# Paraphrasing

## Potential Applications for Plagiarism Detection

Marta Vila

CLiC – Universitat de Barcelona

WIQE 2010      Valencia

September 13, 2010

# Outline

## Outline

## Why paraphrasing?

Paraphrasing stands for (approximate) **sameness of meaning** between **different wordings**.

- He announced that the company wouldn't participate
  He stated that the company wouldn't participate

- My father built the house
  The house was built by my father

## Why paraphrasing?

Paraphrasing stands for (approximate) **sameness of meaning**
between **different wordings**.

- He announced that the company wouldn't participate
  He stated that the company wouldn't participate
- My father built the house
  The house was built by my father

### Why is paraphrasing relevant for plagiarism?

## Why paraphrasing?

Paraphrasing stands for (approximate) **sameness of meaning** between **different wordings**.

- He announced that the company wouldn't participate
  He stated that the company wouldn't participate
- My father built the house
  The house was built by my father

### Why is paraphrasing relevant for plagiarism?

Paraphrasing is a linguistic ability used in plagiarism.
To plagiarize involves, on many occasions, paraphrasing.

- Hamlet's pretense of madness
  Hamlet adopts a pretense of madness[1]

[1] http://www.princeton.edu/pr/pub/integrity/08/plagiarism/

## What is relevant in paraphrasing?

- Paraphrasing typologies
  - → A paraphrasing typology is also a plagiarism phenomena typology

## What is relevant in paraphrasing?

- Paraphrasing typologies
  - $\rightarrow$ A paraphrasing typology is also a plagiarism phenomena typology
- Paraphrasing corpora
- NLP approaches to paraphrasing
  - $\rightarrow$ A number of techniques and corpora can also be useful in plagiarism detection

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
A Paraphrasing Typology

## Outline

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

**Paraphrasing Complexity**
State of the Art
A Paraphrasing Typology

# Outline

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

**Paraphrasing Complexity**
State of the Art
A Paraphrasing Typology

# Paraphrasing Complexity

He announced that the company wouldn't participate
He stated that the company wouldn't participate

My father built the house
The house was built by my father

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

**Paraphrasing Complexity**
State of the Art
A Paraphrasing Typology

# Paraphrasing Complexity

Mary buttered her toast with expensive butter
Mary buttered her toast

Emma cried
Emma burst into tears

They got married last year
They got married in 2004

He announced that the company wouldn't participate
He stated that the company wouldn't participate

Patrick Ewing scored a personal season high of 41 points
Patrick Ewing scored 41 points. It was a personal season high

The pilot was having breakfast
The commander was having breakfast

My father built the house
The house was built by my father

I want some fresh air
Could you open the window?

She used to *only* eat hot dishes
She used to eat *only* hot dishes

Do you like French movies?
Are you interested in French cinema?

Steven made an attempt to stop playing Hearts
Steven attempted to stop playing Hearts

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

**Paraphrasing Complexity**
State of the Art
A Paraphrasing Typology

## Paraphrasing Complexity

Let's put it in order!

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
**State of the Art**
A Paraphrasing Typology

## Outline

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
**State of the Art**
A Paraphrasing Typology

## State of the Art

Paraphrasing typologies:

- Dras (1999): syntactic paraphrases

  $$LV + NP + inf\text{-}VP \;\rightleftharpoons\; V + inf\text{-}VP$$

  (2)    a.   Steven **made an attempt** to stop playing Hearts.
         b.   Steven **attempted** to stop playing Hearts.

- Fuijta (2005): lexical and structural paraphrases

  **Paraphrasing of common nouns to their synonyms** (Fujita and Inui, 2001; Yamamoto, 2002b; Okamoto *et al.*, 2003)

  s. ***kyuryo-ni*** *kinenkan-ga*      *kansei-shi-ta.*
     hill-LOC    a memorial hall-NOM   to build up-PAST
       A memorial hall was completed on the hill.

  t. ***takadai-ni*** *kinenkan-ga*      *kansei-shi-ta.*
     hill-LOC    a memorial hall-NOM   to build up-PAST

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
**State of the Art**
A Paraphrasing Typology

## State of the Art

Paraphrasing typologies:

- Bhagat (2009): lexical and structural paraphrases

  **Semantic implication:** Replacing a word or a phrase denoting an action, event etc. by a word or phrase denoting its possible future effect, in the appropriate context, results in a paraphrase of the original sentence of phrase. This may be accompanied by the addition/deletion of appropriate function words and sentence restructuring. This often generates a quasi-paraphrase.

  *Accompanying structural changes:* Substitution, Addition/Deletion, Permutation.

  *Example:*
  Google *is in talks to buy* YouTube. ⇔ Google *bought* YouTube.
  The Marines are *fighting* the terrorists. ⇔ The Marines are *eliminating* the terrorists.

- Žolkovskij and Mel'čuk (1965, 1966, 1967) and
  Milićević (2007): syntactic and semantic paraphrases

  a republic/a republican state:

  $$C_0 \Leftrightarrow \text{Gener}(C_0) \xrightarrow{\text{ATTR}} A_0/\text{Adv}_1(C_0)$$

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
**State of the Art**
A Paraphrasing Typology

## State of the Art

Problems:

- They do not cover the paraphrasing phenomenon as a whole
- Tied to a specific linguistic theory and formalism
  - $\rightarrow$ Meaning-Text Theory in Žolkovskij and Mel'čuk, and Milićević

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

## Outline

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# A Paraphrasing Typology

- Substitution
  - $\rightarrow$ He **announced** that the company wouldn't participate
    He **stated** that the company wouldn't participate

- Deletion
  - $\rightarrow$ Mary opened the **bottle of** wine
    Mary opened the wine

- Transformation
  - $\rightarrow$ My father built the house
    The house was built by my father

- Splitting
  - $\rightarrow$ Patrick Ewing scored a personal season high of 41 points
    Patrick Ewing scored 41 points. It was a personal season high

- Change of order
  - $\rightarrow$ She used to **only** eat hot dishes
    She used to eat **only** hot dishes[1]

---

[1]Examples extracted from Bhagat (2009), Dras (1999), Fuijta (2005), Pustejovsky (1995) and ourselves.

