# AYLIEN

# Using NLP to understand textual content at scale

**Parsa Ghaffari**, CEO & Founder

**aylien.com** / @_aylien

# Introduction

Documents

Tweets

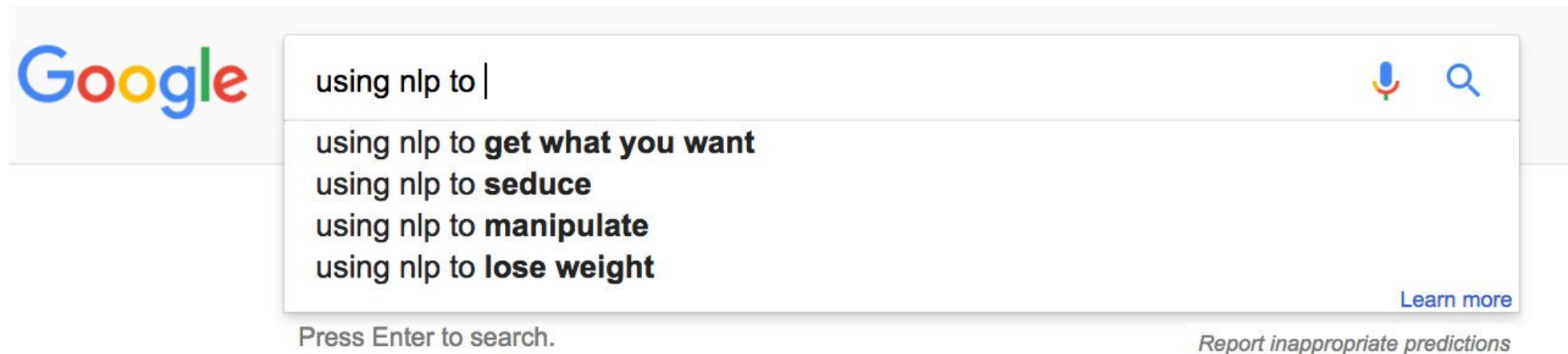Web pages
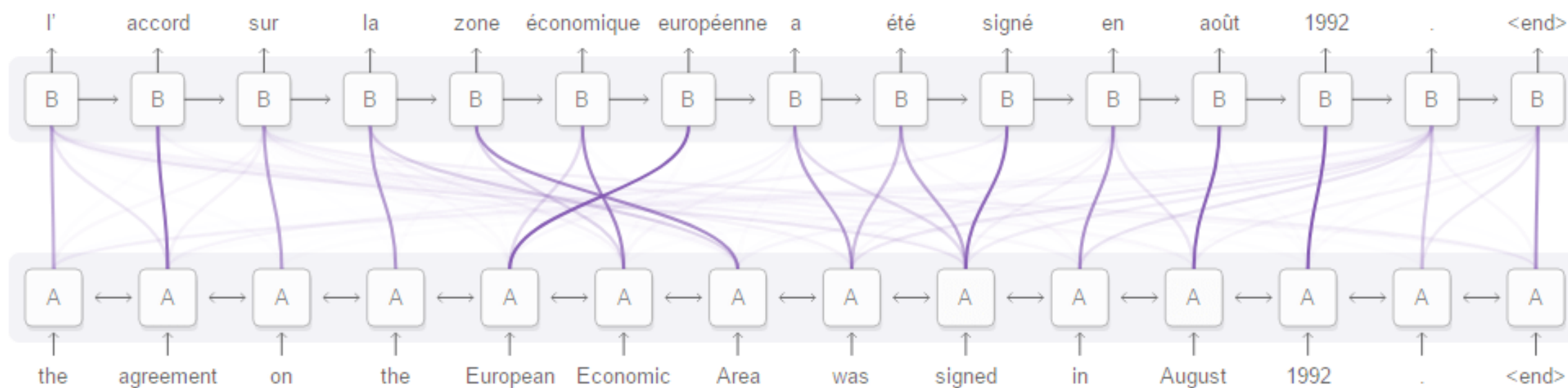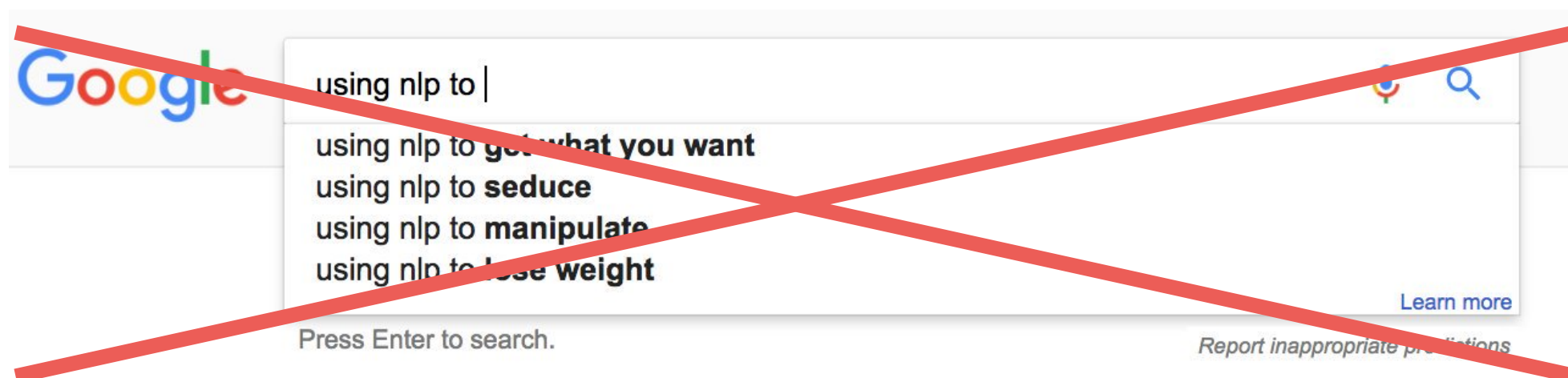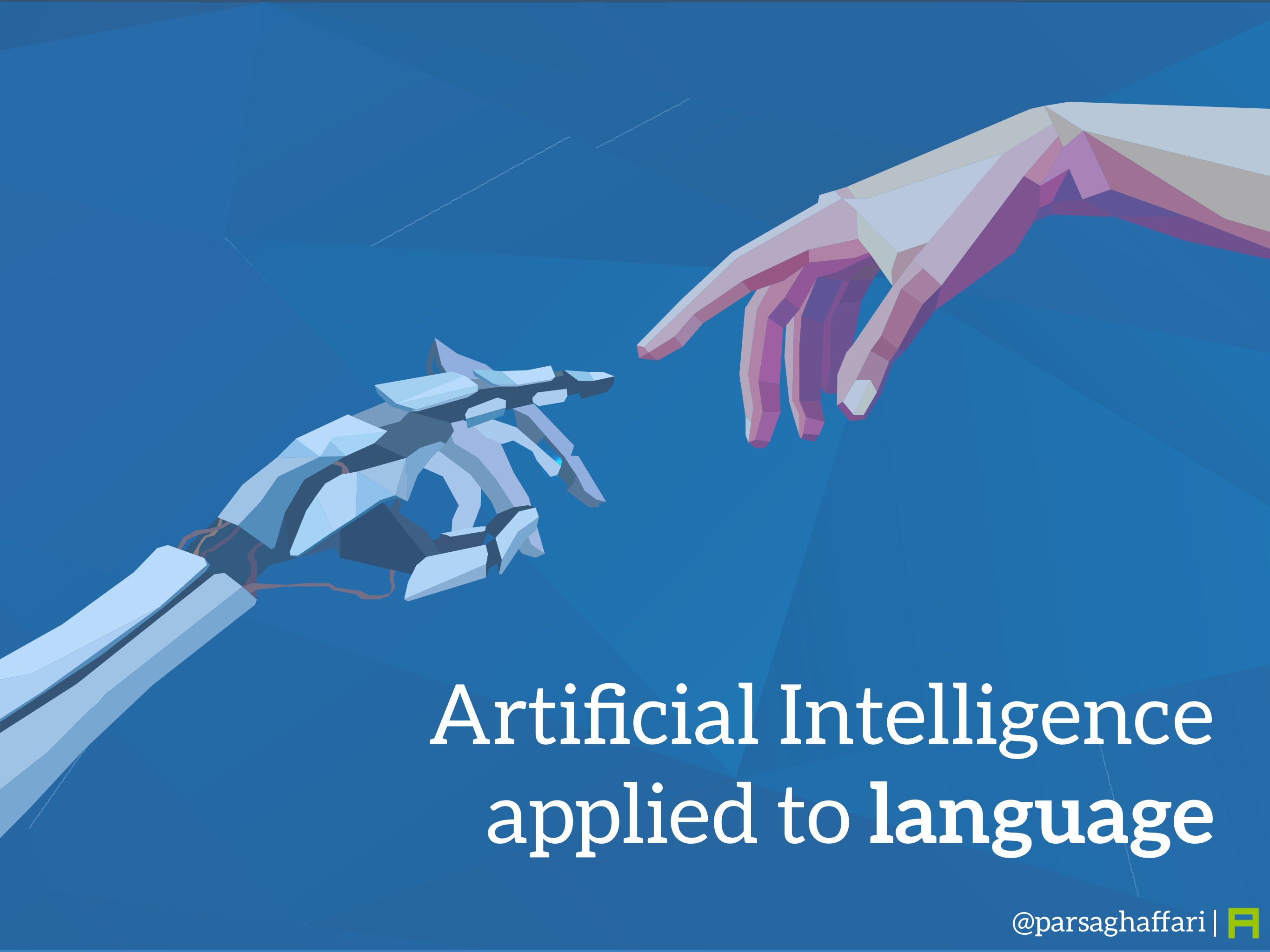
Document

