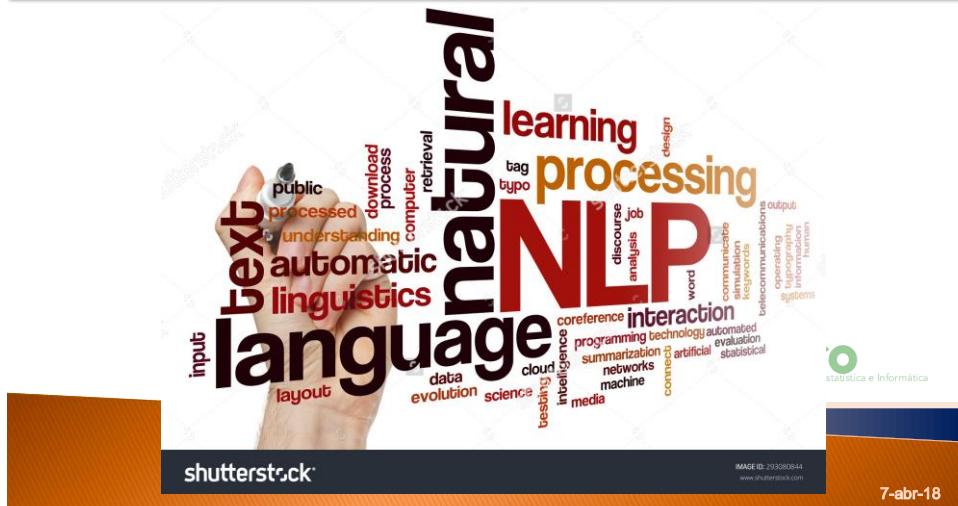
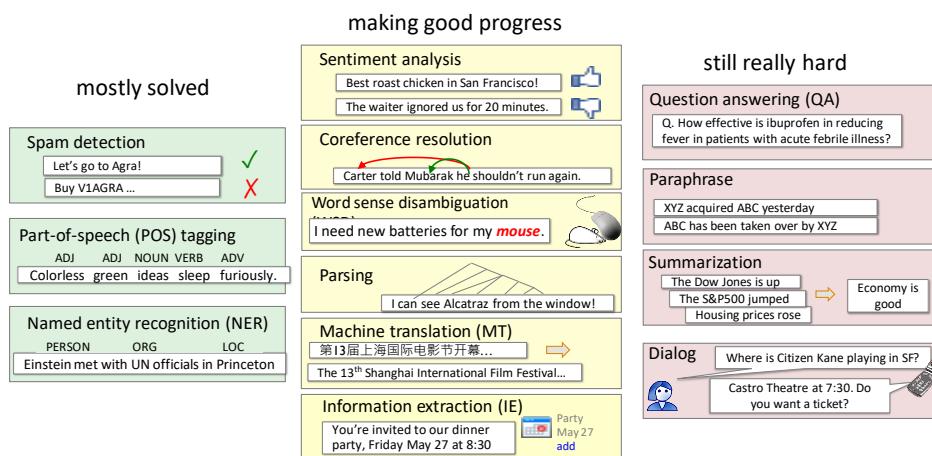


## Last Lecture: Main points to recall



## Emerging Applications in TM



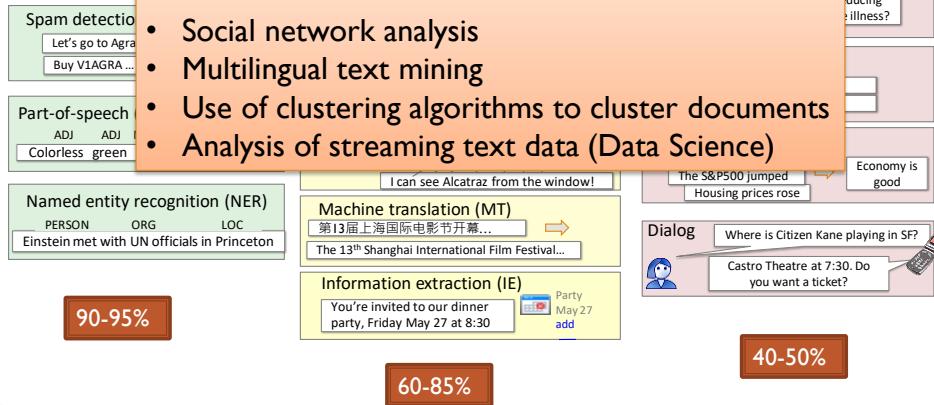
## Emerging Applications in TM



making good progress

### And many others...

- Social network analysis
- Multilingual text mining
- Use of clustering algorithms to cluster documents
- Analysis of streaming text data (Data Science)



3

## Why NL understanding is difficult?



### non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #never say never & you yourself should never give up either♥

### segmentation issues

the New York-New Haven Railroad  
the New York New Haven Railroad

### idioms

dark horse  
get cold feet  
lose face  
throw in the towel

### neologisms

unfriend  
Retweet  
bromance

### world knowledge

Mary and Sue are sisters.  
Mary and Sue are mothers.

### tricky entity names

Where is *A Bug's Life* playing ...  
*Let It Be* was recorded ...  
... a mutation on the *for* gene ...

### Ambiguity

include your children when baking cookies

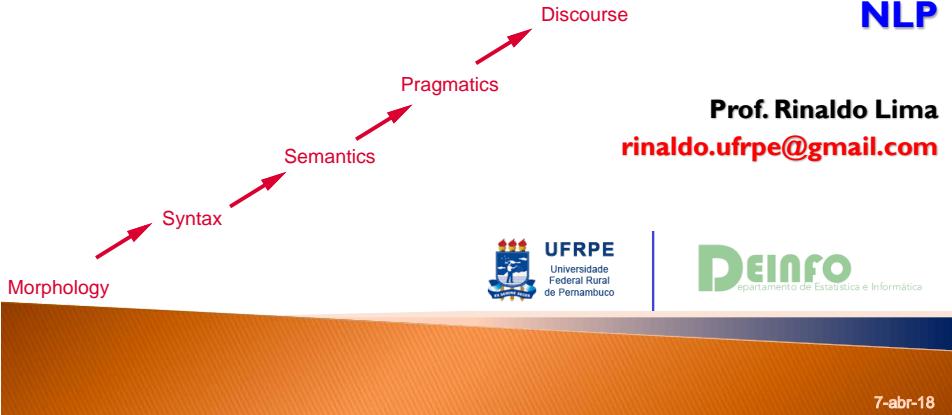
two soviet ships collide, one dies

But that's what makes it fun!

