

# Hypergraphs and Information Fusion for Term Representation Enrichment. Applications to Named Entity Recognition and Word Sense Disambiguation

Ph.D. Thesis Defense

---

Pavel Soriano-Morales

Supervised by Sabine Loudcher and Julien Ah-Pine

February 7th, 2018



UNIVERSITÉ  
LUMIÈRE  
LYON 2



UNIVERSITÉ  
DE LYON

INSTITUT  
DES SCIENCES  
DE L'HOMME



# Introduction

## Why is it useful to us to understand text?






[All](#) [Images](#) [Shopping](#) [Videos](#) [News](#) [More](#) [Settings](#) [Tools](#)

About 520,000 results (0.63 seconds)

### Guido van Rossum

Python was conceived in the late 1980s, and its implementation began in December 1989 by **Guido van Rossum** at Centrum Wiskunde & Informatica (CWI) in the Netherlands as a successor to the ABC language (itself inspired by SETL) capable of exception handling and interfacing with the operating system Amoeba. **Van Rossum** is ...



☒ **Flight to Madrid (IB 8709)** 

When

Tue, October 24, 07:15 – 09:05

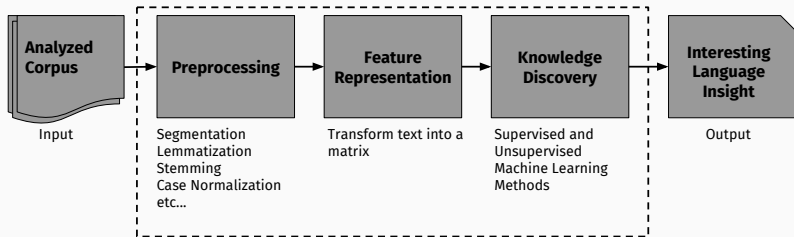
Where **Lyon LYS** [map](#)

This event was automatically created from an email.

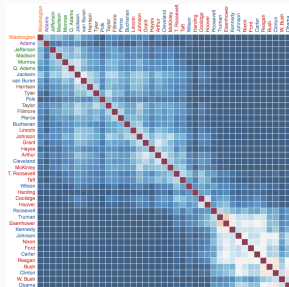
[View confirmation](#) [Remove](#) [More details](#)

# How do we extract meaning from text?

We use **Natural Language Processing** (NLP), a field of computer science interested in making computers comprehend text and obtain useful information from it

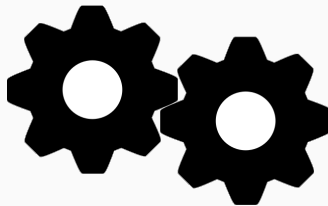


How do we represent text for the machine to understand?



Dealing with data sparsity  
Leveraging heterogeneity

What techniques do we use to discover meaning from text?



