# 4TH INTERNATIONAL FORUM BIG DATA DAY

## - BAKU 2018 -
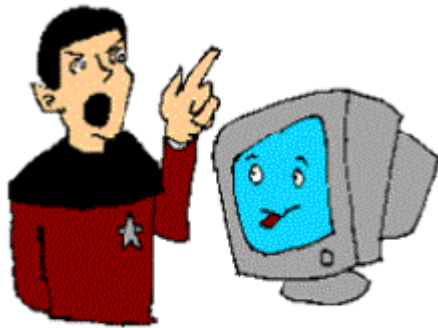
## *Natural Language Processing and its Application*

**Dr., Samir Rustamov,**
Assistant Professor, School of IT & Engineering, ADA University

# What is Natural Language Processing (NLP)?

- Natural Language Processing (NLP) is a field of artificial intelligence that enables computers interact with human in natural language.
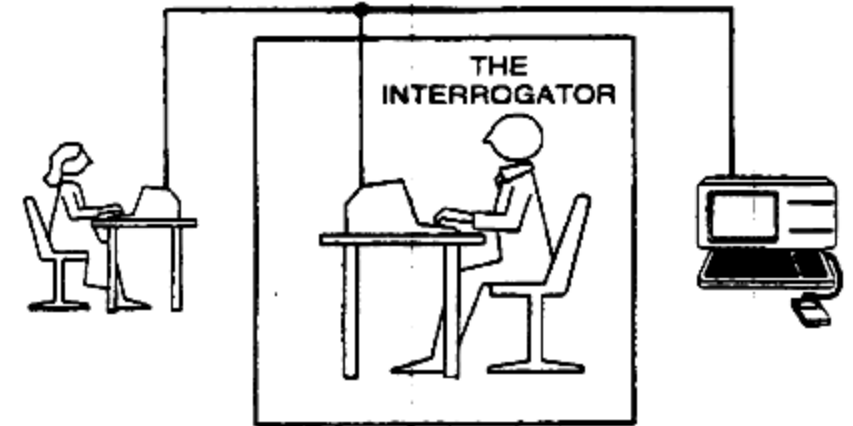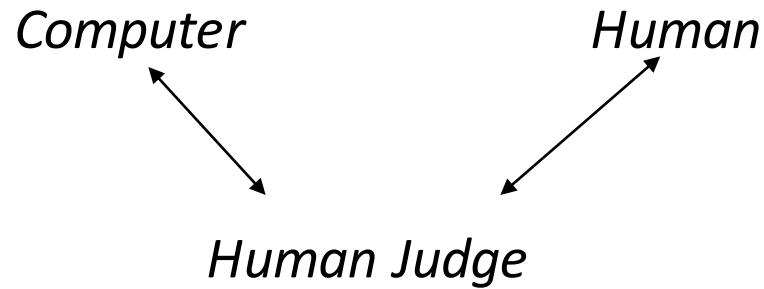


Natural Language Processing

Ultimate goal: Natural human-to-computer communication

# The Turing Test
## ([Can Machine think? A. M. Turing, 1950](#))

THE INTERROGATOR

*Computer*     *Human*

*Human Judge*

- *Human Judge* asks tele-typed questions to *Computer* and *Human.*
- *Computer's* job is to act like a human.
- *Human's* job is to convince Judge that he is not machine.
- *Computer* is judged "intelligent" if it can fool the judge
- Judgment of intelligence is linked to appropriate answers to questions from the system.
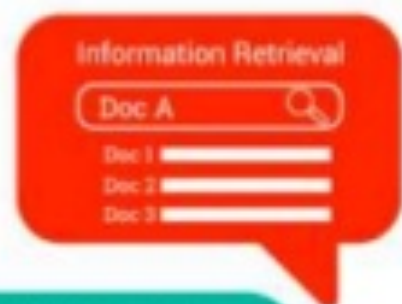
# NLP in the Commercial World

# Natural Language Processing Market to Reach $22.3 Billion by 2025



Natural Language Processing Total Revenue by Segment, World Markets: 2016-2025

Source: Tractica

# Forms of Natural Language

- The input/output of a NLP system can be:
  - **written text**
  - **speech**
- To process written text, we need:
  - **lexical, syntactic, semantic knowledge about the language**
  - **discourse information, real world knowledge**
- To process spoken language, we need everything required to process written text, plus the challenges of speech recognition and speech synthesis.