Sentiment
Classification
Entities
Summary

AYLIEN is a platform for **aggregating** large volumes of textual content and **analyzing** and **understanding** it using **NLP.**

# What is NLP?

Google

using nlp to |

using nlp to **get what you want**
using nlp to **seduce**
using nlp to **manipulate**
using nlp to **lose weight**

Learn more

Press Enter to search.

Report inappropriate predictions

# What is NLP *not?*

Artificial Intelligence
applied to **language**

@parsaghaffari |

# Why NLP?

Language is a proxy to people's thoughts, emotions and feelings, and intentions

↓

We can understand people by understanding language

↓

Computers can help people accomplish things more efficiently, if they understand language

↓

Wildly applicable technology capable of creating $B's in value

@parsaghaffari |

# Challenges of NLP

# Challenge #1
The inherent complexity of language

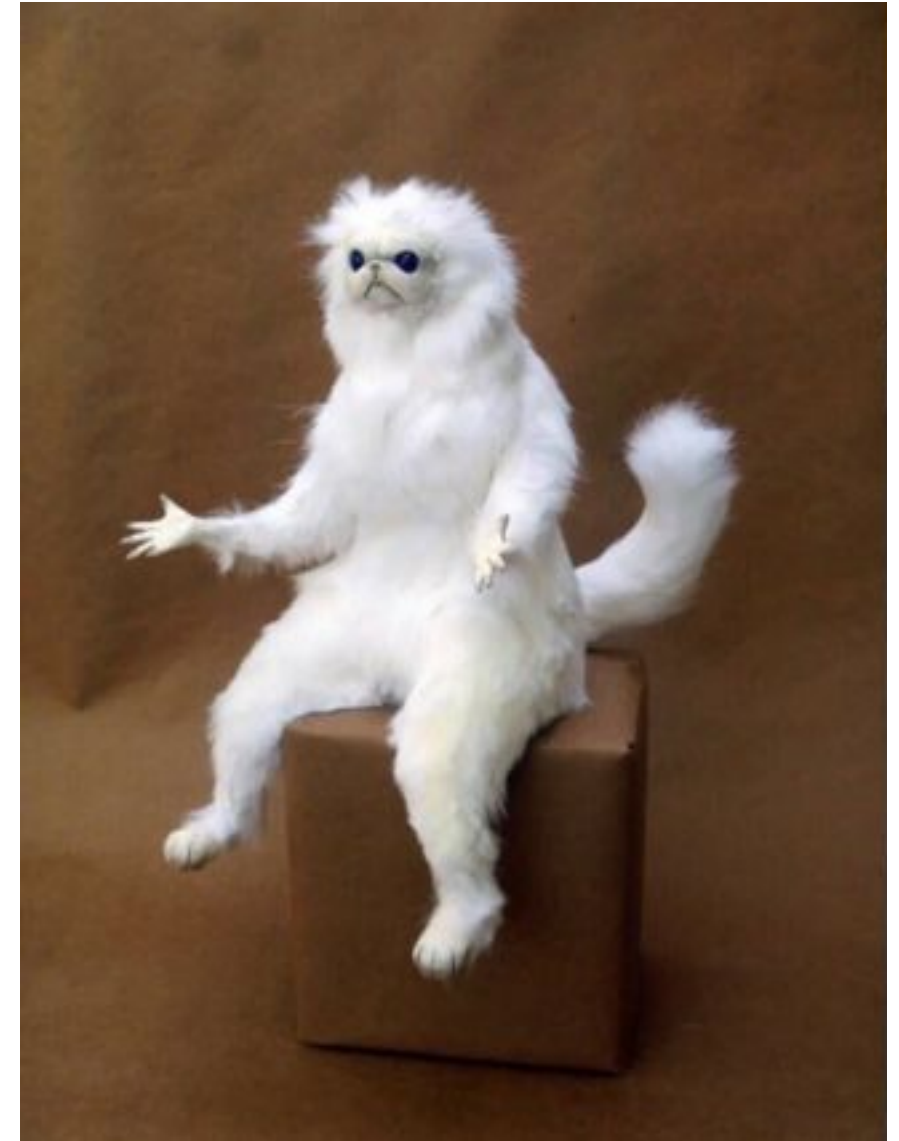What does it mean when someone says:
"I made her duck"?

How many distinct interpretations can you count?

# Challenge #1

The inherent complexity of language

"I made her duck" could be interpreted as:

- I cooked a duck for her
- I cooked a duck belonging to her
- I created a duck for her
- I created a duck that now belongs to her
- I caused her to lower her head
- I turned her into a duck (!)

@parsaghaffari |

# Complexity of Language

Adjectives:
"**useless**",
"**sad**"

Emotion:
**Frustration**

Tone:
**Negative**,
**Subjective**

"Literally ur facebook message app is useless, you only want it to increase profit. Please fix yourself. Its sad @facebook"