Finding semantic communities

- **Common ways to represent text**

- Lexical
- Syntactic
  - Constituency Tree
  - Dependency Tree
- Semantic

- **Common ways to represent text**

- Lexical
- Syntactic
  - Constituency Tree
  - Dependency Tree
- Semantic

- **Example Phrase**

*The report contains copies of the minutes of these meetings*

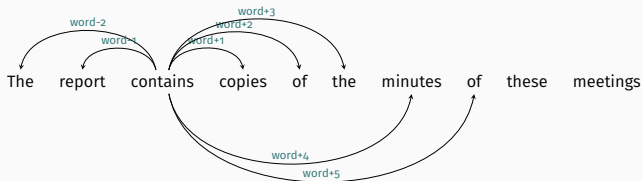
- **Common ways to represent text**

- Lexical
- Syntactic
  - Constituency Tree
  - Dependency Tree
- Semantic

- **Example Phrase**

*The report contains copies of the minutes of these meetings*

### Lexical Representation



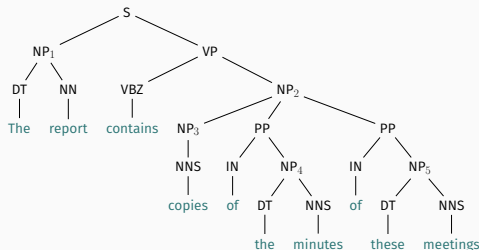
- **Common ways to represent text**

- Lexical
- Syntactic
  - Constituency Tree
  - Dependency Tree
- Semantic

- **Example Phrase**

*The report contains copies of the minutes of these meetings*

### Constituency Tree Representation





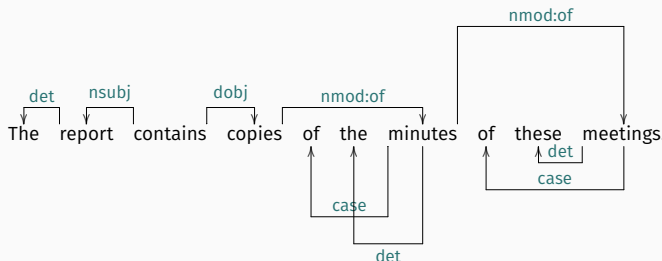
- **Common ways to represent text**

- Lexical
- Syntactic
  - Constituency Tree
  - Dependency Tree
- Semantic

- **Example Phrase**

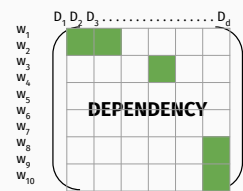
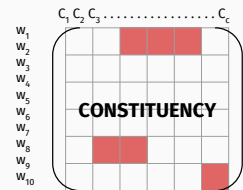
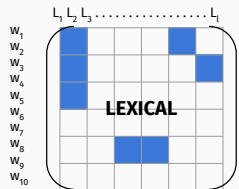
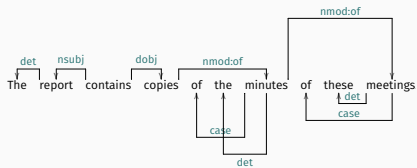
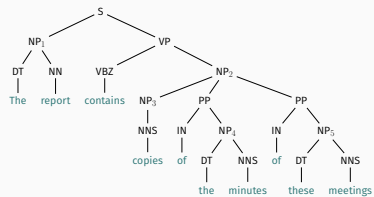
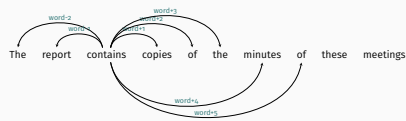
*The report contains copies of the minutes of these meetings*

## Dependency Tree Representation



# Introduction

## Representation Models



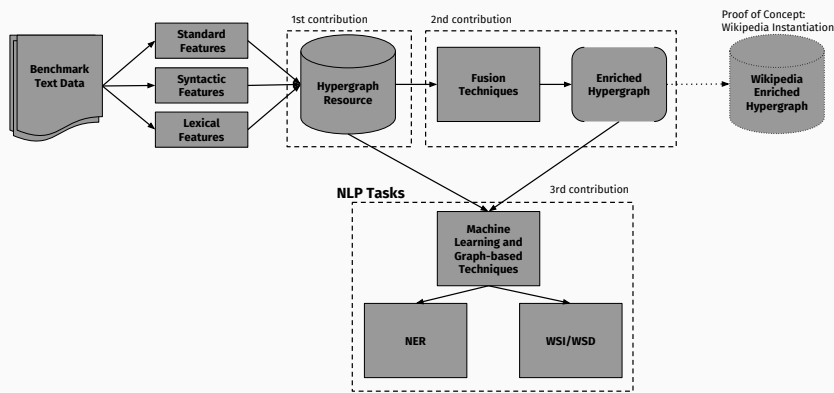
1. What type of model can we employ to represent a corpus **using heterogeneous features**?
  - *Hypergraph model to hold different types of linguistic information*

1. What type of model can we employ to represent a corpus **using heterogeneous features**?
  - *Hypergraph model to hold different types of linguistic information*
2. How can we combine these features while **dealing with feature sparsity**?
  - *Multimedia fusion techniques to combine and densify representation spaces*

