# Natural Language Processing & Applications

# Why Text ?

Source : RECOMND

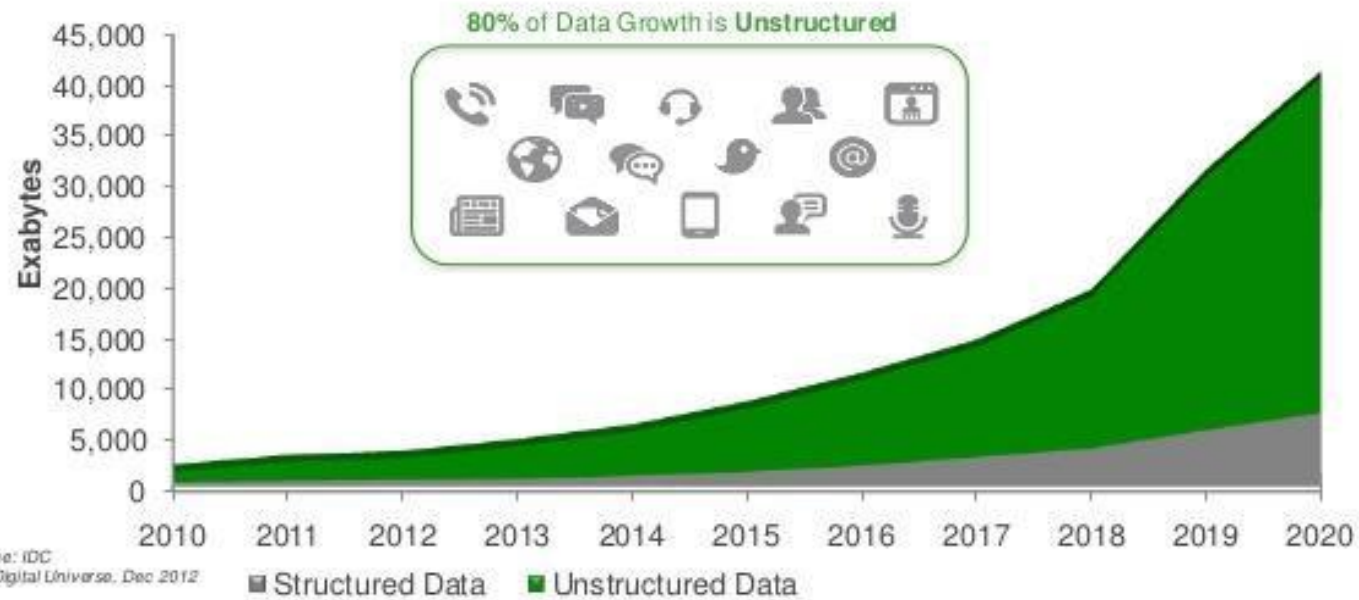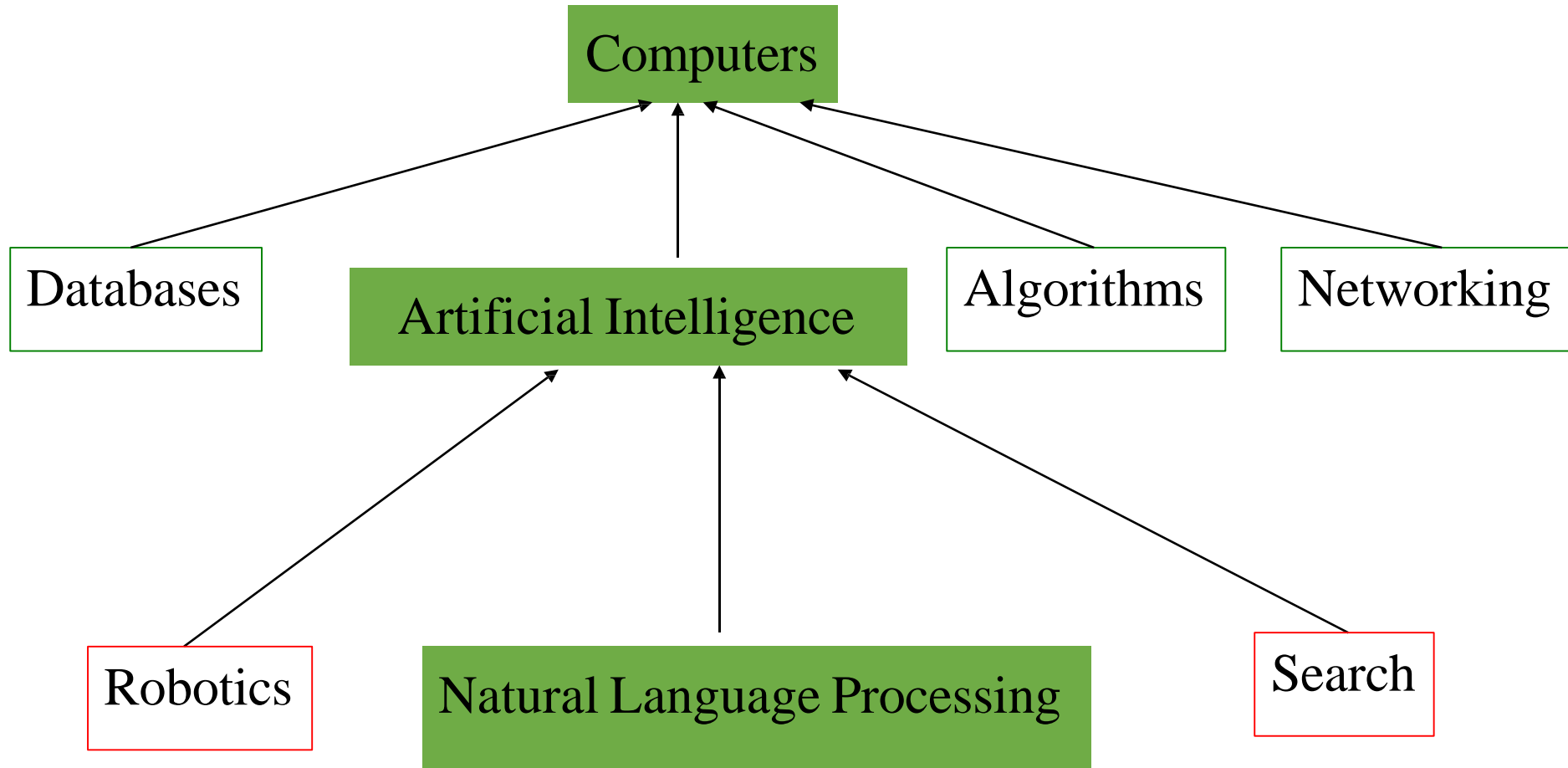# Natural Language Processing