Product:
**Messenger App**

Language:
**English**, **Informal**

Organization:
**Facebook**

# Complexity of Language

And it's not just **English**!

- **German:** "Donaudampfschifffahrtsgesellschaftskapitän" (5 "words")
- **Chinese:** 50,000 different characters (2-3k to read a newspaper)
- **Japanese:** 3 writing systems
- **Thai:** Ambiguous word boundaries and sentence concepts
- **Slavic:** Different word forms depending on gender, case, tense

@parsaghaffari |

# Complexity of Language

**Machine Learning** is the most suitable toolbox available to us for dealing with this complexity.

However, there are still challenges:

1. Variety of tasks
2. Diversity of domains
3. Data preparation
4. Model training
5. Model evaluation
6. Workflow issues

# Variety of Tasks

"Literally ur **facebook message app** is **useless**, you only want it to increase profit. Please fix yourself. Its sad @facebook"

Sentiment Analysis (**negative**)

Emotion Detection (**frustration**)

Language Detection (**English**)

Entity Extraction/Linking (**Facebook**, **Facebook Messenger**)

Machine Translation (**EN -> X**)

Summarization (abstractive/extractive)

Question Answering

Topic modeling/Clustering

...

# Diversity of Domains

**Genres & sources:**

- News articles
- Social media updates
- Reviews
- ...

**Languages & dialects:**

- English:
  - British:
    - Northern:
      - Cheshire
      - ...
    - ...
  - American:
    - ...
  - ...
- ...

**Other types of domains:**

- Industries
- Users
- Time
- ...