1. What type of model can we employ to represent a corpus **using heterogeneous features**?
  - *Hypergraph model to hold different types of linguistic information*
2. How can we combine these features while **dealing with feature sparsity**?
  - *Multimedia fusion techniques to combine and densify representation spaces*
3. How can we **find communities** existing within the language networks?
  - *An alternative network-based algorithm to discover semantically related words within a text*

# Introduction

## Work Overview



# **Contributions in Detail**

## **Hypergraph Linguistic Model**

---

- **Leveraging contexts**

- We extract linguistic information from words based on the **distributional hypothesis** (a word is defined by its surroundings)
- These surroundings are defined as contexts
- Contexts are formed by the interactions a word participates in. These interactions can be lexical or syntactical or other types.



- **Leveraging contexts**

- We extract linguistic information from words based on the **distributional hypothesis** (a word is defined by its surroundings)
- These surroundings are defined as contexts
- Contexts are formed by the interactions a word participates in. These interactions can be lexical or syntactical or other types.

- **We use network models to represent contexts**

- Graphs structures can give us a clearer view into the relations of words within a text
- Allow us to apply methods from graph theory
- Ultimately graphs are transformed to a vectorial representation through the adjacency/incidence matrices

## Example phrase

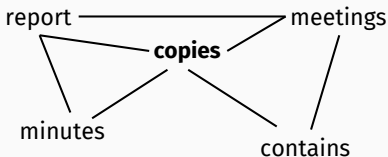
*The report contains copies of the minutes of these meetings*

## Example phrase

*The report contains copies of the minutes of these meetings*

### Lexical Networks

Sentence Level

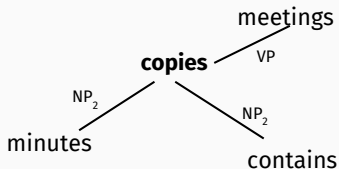


## Example phrase

*The report contains copies of the minutes of these meetings*

### Syntactic Networks

Constituency Tree

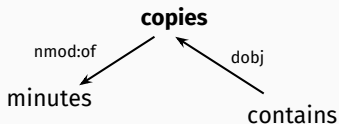


## Example phrase

*The report contains copies of the minutes of these meetings*

### Syntactic Networks

Dependency Tree



- **Limitations of existing representations**
  - Language networks generally employ a single type of textual information
  - The edges of the network relate maximum two words at each time

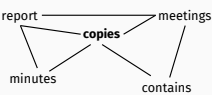
- **Limitations of existing representations**
  - Language networks generally employ a single type of textual information
  - The edges of the network relate maximum two words at each time
- **Proposition**
  - Use a hypergraph model to link together the different types of networks
  - This allows for a semantic overview at three different layers: short range, medium range, and long range at once
  - Relating more than two words at the same time

# Hypergraph Linguistic Model

## Proposed Model

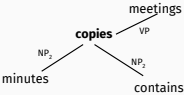
### Lexical Networks

Sentence Level



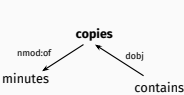
### Syntactic Networks

Constituency Tree



### Syntactic Networks

Dependency Tree



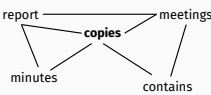


# Hypergraph Linguistic Model

## Proposed Model

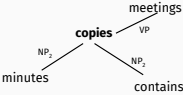
### Lexical Networks

Sentence Level



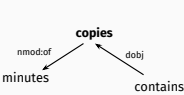
### Syntactic Networks

Constituency Tree



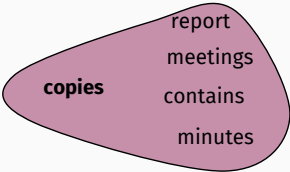
### Syntactic Networks

Dependency Tree



## Hypergraph Model

Lexical

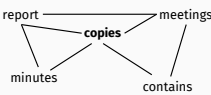


# Hypergraph Linguistic Model

## Proposed Model

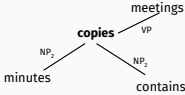
### Lexical Networks

Sentence Level



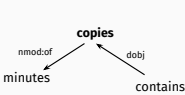
### Syntactic Networks

Constituency Tree

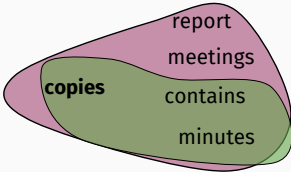


### Syntactic Networks

Dependency Tree



## Hypergraph Model



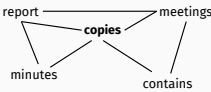
- Lexical
- Constituency (NP<sub>2</sub>)