- A hallmark of human intelligence.

- Natural Language Processing
  - Natural Language Understanding
  - Natural Language Generation
  - Process information contained in natural language text.
  - Computational Linguistics (CL), Human Language Technology (HLT), Natural Language Engineering (NLE)

- Can machines understand human language?

  Ultimate goal

  Analyze, understand and generate human languages just like humans do.

# Fitting in CS taxonomy

# NLP- Tasks

## Natural Language Understanding

Taking some spoken/typed sentence and working out what it means

## Natural Language Generation

Taking some formal representation of what you want to say and working out a way to express it in a natural (human) language (e.g., English)

# Working towards

- Applying computational techniques to language domain.

- Use the theories to build systems that can be of social use.

- Make computers learn our language rather than we learn theirs.

# Natural language understanding

Raw speech signal /Raw Text

    ↓   **• Speech recognition**

Sequence of words spoken /written

    ↓   **• Syntactic analysis**

Structure of the sentence

    ↓   **• Semantic analysis**

Partial representation of meaning of sentence

    ↓   **• Discourse & Pragmatic analysis**

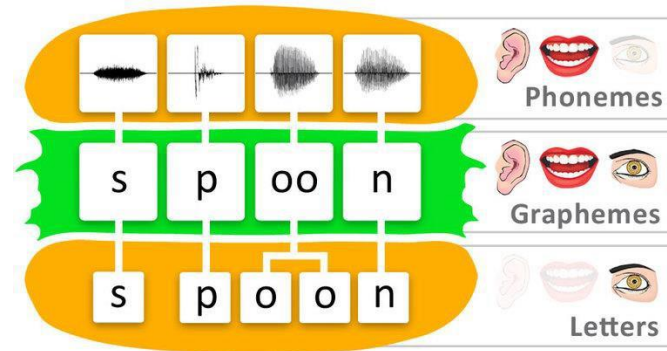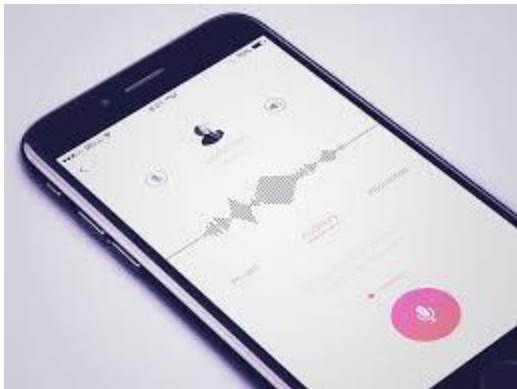Final representation of meaning of sentence

# Aspects of Language Processing

- Phonology – Speech processing

- Word, lexicon: lexical analysis
  - Morphology, word segmentation

- Syntax
  - Sentence structure, phrase, grammar, …

- Semantics
  - Meaning

- Discourse analysis
  - Meaning of a text
  - Relationship between sentences

- Pragmatics
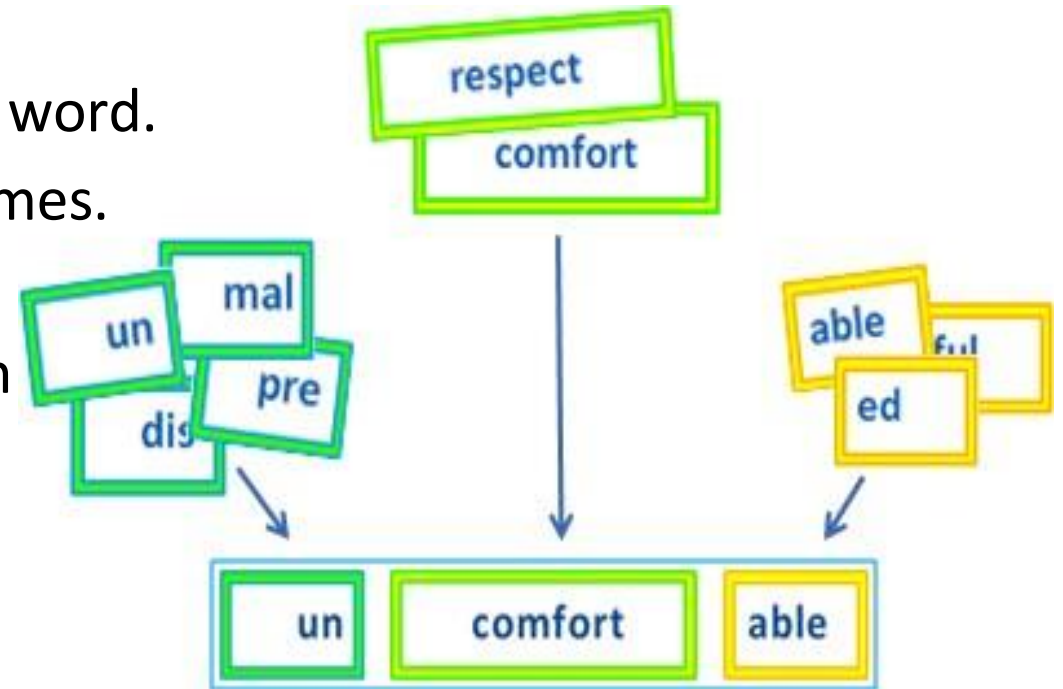  The study of meaning in different contexts of use

# Phonology



## Speech processing

- Humans process speech remarkably well.
- Speech interface can replace keyboards and monitors.
- Convert Acoustic signals to Text.
- Phonemes are the smallest recognizable speech unit in a language.





