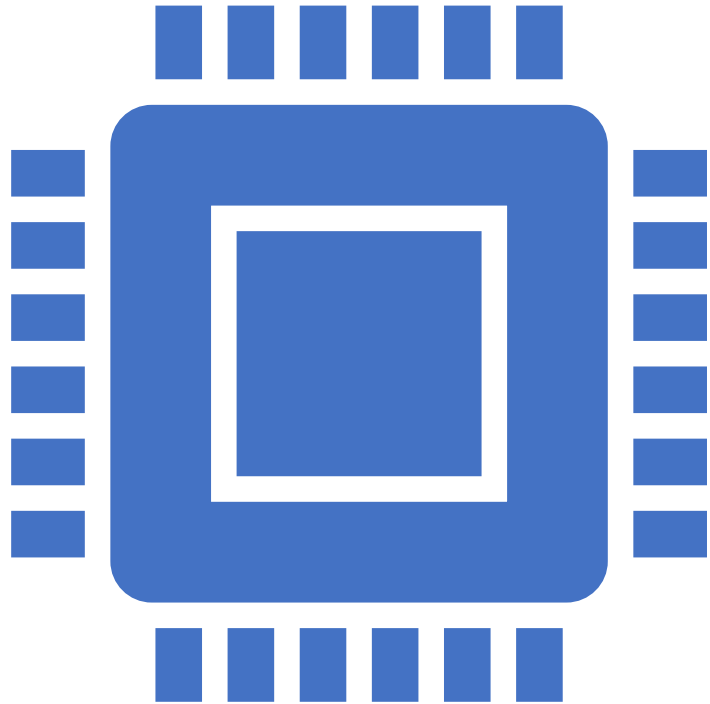# NLP - Basic units of processing

Liad Magen – Asigmo Spring 2022

Basic units of
processing

# What do we process?

- A Collection of Documents
    - By date or domain
- Document
- Section
    - Subsection
- Paragraph
- Sentence
    - Phrases
- Word
    - Characters / Morphemes

# What do we process?

- A Collection of Documents
    - By date or domain
- **Document**
- Section
    - Subsection
- Paragraph
- Sentence
    - Phrases
- Word
    - Characters / Morphemes

# Document Processing

# What do we process?

- A Collection of Documents
    - By date or domain
- Document
- Section
    - Subsection
- Paragraph
- **Sentence**
    - Phrases
- Word
    - Characters / Morphemes

# What is a sentence?

Python:

```
sentences = text.split(".")
```

Sentence can also end with: **! : ; ? …**

Numbers: 0.99 // 1.234,56

Date : 28.02.2022

A '.' Can also represent part of a sentence:

- www.url.com
- I.B.M // E.E.G

Full stop, period.

# What is a word?

Sequence of characters?

Basic unit of meaning?

White-space **tokenization**?

# What is a word? (Clitics)

There's

Doesn't

John's

unlucky

# Chinese has no spaces!

# Compounds & Orthography

considered as one word

- Ice cream
- Web site
- New York
- New York-based

# Phrases

- Gave up
- Gave away
- Took a picture
- Took a brake
- Took a risk
- Took the toy apart

- German: Schalten Sie es ein (➜ einschalten)

# Reflexive verbs // Verb conjunction

- **French**:
  Tu te llaves

- **Portuguese**:
  Sente-se

- **Spanish**:
  Recuérdame / le dijeron

- **Albanian**:
  Laj duart

  Noun conjunction: duar -> duart

Case study:
Sentence Boundary
Detection

# Question:

When do we need (or don't need) to divide a document into sentences?

Not needed:

if it's already divided; Search Query; Working on a document-level (e.g., document classification)

What about translation – needed or not?

# Anaphora / Cataphora / Co-reference resolution

*Human*: What do we want?

*Computer*: Natural-Language Processing!

*Human*: When do we want it?

*Computer*: When do we want what?

# Sentence Boundary Detection

Python:

```python
sentences = text.split(".")
```

# Sentence Boundary Detection

However, if you picked up a certain book from the library, you can opt to start up your own biplane instead of boarding the zeppelin and forgo the hassle of having to navigate around the blimp's infrastructure. Flying the plane requires you to utilize the number pad (in a slightly less frustrating way than the fights) to dodge the Luftwaffe assault while Henry blasts the enemy planes. Depending on how many planes you destroy before you get shot down, you will have to face fewer roadblocks on the way to Alexandretta when you continue the rest of the journey by car. The guards you stop you at each roadblock act as in Brunwald; i.e., find out a way to talk or bribe your way around them, or just get to punching their faces in. A certain item can even help you skip over every single guard just by showing it to them!