# Tasks, Domains & Languages

# Data Preparation

- Gathering data for training/testing
- Defining labels/a taxonomy
- Cleaning up the data
- Getting the data in the right format
- Annotating the data
- Splitting the data

@parsaghaffari |

# Model Training

- Selecting the right algorithm/NN architecture
- Picking the right representations (e.g. pre-trained word embeddings)
- Hyper-parameter tuning
- Pre-training on other data (e.g. pre-training on sentiment data before emotion detection)
- Domain adaptation/transfer learning

# Model Evaluation

- Picking/defining the right metrics
- Creating a test set to evaluate models on
- Comparing models
- Qualitative evaluation
- Explaining predictions

# Workflow-related Issues

- Continuously updating models and getting feedback
- Active learning
- Domain adaptation
- Dataset visualization
- Consuming trained models

@parsaghaffari |

# Challenge #2

The scale and production rate of unstructured data

## The Rise of Unstructured Data

Merrill Lynch: Unstructured data accounts for 80%+ of all data in organizations and 95% of data generated daily online.

Exabytes
1800

- Structured data
- Unstructured data

0

2005                                          2015

Year

Research from International Data Corporation (IDC) shows that unstructured content accounts for 90% of all digital information.

"More content was uploaded yesterday than any one human could ever consume in their entire life." – Condé Nast

# Our Solution

# Text Analysis Platform (TAP) ß

An in-browser environment for building **custom NLP models**

Entities

Your Datasets

Categorization

Sentiment

CSV

TEXT ANALYSIS PLATFORM

# Build Datasets

From **CSV** files, using our **Knowledge Base** or by forking other datasets

---



**TAP**    🌐 Explore    🗄 My Datasets    🪄 My Models      ⚙ Jobs 0    👤 parsa ▾    ❓

## Dataset Preview

The following table shows a few sample rows of your Dataset.

Please choose which columns should be treated as Labels and which should be treated as Documents. Once you're done, hit "Convert to Dataset" and we'll create your new Dataset.

**Quick tip:** You also need to tell TAP whether or not it should ignore the header row in your file using the header row switch.

| Column Role | ⚪ Attribute | 🔴 Document | 🔵 Collection | ⚪ Attribute | 🟢 Label | ⚪ Attribu |
|---|---|---|---|---|---|---|
| Column Data Type | 99 String | 99 String | 99 String | 99 String | 99 String | 99 String |

**Header row?**

| other_topic | resolution_topics | gender | name | Resolution_Category | retweet_ |
|---|---|---|---|---|---|
| Read moore books, read less facebook. | Eat healthier | female | Dena_Marina | Health & Fitness | 0 |
| | Humor about Personal Growth and Interests Resolutions | female | ninjagirl325 | Humor | 1 |
| | Be More Confident | male | RickyDelReyy | Personal Growth | 0 |
| Help Morespread pet | | | | | |

# Build Datasets

From **CSV** files, using our **Knowledge Base** or by forking other datasets

# Build Datasets

From **CSV** files, using our **Knowledge Base** or by forking other datasets

## Knowledge Base

Amazon reviews ▼ | camera AND lens AND overall:5 | 🔍 Search

reviewer_id [string] | text [string] | overall [long] | summary [string] | class_labels [string] | version [string] | review_time [date] | id [long]

You can import all results or select some specific rows to import. Select rows by clicking on each row.

| Text | Reviewer Id | Overall | Summary | Class Labels | Version | Review Time |
|---|---|---|---|---|---|---|
| Best *lens* for the money. I keep this *lens* on my *camera* 90% of the time. Great *lens* for starting photography | A3USKITGRX OIVG | 5 | Amazing *lens* | Electronics | 1.0 | 2012-11-28T00:00:00... |
| Fotodiox *Lens* Mount Adapter -- Nikon *Lens* to Sony NEX E-Series *Camera* arrived on time as advertised ... | ASJ7QB66JU 4RL | 5 | Fotodiox *Lens* | Electronics | 1.0 | 2013-08-02T00:00:00... |
| *Lens* is more compact than I expected. This is a great all around *camera lens*! | A1Z35O63NE 2BZE | 5 | This is a great all ... | Electronics | 1.0 | 2014-06-27T00:00:00... |
| My first time have a *lens* of this quality. Great *lens* to keep dirt and other items from *camera lens* ... | A1EMTRTDM XLBVB | 5 | Great *Lens*!... | Electronics | 1.0 | 2013-07-15T00:00:00... |
| The *camera lens* is just perfect. It replaced a tamaron *lens* of the same specifications. I | A3VBPLKSRH | 5 | *Camera* ... | Electronics | 1.0 | 2011-12-24T00:00:00 |