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# A Paraphrasing Typology

- Substitution
    - $\rightarrow$ He **announced** that the company wouldn't participate
      He **stated** that the company wouldn't participate

- Deletion (Addition)
    - $\rightarrow$ Mary opened the **bottle of** wine
      Mary opened the wine

- Transformation
    - $\rightarrow$ My father built the house
      The house was built by my father

- Splitting (Combining)
    - $\rightarrow$ Patrick Ewing scored a personal season high of 41 points
      Patrick Ewing scored 41 points. It was a personal season high

- Change of order
    - $\rightarrow$ She used to **only** eat hot dishes
      She used to eat **only** hot dishes

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# A Paraphrasing Typology

- Substitution
  - $\rightarrow$ He **announced** that the company wouldn't participate
    He **stated** that the company wouldn't participate

- Deletion
  - $\rightarrow$ Mary opened the **bottle of** wine
    Mary opened the wine

- Transformation
  - $\rightarrow$ My father built the house
    The house was built by my father

- Splitting
  - $\rightarrow$ Patrick Ewing scored a personal season high of 41 points
    Patrick Ewing scored 41 points. It was a personal season high

- Change of order
  - $\rightarrow$ She used to **only** eat hot dishes
    She used to eat **only** hot dishes

### Usually combined!

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# A Paraphrasing Typology

#### Substitution

$\rightarrow$ He **announced** that the company wouldn't participate
He **stated** that the company wouldn't participate

- Deletion
  - $\rightarrow$ Mary opened the **bottle of** wine
    Mary opened the wine

- Transformation
  - $\rightarrow$ My father built the house
    The house was built by my father

- Splitting
  - $\rightarrow$ Patrick Ewing scored a personal season high of 41 points
    Patrick Ewing scored 41 points. It was a personal season high

- Change of order
  - $\rightarrow$ She used to **only** eat hot dishes
    She used to eat **only** hot dishes

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

## Substitution

- Synonymy
  - $\rightarrow$ He **announced** that the company wouldn't participate
    He **stated** that the company wouldn't participate
    ✔ WordNet, dictionary of synonyms, thesaurus

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# Substitution

- Synonymy
  - → He **announced** that the company wouldn't participate
    He **stated** that the company wouldn't participate
    ✔WordNet, dictionary of synonyms, thesaurus
  - → **Since** it was sunny yesterday, the laundry dried well
    **As** it was sunny yesterday, the laundry dried well
    ✔Dictionary of synonyms, thesaurus

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
A Paraphrasing Typology

# Substitution

- Synonymy
  - $\rightarrow$ He **announced** that the company wouldn't participate
    He **stated** that the company wouldn't participate
    ✔ WordNet, dictionary of synonyms, thesaurus
  - $\rightarrow$ **Since** it was sunny yesterday, the laundry dried well
    **As** it was sunny yesterday, the laundry dried well
    ✔ Dictionary of synonyms, thesaurus
- Generalization
  - $\rightarrow$ I have been studying the reproduction of **cats** for ten years
    I have been studying the reproduction of **felines** for ten years
    ✔ WordNet, taxonomies

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# Substitution

- Synonymy
  - → He **announced** that the company wouldn't participate
    He **stated** that the company wouldn't participate
    ✔ WordNet, dictionary of synonyms, thesaurus
  - → **Since** it was sunny yesterday, the laundry dried well
    **As** it was sunny yesterday, the laundry dried well
    ✔ Dictionary of synonyms, thesaurus

- Generalization
  - → I have been studying the reproduction of **cats** for ten years
    I have been studying the reproduction of **felines** for ten years
    ✔ WordNet, taxonomies
  - → Mary **has worn** soft contact lenses since college
    Mary **has used** soft contact lenses since college
    ? WordNet, taxonomies

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

## Substitution

- Antonymy
  - → The neighboring town is **poorer** in forest resources than our town
    
    Our town is **richer** in forest resources that the neighboring town
    
    ✔WordNet, dictionary of antonyms, thesaurus
    ✔Matching

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
**NLP Approaches**
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# Substitution

- Antonymy
  - → The neighboring town is **poorer** in forest resources than our town
    Our town is **richer** in forest resources that the neighboring town
    ✔WordNet, dictionary of antonyms, thesaurus
    ✔Matching
  - → I **lost interest** in the endeavor
    I **developed disinterest** in the endeavor
    ? WordNet, dictionary of antonyms, thesaurus

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

## Substitution

- Actant-Action Substitution
  - → I dislike rash **drivers**
    I dislike rash **driving**
    ✔Lemmatizer

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
A Paraphrasing Typology

# Substitution

- Actant-Action Substitution
  - → I dislike rash **drivers**
    I dislike rash **driving**
    ✔Lemmatizer
  - → Mary is John's **student**
    John **teaches** Mary
    ✔Argument structure

Use of the GL lexical entries (Pustejovsky, 2005) for paraphrasing

$$
\begin{bmatrix}
\textbf{teach} \\
\text{ARGSTR} = \begin{bmatrix} \text{ARG}_1 = \boxed{1}\,[\textbf{teacher}] \\ \text{ARG}_2 = \boxed{2}\,[\textbf{student}] \end{bmatrix} \\
\\
\text{QUALIA} = \begin{bmatrix} \text{AGENT} = \textbf{teach\_act}\,(e_1, \boxed{1}, \boxed{2}) \end{bmatrix}
\end{bmatrix}
$$

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# Substitution

- Word-Definition Substitution
  - → Temporary space for rubble and **scrap** wood
    Temporary space for rubble and wood **that became unnecessary**
    ✔MRD, WordNet

