# Information Retrieval: between Natural Language, Texts and Meaning

**Roberto Basili**

(Università di Roma, Tor Vergata, basili@info.uniroma2.it)

Some slider borrowed from the tutorial «Natural Language Understanding: Foundations and State-of-the-Art", by Percy Liang (Stanford University).

Web Mining & Retrieval, a.a. 2020-21

# Overview

- Documents in Information Retrieval

  - Information, Representation, (re)current challenges, success(and unsuccess)ful stories

- Information and Content

  - Natural Language Processing: introduction to the linguistic background
    - Natural Language and Content
    - NL Syntax
    - NL Semantics

  - Document Representation and IR models

- Summary

# Semantics, Open Data and Natural Language

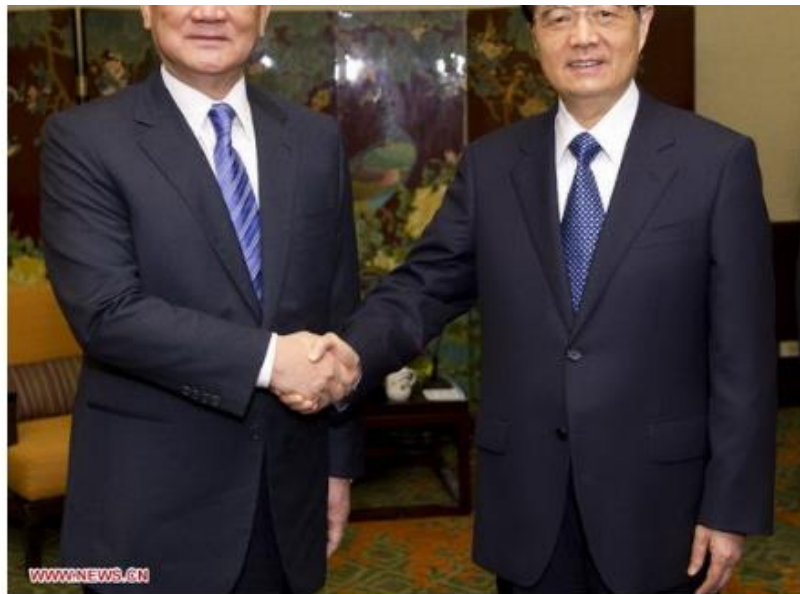- **Web contents, characterized by rich multimedia information, are mostly opaque from a semantic standpoint**

Chinese President Hu Jintao (R) shakes hands with Honorary Chairman of the Chinese Kuomintang (KMT) Lien Chan, in Honolulu, Hawaii, the U.S., Nov. 11, 2011. (Xinhua/Huang Jingwen)

HONOLULU, United States, Nov. 11 (Xinhua) -- Hu Jintao, general secretary of the Central

*Who is Hu Jintao?*

# Content Semantics and Natural Language

- Human languages are the main carrier of the information involved in processes such as retrieval, publication and exchange of knowledge as it is associated to the open Web contents

- Words and NL syntactic structures express concepts, activities, events, abstractions and conceptual relations we usually share through data

- "Language is parasitic to knowledge representation languages but the viceversa is not true" (Wilks, 2001)

- From Learning to Read to Knowledge Distillation as a(n integrated pool of) Semantic interpretation Task(s)

# Texts, Information & Document Structures

## What is a document?

Sailing in Greece

B. Smith

content structure

external attributes
author = 'B. Smith'
crdate = '25.05.98'
ladate = '30.06.99'

layout structure

logical structure

Sailing
Greece
Mediterenean
Fish
Sunset

head
title
author
chapter
section
section
chapter

Home

# Information Retrieval Models

- An IR model must specify (at least) :

    - A representation of the document

    - A rapresentation of individual queries

    - The retrieval function

- The model determines a specifici notion of relevance.

- Relevance can be discrete (e.g. binary) or continuous (i.e. rank or relevance order).