# Components of NLP

**Natural Language Understanding**

Mapping the given input in the natural language into a useful representation.

**Natural Language Generation**

Producing output in the natural language from some internal representation.

NL Understanding is much harder than NL Generation. But, still both of them are hard.

# Why NL Understanding is hard?

- Natural language is extremely rich in form and structure, and **very ambiguous**.
- One input can mean many different things. Ambiguity can be at different levels.
  - Lexical (word level) ambiguity -- different meanings of words
  - Syntactic ambiguity -- different ways to parse the sentence
  - Interpreting partial information -- how to interpret pronouns
- Many input can mean the same thing.
- Interaction among components of the input is not clear.

# Example of ambiguity

- Some interpretations of:        **Adamı gördüm.**
  1. I saw the man.
  2. I saw Adam
  3. I saw my island.
  4. I visited my island.
  5. I saw my ADA
  6. I visited my ADA
  7. I bribed the man.

- Semantic Ambiguity:
  - gör          to see
  - gör          to visit
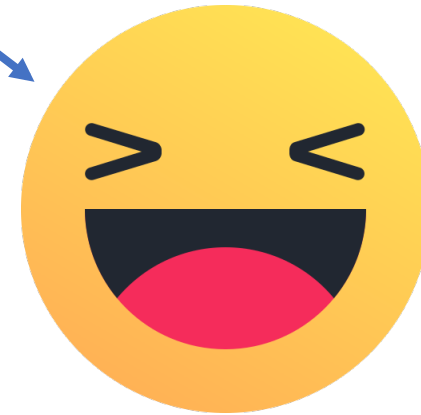  - gör          to bribe

# Resolve Ambiguities

- **lexical disambiguation** -- Resolution of part-of-speech and word-sense ambiguities are two important kinds of lexical disambiguation.

- **syntactic ambiguity** -- can be addressed by probabilistic parsing.
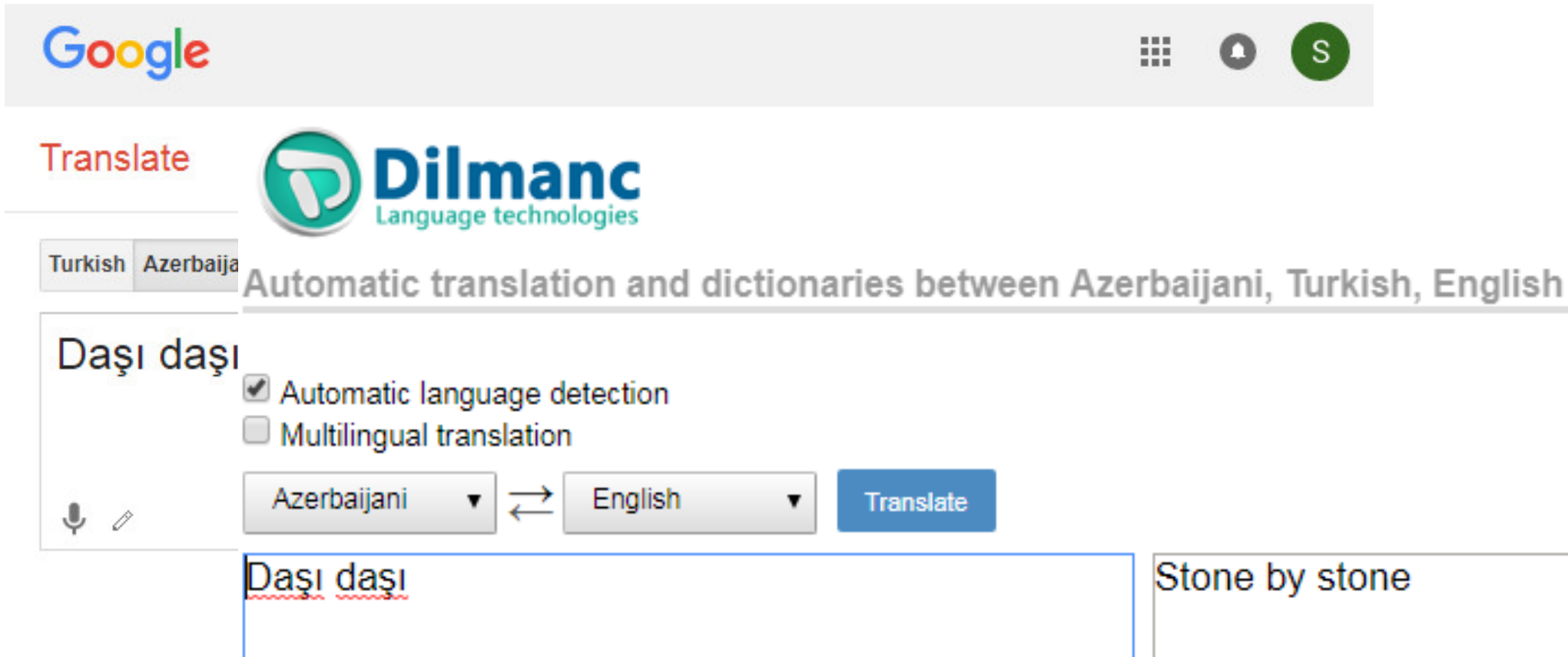
# What is PoS tagging? Why is it important?