# Hypergraph Linguistic Model

## Proposed Model

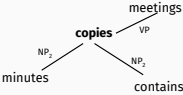
### Lexical Networks

Sentence Level



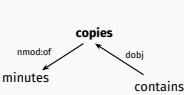
### Syntactic Networks

Constituency Tree

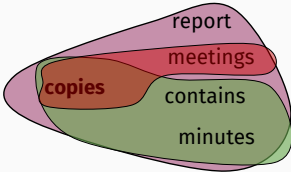


### Syntactic Networks

Dependency Tree



## Hypergraph Model



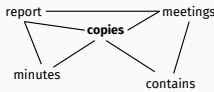
- Lexical
- Constituency (NP<sub>2</sub>)
- Constituency (VP)

# Hypergraph Linguistic Model

## Proposed Model

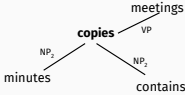
### Lexical Networks

Sentence Level



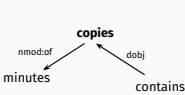
### Syntactic Networks

Constituency Tree

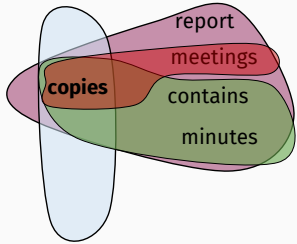


### Syntactic Networks

Dependency Tree



## Hypergraph Model



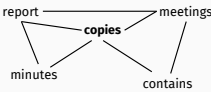
- Lexical
- Constituency (NP<sub>2</sub>)
- Constituency (VP)
- Dependency (dobj:contains)

# Hypergraph Linguistic Model

## Proposed Model

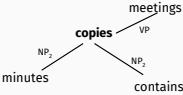
### Lexical Networks

Sentence Level



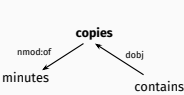
### Syntactic Networks

Constituency Tree

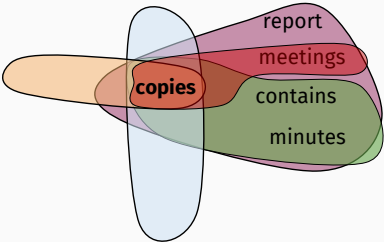


### Syntactic Networks

Dependency Tree



## Hypergraph Model



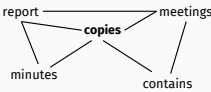
- Lexical
- Constituency (NP<sub>2</sub>)
- Constituency (VP)
- Dependency (dobj:contains)
- Dependency (nmod:of)

# Hypergraph Linguistic Model

## Proposed Model

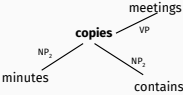
### Lexical Networks

Sentence Level



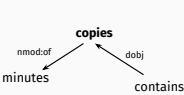
### Syntactic Networks

Constituency Tree

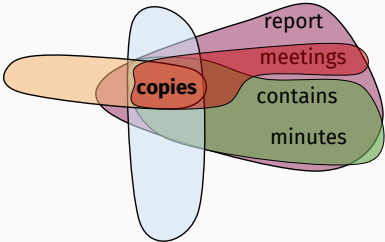


### Syntactic Networks

Dependency Tree



## Hypergraph Model



- Lexical
- Constituency (NP<sub>2</sub>)
- Constituency (VP)
- Dependency (dobj:contains)
- Dependency (nmod:of)

How to combine these heterogeneous networks into a single representation?