- It is a perfect example of learnable function through induction from examples (see Google)

# IR models (2)

Set Theoretic

Fuzzy
Extended Boolean

Classic Models

boolean
vector
probabilistic

Algebraic

Generalized Vector
Lat. Semantic Index
Neural Networks

User Task

Retrieval:
Adhoc
Filtering

Probabilistic

Inference Network
Belief Network

Structured Models

Non-Overlapping Lists
Proximal Nodes

Browsing

Browsing

Flat
Structure Guided
Hypertext

10

# Model Families for IR

- Boolean Models (set theoretic)
  - Standard boolean
  - Extended Boolean

- Vector Models (algebraic)
  - Generalized Vector Space
  - Latent Semantic Indexing
  - Neural models

- Probabilistic Models
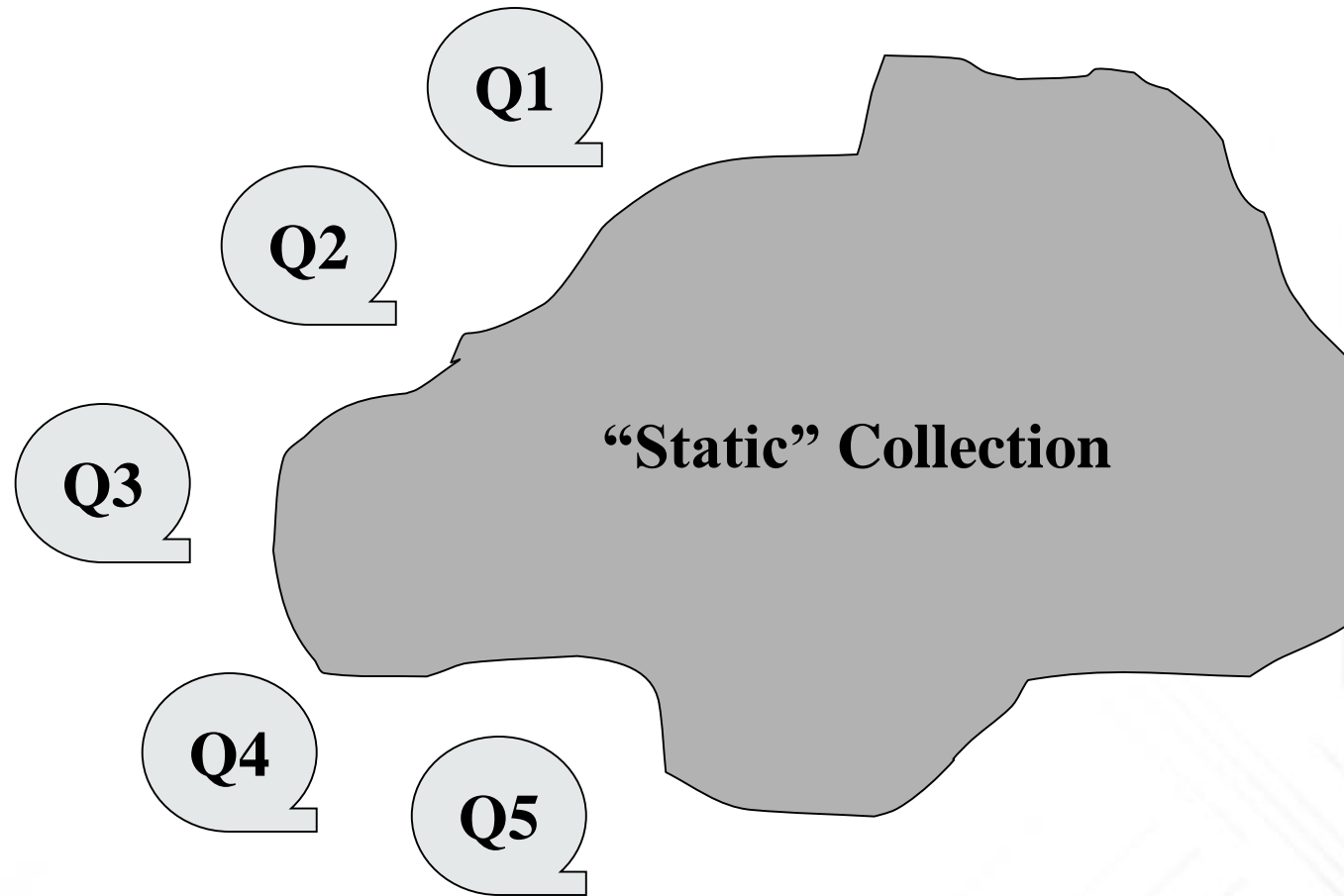
# Other classification Dimensions