Phonemes

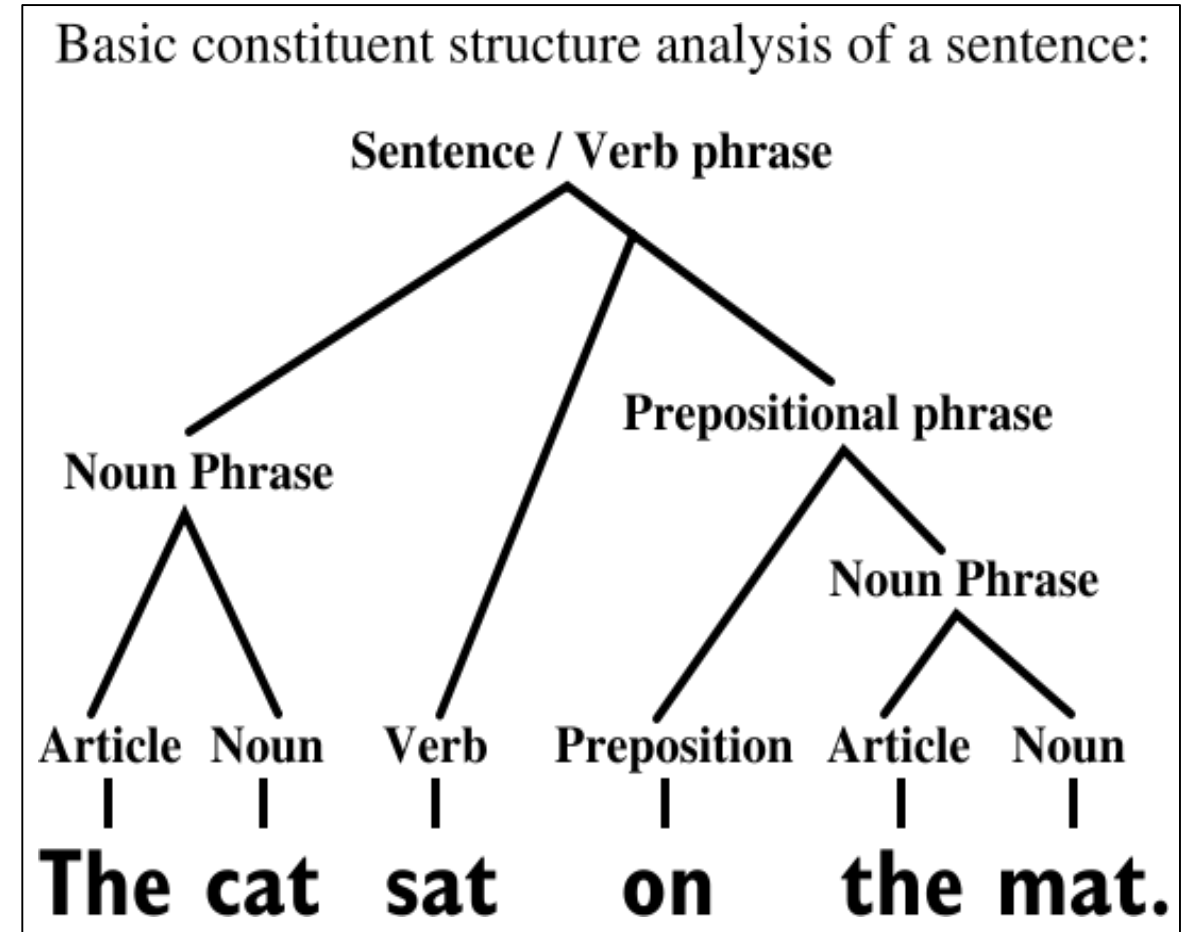s p oo n — Graphemes

s p o o n — Letters

# Morphology

- Structures and patterns in words
- Words are a sequence of Morphemes.
  - Morpheme – smallest meaningful unit in a word.
- Analyses how words are formed from morphemes.