4

# Processamento de Linguagem Natural

## Aula 02: Introduction to NLP



## Contents



- ▶ **NLP Corpora**
- ▶ **NLP subtasks**
  1. Tokenization
  2. Sentence Splitting
  3. Lemmatization/Stemming
  4. POS tagging
  5. NER
  6. Parsing
  7. Chunking
  8. Coreference Resolution
  9. Semantic Roles
- ▶ **Applying NLP subtasks in many TM applications**

# NLP CORPORA (DATASETS)



**UFRPE**  
Universidade  
Federal Rural  
de Pernambuco

**DeInfo**  
Departamento de Estatística e Informática

7-abr-18

## What is a corpus?



A corpus is a finite machine-readable body of naturally occurring text, selected according to specified criteria, eg:

- ▶ **Language and type:** English/German/Arabic/..., dialects v. “standard”
- ▶ **Genre and Domain:** 18th century novels, newspaper text, software manuals...
- ▶ **Web as Corpus:** URL “domain” = country:.uk .ar
- ▶ **Media:** “Written” Text, Audio, Transcriptions, Video.
- ▶ **Size:** 1000 words, 50K words, 1M words, 100M words,

## Corpora Examples : Brown and LOB



**Brown:** The most famous corpus! (well, first widely-used corpus)

- ▶ A balanced corpus: representative of a whole language
- ▶ Brown: published American English from 1960s (newspapers, books, etc)
- ▶ 1 million words, POS tagged.

It will be used in our assignments

**LOB:** Lancaster-Oslo Bergen corpus

- ▶ British English version
- ▶ published British English text from 1960s

10

## Some recent corpora



**Corpus features: Size, Domain, Language**

- ▶ **British National Corpus:** 100M words, balanced British English
- ▶ **Newswire Corpus:** 600M words, newswire, American English
- ▶ **UN or EU proceedings:** 20M+ words, legal, 10 language pairs
- ▶ **Penn Treebank:** 2M words, newswire American English
- ▶ **Web:** 8 billion(?) words, many domains and languages
- ▶ **Web-as-Corpus:**
  - harvest your own corpus from WWW,
  - via “seed terms” → Google API → web-pages → Corpus!

11

## Corpus Annotation



**Annotation** is a process in which linguistics experts add (linguistic) information to the corpus that is not explicitly there (increases utility of a corpus), e.g.:

- **Text Headers**: meta-data for each text: author, date, type,...
- **Part-of-speech tag** for each word (very common!).
- **Syntactic structure**: parse-tree for each sentence
- **Word Sense** label for each word
- **Prosodic information**: pauses, rise and fall in pitch, etc.

12

## Annotation with XML: Text Centric vs Data Centric XML



Listing 2.1 Purchase Order in XML

```
<po id="43871" submitted="2001-10-05">
  <billTo>
    <company>The Skateboard Warehouse</company>
    <street>One Warehouse Park</street>
    <street>Building 17</street>
    <city>Boston</city>
    <state>MA</state>
    <postalCode>01775</postalCode>
  </billTo>
  <shipTo>
    <company>The Skateboard Warehouse</company>
    <street>One Warehouse Park</street>
    <street>Building 17</street>
    <city>Boston</city>
    <state>MA</state>
    <postalCode>01775</postalCode>
  </shipTo>
</order>

      <H1>Skateboard Usage Requirements</H1>
      <P>In order to use the <B>FastGlide</B> skateboard you have to
      have:</P>
      <LIST>
        <ITEM>A strong pair of legs.</ITEM>
        <ITEM>A reasonably long stretch of smooth road surface.</ITEM>
        <ITEM>The impulse to impress others.</ITEM>
      </LIST>
      <P>If you have all of the above, you can proceed to <LINK
      HREF="Chapter2.xml">Getting on the Board</LINK>.</P>
```

DATA CENTRIC XML  
DOC

TEXT CENTRIC XML  
DOC

13



## Annotation Example: POS tagging

Some texts are annotated with Part-of-speech (POS) tags.

- ▶ POS tags encode simple grammatical functions.

```
<s><w pos=RN> Here </w> <w pos=BEZ> is </w> <w pos=IN>
a </w><w pos>NN> sentence </w>.</s>
```

**Several tag sets are available:**

- ▶ Brown tag set (**87 tags**) in Brown corpus
  - Penn tag set (45 tags) in Penn Treebank
- ▶ CLAWS / LOB tag set (**132 tags**) in LOB corpus
  - CLAWS c5 tag set (62 tags) in BNC (British National Corpus)
- ▶ Tagging is usually done automatically (then proofread and corrected by humans)

14



## Corpora Annotation: Stand-off vs Inline

### ▶ **Inline:**

- Data and metadata (annotation or markup) are combined together into a single file.

### ▶ **Stand-off:**

- Metadata is stored in a separate document, using reference **anchors**
- Alignment based on token or **character offsets**
- Primary data is left **untouched**

15

## Inline Annotation Example - POS tagging



John went to Paris yesterday. He loved the excursion.

```
John_NNP went_VBD to_TO
Paris_NNP yesterday_NN ._
He_PRP loved_VBD the_DT
excursion_NN ._.
```

Horizontal format

```
John_NNP
went_VBD
to_TO
Paris_NNP
yesterday_NN
. .
He_PRP
loved_VBD
the_DT
excursion_NN
. .
```

Vertical format

16

## Stand-off Annotation Example - POS tagging



1234567890123456789012345678901234567890123  
 1            2            3            4            5  
 John went to Paris yesterday. He loved the excursion.