- Document Logical Model
  - Type of Indexes
    - Structures vs. Content
    - Metadata vs. Content
  - Full text as a model of the content
  - Full text viz. Document (Hypertextual) Structure
    - Declarative vs. operational semantics

- The role of user
  - Subjective vs. Objective forms of relevance
  - Operational environment
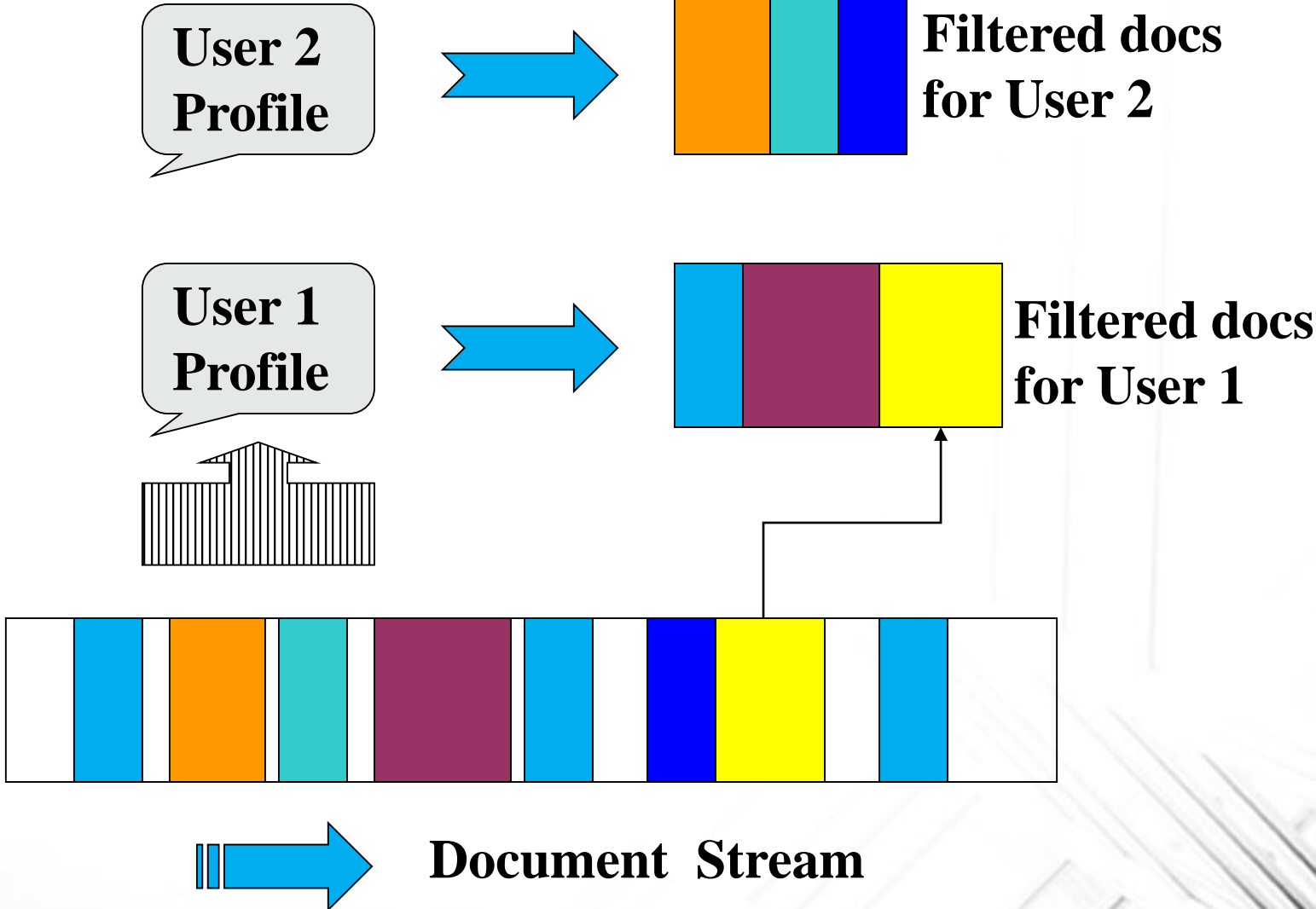    - Search vs. Browsing