  e.g.,        dogs= dog+s.
- Inflectional Morphology – Same Part of Speech
  - Buses = Bus + es
  - Carried = Carry + ed
- Derivational Morphology – Change PoS.
  - Destruct + ion = Destruction (Noun)
  - Beauty + ful = Beautiful (Adjective)
- Affixes – Prefixes, Suffixes Rules govern the fusion.

respect
comfort

un   mal
dis   pre

able   ful
ed

un   comfort   able

# Syntax

- Words when put together they convey more.
- Syntax is the grammatical structure of the sentence.
- Syntactic Analysis (Parsing)

  Process of assigning a parse tree to a sentence.



Basic constituent structure analysis of a sentence:

# Syntactic Analysis - Grammar

```
sentence -> noun_phrase, verb_phrase

noun_phrase -> proper_noun

noun_phrase -> determiner, noun

verb_phrase -> verb,
               noun_phrase

proper_noun -> [mary]

noun -> [apple]

verb -> [ate]
determiner -> [the]
```
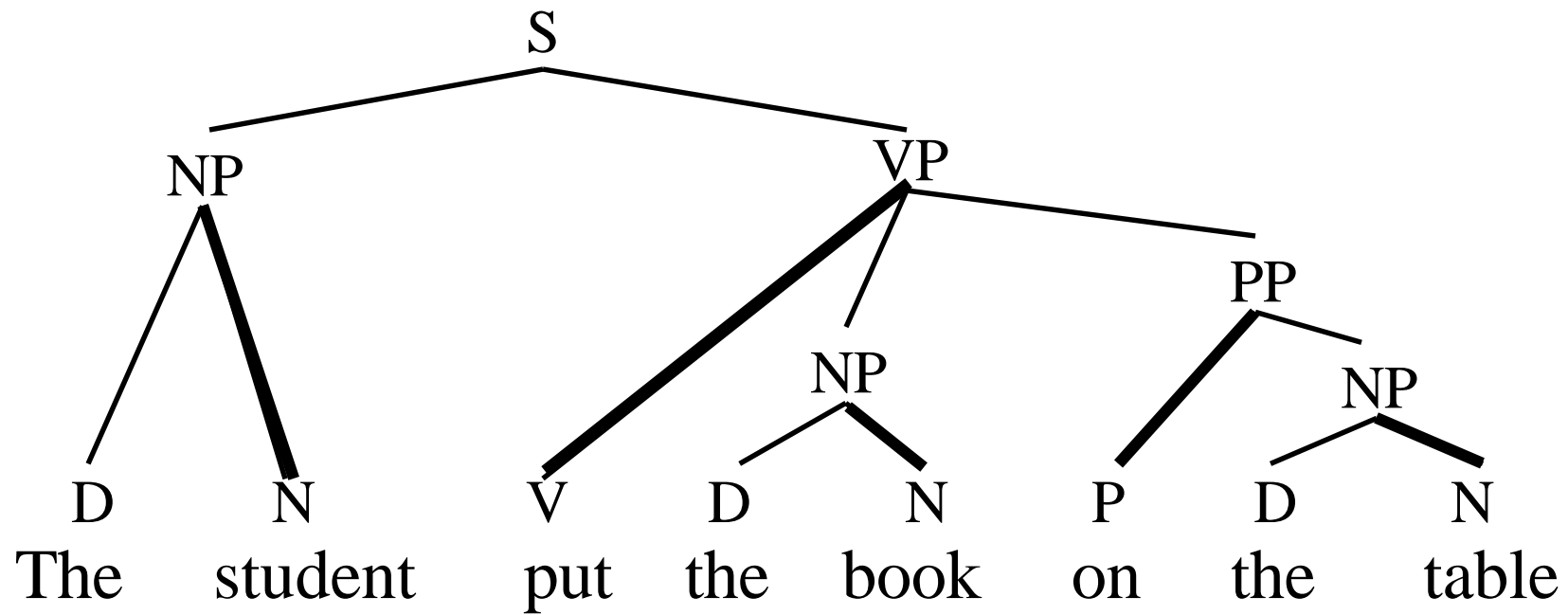
# Parsing

- Analyze the structure of a sentence

# Semantic Analysis

- What do you mean..?