Character Offsets

```
1 4 NNP
6 9 VBD
11 12 TO
14 18 NNP
20 28 NN
29 29 .
31 32 PRP
34 38 VBD
40 42 DT
44 52 NN
53 53 .
```

17

## Stand-off Annotation: a more complex example



text.xml:

```
...</header>
<body>Fürchtet euch nicht ! Die einstige Fußball-Weltmacht zittert vor einem Winzling
. Mit seinem Tor zum 1:0 für die Ukraine stürzte der 1,62 Meter große Gennadi Subow die
deutsche Nationalelf...</body>
```

tok.xml:

```
<markList type="token" xml:base="text.xml">
<mark id="t1" xlink:href="#xpmarker(string-range(/body,'',0,8))"/> <!-- Fürchtet -->
<mark id="t2" xlink:href="#xpmarker(string-range(/body,'',9,4))"/> <!-- euch -->
<mark id="t3" xlink:href="#xpmarker(string-range(/body,'',14,5))"/> <!-- nicht -->
<mark id="t4" xlink:href="#xpmarker(string-range(/body,'',20,1))"/> <!-- ! -->
<mark id="t5" xlink:href="#xpmarker(string-range(/body,'',22,3))"/> <!-- Die -->
...
```

**"Layers"** of annotations in several separated XML files

infStat.xml:

```
<featList type="information_status" xml:base="tok.xml">
...
<!-- euch: new -->
<feat xlink:href="#xpmarker(id('t2'))" value="type.infStat.xml#new"/>
<!-- Die einstige Fußball-Weltmacht: accessible -->
<feat xlink:href="#xpmarker(id('t5')/range-to(id('t7')))" value="type.infStat.xml#acc"/>
...

```

type.infStat.xml:

```
<typeList type="information_status">
<type id="giv" name="giv" descr="The referent is given in the discourse."/>
<type id="new" name="new" descr="The referent is new in the discourse."/>
<type id="acc" name="acc" descr="The referent is accessible."/>
...
</type>
...
```

18

## Corpora Annotation Tools



brat rapid annotation tool

<http://brat.nlplab.org/>

online environment for collaborative text annotation

The screenshot shows the Brat interface with a text document. Annotations are made using colored boxes (Org, Money) and arrows. A tooltip "Connect by drag and drop" is visible.

Learn more:

- [What is it?](#)
- [What can you do with it?](#)
- [What does it do?](#)
- [What do I need to run it?](#)

[Try brat online](#)

(username: "crunchy", password: "frog")

Take a tutorial: [news](#) ([reset](#)), [bio](#) ([reset](#)).

Runs in your browser: no installation required

Intuitive annotation visualization and editing.

Create your own local brat installation:

[Download v1.3](#)[\(MD5, SHA512, Repository \(GitHub\), Older versions\)](#)

Manage your own annotation effort

Easy to set up: [Installation instructions](#)[Instructions for upgrading to v1.3 \(Crunchy Frog\)](#)Open source ([MIT License](#))

Current version: v1.3 Crunchy Frog (2012-11-08).

19

## Corpora Annotation Tools



### Annotation & Text-Processing Tools

#### Text Coding/(Manual) Annotation Programs/Text-analysis Tools & Search Engines

Please note that some of these programs produce XML files in *standoff* format, which separates the text into different linked levels. The advantage of this type is possible to link various types of annotation to the same set of data, but the disadvantage is that it's usually not possible to 'interact' directly with that data unless interface it's been created with. In other words, creating standoff annotations usually ties one into specific programs and the functionality they provide for merging

DART (Dialogue Annotation & Research Tool)

An annotation & analysis tool designed for the semi-automatic annotation of spoken (transcribed) dialogues on the levels of syntax, pragmatics (surface) polarity, semantics (topics), & semantic-pragmatics (modes, IFIDs).

DART produces annotations in what I refer to as 'Simple XML', a highly readable format that still allows the corpus user to 'interact' with the data by performing corrections, add additional annotations, etc.

As a research tool, DART also offers facilities for the creation of dialogue corpora & their associated analysis resources, a built-in corpus editor, analysis, as well as speech-act statistics.

For a detailed description, see my recent article in [Corpus Linguistics and Linguistic Theory](#) about version 1.

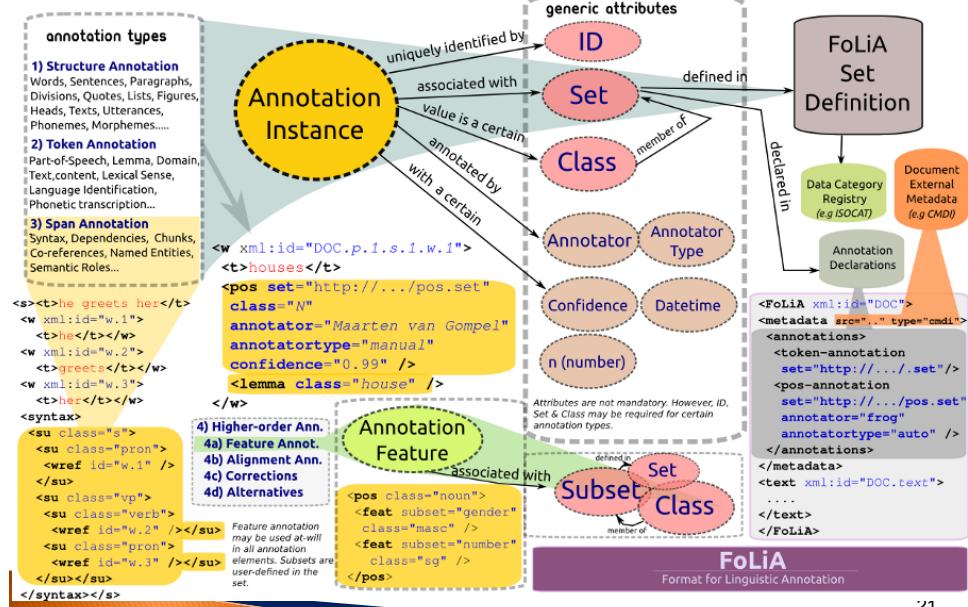
Version 2, as well as DART's 'big brother', the Text Analysis & Research Tool (TART), are currently under development.