# Retrieval Tasks

- **Ad hoc retrieval**

  - DEF. Relatively stable document collection vs. highly variable queries.

- **Information Filtering:**

  - DEF. Fixed Queries and continuous streams of documents

  - Type of Filtering
    - User Filter: static model of the subjective preferences
    - Category based filtering: static model of categories as domain preferences

  - Target Function: binary decision, in general

- **Information Routing:**

  - DEF. When filters define dynamic e non binary models of preference.

13

# Ad Hoc Retrieval



Q1

Q2

Q3

"Static" Collection

Q4

Q5

# Filtering



User 2 Profile → Filtered docs for User 2

User 1 Profile → Filtered docs for User 1

Document Stream

# Learning and IR

- The task in IR and the need of modeling either documents and queries are strongly related to Machine Learning

- First, **no analytical function is available** for every domain, document collection, user and query is available

- Second, **unstructured data** (as much frequently occurring in Web applications) are hard to be modeled without resorting to a reference notion of **content**

- CHALLENGE: How to deal in an **efficient** manner with the tasks of **representing**, **querying**, **matching**, **filtering** and **sorting** the complex contents characterizing the arbitrarily **distributed and unstructured Web data**?

  - In the case of textual document: **how can we learn to formalize the vague notion of content for a document?**

# Semantics, Natural Language & Learning: tasks

- In order to make contents explicit in an IR process they must be recognized in the contexts of their use

- All these process (also called **Learning to Read** or **Knowledge Distillation**) proceed as a (integrated pool of) **Semantic interpretation Task(s)**

  - **Information Extraction** (from text to machine readable concepts)
    - Entity Recognition and Classification
    - Relation Extraction
    - Semantic Role Labeling (Shallow Semantic Parsing)

  - **Estimation of Text Similarity** (from text to quantitative semantic measures)
    - Structured Text Similarity/Textual Entailment Recognition
    - Sense disambiguation