# Sentence Boundary Detection

However, if you picked up a certain book from the library, you can opt to start up your own biplane instead of boarding the zeppelin and forgo the hassle of having to navigate around the blimp's infrastructure. Flying the plane requires you to utilize the number pad (in a slightly less frustrating way than the fights) to dodge the Luftwaffe assault while Henry blasts the enemy planes. Depending on how many planes you destroy before you get shot down, you will have to face fewer roadblocks on the way to Alexandretta when you continue the rest of the journey by car. The guards you stop you at each roadblock act as in Brunwald**; i.e.,** find out a way to talk or bribe your way around them, or just get to punching their faces in. A certain item can even help you skip over every single guard just by showing it to them!

# Sentence Boundary Detection

Python:

```python
sentences = text.split(". ")
```

# Sentence Boundary Detection

WASHINGTON — Former President George W. Bush called on Americans on Saturday to put aside partisan differences, heed the guidance of medical professionals and show empathy for those stricken by the coronavirus and the resulting economic devastation.

In a three-minute video message, Mr. Bush, who rarely speaks out on current events, struck a tone of unity that contrasted with the more combative approach taken at times by President Trump as the former president evoked the sense of national solidarity in the wake of the attacks of Sept. 11, 2001.

# Sentence Boundary Detection

Learning system, that **learns**, according to text features, which ones are attributed to sentence ends.

- Q: How to get annotated data?
- Q: Which kind of data to use?
- Q: Which clues could hint of end-of-sentences?
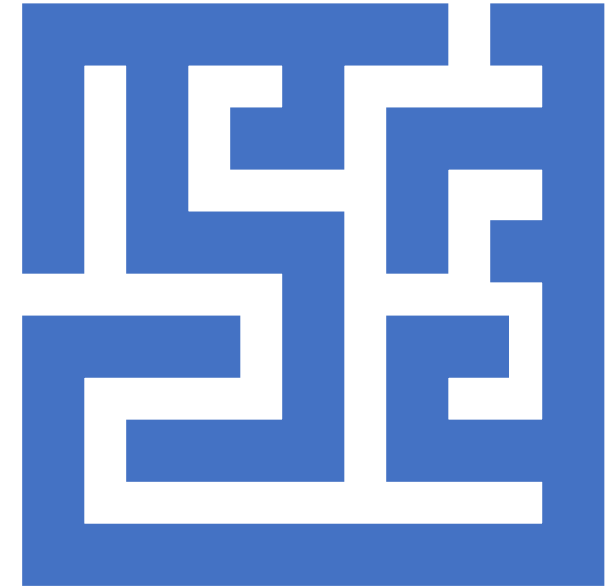
# Clues to identify end-of-sentence

- A dictionary of safe-words:
  *Mr. Dr. Ltd. Inc. U.S. e.g. i.e.*

- Hinting words in a windowed neighborhood

- An ensemble or a combination of clues:
  - *Capital letter after ". "*
    Though not perfect: *George W. Bush*

- Word length: short words are probably part of a word


ML Task: Weighting the hints as a learnable solution.

- Collect the data

- Annotate

- Define the 'hints' (a.k.a. *feature engineering*)

- Identify the problem: Binary classification
  - Decide which of the hints to use to determine end-of-sentence

- Train

- Evaluate

# Weighting the hints as a learnable solution

# Sentence Boundary Detection and the Problem with the U.S.

## Dan Gillick

Computer Science Division
University of California, Berkeley
dgillick@cs.berkeley.edu

## Worth Reading

# SOTA is still not 100% solved

It's not a good time. (It's never a good time).

In a quiet voice, he said "this will not work. I am quitting", and then he left the room.

# CODE REVIEW

Home

**Questions**

Tags

Users

Unanswered

Jobs

## Retrieving a substring from an exponentially growing string

Asked today    Active today    Viewed 68 times

8

1

For integer A, the base level of 1 is "XYZ" For subsequent levels, the levels become "X" + Level (A - 1) + "Y" + Level (A - 1) + "Z". So for level 2, the string would be "XXYZYXYZZ"

The objective is to return the substring from the Level string using the Start and End Index.

Example: If entering 2 3 7, it would Level 2, 3rd character to the 7th character and the result would be "YZYXY" from "XXYZYXYZZ"

The following constraints are given:

- $1 \leq$ Level K $\leq 50$,

- $1 \leq$ start $\leq$ end $\leq$ length of Level K String,

- $1 \leq$ end - start + 1 $\leq 100$.