[http://martinweisser.org/corpora\\_site/annotation\\_tools.html](http://martinweisser.org/corpora_site/annotation_tools.html)

20

## FOLIA – Format for Linguistic Annotation (2012)

Radboud University - <http://proycon.github.io/folia/>



21

## NLP Subtasks



**DeInfo**  
Departamento de Estatística e Informática

7-abr-18

### Levels of NLP Analysis



- ▶ Morphology: Concerns the way words are built up from smaller meaning bearing units. (come(s), walk(ed))
- ▶ Syntax: concerns how words are put together to form correct sentences and what structural role each word has.
- ▶ Semantics: concerns what words mean and how these meanings combine in sentences to form sentence meanings.
- ▶ Pragmatics: concerns how sentences are used in different situations and how its use affects the interpretation of the sentence.
- ▶ Discourse: concerns how the immediately preceding sentences affect the interpretation of the next sentence.

23

## Tokenization (Word Segmentation)



**Tokenization** is a processing step where the input text is automatically divided into units called tokens where each is either a word or a number or a punctuation mark

- ▶ In some written languages (e.g. Chinese) words are not separated by spaces.
- ▶ Even in English, characters other than white-space can be used to separate words [e.g., ; . - : ( ) ]

24

## Issues in Tokenization



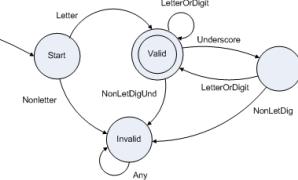
- ▶ Finland's capital → Finland Finlands Finland's ?
- ▶ what're, I'm, isn't → What are, I am, is not
- ▶ Hewlett-Packard → Hewlett Packard ?
- ▶ state-of-the-art → state of the art ?
- ▶ Lowercase → lower-case lowercase lower case ?
- ▶ San Francisco → one token or two?
- ▶ m.p.h., PhD. → ??

25

## Tokenization Algorithm for English



### A finite automata



```

Initialize:
  Set Stream to the input text string
  Set currentPosition to 0 and internalQuoteFlag to false
  Set delimiterSet to ', . ; ! ? () <>+ \n \t space'
  Set whiteSpace to \t\n space

Procedure getNextToken()
  L1: cursor := currentPosition; ch := charAt(cursor);
  If ch = endOfStream then return null; endif
  L2: while ch is not endOfStream nor instanceOf(delimiterSet) do
    increment cursor by 1; ch := charAt(cursor);
  endwhile
  If ch = endOfStream then
    If cursor = currentPosition then return null; endif
  endif
  If ch is whiteSpace then
    If currentPosition = cursor then
      increment currentPosition by 1 and goto L1;
    else
      Token := substring(Stream,currentPosition,cursor-1);
      currentPosition := cursor+1; return Token;
    endif
  endif
  If ch = '*' then
    If charAt(cursor-1) = instanceOf(delimiterSet) then
      internalQuoteFlag := true; increment currentPosition by 1; goto L1;
    endif
    If charAt(cursor+1) != instanceOf(delimiterSet) then
      increment cursor by 1; ch := charAt(cursor); goto L2;
    else
      internalQuoteFlag = true then
        Token := substring(Stream,currentPosition,cursor-1);
        internalQuoteFlag := false;
      else
        Token := substring(Stream,currentPosition,cursor);
      endif
      currentPosition := cursor+1; return Token;
    endif
  if cursor = currentPosition then
    Token := ch; currentPosition := cursor+1;
  else
    Token := substring(Stream,currentPosition,cursor-1);
    currentPosition := cursor;
  endif
  return Token;
endprocedure
  
```



26

## Sentence Segmentation (Splitting)



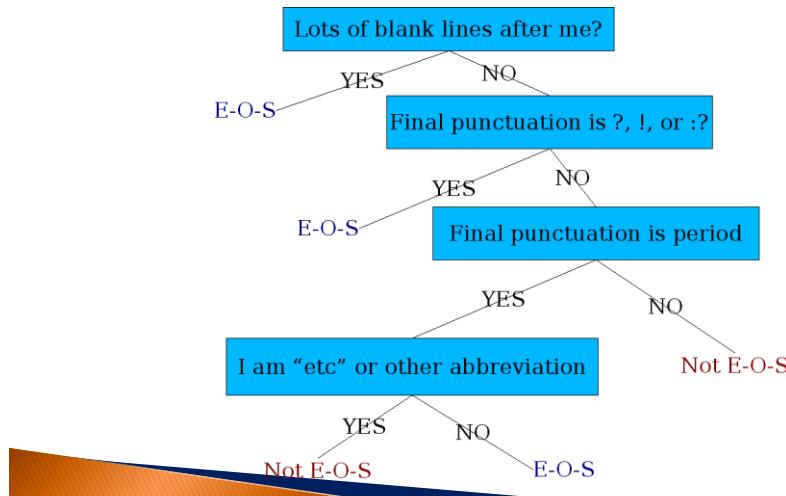
- ▶ !, ? are relatively unambiguous
- ▶ Period “.” is quite ambiguous
  - Sentence boundary
  - Abbreviations like Inc. or Dr.
  - Numbers like .02% or 4.3
- ▶ Build a **binary classifier**
  - Looks at a “.”
  - Decides EndOfSentence/NotEndOfSentence
  - Classifiers:
    - hand-written rules,
    - regular expressions, or
    - machine-learning

27

## Sentence Segmentation (Splitting)



### Determining if a word is end-of-sentence: a Decision Tree



28

**Input:** a text with periods  
**Output:** same text with End-of-Sentence (EOS) periods identified

### Sentence Splitting Algorithm for English



### Ruled-based

#### Rules:

- All ? ! are EOS
- If " or ' appears before period, it is EOS
- If the following character is not white space, it is not EOS
- If ) ] before period, it is EOS
- If the token to which the period is attached is capitalized and is < 5 characters and the next token begins uppercase, it is not EOS
- If the token to which the period is attached has other periods, it is not EOS
- If the token to which the period is attached begins with a lowercase letter and the next token following whitespace is uppercase, it is EOS
- If the token to which the period is attached has < 2 characters, it is not EOS
- If the next token following whitespace begins with \$ { [ " ' it is EOS
- Otherwise, the period is not EOS

29

## Morphological Analysis



- ▶ **Morphology** is the field of linguistics that studies the internal structure of words
- ▶ A **morpheme** is the smallest linguistic unit that has semantic meaning
  - e.g. “carry”, “pre”, “ed”, “ly”, “s”
- ▶ Morphological analysis is the task of segmenting a word into its morphemes:
  - carried ⇒ carry + ed (past tense)
  - independently ⇒ in + (depend + ent) + ly
  - Googlers ⇒ (Google + er) + s (plural)
  - unlockable ⇒ un + (lock + able) ?
    - ⇒ (un + lock) + able ?

