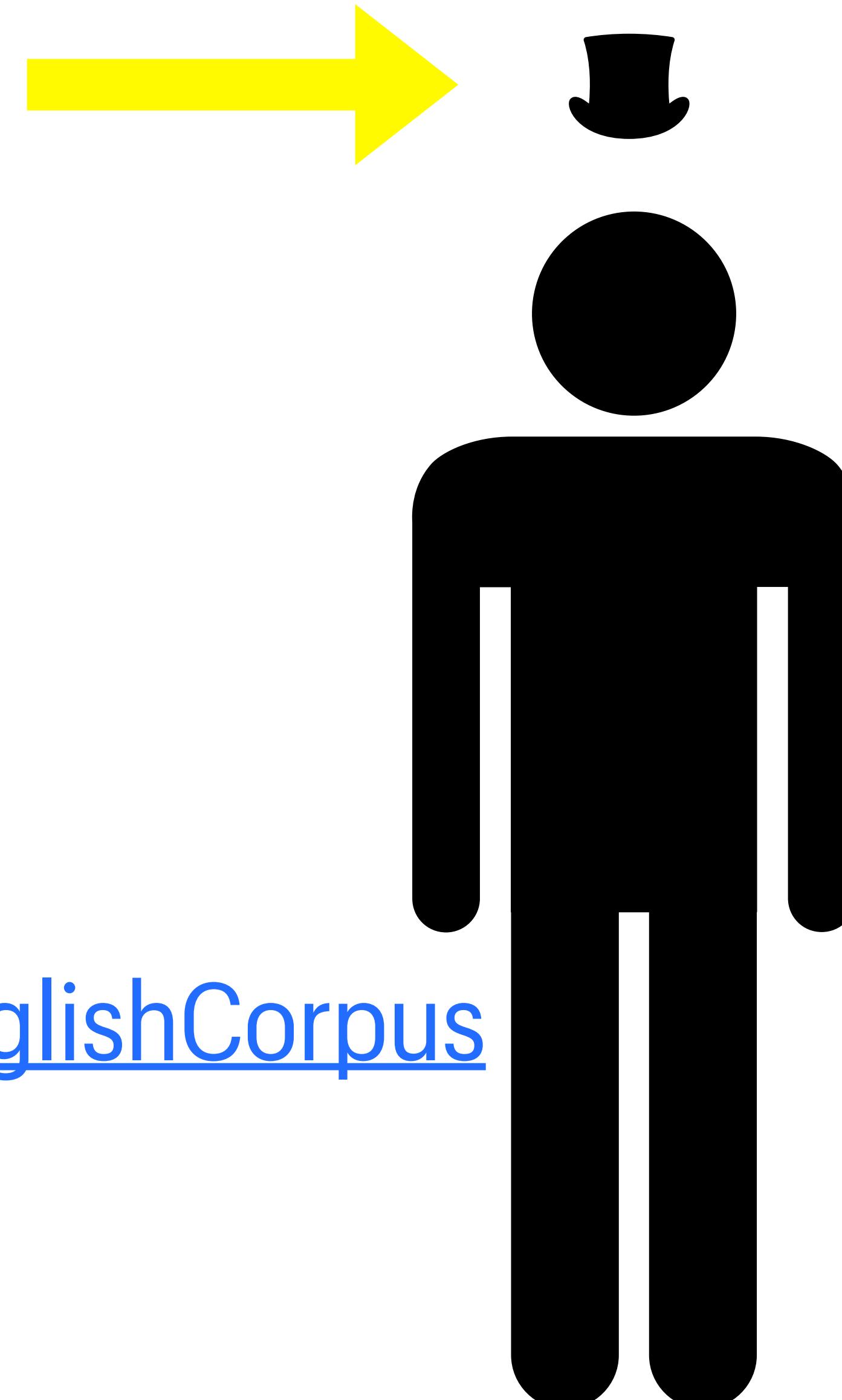
K Keras

Simple. Flexible. Powerful.