Each word has a part-of-speech tag to describe its category. POS Taggers try to find POS tags for the words.

"gül" - ?

# Why it matters? Applications

► Machine translation – "Daşı daşı"

# Language Processing



Discourse ( Connected sentence processing in a larger body of text)

Pragmatics (Meaning in context & for a purpose )

Semantics (Meaning of sentences)

Syntax, Parsing (Structure of sentences )

Morphology, Lexicon Words & their forms (Words & their forms )

Phonetics & Phonology (Speech sound)

# Applications for spelling correction

## Word processing

Spell checking is a componant of

**Spelling and Grammar: English (US)**

Not in dictionary:

Spell checking is a **componant** of

Ignore
Ignore All
Add

Suggestions:

component

Change
Change All
AutoCorrect

## Phones

**New iMessage**    Cancel

To:    Dan Jurafsky

late ×

Sorry, running layr    **Send**

Q W E R T Y U I O P
A S D F G H J K L
Z X C V B N M
123    space    return

## Web search

ploogle    natural langage processing

Showing results for natural *language* processing
Search instead for natural langage processing

# Spelling correction

$$f\left(\begin{array}{l}\text{- Cin seddi dunyanın yedi mocusəsindən birdir.} \\ \text{- Niye?} \\ \text{- Cunki duzəldikleti en uzunomurlu sheydi.}\end{array}\right) = \begin{array}{l}\text{-Çin səddi dünyanın yeddi möcüzəsindən biridir.} \\ \text{- Niyə?} \\ \text{- Çünki düzəltdikləti ən uzunömürlü şeydir.}\end{array}$$

# Machine translation

When the space organization NASA first started sending up astronaunts they discovered ballpoint pens would not work in zero gravity. To solve the problem, NASA scientists spent ten years and 12 billion to develop a pen that would write in zero gravity, upside down, under water, on all types of surface, and at temperatures ranging from below freezing to 300C. Russians used a pencil.