▶ **Morphemes:**

- The small meaningful units that make up words
- **Stems:** The core meaning-bearing units
- **Affixes:** pieces that adhere to stems

30

## Tokenization: language issues



- ▶ **French**
  - L'ensemble → one token or two?
    - L?    L'?    Le ?
- ▶ **German noun compounds are not segmented**
  - Lebensversicherungsgesellschaftsanstellter
    - ‘life insurance company employee’
  - German information retrieval needs compound splitter

31

## Lemmatization



- ▶ Reduce inflections or variant forms to base form.

### Dictionary form of a word.

- *am, are, is* → *be*
- *car, cars, car's, cars'* → *car*

- ▶ *the boy's cars are different colors* → ***the boy car be different color***

- ▶ Most used in Machine translation

- Spanish *quiero* ('I want'), *quieres* ('you want') same lemma as *querer* 'want'

32

## Stemming



- ▶ Reduce terms to their stems or roots

- ▶ **Stemming is the crude chopping of affixes**

- language dependent
- e.g., *automate(s), automatic, automation* all reduced to *automat.*

for example *compressed* and *compression* are both accepted as equivalent to *compress*.



for exempl **compress** and **compress** ar both accept as equival to compress

33

## Porter's algorithm: the most used English stemmer



### Step 1a

sses → ss	caresses → caress
ies → i	ponies → poni
ss → ss	caress → caress
s → ø	cats → cat

### Step 1b

(*v*)ing → ø	walking → walk
	sing → sing
(*v*)ed → ø	plastered → plaster
...	

### Step 2 (for long stems)

ational → ate	relational → relate
izer → ize	digitizer → digitize
ator → ate	operator → operate
...	

### Step 3 (for longer stems)

al → ø	revival → reviv
able → ø	adjustable → adjust
ate → ø	activate → activ
...	

## Why only strip -ing if there is a vowel?

(\*v\*)ing → ø walking → walk  
sing → sing

34

**Input:** a text token and a dictionary  
**Doubling consonants:** b d g k m n p r l t



**Rules:**

- If token length < 4 **return** token
- If token is number **return** token
- If token is acronym **return** token
- If token in dictionary **return** the stored stem
- If token ends in s'  
strip the ' and **return** stripped token
- If token ends in 's'  
strip the 's' and **return** stripped token
- If token ends in "is", "us", or "ss" **return** token
- If token ends in s'  
strip s, check in dictionary, and **return** stripped token if there
- If token ends with es  
strip es, check in dictionary, and **return** stripped token if there
- If token ends in les  
replace ies by y and **return** changed token
- If token ends in s'  
strip s and **return** stripped token
- If token doesn't end with ed or ing and **return** token
- If token ends with ed  
strip ed, check in dictionary and **return** stripped token if there
- If token ends in led  
replace led by y and **return** changed token
- If token ends in eed  
remove d and **return** stripped token if in dictionary
- If token ends with ing  
strip ing (if length > 5) and **return** stripped token if in dictionary
- If token ends with ing and length ≤ 5 **return** token
- // Now we have SS, the stripped stem, without ed or ing and it's not in the dictionary (otherwise algorithm would terminate)
- If SS ends in doubling consonant  
strip final consonant and **return** the changed SS if in dictionary
- If doubling consonant was l **return** original SS
- If no doubled consonants in SS  
add e and **return** changed SS if in dictionary
- If SS ends in c or z, or there is a g or l before the final doubling consonant  
add e and **return** changed SS
- If SS ends in any consonant that is preceded by a single vowel  
add e and **return** changed SS

**return** SS

## Stemming Algorithm for English



## Ruled-based

35

## Part-Of-Speech (POS) Tagging



- ▶ Annotate each word in a sentence with a part-of-speech tag

I ate the spaghetti with meatballs.

Pro V Det N Prep N

John saw the saw and decided to take it to the table.

PN V Det N Con V Part V Pro Prep Det N

- ▶ Useful for subsequent syntactic parsing and word sense disambiguation

36

## Part Of Speech (POS) Tagging



Alphabetical list of part-of-speech tags used in the Penn Treebank Project:

POS labels  
were defined  
by the  
Penn Treebank  
project

Number	Tag	Description
1.	CC	Coordinating conjunction
2.	CD	Cardinal number
3.	DT	Determiner
4.	EX	Existential <i>there</i>
5.	FW	Foreign word
6.	IN	Preposition or subordinating conjunction
7.	JJ	Adjective
8.	JJR	Adjective, comparative
9.	JJS	Adjective, superlative
10.	LS	List item marker
11.	MD	Modal
12.	NN	Noun, singular or mass
13.	NNS	Noun, plural
14.	NNP	Proper noun, singular
15.	NNPS	Proper noun, plural
16.	PDT	Predeterminer
17.	POS	Possessive ending
18.	PRP	Personal pronoun

37

## Phrase Chunking



- ▶ Find all non-recursive noun phrases (NPs), verb phrases (VPs), and Prepositional phrases in a sentence.
  - [NP I] [VP ate] [NP the spaghetti] [PP with] [NP meatballs].
  - [NP He] [VP reckons] [NP the current account deficit] [VP will narrow] [PP to] [NP only # 1.8 billion] [PP in] [NP September]

38

## Named Entity Recognition (NER)



A **Named Entity** denotes

The most common  
provided by several tools

- ▶ named (PERSON, LOCATION, ORGANIZATION, MISC),
- ▶ numerical (MONEY, NUMBER, ORDINAL, PERCENT), and
- ▶ temporal (DATE, TIME, DURATION, SET) entities.
- ▶ Named entities are **recognized** using **supervised machine learning algorithms**
- ▶ Numerical entities are recognized using a rule-based system.
- ▶ Numerical entities that require normalization, e.g., dates.
- ▶ It can also be defined by Regex (regular expressions)

39

## Gazeteers – NER by lists



- ▶ A gazetteer consists in a set of lists containing names of entities such as **cities, organizations, days of the week**, etc.
- ▶ These lists are used to find occurrences of these names in text, e.g. for the NER task .
- ▶ The word ‘gazetteer’ is often used interchangeably for both the set of entity lists and for the processing resource that makes use of those lists

40

## Gazeteers - Examples



Below is a section of the list for units of currency:

Ecu  
 European Currency Units  
 FFr  
 Fr  
 German mark  
 German marks  
 New Taiwan dollar  
 New Taiwan dollars  
 NT dollar  
 NT dollars  
 Real  
 ...