# Modeling

# Data Sets:

- Cambridge Readability Data Set



<https://ilexir.co.uk/datasets/index.html>

- One Stop English Corpus

<https://github.com/nishkalavallabhi/OneStopEnglishCorpus>

# grammar

A circular word cloud centered around the word "complexity". The words are arranged in a circle, with some words having associated smaller text or arrows pointing towards them. The words include:

- Complexity (large, central)
- Learning
- Neural
- Hets
- GFR
- Words
- Win
- Machine
- Notes
- Vocabulary
- Five
- Levels
- Three
- Robust
- Wanted
- Experiment
- Level
- Capture
- Text
- Machine
- Networks
- Fail
- Make
- Five
- Levels
- Three
- Robust
- Wanted
- Experiment
- Level
- Capture
- Text
- Machine
- Networks
- Fail
- Make
- Success
- List
- Determine
- Use
- Go
- Laugh
- Classical
- Model

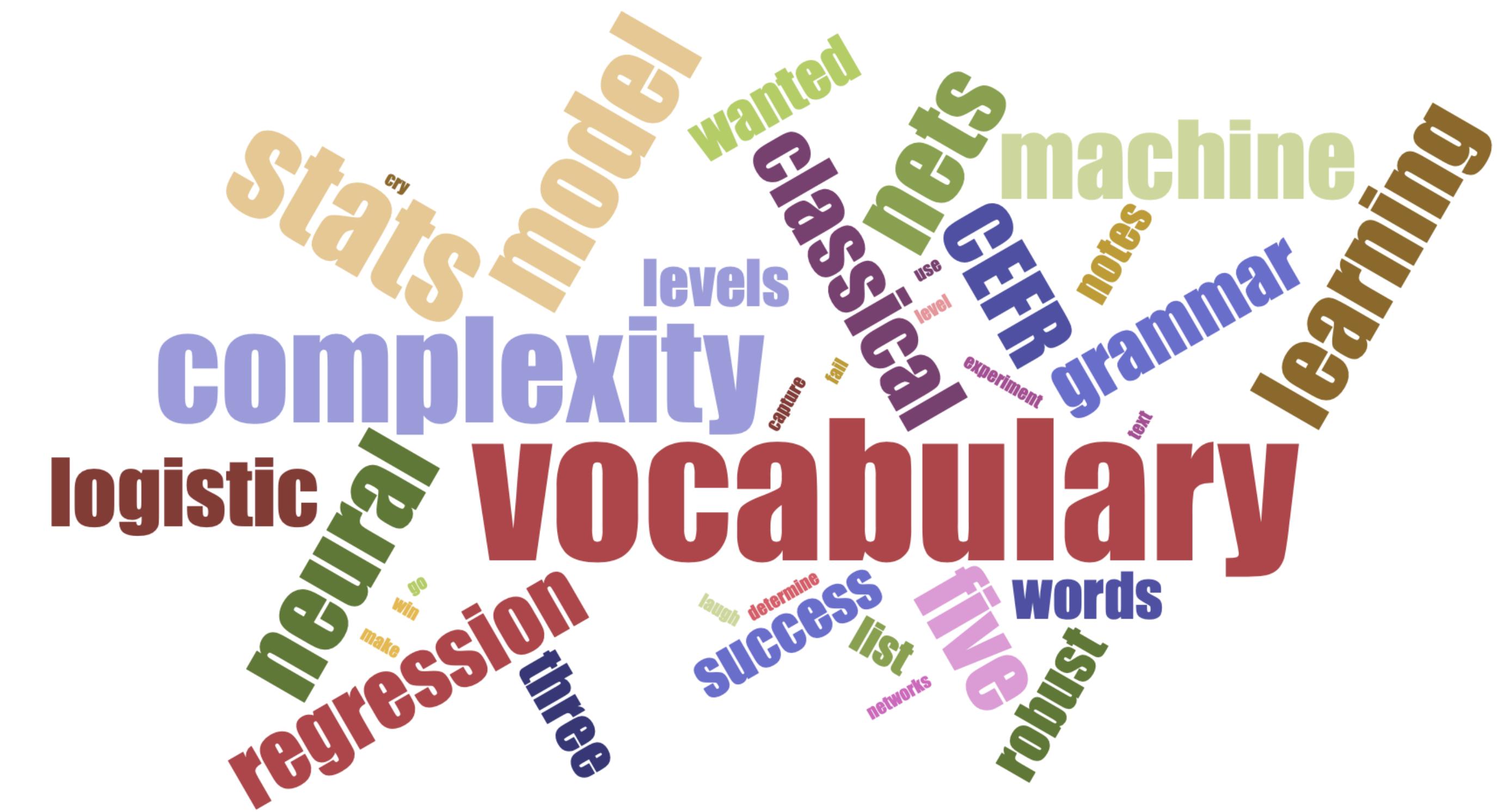
A circular word cloud centered on the word "vocabulary". The words are arranged in a circle, with their sizes and colors varying. The central word is "vocabulary" in large red font. Other prominent words include "logistic" (dark red), "neural" (green), "regression" (red), "three" (blue), "success" (purple), "list" (green), "five" (pink), "words" (blue), "robust" (green), "text" (pink), "grammar" (blue), "notes" (yellow), "CEFR" (blue), "levels" (purple), "capture" (pink), "tail" (pink), "experiment" (purple), "use" (pink), "wanted" (green), "classical" (purple), "nets" (green), "machine" (green), "learning" (brown), "complexity" (purple), "stats" (yellow), "cry" (pink), and "model" (yellow).