NASA kosmik təşkilatı ilk dəfə astronavtların göndərilməsinə başladıqları zaman, ballpoint qələmləri sıfır çəkisi ilə işləməyəcəkdi. Problemi həll etmək üçün, NASA alimləri sıfır ağırlıq, baş aşağı, su altında, bütün səthlərdə və aşağıda dondurmadan 300 dərəcə qədər dəyişən temperaturda yazacaq bir qələm hazırlamaq üçün on il və 12 milyard dollar sərf etmişdir. Ruslar bir qələm istifadə edirdi.
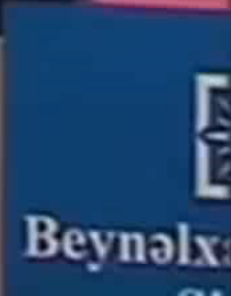
Kosmos təşkilatı NASA ilk astronaunts qaldıraraq başlayanda onlar kəşf diyircəkli qələmlər çəkisizlik şəraitində işləməyəcək. Problemin həlli, NASA alimləri çəkisizlik yazacaq qələm inkişaf etdirmək, on il və 12 milyard xərcləyib, tərsinə, su altında, bütün səthinin növləri, 300c-a şaxta tutmuş və temperaturda . karandaş istifadə olunan Ruslar.
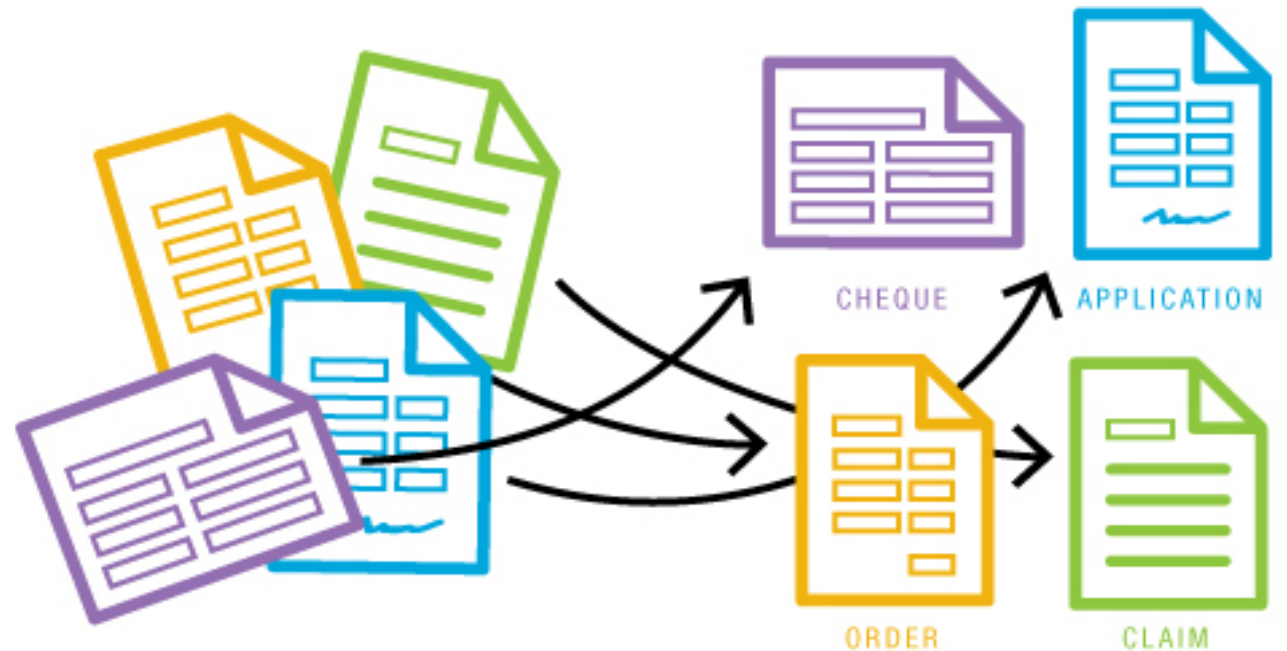
# Why is machine translation hard?

- Requires both understanding the "from" language and generating the "to" language.

- How can we teach a computer a "second language" when it doesn't even really have a first language?

- Can we do machine translation without solving *natural language understanding* and *natural language generation* first?

# Text Classification

- Assigning subject categories, topics, or genres
- Spam detection
- Authorship identification
- Age/gender identification
- Language Identification
- Sentiment analysis
- …