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# A Paraphrasing Typology

- Substitution
    - → He **announced** that the company wouldn't participate
      He **stated** that the company wouldn't participate

Deletion

→ Mary opened the **bottle of** wine
  Mary opened the wine

- Transformation
    - → My father built the house
      The house was built by my father

- Splitting
    - → Patrick Ewing scored a personal season high of 41 points
      Patrick Ewing scored 41 points. It was a personal season high

- Change of order
    - → She used to **only** eat hot dishes
      She used to eat **only** hot dishes

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
A Paraphrasing Typology

# Deletion

- Non-propositional content deletion
  → Steven **made an attempt** to stop playing Hearts
  Steven **attempted** to stop paying Hearts
  ✔List of light verbs

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# Deletion

- Non-propositional content deletion
  - → Steven **made an attempt** to stop playing Hearts
    Steven **attempted** to stop paying Hearts
    ✔List of light verbs
- Argument deletion
  - → **My father** built the house
    The house was built
    ✔Argument structure, diathesis alternations

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
**NLP Approaches**
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# Deletion

- Non-propositional content deletion
  - → Steven **made an attempt** to stop playing Hearts
    Steven **attempted** to stop paying Hearts
    ✔List of light verbs

- Argument deletion
  - → **My father** built the house
    The house was built
    ✔Argument structure, diathesis alternations

- Adjunct deletion
  - → John ran home **at noon**
    John ran home
    **?** Argument structure, diathesis alternations

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

## Deletion

- Deletion recoverable by coercion
  - → John began **reading** a book
    John began a book
    ✔Qualia

$$
\begin{bmatrix}
\textbf{book} \\
\text{ARGSTR} = \begin{bmatrix} \text{ARG}_1 = \textbf{x: information} \\ \text{ARG}_2 = \textbf{y : phys\_obj} \end{bmatrix} \\[2em]
\text{QUALIA} = \begin{bmatrix} \textbf{information·phys\_obj\_lcp} \\ \text{FORMAL} = \textbf{hold (y,x)} \\ \text{TELIC} = \textbf{read (e,w,x·y)} \\ \text{AGENT} = \textbf{write (e',v,w·y)} \end{bmatrix}
\end{bmatrix}
$$

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
A Paraphrasing Typology

# A Paraphrasing Typology

- Substitution
    - $\rightarrow$ He **announced** that the company wouldn't participate
      He **stated** that the company wouldn't participate

- Deletion
    - $\rightarrow$ Mary opened the **bottle of** wine
      Mary opened the wine

Transformation

$\rightarrow$ My father built the house
  The house was built by my father

- Splitting
    - $\rightarrow$ Patrick Ewing scored a personal season high of 41 points
      Patrick Ewing scored 41 points. It was a personal season high

- Change of order
    - $\rightarrow$ She used to **only** eat hot dishes
      She used to eat **only** hot dishes

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
**NLP Approaches**
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

## Transformation

- Syntactic transformation
  - → My father built the house
    The house was built by my father
  - → The laundry sways in the breeze
    The breeze makes the laundry sway
    ✔ Diathesis alternations, Dras (1999)

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# A Paraphrasing Typology

- **Substitution**
  - → He **announced** that the company wouldn't participate
    He **stated** that the company wouldn't participate

- **Deletion**
  - → Mary opened the **bottle of** wine
    Mary opened the wine

- **Transformation**
  - → My father built the house
    The house was built by my father

Splitting

→ Patrick Ewing scored a personal season high of 41 points
  Patrick Ewing scored 41 points. It was a personal season high

- **Change of order**
  - → She used to **only** eat hot dishes
    She used to eat **only** hot dishes

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# Splitting

- Splitting
  - → Patrick Ewing scored a personal season high of 41 points
    Patrick Ewing scored 41 points. It was a personal season high
    ✔ Matching

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
A Paraphrasing Typology

# A Paraphrasing Typology

- Substitution
    - → He **announced** that the company wouldn't participate
      He **stated** that the company wouldn't participate

- Deletion
    - → Mary opened the **bottle of** wine
      Mary opened the wine

- Transformation
    - → My father built the house
      The house was built by my father

- Splitting
    - → Patrick Ewing scored a personal season high of 41 points
      Patrick Ewing scored 41 points. It was a personal season high

Change of order

→ She used to **only** eat hot dishes
  She used to eat **only** hot dishes

Paraphrasing and Plagiarism Detection
**Typologies**
Corpora
NLP Approaches
The WRPA System

Paraphrasing Complexity
State of the Art
**A Paraphrasing Typology**

# Change of order

- Change of order
  - → She used to **only** eat hot dishes
    She used to eat **only** hot dishes
  - → The student copied the critical diagrams **before returning the book**
    **Before returning the book,** the student copied the critical diagrams
    ✔Matching

Paraphrasing and Plagiarism Detection
Typologies
**Corpora**
NLP Approaches
The WRPA System

State of the Art
CoCo Interface

## Outline

Paraphrasing and Plagiarism Detection
Typologies
**Corpora**
NLP Approaches
The WRPA System

**State of the Art**
CoCo Interface

## Outline

Paraphrasing and Plagiarism Detection
Typologies
**Corpora**
NLP Approaches
The WRPA System

State of the Art
CoCo Interface

## Corpora

Corpora used for paraphrasing treatment:

- Parallel corpora (multiple parallel translations, source)

  Barzilay and McKeown (2001)

- Comparable corpora (newspaper articles, same event)

  Barzilay, McKeown and Elhadad (1999)

- Bilingual parallel corpora (source – translation)

  Zhao et al. (2009)

- Monolingual corpus

  Bhagat and Ravichandran (2008)

- Wikipedia

  Vila, Rodríguez and Martí (2010)

- Web

  Dolan, Quirk, Brockett (2004)