# grammar

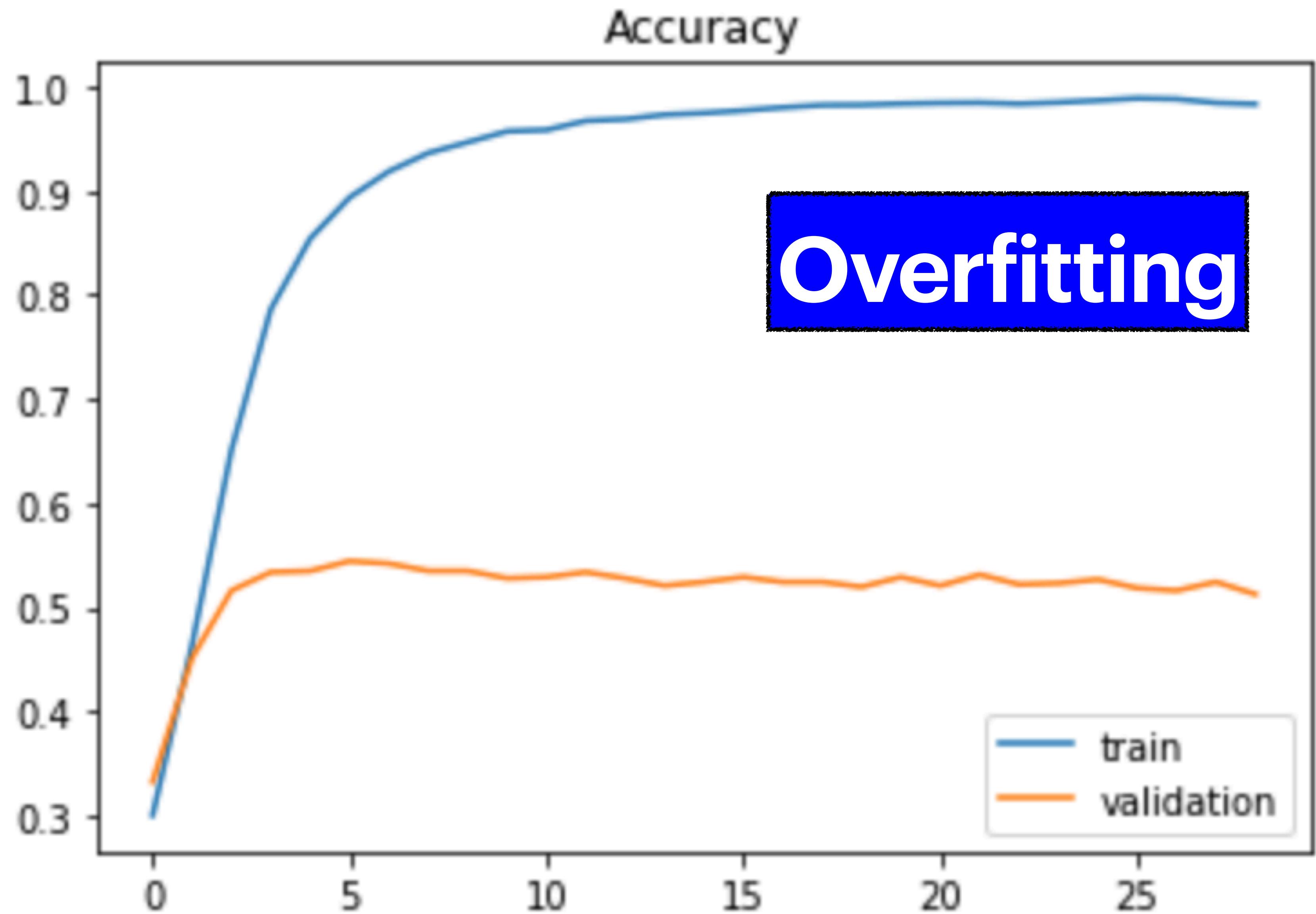


Neural Networks  
with tokenized  
vectorization