# SENTIMENT ANALYSIS

Discovering people opinions, emotions and feelings about a product or service

# Information retrieval

- Information retrieval is the activity of obtaining information resources relevant to an information need from a collection of information resources.
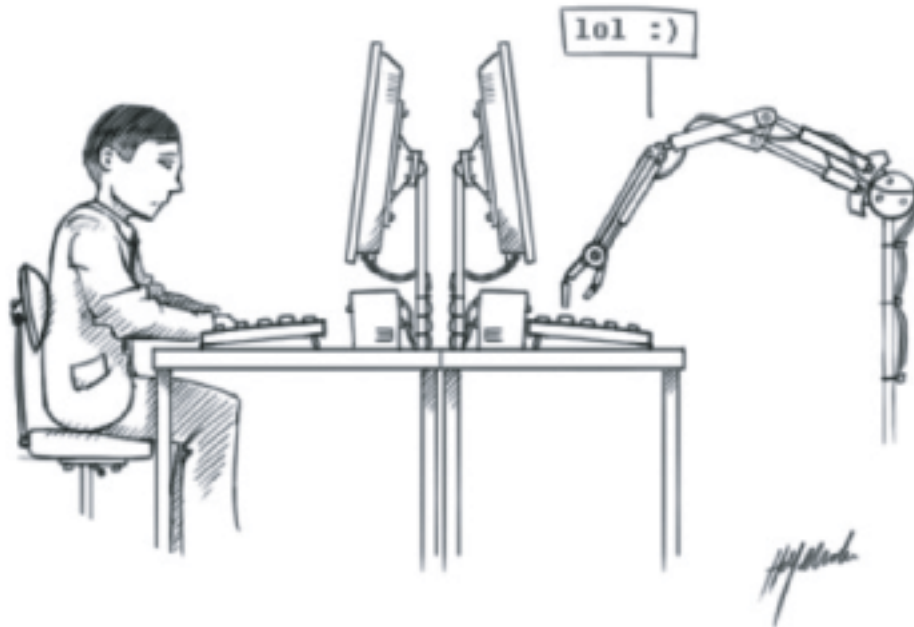
# Text Summarization

- **Goal**: produce an abridged version of a text that contains information that is important or relevant to a user.



**Summary:** In many cases, the zombie machines are used to send out spam, to perform click fraud, to aid in identity theft, or are directed to attack another web server on the internet, as was recently seen with the Twitter/Facebook/LiveJournal attacks of last month... Using a private newsgroup, the trojan executes a command which logs it into the newsgroup and requests a specific page... The trojan is attempting to remain discreet and undetected, being used to subtly gather information and potentially determine its future attack targets...