Paraphrasing and Plagiarism Detection
Typologies
**Corpora**
NLP Approaches
The WRPA System

State of the Art
CoCo Interface

Paraphrasing corpora:

- Microsoft Research Paraphrase Corpus

  Dolan, Brockett and Quirk (2005)

- Barzilay and Lee (2003)

  http://www.cs.cornell.edu/Info/Projects/NLP/statpar.html

- CoCo interface

  España et al. (2009)

Paraphrasing and Plagiarism Detection
Typologies
**Corpora**
NLP Approaches
The WRPA System

State of the Art
CoCo Interface

## Outline

Paraphrasing and Plagiarism Detection
Typologies
**Corpora**
NLP Approaches
The WRPA System

State of the Art
CoCo Interface

# COCO *The </Text-Mess> COrpora COmpilation*

Home | Goals | Statistics | Tasks

Current Task: Paraphrasing

| Evaluate a pair | Generate a pair | Complete a pair | Generate templates | > Search < | > Modify < |

**Evaluate an existing pair**

Choose Sentence 1 from [Microsoft corpus ▼] and Sentence 2 from [Users generated (EN) ▼]

Select sentences:  ○ Randomly
○ Sequentially
○ Filtering

[ START ▷ ]

Sentence 1 | Other, more traditional tests are also available.

Sentence 2 | You can find more traditional tests as well.

Are they paraphrases?   ○ YES   ○ NO

Why? _____

[ SUBMIT ]
[ * FINISH * ]

Paraphrasing and Plagiarism Detection
Typologies
**Corpora**
NLP Approaches
The WRPA System

State of the Art
CoCo Interface

Paraphrasing and Plagiarism Detection
Typologies
Corpora
**NLP Approaches**
The WRPA System

State of the Art

## Outline

Paraphrasing and Plagiarism Detection
Typologies
Corpora
**NLP Approaches**
The WRPA System

State of the Art

## Outline

Paraphrasing and Plagiarism Detection
Typologies
Corpora
**NLP Approaches**
The WRPA System

State of the Art

# State of the Art

## Distributional Hypothesis (Harris, 1954)

Words that occur in the same contexts tend to have similar meanings.

| The students | solved | the problem |
| The students | found a solution to | the problem |

Bhagat and Ravichandran (2008)

## Extended Distributional Hypothesis (Lin and Pantel, 2001)

Paths that link the same sets of words tend to have similar meanings.

The students ←N:subj:V←solved→V:obj:N→the problem
The students ← N:subj:V←find→V:obj:N→solution→N:to:N→the problem

Paraphrasing and Plagiarism Detection
Typologies
Corpora
**NLP Approaches**
The WRPA System

State of the Art

# State of the Art

## Matching

Strings that share a high number of units tend to have similar meanings.

- Bag of words (NEs)
- N-grams

DATE: NUM1 are killed and around NUM2 injured when suicide bomber blows up his explosive-packed belt at X1 in X2.

palestinian suicide bomber blew himself up at X1 in X2 DATE, killing NUM1 and wounding NUM2 police said.

Barzilay and Lee (2003)

Paraphrasing and Plagiarism Detection
Typologies
Corpora
**NLP Approaches**
The WRPA System

State of the Art

# State of the Art

### Edit distance (Levenshtein Distance)

Strings separated by a small edit distance tend to have similar meanings.

Dolan, Quirk and Brockett (2004)

Paraphrasing and Plagiarism Detection
Typologies
Corpora
**NLP Approaches**
The WRPA System

State of the Art

## State of the Art

### Multiple translation

Strings resulting from the translation of another string in another language are paraphrases.

Language 1                    String
                    ↙         ↓         ↘
Language 2    String 1    String 2    String 3    = **Paraphrases**

Barzilay and McKeown (2001)

Zhao et al. (2009)

Paraphrasing and Plagiarism Detection
Typologies
Corpora
**NLP Approaches**
The WRPA System

State of the Art

# State of the Art

### Rule application

To check if strings satisfy a set of (manually) created paraphrasing rules.

**Head omission**: group of students/students
**Ordering of sentence components**: Tuesday they met.../They met ... Tuesday

Barzilay, McKeown and Elhadad (1999)

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Outline

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Outline

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# WRPA
## A System for Relational Paraphrase Acquisition from Wikipedia

Vila, Rodríguez and Martí (2010)

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Relational Paraphrases

Paraphrases expressing a relation between two entities <*Source*, Target>.

- Cervantes **wrote** El Quijote/El Quijote **by** Cervantes
- Joan Ponç **was born in** Barcelona/Barcelona, Joan Ponç**'s home town**

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
**The WRPA System**

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Hypothesis

### Distributional Hypothesis (Harris, 1954)

Words that occur in the same contexts tend to have similar meanings.

| The students | **solved** | the problem |
|---|---|---|
| The students | **found a solution to** | the problem |

### Our hypothesis:

same *Sources* and *Targets* ↔ paraphrase candidates

| | AUTHOR | **wrote** | WORK |
|---|---|---|---|
| **The designer of** | WORK | **was** | AUTHOR |
| | WORK | **was created by** | AUTHOR |
| | AUTHOR | **, inventor of** | WORK |
| | WORK | **by** | AUTHOR |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
**The WRPA System**

**Presentation**
Methodology
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia



WIKIPEDIA
*The Free Encyclopedia*

$\longrightarrow$

Wu y Weld (2007, 2010)
Wu, Hoffmann y Weld (2008)

| **WRPA** | **Others** |
|---|---|
| Complex relations | Simple relations |
| + Other sections | + Infoboxes |
| Paraphrasing | Information Extraction |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Outline

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia

Structured information
↓                              ↓
{Source}              {Target}
↓                              ↓
NON-structured information
↓
{Source **paraphrase** Target}