I have written a brute force approach for this problem in Python as follows

```python
def my_substring():
    level = int(input())
    start_index = int(input()) - 1
    end_index = int(input())
    strings_list = [None] * level
```
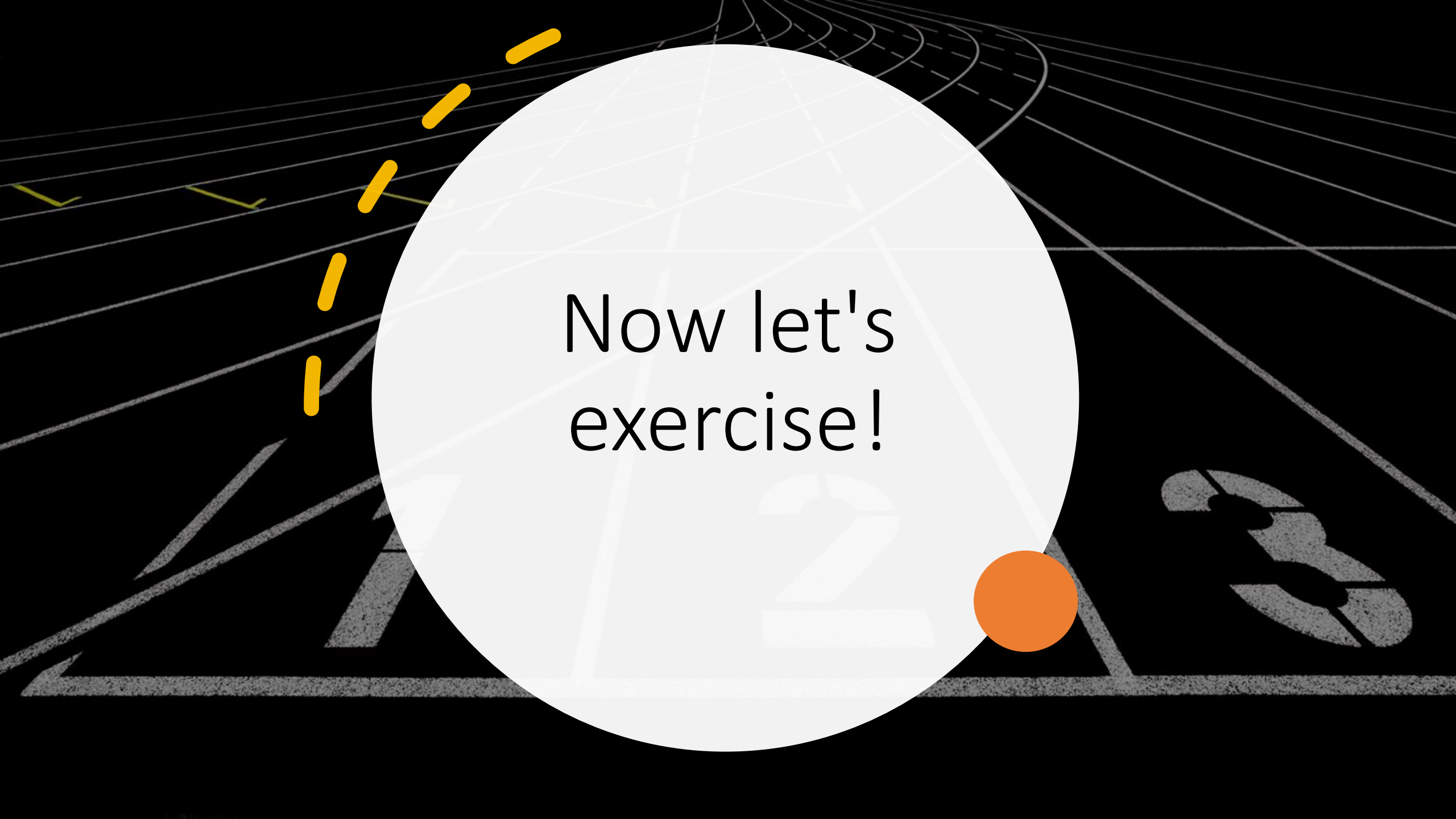
Real world example

# Pop quiz

Name the 5 levels of Layers of
Linguistic Knowledge
(and what do they mean)

- Phonology
- Morphology
- Syntax
- Semantics
- Pragmatics

Now let's exercise!