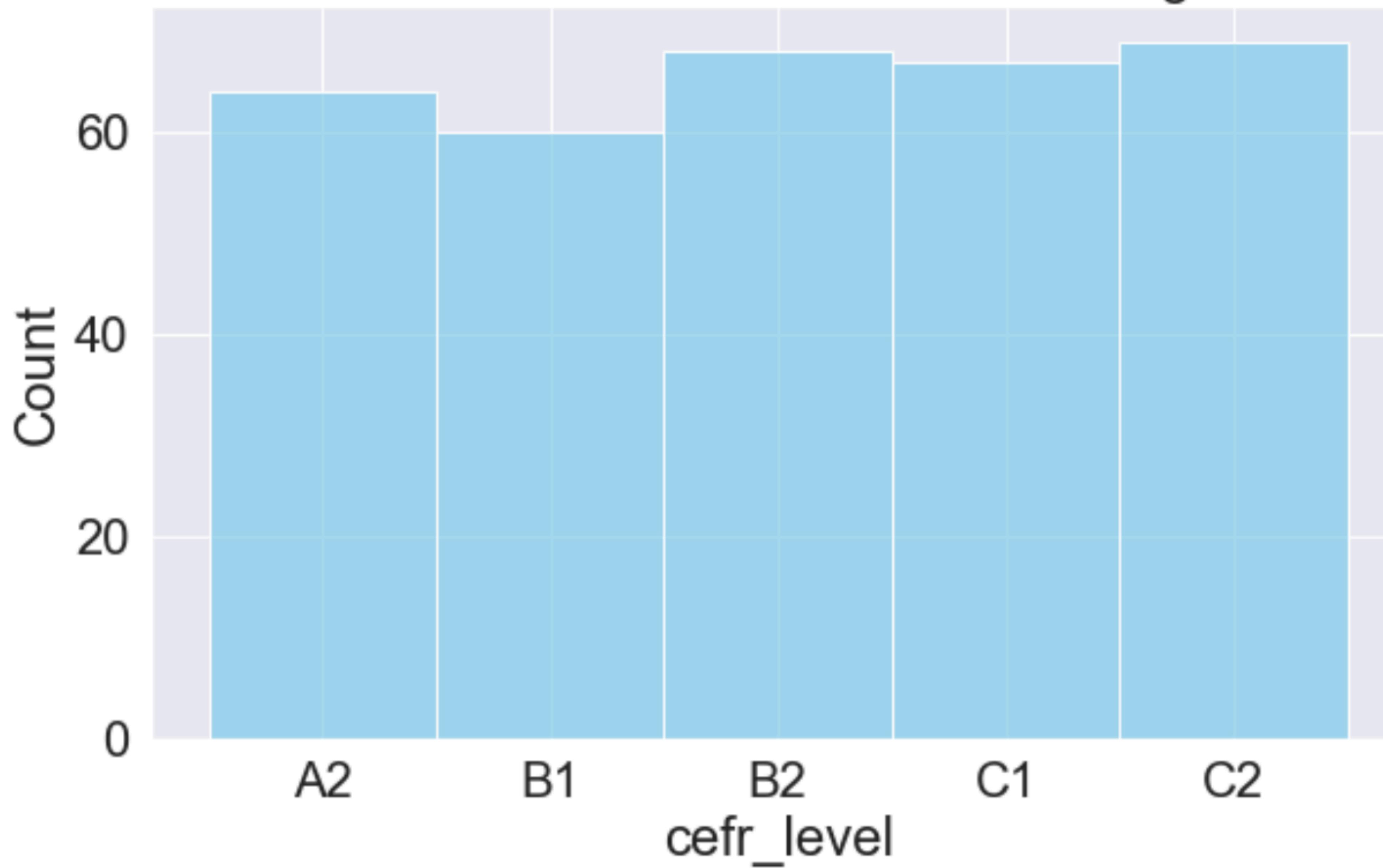
Classical Machine  
Learning with statistical  
vectorization