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
**The WRPA System**

Presentation
**Methodology**
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia

$$
\begin{array}{ccc}
 & \text{Structured information} & \\
\downarrow & & \downarrow \\
\{\text{Source}\} & & \{\text{Target}\} \\
\downarrow & & \downarrow \\
 & \text{NON-structured information} & \\
 & \downarrow & \\
 & \{\text{Source } \mathbf{paraphrase} \text{ Target}\} & \\
\downarrow & & \downarrow \\
\text{anchor} & & \text{anchor} \\
\searrow & & \swarrow \\
 & \text{PATTERN} & \\
\end{array}
$$

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia

Structured information

↓                           ↓                           ↘

{Source}                    {Target}                    CFG

↓                           ↓                           ↗

NON-structured information

↓

{Source **paraphrase** Target}

↓                           ↓

anchor                      anchor

↘                   ↙

PATTERN

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia

{**text**}     [X]     {**text**}     Y     {**text**}     [Z]     {**text**}

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia

| {**text**} | [X] | {**text**} | Y | {**text**} | [Z] | {**text**} |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| **paraphr.** | | **paraphr.** | | **paraphr.** | | **paraphr.** |
| | Source | | Target | | Compl. info | |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia

| {**text**} | [X] | {**text**} | Y | {**text**} | [Z] | {**text**} |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| **paraphr.** | | **paraphr.** | | **paraphr.** | | **paraphr.** |
| | Source | | Target | | Compl. info | |

|  | Person | | Date of birth Place of birth Date of death | | – | |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia

| {**text**} | [X] | {**text**} | Y | {**text**} | [Z] | {**text**} |
|---|---|---|---|---|---|---|
| ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| **paraphr.** | | **paraphr.** | | **paraphr.** | | **paraphr.** |
| | Source | | Target | | Compl. info | |
| | Person | | Date of birth Place of birth Date of death | | – | |
| | Author | | Work | | Date | |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
**The WRPA System**

Presentation
**Methodology**
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia

| {**text**} | [X] | {**text**} | Y | {**text**} | [Z] | {**text**} |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| **paraphr.** | | **paraphr.** | | **paraphr.** | | **paraphr.** |
| | Source | | Target | | Compl. info | |
| | Person | Date of birth Place of birth Date of death | | | – | (English) |
| | Author | | Work | | Date | (Spanish) |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# WRPA–Relational Paraphrase Aquisition from Wikipedia

## Authorship

- Complex relation (painters, sculptors, architects, writers, composers, singer-songwriters, directors, philosophers, inventors and scientists).
- Rich casuistic of manifestations.
  - $\rightarrow$ Appropriateness of this relation for research in paraphrasing.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

Wikipedia components and sections relevant to our work:

- Categories
- Infoboxes
- Work section
- Work page

Paraphrasing and Plagiarism Detection
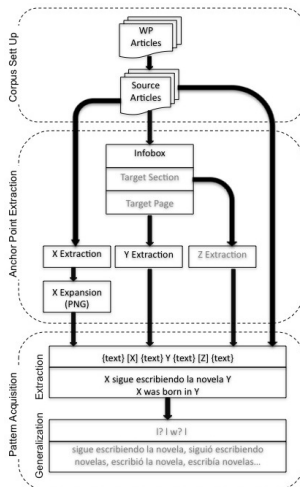Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

Wikipedia components and sections relevant to our work:

- Categories
- Infoboxes
- Work section
- Work page

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## System scheme

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Corpus Set Up

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Corpus Set Up

Date of birth
Place of birth      →      Persons
Date of death
Authorship          →      Authors

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
**The WRPA System**

Presentation
**Methodology**
Evaluation and Results
Conclusions and Future Work

## Corpus Set Up

Date of birth
Place of birth    →    Persons
Date of death
Authorship        →    Authors

\*A category of authors does not exist in the Spanish Wikipedia →
done manually

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Anchor Point Extraction

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Anchor Point Extraction

### X Extraction

- Titles
- Redirection pages

### Y and Z Extraction

**Person:**

- 'Date of birth', 'place of birth' and 'date of death' attributes in infoboxes.

**Authorship:**

- 'Work' attribute in infoboxes.
- Work section
- Work page

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Anchor Point Extraction

### X Extraction

- Titles
- Redirection pages

### Y and Z Extraction

**Person:**
- 'Date of birth', 'place of birth' and 'date of death' attributes in infoboxes.

**Authorship:**
- 'Work' attribute in infoboxes.
- Work section
- Work page

*Attributes are expressed in different ways $\rightarrow$ done manually/KBP track (TAC 2010)

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Anchor Point Extraction

## X Extraction

- Titles
- Redirection pages

## Y and Z Extraction

**Person:**

- 'Date of birth', 'place of birth' and 'date of death' attribute in infoboxes.

**Authorship:**

- 'Work' attribute in infoboxes.
- Work sections
- Work pages

---

* Attributes are univalued.
* The English Wikipedia is very extensive.
* There is a relatively high number of infoboxes in person pages (34 %).
* The corresponding attributes generally appear in the infoboxes.

* Attributes are multivalued.
* The Spanish Wikipedia is smaller.
* Most author pages lack an infobox.
* Most infoboxes lack a work attribute.
* Infoboxes only contain the most important works (2 approx.).

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Anchor Point Extraction

## X Extraction

- Titles
- Redirection pages

## Y and Z Extraction

**Person:**

- 'Date of birth', 'place of birth' and 'date of death' attribute in infoboxes.

**Authorship:**

- 'Work' attribute in infoboxes.
- Work sections
- Work pages

* Attributes are univalued.
* The English Wikipedia is very extensive.
* There is a relatively high number of infoboxes in person pages (34 %).
* The corresponding attributes generally appear in the infoboxes.

* Attributes are multivalued.
* The Spanish Wikipedia is smaller.
* Most author pages lack an infobox.
* Most infoboxes lack a work attribute.
* Infoboxes only contain the most important works (2 approx.).