Usando a Web para construir gazetters.



 **GeoNames**

The GeoNames geographical database covers all countries and contains over eleven million placenames that are available for download free of charge.



A screenshot of the GeoNames search interface. It features a search bar with the placeholder "enter a location name, ex: "Paris", "Mount Everest", "New York"" and a dropdown menu set to "all countries". Below the search bar are three buttons: "search", "show on map", and "[advanced search]".

41

# PARSING

## I. Constituent



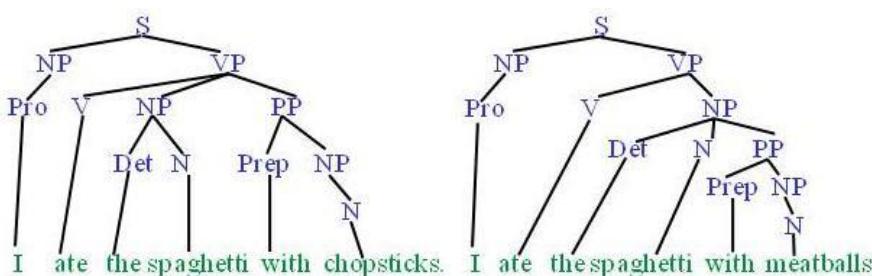
**DeInfo**  
Departamento de Estatística e Informática

7-abr-18

### Syntactic Parsing (Constituent)



- It produces the **most probable syntactic parse tree** for a sentence
- It is usually based on a **context-free grammar**
- It organizes the words in a sentence in a **hierarchical structured** according to their internal phrases



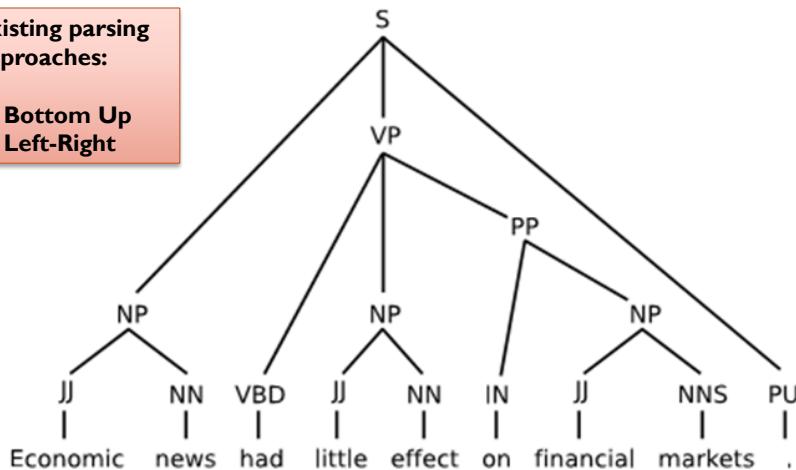
43

## Constituent Parsing Tree - Example



**Existing parsing approaches:**

- Bottom Up
- Left-Right



44

## PARSING 2. Dependency

## Dependency Parsing



**Dependency Parsing** is an approach based on the linguistic theory of **dependency grammar**.

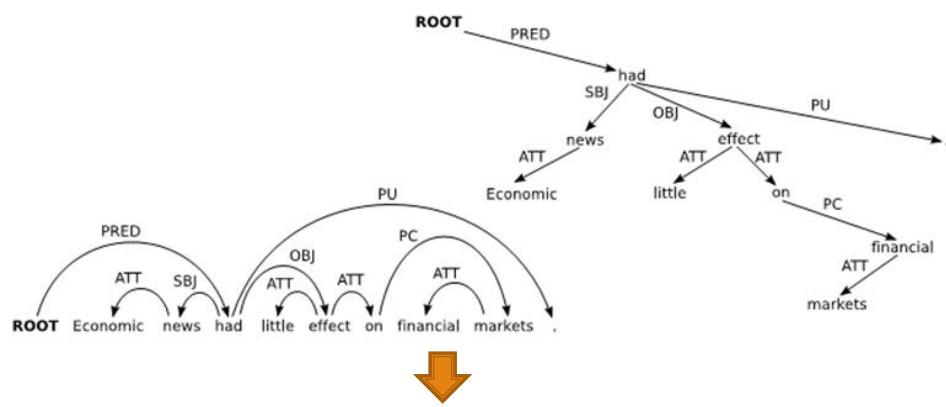
The basic assumption:

syntactic structure consists essentially of words connected by **asymmetric binary relations** of dependency relationships between them

•46

46

## Dependency Parsing Example



### Generated Relations:

- amod ( news , Economic )
- nsubj ( had, news )
- amod ( effect , little)
- dobj ( had , effect )
- amod ( markets , financial )
- prep\_on ( effect , markets )

The verb "had" has the subject "news"

47

## Coreference Resolution (Anaphora)



- ▶ Determine which phrases in a document refer to the same underlying entity.
  - John put the **carrot** on the **plate** and ate **it**.
  - **Bush** started the war in Iraq. But the **president** needed the consent of Congress.
- ▶ It is also called **anaphora resolution**
- ▶ CR types:
  - Pronominal
  - nominal
- ▶ Some cases require difficult reasoning

48

## Semantic Role Labeling



- ▶ Semantic Role Labeling (**SRL**) aims at automatically assigning semantic roles for each predicate in a sentence.
- ▶ SRL is also called **thematic role labeling**, or **shallow semantic parsing**.
- ▶ It determine which constituents in a sentence are **semantic arguments for a given predicate**, and then determining the appropriate role for each of those arguments.
- ▶ SRL has the potential to improve performance in any language understanding task
- ▶ The state-of-the-art approaches to SRL are based on supervised machine learning and **Semantic Roles** acquired from a corpus

49

## Semantic Role Labeling with FrameNet



FrameNet is a computational resource that provides **Frame Semantics** for NLP

- ▶ The state-of-the-art approaches to SRL are based on supervised machine learning
- ▶ The central idea of **Frame Semantics** is that a word meanings must be described in relation to semantic frames (**syntactic relations**)
- ▶ **FrameNet** is a collection of facts that specify "characteristic features, attributes, and functions of a word, and its characteristic interactions with things necessarily or typically associated with it."