‹ Previous | Results from 1 to 10 of 65,637 | › Next

Close | Import documents into 🏷 positive | Start | 10000 | ⬆ Import 10,000 documents

# Train Models

**Compare** various models/parameters and pick the best one

# Evaluate **Models**

**Understand** different aspects of each model

# Evaluate **Models**

**Explain** predictions made by a model

# Use Models

**Deploy** the best version of your model to production, and use it as an **API**

---

TAP | Explore | My Datasets | My Models | Jobs 0 | parsa ▾ | ?

## parsa / sentiment_fork (1)_model1

View Evaluation | Retrain | ...

Dataset | sentiment_fork (1) | API calls | 1
Private | Make it public

This is a sentiment analysis model for tweets. ✎

sentiment-analysis | Manage tags

### Usage details

**Endpoint**

http://api.tap.aylien.com/v1/models/2ee632e2-81b9-43fc-bf38-662c772cdf2d | Copy

**API Key**

814c6cb9f0bf4702aebf61f3c1596171

**Example request:**

```
curl -X POST \
    -H "x-aylien-tap-application-key: 814c6cb9f0bf4702aebf61f3c1596171" \
    --data-urlencode "text=RT @vancitynomad: Thanks SD Carl and crew for the great flight ac3 to Narita sept11 @MomonaKomagata @AirC
    http://api.tap.aylien.com/v1/models/2ee632e2-81b9-43fc-bf38-662c772cdf2d
```

### Model performance

# Leverage the **Marketplace**

**Share** your **models** and **datasets** with other users, or leverage theirs



TAP    🌐 Explore    🗄 My Datasets    🪄 My Models              ⚙ Jobs ② | 👤 parsa ▾ | ❓

## 🗄 Datasets

Here you can see a list of Datasets that other users have made public.
Take a closer look at any Dataset by clicking on it or copy any Dataset for your own use using the **Fork** button.

### Emojis
robson@aylien.com

Various emojis and messages associated with them.

🕐 Updated 5 hours ago

👁 Inspect | ⌥ Fork  ⟨ 6

### JIRA Tickets
robson@aylien.com

A dataset consisting of different categories of JIRA tickets (bugs, feature requests, etc) along with thousands of tickets for each category.

🕐 Updated 8 days ago

👁 Inspect | ⌥ Fork  ⟨ 8

### Offensive Speech
robson@aylien.com

Social media posts and messages categorized based on their offensiveness.

🕐 Updated 24 days ago

👁 Inspect | ⌥ Fork  ⟨ 5

### Large Movie Review
robson@aylien.com

A large collection of positive and negative movie reviews.

🕐 Updated a month ago

👁 Inspect | ⌥ Fork  ⟨ 7

### Rotten Tomatoes Phrases
robson@aylien.com

A collection of short phrases about movies, with positive, neutral and negative annotations.

🕐 Updated a month ago

👁 Inspect | ⌥ Fork  ⟨ 1

### Political Ideology
robson@aylien.com

A dataset consisting of 4,062 sentences annotated for political ideology–2025 liberal sentences, 1701 conservative sentences, and 600 neutral sentences.

🕐 Updated a month ago

👁 Inspect | ⌥ Fork  ⟨ 4

‹ Previous | Next ›

# Key takeaways

- NLP brings us closer to people's thoughts, emotions and intentions => It has the potential to impact billions of people;
- We're not doing a great job at modeling the processes of which language is an output => We face the problem of too many <task, domain, language>s => Infinite room for research in this area;
- NLP(/ML) is not a one-size-fits-all problem => We need to rely on good engineering and products to fill the gap;

@parsaghaffari |