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Anchor Point Extraction

**Paul Auster:**

| | |
|---|---|
| Infobox | 0 |
| Work section | 33 (25 correct) |
| Work page | 18 (all correct) |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Anchor Point Extraction

| X | Y(Z) |
|---|---|

### Authorship

| | |
|---|---|
| Canaletto | La Riva degli Schiavoni (1730-31) |
| Edgar Allan Poe | [[Eureka]] |
| Luis Eduardo Aute | Templo de carne (1986) |

### Date of birth

| | |
|---|---|
| David Kaye | [[October 14]], [[1964]] |
| Sara Rue | 1979|01|26 |
| Joan of Arc | c. [[1412]] |

### Place of birth

| | |
|---|---|
| Giovanni Branca | [[San Angelo]] in [[Lizzola]], [[Pesaro]] |
| Grigore Preoteasa | [[Bucharest]], [[Romania]] |
| Tomas Plekanec | [[Kladno]], [[Czech Republic |CZE]] |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Pattern Extraction

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Pattern Extraction

. {**text**} [X] {**text**} Y {**text**} [Z] {**text**} .

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Pattern Extraction

En contra de la guerra, con ocasión de su adhesión al "Consejo Mundial de la Paz" pintó el famoso "YYYYYYY " en ( ZZZZZZZ).

En los años siguientes XXXXXX escribió prolíficamente: "YYYYYY" apareció en **z4**, "**y5**" en **z5**, "**y6**" en **z6**

En los años siguientes XXXXXX escribió prolíficamente: "**y4**" apareció en **z4**, "YYYYYY" en **z5**, "**y6**" en **z6**

En los años siguientes XXXXXX escribió prolíficamente: "**y4**" apareció en **z4**, "**y5**" en **z5**, "YYYYYY" en **z6**

La última obra de XXXXXX es un poema: "YYYYYYY" (**z35**).

En Mijáilovskoye, tras la reprimenda paterna y acogido por su amada aya, XXXXXX compuso seis capítulos de "YYYYYYY", el dra

En Mijáilovskoye, tras la reprimenda paterna y acogido por su amada aya, XXXXXX compuso seis capítulos de "**y15**", el dra

:*Ágora - "YYYYYYY"

:*YYYYYY de l"Assut de l"Or

XXXXXX escribió dos libros sobre la reina Isabel I de Inglaterra, "YYYYYY" («Trompetas para Isabel», **z3**) y "**y4**" («Las

XXXXXX escribió dos libros sobre la reina Isabel I de Inglaterra, "**y3**" («Trompetas para Isabel», **z3**) y "**y4**" («Las

En enero de 1845, publicó un poema que le haría célebre: "YYYYYY".

Algunas de estas leyendas inspirarían en su momento una de sus obras fundamentales: "YYYYYY".

Sobre "YYYYYY", XXXXXX dijo a la agencia española EFE el 25 de junio de ZZZZZZ: "No soy socióloga, ni psicóloga ni he querid

Las primeras obras poéticas compuestas por XXXXXX son "L"Allegro" e "YYYYYY" (**z2**), dos pastorales, que reflejan el disfr

Las primeras obras poéticas compuestas por XXXXXX son "L"Allegro" e "**y5**" (**z2**), dos pastorales, que reflejan el disfr

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Pattern Extraction

| | | | |
|---|---|---|---|
| En contra de la guerra, con ocasión de su adhesión al "Consejo Mundial de la Paz" pintó el famoso "YYYYYYY " en ( ZZZZZZZ). | | | |
| En los años siguientes XXXXXX escribió prolíficamente: "YYYYYYY" apareció en **z4**, "**y5**" en **z5**, "**y6**" en **z6* | | | |
| En los años siguientes XXXXXX escribió prolíficamente: "**y4**" apareció en **z4**, "YYYYYYY" en **z5**, "**y6**" en **z6* | | | |
| En los años siguientes XXXXXX escribió prolíficamente: "**y4**" apareció en **z4**, "**y5**" en **z5**, "YYYYYYY" en **z6* | | | |
| **La última obra de XXXXXX es un poema: "YYYYYYY" (**z5**).** | | | |
| En Mijáilovskoye, tras la reprimenda paterna y acogido por su amada aya, XXXXXX compuso seis capítulos de "YYYYYYY", el dra | | | |
| En Mijáilovskoye, tras la reprimenda paterna y acogido por su amada aya, XXXXXX compuso seis capítulos de "**y15**", el dra | | | |
| :*Ágora - "YYYYYYY" | | | |
| :*YYYYYY de l"Assut de l"Or | | | |
| XXXXXX escribió dos libros sobre la reina Isabel I de Inglaterra, "YYYYYY" («Trompetas para Isabel», **z3**) y "**y4**" («Las | | | |
| XXXXXX escribió dos libros sobre la reina Isabel I de Inglaterra, "**y3**" («Trompetas para Isabel», **z3**) y "**y4**" («Las | | | |
| En enero de 1845, publicó un poema que le haría célebre: "YYYYYY". | | | |
| Algunas de estas leyendas inspirarían en su momento una de sus obras fundamentales: "YYYYYY". | | | |
| Sobre "YYYYYY", XXXXXX dijo a la agencia española EFE el 25 de junio de ZZZZZZ: "No soy socióloga, ni psicóloga ni he querid | | | |
| Las primeras obras poéticas compuestas por XXXXXX son "L"Allegro" e "YYYYYY" (**z2**), dos pastorales, que reflejan el disfru | | | |
| Las primeras obras poéticas compuestas por XXXXXX son "L"Allegro" e "**y5**" (**z2**), dos pastorales, que reflejan el disfru | | | |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Pattern Extraction

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Pattern Extraction

En contra de la guerra, con ocasión de su adhesión al "Consejo Mundial de la Paz" pintó el famoso "YYYYYYY " en ( ZZZZZZZ).