## Framenet Original Project (English)



About FrameNet ▾ Documentation ▾ FrameNet Data ▾ Related Projects ▾ Bibliography

*FrameNet maps meaning to form in contemporary English through the theory of Frame Semantics.*

### Recent News

- Multilingual FrameNet Project
- Website makeover
- Release 1.7 data included in NLTK
- FrameNet at the Linguistic Society of America 2017

[GET THE DATA](#)

<https://framenet.icsi.berkeley.edu/fndrupal/>

## FrameNet Example



### Apply\_heat: Frame Elements

Cooker

Food

Temperature\_setting

Duration

Heating\_instrument

Medium

Lila **FRIED** the eggs in a copper pan.

52

## Semantic Role Labeling (SRL): Example



- ▶ For each clause, determine the semantic role played by each noun phrase that is an argument to the verb.

agent patient source destination instrument

- John drove Mary from Austin to Dallas in his Toyota Prius.
- The hammer broke the window.

53

53

# FrameNet Brazilian Portuguese



[Inicial](#) [Projetos](#) [Dados](#) [Pessoas](#) [Publicações](#) [m.knob](#) [FN-Br pelo Mundo](#) [Oportunidades](#) [Dicionário da Copa do Mundo](#)

A FrameNet Brasil é um laboratório de Linguística Computacional sediado na Universidade Federal de Juiz de Fora. Nossa missão é a de desenvolver soluções computacionais para a Compreensão de Língua Natural com base na Semântica de Frames e na Gramática das Construções.



m.knob.com  
Blend in

**FrameNet Brasil lança guia turístico virtual**

O m.knob, Multilingual Knowledge Base (Base de Conhecimentos Multilingüe), é um guia virtual que fornece recomendações personalizadas de atrações e eventos turísticos para seus usuários. Fruto de um projeto de pesquisa interdisciplinar conduzido, no Laboratório FrameNet Brasil, por pesquisadores e alunos de graduação e pós dos departamentos de Letras e Ciência da Computação da Universidade Federal de Juiz de Fora, o m.knob alia processamento [...] [Leia mais](#)

**Publicado volume temático da Revista Veredas sobre Semântica de Frames e Suas Aplicações Tecnológicas**

Acaba de ser publicado o volume 17 (1) da Revista Veredas, o qual tem como tema "Semântica de Frames e Suas Aplicações Tecnológicas". Além de trazer uma qualificada amostra dos trabalhos desenvolvidos na área por pesquisadores de diversas instituições localizadas no Japão, Suécia, Alemanha, Estados Unidos e Brasil, o volume é uma homenagem à contribuição do linguista

[youtube.com/framenetbrasil](http://youtube.com/framenetbrasil)

**Thomas Herbst - Restricting t...**



International conference on construction grammar

<http://www.ufjf.br/framenetbr/>

54

# FrameNet Brazilian Portuguese: Example



**Frames**

Search Frame

- Abandono
- Absorção\_de\_calor
- Abundância
- Abundar\_com
- Abusar
- Acetitar\_ou\_recusar\_a\_agir
- Acessórios Vestuário
- Acidente
- Acomodação
- Acompanhamento
- Acordar
- Adição
- Adjacency
- Adornar
- Adotar\_seleção
- Adquirir
- Afirmar\_ou\_negar
- Agir\_intencionalmente
- Agregado
- Agricultura
- Agrupar

**Abandono**

**Definição**  
Um Agente deixa pra trás um Tema, de modo que este não esteja mais sob seu controle ou sua propriedade. Carolina abandona a criança é considerado um crime sério.

**Exemplo(s)**

**Elementos de Frame Nucleares**

**Agente [agent]** O Agente é a pessoa que deixa para trás o **Tema**.

**Tema [theme]** O Tema é a entidade que é abandonada pelo Agente.

**Elementos de Frame Não-Nucleares**

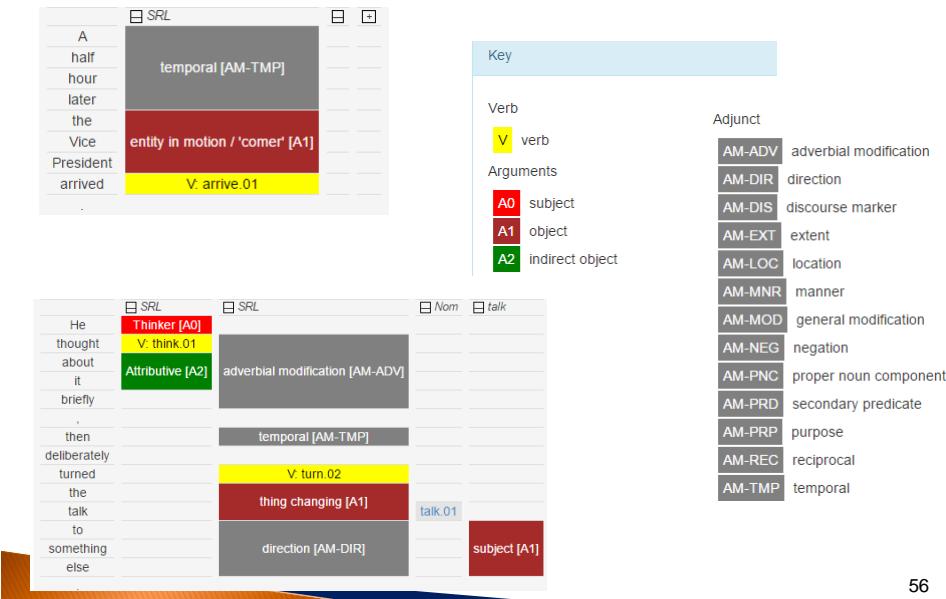
**Relações**

**Unidades Lexicais**

[abandonar.v](#) [abandono.n](#) [deixar.v](#)