  - **Semantic Search**, **Question Classification** and **Answer Ranking**

  - **Knowledge Acquisition**, e.g. ontology learning

  - **Social Network Analysis**: Opinion Mining, Recommending
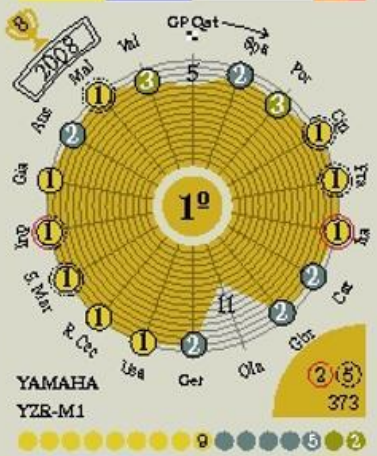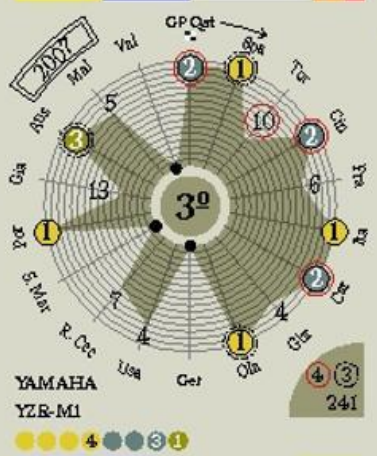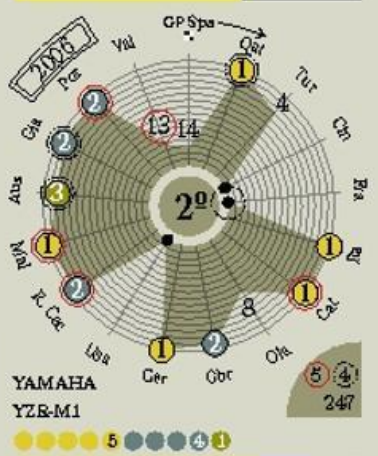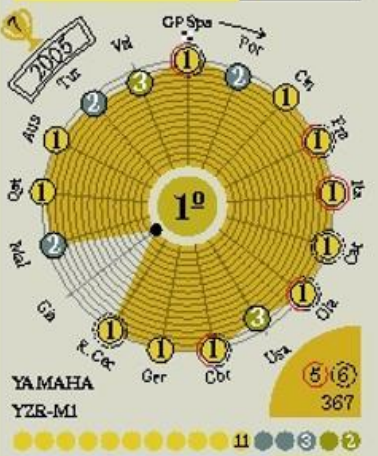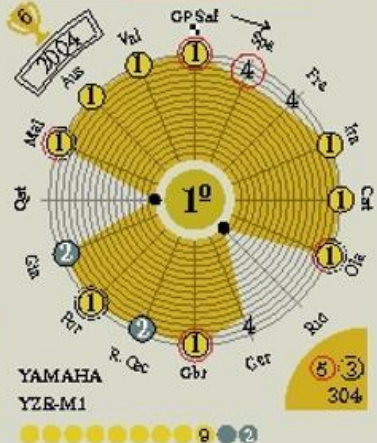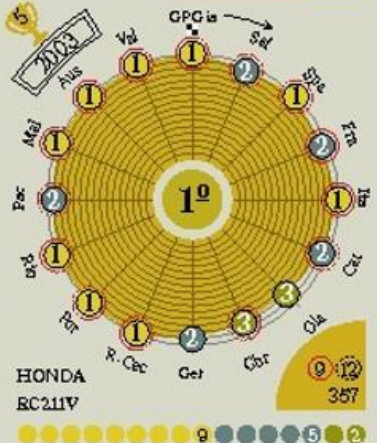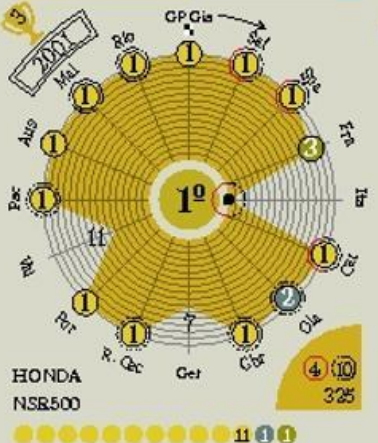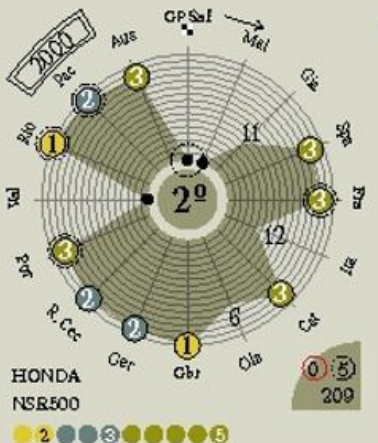
# Two major objectives

- Discuss the **nature of content** in unstructured data within a semantic perspective over natural language

  - What constitute a useful notion of content within unstructured data collections (that are largely made of linguistic information, e.g. Web pages or infographics)

  - What is **natural langauge semantics** and how can we model it formally?

  - What is the **meaning of a linguistic expression**?

- What is **the notion of document** that we can use within IR processes

  - Nature and role of document information

  - Relationship between a declarative view on content wrt an operational view of content

  - How this has to do with IR and ML?