En los años siguientes XXXXXX escribió prolíficamente: "YYYYYY" apareció en **z4**, "**y5**" en **z5**, "**y6**" en **z6**

En los años siguientes XXXXXX escribió prolíficamente: "**y4**" apareció en **z4**, "YYYYYY" en **z5**, "**y6**" en **z6**

En los años siguientes XXXXXX escribió prolíficamente: "**y4**" apareció en **z4**, "**y5**" en **z5**, "YYYYYY" en **z6**

La última obra de XXXXXX es un poema: "YYYYYYY" (**z35**).

En Mijáilovskoye, tras la reprimenda paterna y acogido por su amada aya, XXXXXX compuso seis capítulos de "YYYYYYY", el dra

En Mijáilovskoye, tras la reprimenda paterna y acogido por su amada aya, XXXXXX compuso seis capítulos de "**y15**", el dra

:*Ágora - "YYYYYYY"

:*YYYYYYY de l"Assut de l"Or

XXXXXX escribió dos libros sobre la reina Isabel I de Inglaterra, "YYYYYY" («Trompetas para Isabel», **z3**) y "**y4**" («Las

XXXXXX escribió dos libros sobre la reina Isabel I de Inglaterra, "**y3**" («Trompetas para Isabel», **z3**) y "YYYYYY" («Las

En enero de 1845, publicó un poema que le haría célebre: "YYYYY".

Algunas de estas leyendas inspirarían en su momento una de sus obras fundamentales: "YYYYYY".

Sobre "YYYYYY", XXXXXX dijo a la agencia española EFE el 25 de junio de ZZZZZZ: "No soy socióloga, ni psicóloga ni he querid

Las primeras obras poéticas compuestas por XXXXXX son "L"Allegro" e "YYYYYY" (**z2**), dos pastorales, que reflejan el disfr

Las primeras obras poéticas compuestas por XXXXXX son "L"Allegro" e "**y5**" (**z2**), dos pastorales, que reflejan el disfr

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Pattern Extraction

. {**text**} [X] {**text**} Y {**text**} [Z] {**text**} .
↓
. X {**text**} Y .

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Pattern Extraction

### Authorship

1  X sigue escribiendo la novela Y
2  X comenzó a grabar su álbum debut, "Y"
3  X dirigió 'Educando a Rita"(YEAR) y "Y"

### Date of birth

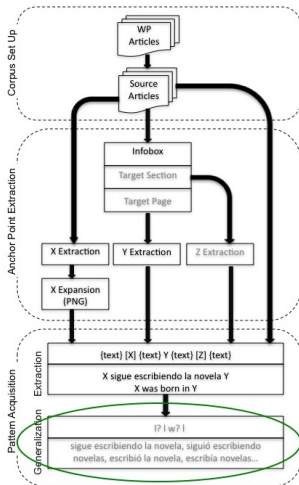4  X known as PERSON was born in Y
5  X was born in Y
6  X was born in PLACE on Y

### Date of death

7  X DATE-Y
8  X DATE to Y
9  X PERSON DATE - Y

### Place of death

10  X born DATE in Y
11  X DATE in Y
12  X DATE Y

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Pattern Generalitzation

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
**The WRPA System**

Presentation
**Methodology**
Evaluation and Results
Conclusions and Future Work

## Pattern Generalitzation

To gather the variants of a generic pattern (only for authorship).

- To lemmatize and PoS tag the patterns (Freeling).
- To represent each pattern as a sequence of $<$word, lemma, PoS$>$ tuples.

| Word | sigue | escribiendo | la | novela |
|------|-------|-------------|-----|--------|
| **Lemma** | seguir | escribir | el | novela |
| **PoS** | v | v | da | nc |

- To use an A\* approach until achieving *n* matches with other patterns.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Pattern Generalitzation

- Each state is represented as a sequence of tuples consisting of a token (word, lemma or PoS) and a condition (obligatory, skippable or omitted):

| Initial state | \<sigue:w\>  \<escribiendo:w\>  \<la:w\>  \<novela:w\> |
|---|---|
| Matching pattern | sigue escribiendo la novela |
| Generalization | \<sigue:l?\>  \<escribiendo:l\>  \<la:w?\>  \<novela:l\> |
| Matching patterns | 1) sigue escribiendo la novela<br>2) siguió escribiendo novelas<br>3) escribió la novela<br>4) escribía novelas |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
**Evaluation and Results**
Conclusions and Future Work

## Outline

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
**Evaluation and Results**
Conclusions and Future Work

## Evaluation

### Precision

To apply each pattern to its original page and verify whether the output is a Y (or a variant) of the corresponding infobox (correct Ys).

$$\frac{Num.\ of\ correct\ Ys}{Num.\ of\ extracted\ Ys}$$

### Recall

**Person:** To apply the pattern to its original page and verify whether the output is a Y (it can be applied to that page).
Conservative assumption: all the pages contain a Y.

$$\frac{Num.\ of\ extracted\ Ys}{Num.\ of\ pages\ of\ the\ corpus}$$

### Baseline

**Person:** The most frequent pattern.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
**Evaluation and Results**
Conclusions and Future Work