- Words – Lexical Semantics

- Sentences – Compositional Semantics

- Converting the syntactic structures to semantic format – <span style="color:red">meaning representation.</span>

- Semantics: the meaning of a word or phrase within a sentence

e.g., "I saw the prudential building flying into Boston"

# Semantic Representations

- Meaning representation of the sentence from its syntactic structure(s)
- Ways of meaning representing the sentence:
  - Logical forms

    **Representation:** $\Xi$x man(x) & tall(x) & plays(x, basketball)
  - Semantic role labelling

    **Representation:** Kill( Agent : X, Victim: 3 people, Instrument : Gun)

# Discourse Analysis

- The meaning of an individual sentence may *depend on the sentences that precede* it and may *influence the meaning of the sentence that follow it*.

- Issues related to discourse Integration
  - Anaphora

    Resolving the pronoun's reference. Co-reference resolution

    *Coreference resolution is the task of finding all expressions that refer to the same entity in a text.*
  - Ellipsis

    Incomplete sentences

# Anaphora Resolution



Anaphoric Pronoun Resolution

**John** likes **ice-cream, he** eats **it** every day.

**Finding links**
- Pronoun to antecedent

**Enriching text**
- Input: preprocessed document
- Output: All found anaphoric pronoun references to words/phrases

# Discourse Structures- Ellipsis

- Ellipsis – Incomplete sentences

  The second sentence is not complete, but what it means can be inferred from the first one.

# Pragmatics

- Uses context of utterance

  - Where, by who, to whom, why, when it was said

  - *Intentions:* *inform, request, promise, criticize, …*

# Syntax Ambiguity

# Semantic Ambiguity

Semantic ambiguity: **"I saw the prudential building flying into Boston"**



**Semantic Restriction, Domain Knowledge - Ontology**

# Pragmatics Ambiguity

Pragmatic ambiguity: **"you're late"**

What's the speaker's intention: informing or criticizing?

# Computing Techniques

- Stemming
  - Reduce *words* to base form.
- POS Tagging
  - Determine for each word whether it is a noun, adjective, verb, …..
- Parsing
  - sentence $\Rightarrow$ parse tree
- Word Sense Disambiguation
  - orange juice vs. orange coat
- Semantic analysis
  - Semantic representations and Evaluation

# Application Areas

➢ Machine Translation

➢ Information Retrieval

     Selecting from a set of documents the ones that are relevant to a query

➢ Text Categorization

     Sorting text into fixed topic categories

➢ Question Answering

➢ Information Extraction

   Converting unstructured text into structured data

# Application Areas (cont..)

- Spoken language control systems

- Spelling and grammar checkers

- Sentiment Analysis

- Text-to-Speech & Speech recognition

- Natural Language Dialogue Interfaces to Databases

- Plagiarism detection

# Information Retrieval

# Email Spam Filtering/Categorizing

# Machine Translation

- Multilingual Usage
- Machine-assisted human Translation
- Scope

    Creating Language resources.

# Information Extraction

Information extraction systems
- Find and understand relevant parts of text.
- Produce a structured representation of the relevant information from text, in the form of :
  - entities,
  - relations between entities ,
  - events in which the entities are involved.
- Produce a structured representation of the relevant information- *relations/events*

# Plagiarism Detection

# Current Technology Giants

- IBM Watson
- Google
- Microsoft Oxford
- Wolfram Natural Language Understanding System

# TOOLS

# Conclusion

- Complete human-level natural language understanding is still a distant goal

- Develop Algorithms for each level.

- Find appropriate match between application domain and the available methods

# References

Books
- Allen.J, "Natural Language Understanding ", second edition, Benjamin Cummings Publishing Co .
- Manning.C.D and Schutze.H, "Foundations of Statistical Natural Language Processing ", MIT Press.

Web references
- research.microsoft.com/en-us/groups/nlp/
- nlp.stanford.edu/