# *Dialogue systems*

- A **dialog system** or **conversational agent** (**CA**) is a computer system intended to converse with a human, with a coherent structure.
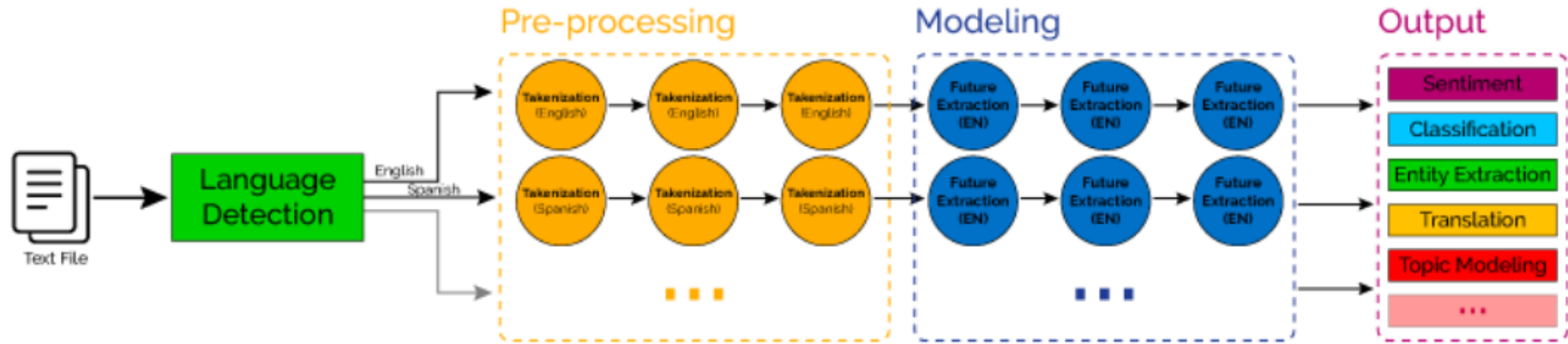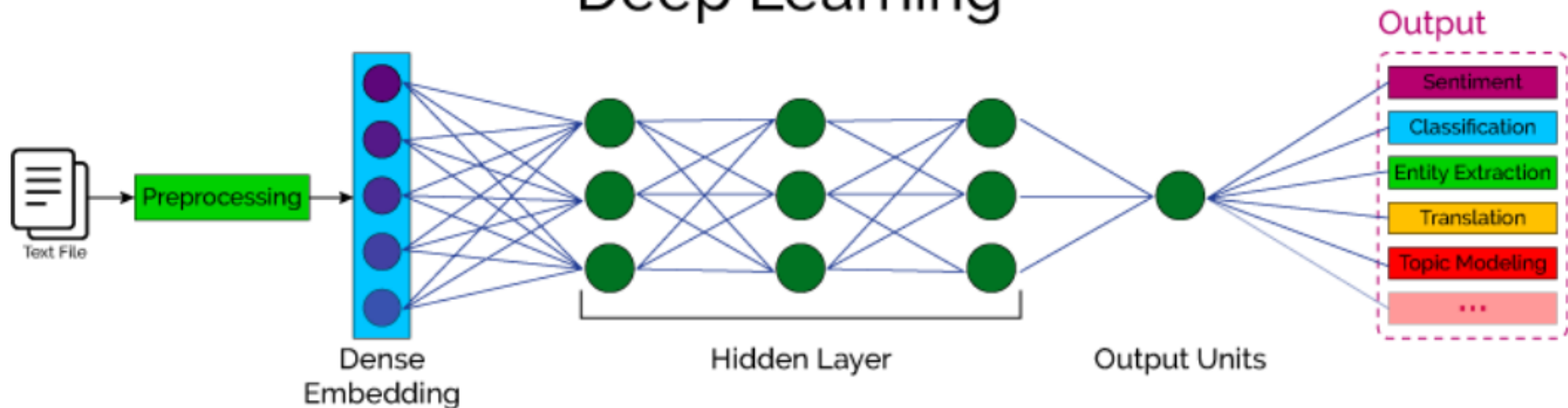
# Question Answering:

| Deep Learning Algorithms | NLP Usage |
|---|---|
| Neural Network – NN (feed) | •Part-of-speech Tagging<br>•Tokenization<br>•Named Entity Recognition<br>•Intent Extraction |
| Recurrent Neural Networks -(RNN) | •Machine Translation<br>•Question Answering System<br>•Image Captioning |
| Recursive Neural Networks | •Parsing sentences<br>•Sentiment Analysis<br>•Paraphrase detection<br>•Relation Classification<br>•Object detection |
| Convolutional Neural Network -(CNN) | •Sentence/ Text classification<br>•Relation extraction and classification<br>•Spam detection<br>•Categorization of search queries<br>•Semantic relation extraction |

# Difference Between Classical NLP & Deep Learning NLP

## Classical NLP

**Pre-processing**

**Modeling**

**Output**

Text File → Language Detection

English / Spanish

Tokenization (English) → Tokenization (English) → Tokenization (English)

Tokenization (Spanish) → Tokenization (Spanish) → Tokenization (Spanish)

Future Extraction (EN) → Future Extraction (EN) → Future Extraction (EN)

Future Extraction (EN) → Future Extraction (EN) → Future Extraction (EN)

Sentiment
Classification
Entity Extraction
Translation
Topic Modeling
...

## Deep Learning

**Output**

Text File → Preprocessing → Dense Embedding → Hidden Layer → Output Units

Sentiment
Classification
Entity Extraction
Translation
Topic Modeling
...

# References:

- 1. Overview of Artificial Intelligence and Natural Language Processing. NAVDEEP SINGH GILL. https://www.upwork.com/hiring/for-clients/artificial-intelligence-and-natural-language-processing-in-big-data/

2. **Natural Language Processing Market to Reach $22.3 Billion by 2025**

https://www.tractica.com/newsroom/press-releases/natural-language-processing-market-to-reach-22-3-billion-by-2025/

3. Dan Jurafsky. Natural Language Processing Lectures.

4. BİL711  Natural Language Processing. Prof. Dr. İlyas Çiçekli.