## Results

|  |  |  | Precision | Recall | F1 |
|---|---|---|---|---|---|
| Person | Date of birth | X born Y | 0.95 | 0.57 | 0.71 |
|  |  | X Y | 0.80 | 0.12 | 0.21 |
|  |  | Top 8 patterns | 0.92 | 0.75 | 0.83 |
|  |  | Baseline | 0.95 | 0.57 | 0.71 |
|  | Date of death | X DATE-Y | 0.95 | 0.21 | 0.34 |
|  |  | X PERSON DATE-Y | 0.96 | 0.10 | 0.18 |
|  |  | Top 3 patterns | 0.95 | 0.42 | 0.58 |
|  |  | Baseline | 0.95 | 0.21 | 0.34 |
|  | Place of birth | X born DATE in Y | 0.98 | 0.13 | 0.23 |
|  |  | X DATE in Y | 0.94 | 0.10 | 0.18 |
|  |  | Top 3 patterns | 0.92 | 0.26 | 0.41 |
|  |  | Baseline | 0.98 | 0.13 | 0.23 |
| Authorship | <pintó:w><su:w?><cuadro:w?> | | 0.46 | – | – |
|  | <pintó:w><su:w?><primera:l?><obra:w?> | | 0.49 | – | – |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Results

| | | | Precision | Recall | F1 |
|---|---|---|---|---|---|
| **Person** | **Date of birth** | X born Y | 0.95 | 0.57 | 0.71 |
| | | X Y | 0.80 | 0.12 | 0.21 |
| | | Top 8 patterns | 0.92 | 0.75 | 0.83 |
| | | Baseline | 0.95 | 0.57 | 0.71 |
| | **Date of death** | X DATE-Y | 0.95 | 0.21 | 0.34 |
| | | X PERSON DATE-Y | 0.96 | 0.10 | 0.18 |
| | | Top 3 patterns | 0.95 | 0.42 | 0.58 |
| | | Baseline | 0.95 | 0.21 | 0.34 |
| | **Place of birth** | X born DATE in Y | 0.98 | 0.13 | 0.23 |
| | | X DATE in Y | 0.94 | 0.10 | 0.18 |
| | | Top 3 patterns | 0.92 | 0.26 | 0.41 |
| | | Baseline | 0.98 | 0.13 | 0.23 |
| **Authorship** | <pintó:w><su:w?><cuadro:w?> | | 0.46 | – | – |
| | <pintó:w><su:w?><primera:l?><obra:w?> | | 0.49 | – | – |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
**Evaluation and Results**
Conclusions and Future Work

# Results

|  |  |  | Precision | Recall | F1 |
|---|---|---|---|---|---|
| **Person** | **Date of birth** | X born Y | 0.95 | 0.57 | 0.71 |
|  |  | X Y | 0.80 | 0.12 | 0.21 |
|  |  | Top 8 patterns | 0.92 | 0.75 | 0.83 |
|  |  | Baseline | 0.95 | 0.57 | 0.71 |
|  | **Date of death** | X DATE-Y | 0.95 | 0.21 | 0.34 |
|  |  | X PERSON DATE-Y | 0.96 | 0.10 | 0.18 |
|  |  | Top 3 patterns | 0.95 | 0.42 | 0.58 |
|  |  | Baseline | 0.95 | 0.21 | 0.34 |
|  | **Place of birth** | X born DATE in Y | 0.98 | 0.13 | 0.23 |
|  |  | X DATE in Y | 0.94 | 0.10 | 0.18 |
|  |  | Top 3 patterns | 0.92 | 0.26 | 0.41 |
|  |  | Baseline | 0.98 | 0.13 | 0.23 |
| **Authorship** | <pintó:w><su:w?><cuadro:w?> | | 0.46 | – | – |
|  | <pintó:w><su:w?><primera:l?><obra:w?> | | 0.49 | – | – |

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Outline

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Conclusions

Finding good and numerous anchor points is essential in systems based on the Distributional Hypothesis.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

# Conclusions

Finding good and numerous anchor points is essential in systems based on the Distributional Hypothesis.

WRPA...

- ...guarantees the **quality** of these anchor points as they are directly extracted from structured and semantically labelled data.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Conclusions

Finding good and numerous anchor points is essential in systems based on the Distributional Hypothesis.

WRPA...

- ...guarantees the **quality** of these anchor points as they are directly extracted from structured and semantically labelled data.

- ...guarantees a large **quantity** of anchor points as they are not only extracted from infoboxes but also from section and work pages.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
**The WRPA System**

Presentation
Methodology
Evaluation and Results
**Conclusions and Future Work**

## Conclusions

Finding good and numerous anchor points is essential in systems based on the Distributional Hypothesis.

WRPA...

- ...guarantees the **quality** of these anchor points as they are directly extracted from structured and semantically labelled data.
- ...guarantees a large **quantity** of anchor points as they are not only extracted from infoboxes but also from section and work pages.
- ...can deal with **complex relationships**.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Conclusions

Finding good and numerous anchor points is essential in systems based on the Distributional Hypothesis.

WRPA...

- ...guarantees the **quality** of these anchor points as they are directly extracted from structured and semantically labelled data.
- ...guarantees a large **quantity** of anchor points as they are not only extracted from infoboxes but also from section and work pages.
- ...can deal with **complex relationships**.
- ...is **language independent** and also independent of the **relation**.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
**The WRPA System**

Presentation
Methodology
Evaluation and Results
**Conclusions and Future Work**

## Conclusions

Finding good and numerous anchor points is essential in systems based on the Distributional Hypothesis.

WRPA...

- ...guarantees the **quality** of these anchor points as they are directly extracted from structured and semantically labelled data.
- ...guarantees a large **quantity** of anchor points as they are not only extracted from infoboxes but also from section and work pages.
- ...can deal with **complex relationships**.
- ...is **language independent** and also independent of the **relation**.
- ...only relies on **Wikipedia** and **shallow linguistic processing**.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Future Work

- Recall and baseline for authorship.

Paraphrasing and Plagiarism Detection
Typologies
Corpora
NLP Approaches
The WRPA System

Presentation
Methodology
Evaluation and Results
Conclusions and Future Work

## Future Work

- Recall and baseline for authorship.
- Manual validation of paraphrase candidates: CoCo (España et al., 2009)
  - Do they express an authorship relation?
  - Establishement of paraphrasing relations between them in the framework of a typology.

# Outline

# Thank you!

Questions?

marta.vila@ub.edu