# Overview

- Documents in Information Retrieval

  - Information, Representation, (re)current challenges, success(and unsuccess)ful stories

- Information and Content

  - Natural Language Processing: introduction to the linguistic background

    - Natural Language and Content

    - NL Syntax

    - NL Semantics

  - Document Representation and IR models

- Summary

# Content in unstructured data

- Natural Language

    - Structure

    - Semantics

    - Types of semantics

    - Relationship with Machine Learning

- Examples:

    - NLU: natural language as a logic language

    - Providing more structure: Frame semantics
        - Logic, Frames and Scripts
        - The relationships between syntax and semantics

    - Semantic role labeling

# Natural Language & Ambiguity

# Ambiguità

- "*Dogs must be carried on this escalator*"

can be interpreted in a number of ways:

- *All dogs should have a chance to go on this wonderful escalator ride*

- *This escalator is for dog-holders only*

- *You can't carry your pet on the other escalators*

- *When riding with a pet, carry it*

# The NLP chain

## Levels of linguistic analyses

Pragmatics: what does it do?

Semantics: what does it mean?

Syntax: what is grammatical?

*natural language utterance*

# Analogy with artificial languages

Syntax: no compiler errors

Semantics: no implementation bugs

Pragmatics: implemented the right algorithm

Different syntax, same semantics (5):

$$2 + 3 \Leftrightarrow 3 + 2$$

Same syntax, different semantics (1 and 1.5):

$$3 \ / \ 2 \ (\text{Python 2.7}) \ \nLeftrightarrow \ 3 \ / \ 2 \ (\text{Python 3})$$

Good semantics, bad pragmatics:

correct implementation of deep neural network
for estimating coin flip prob.

# Ambiguity and Linguistic Levels

- Semantics
- Syntax
- Morphology
- Phonology

can/can

eat cake with fork

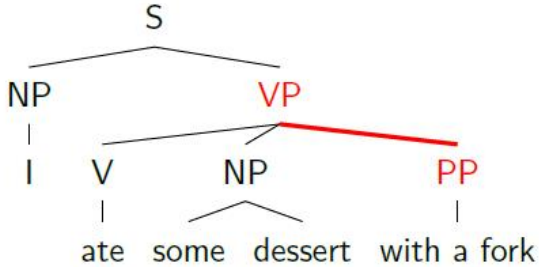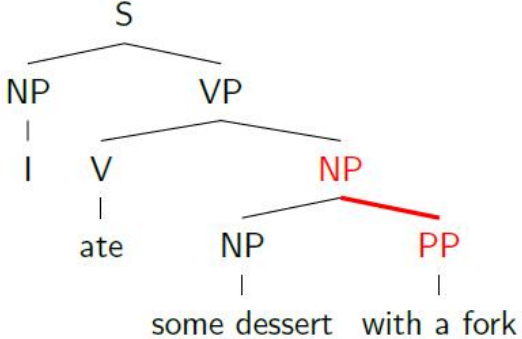earth observation satellite
Eco's book

del (pane)
/del (libro)

compro la borsa
in pelle

il timore dei manager

# Grammars & Ambiguity

# Summary

- IR models necessary in Web mining depend on the ways unstructured data can be made avilable for filtering, classification, retrieval and ranking tasks

- A semantic model for the content of unstructured data is strongly dependent on the linguistic nature of these latter
  - Facts, Entities, Relations, Thematic areas, Subjective information are always rooted in a form of rather free linguistic description

- Studies in Linguistics have provided the basic notion for dealing with the meaning of Natural Language expressions
  - Levels
  - Basic paradigms: lexical description, grammars, logic as a meaning representation language

# References

- AI & Robotics. «Robot Futures», Ilah Reza Nourbakhsh, MIT Press, 2013

- NLP & ML:

  - «Statistical Methods for Speech Recognition», F. Jelinek, MIT Press, 1998

  - «Speech and Language Processing", D. Jurafsky and J. H .Martin, Prentice-Hall, 2009.

  - "Foundations of Statistical Natural Language Processing, Manning & Schutze, MIT Press 2001.

- Sitografia:

  - SAG, Univ. Roma Tor Vergata: http://sag.art.uniroma2.it/

  - Reveal s.r.l.: http://www.revealsrl.it/