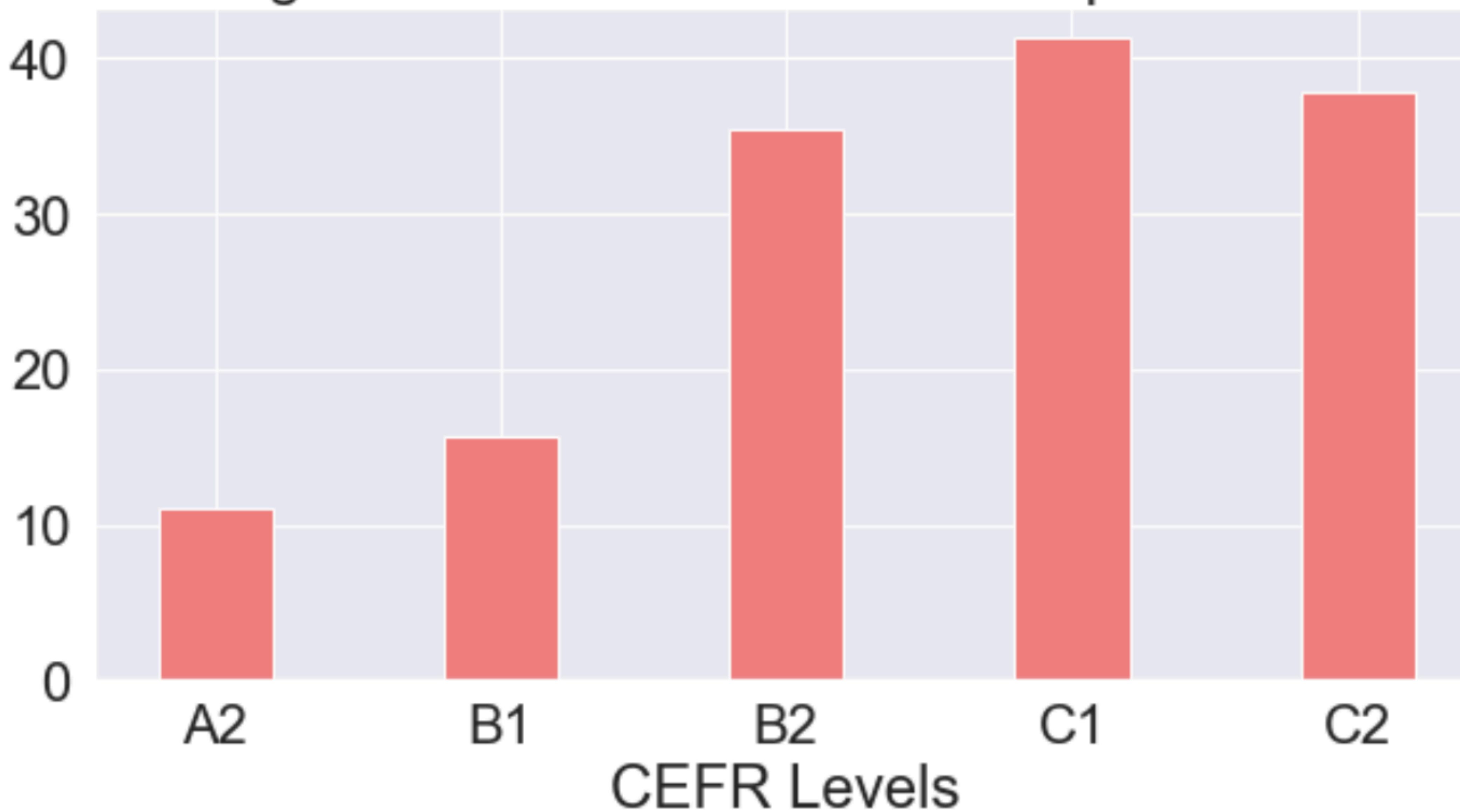
**LSTM**  
**Model**  
with texts  
broken up  
piece-  
wise by  
sentence