<http://webtool.framenetbr.ufjf.br/index.php/fnbr/report/frame/main>

## SRL Example



56

## NLP Tools in English



- ▶ **Stanford CoreNLP** demos
    - <http://nlp.stanford.edu/software>
  - ▶ **GATE** (java)
    - <http://gate.ac.uk>
  - ▶ **OpenNLP**
    - <http://opennlp.apache.org>
  - **ILLINOIS Curator** demos
    - <http://cogcomp.cs.illinois.edu>
  - **LingPipe**
    - <http://alias-i.com/lingpipe-3.9.3/>
  - **SEMAFOR** // demos
    - <http://www.cs.cmu.edu/~ark/SEMAFOR/>
  - **FreeLing**
    - <http://nlp.lsi.upc.edu/freeling/node/30>
- JAVA
PYTHON
- **NLTK** - <http://www.nltk.org/>
  - **Spacy** <https://spacy.io/docs/use-word-vectors-similarities>
  - **CLIPS** <http://www.clips.ua.ac.be/pages/MBSP>
  - **TextBlob** -
    - <https://textblob.readthedocs.io/en/dev/index.html>

57

**CoreNLP**

version 3.7.0

- Overview
- Usage
- Annotators**
- Summary
- Annotator dependencies
- Tokenization
- Sentence Splitting
- Lemmatization
- Parts of Speech
- Named Entity Recognition
- RegexNER (Named Entity Recognition)

**NLP subtasks**

Table of Contents

- About
- Download
- Human languages supported
- Programming languages and operating systems
- License
- Citing Stanford CoreNLP in papers

About

Stanford CoreNLP provides a set of natural language analysis tools. It can give the base forms of words, their parts of speech, whether they are names of companies, people, etc., normalize dates, times, and numeric quantities, mark up the structure of sentences in terms of phrases and word dependencies, indicate which noun phrases refer to the same entities, indicate sentiment, extract particular or open-class relations between entity mentions, get quotes people said, etc.

Tokenization
Sentence Splitting
Lemmatization
Parts of Speech
Named Entity Recognition
RegexNER (Named Entity Recognition)
Constituency Parsing
Dependency Parsing
Coreference Resolution

58

**CoreNLP Pipelines**



## Annotator dependencies

PROPERTY NAME	ANNOTATOR CLASS NAME	REQUIREMENTS
tokenize	TokenizerAnnotator	None
cleanxml	CleanXmlAnnotator	tokenize
ssplit	WordsToSentenceAnnotator	tokenize
pos	POSTaggerAnnotator	tokenize, ssplit
lemma	MorphaAnnotator	tokenize, ssplit, pos
ner	NERClassifierCombiner	tokenize, ssplit, pos, lemma
regexner	RegexNERAnnotator	?
sentiment	SentimentAnnotator	?
parse	ParserAnnotator	tokenize, ssplit
depparse	DependencyParseAnnotator	tokenize, ssplit, pos
dcoref	DeterministicCorefAnnotator	tokenize, ssplit, pos, lemma, ner, parse

59

## CoreNLP Tools Demo

**Stanford CoreNLP**

Output format:  ▾

Please enter your text here:

```
Economic news had little effect on financial markets.
```

**Part-of-Speech:**

```
1 [JJ NN VBD JJ NN IN JJ NNS ]
1 Economic news had little effect on financial markets.
```

**Named Entity Recognition:**

```
1 Economic news had little effect on financial markets.
```

**Coreference:**

```
1 Economic news had little effect on financial markets.
```

**Basic Dependencies:**

```
1 [JJ amod NN nsubj VBD dobj JJ amod NN IN JJ amod NNS ]
1 Economic news had little effect on financial markets.
```

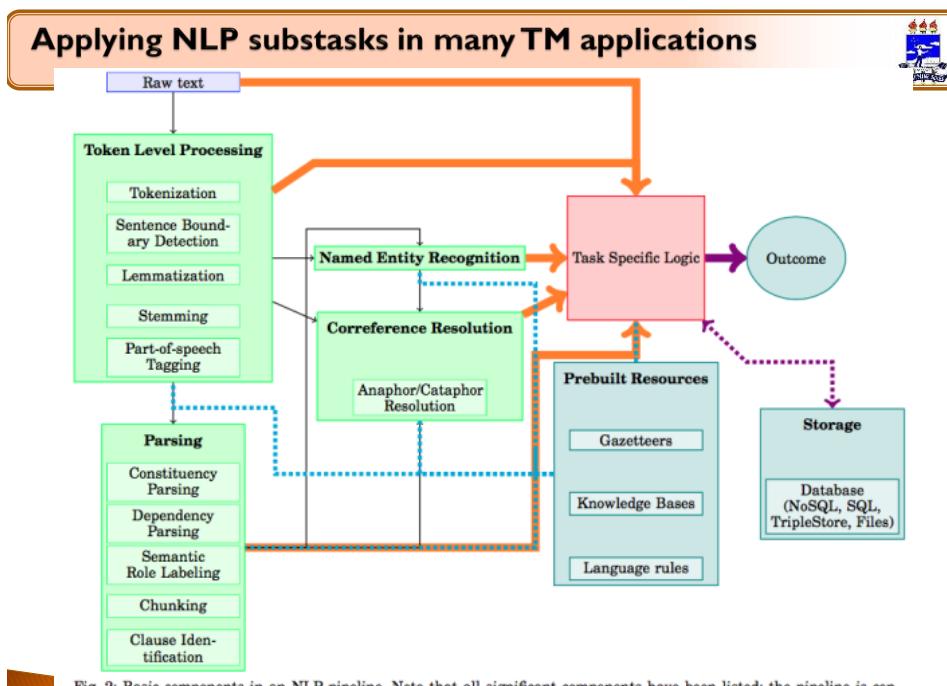


Fig. 2: Basic components in an NLP pipeline. Note that all significant components have been listed; the pipeline is constructed by customizing, adding, and/or removing these components

**Having Fun in the Weekend!**



## I<sup>a</sup>. Lista de Exercícios

62