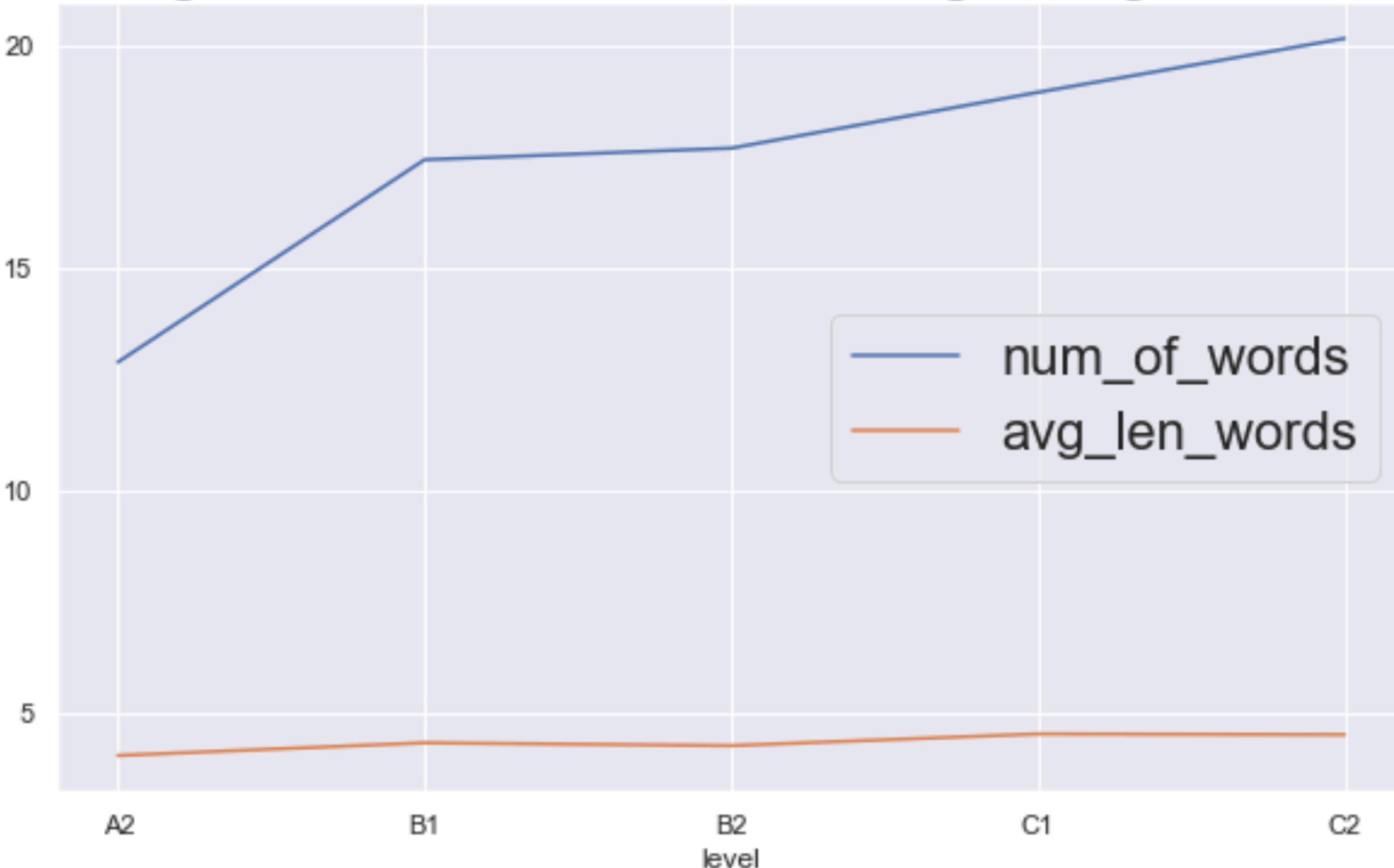
# Distribution of CEFR Level Readings



# Average Total Number of Sentences per Document



# Average Number of Words vs Average Length of Words

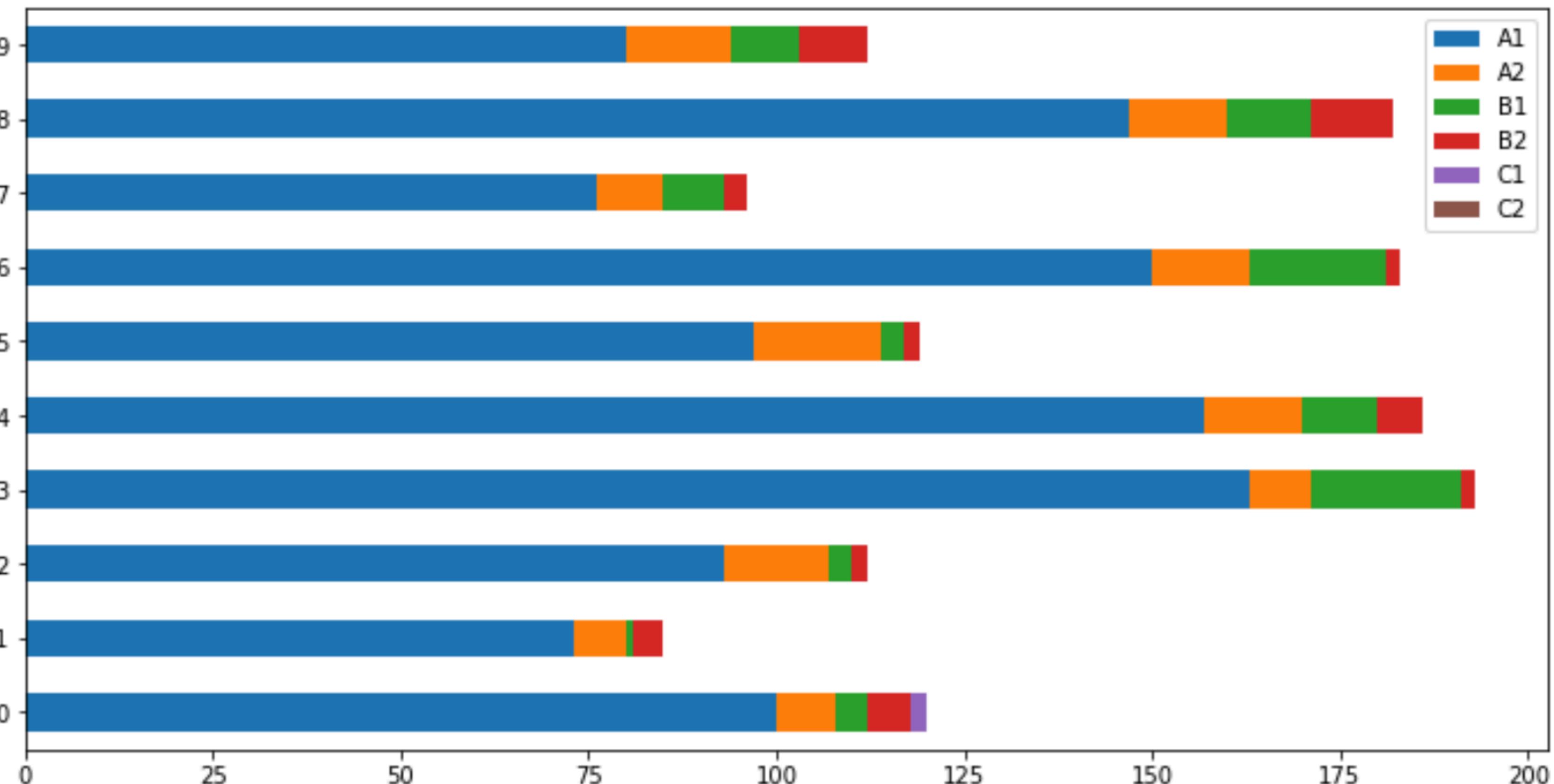


# Two Other Features

**CEFR WORD  
LEVEL  
COUNTS**

**PARTS OF  
SPEECH  
COUNTS**

CEFR Word Level Count Distribution for First Ten Documents



# Models: Logistic Regression: One versus All /Gradient Boasting Classifier

Training

Accuracy:

**0.93**

Validation

Accuracy:

**0.79**

**Overfitting**

Training

Accuracy:

**1.0**

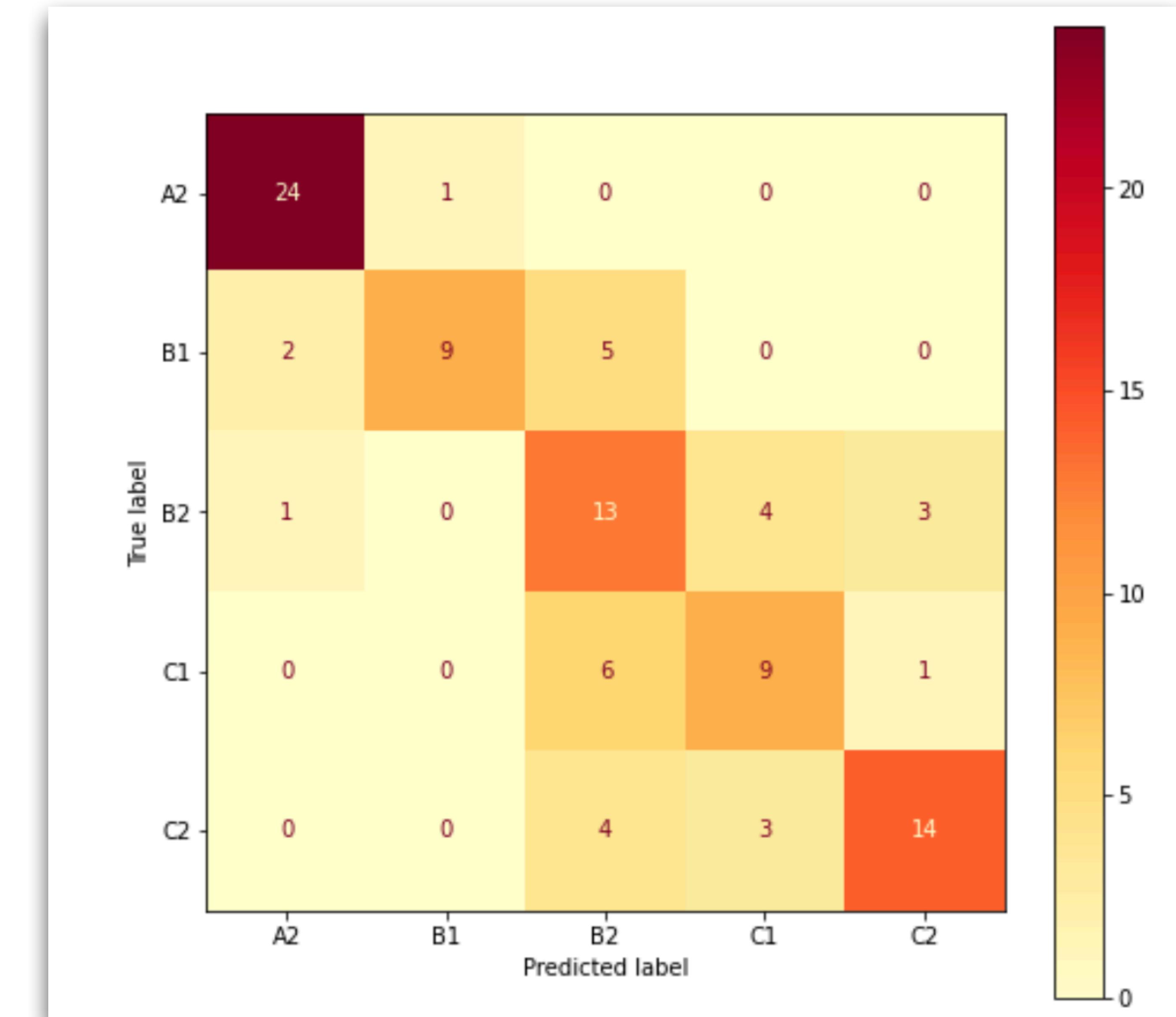
Validation

Accuracy:

**0.8**

# Best Model: Multinomial Naive Bayes

Training  
Accuracy:  
**0.72**  
Validation  
Accuracy:  
**0.70**

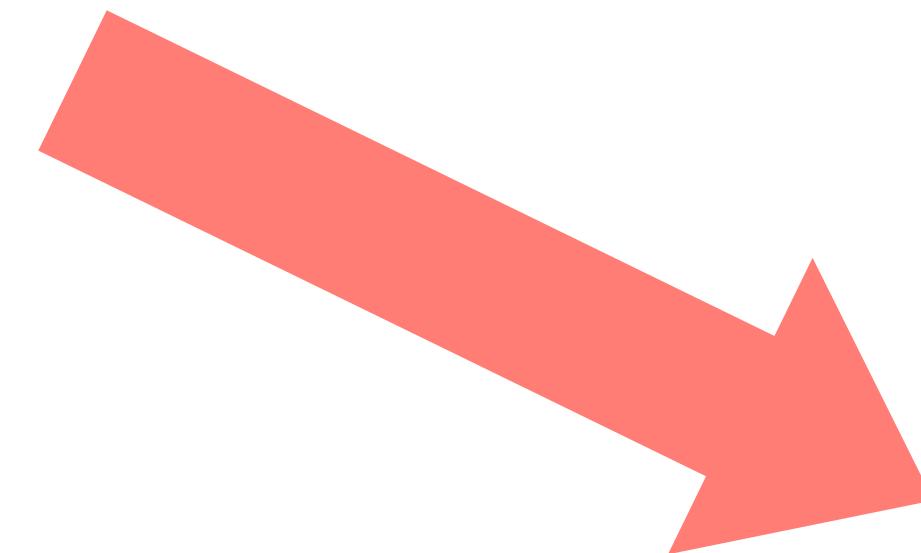


# Challenges:

## Lack of Data

## Use Neural Networks

## Build an App



```
pipe_nb.fit(X_train, y_train)
print("Training Acc.: ", pipe_nb.score(X_train, y_train))
print("Valid Acc.: ", pipe_nb.score(X_test, y_test))
```

```
Training Acc.: 0.7205240174672489
Valid Acc.: 0.6969696969696969
```

```
test_pred = pipe_nb.predict(X_test)
```

```
list_of_keys = cefr_pos_df.columns.tolist()
twilight_sample = pp.process_sample(twilight, list_of_keys,
```

```
pipe_nb.predict([twilight_sample])
```

```
array(['B2'], dtype='<U2')
```

A2

B1

B2

C1

C2

## CEFR Level: CEF Level B2

(Intermediate)  
IELTS Level 5-6

### Suggested vocabulary:

deep-voiced  
aback  
implied  
more—that  
prohibited  
manner  
to ignore

number of words: 73

average sentence length: 15

average word length: 4.6

word complexity: 1538

[DEFINITIONS](#)

[NEW TEXT](#)

[WORD STATS](#)

---

**THANK YOU FOR  
LISTENING**

**I HOPE TO HEAR FROM YOU SOON**