# **Contributions in Detail**

**Combining Features and Dealing with  
Sparsity**

---

- **Definition**

- Used in multimedia analysis tasks to integrate multiple media
- We adapt them to combine textual information
- The goal is to obtain rich insights about the data being treated
- By creating a single representation from heterogeneous information



- **Definition**

- Used in multimedia analysis tasks to integrate multiple media
- We adapt them to combine textual information
- The goal is to obtain rich insights about the data being treated
- By creating a single representation from heterogeneous information

- **Main fusion operators:**

- Early Fusion  $E_{\alpha}(\cdot)$ ,
- Late Fusion  $L_{\beta}(\cdot)$ ,
- Cross Fusion  $X_{\gamma}(\cdot)$

# Combining Features and Dealing with Sparsity

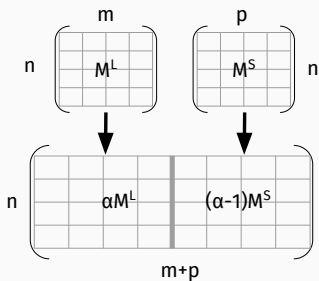
## Early and Late Fusion

### DEFINITIONS

$M^L$	Lexical features	$M^S$	Syntactic features
$S^L$	Lexical similarities	$S^S$	Syntactic similarities

### EARLY FUSION

Matrices  $M^L$  and  $M^S$  have the same number of rows

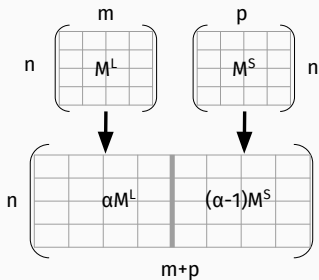


**DEFINITIONS**

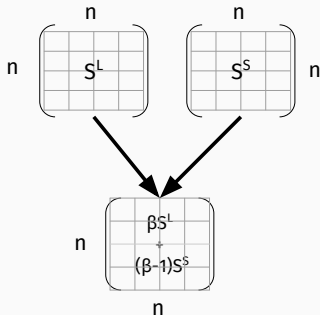
$M^L$	Lexical features	$M^S$	Syntactic features
$S^L$	Lexical similarities	$S^S$	Syntactic similarities

**EARLY FUSION**

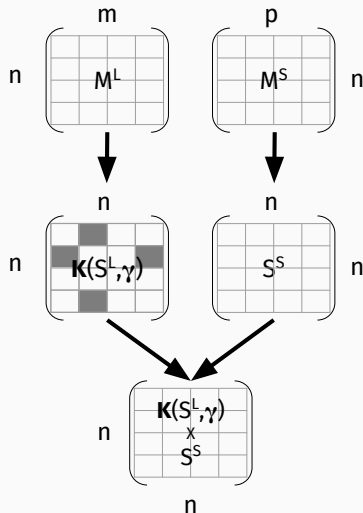
Matrices  $M^L$  and  $M^S$  have the same number of rows

**LATE FUSION: SIMILARITY FUSION**

Matrices  $S^L$  and  $S^S$  have the same size



## CROSS FUSION



- **Combining fusion operators**
  - Applying one function to the result of another to produce a new fusion function

- **Combining fusion operators**

- Applying one function to the result of another to produce a new fusion function

- **First Degree**

- $E(M^L, M^S), L(S^S, M^L)$
- **Cross Feature Fusion:**  $X_F(S^S, M^L)$
- **Cross Similarity Fusion:**  $X_S(S^S, S^L)$

$X_F(S^L, M^S)$

Cross Feature Fusion

$$\begin{matrix} w_1 & w_2 & w_3 \\ \begin{pmatrix} S^L \end{pmatrix} \end{matrix} \times \begin{matrix} f_{S1} & f_{S2} & f_{S3} \\ \begin{pmatrix} M^S \end{pmatrix} \end{matrix} = \begin{matrix} f_{S1} & f_{S2} & f_{S3} \\ \begin{pmatrix} X_F(S^L, M^S) \end{pmatrix} \end{matrix}$$

$X_S(S^L, S^S)$

Cross Similarity Fusion

$$\begin{matrix} w_1 & w_2 & w_3 \\ \begin{pmatrix} S^L \end{pmatrix} \end{matrix} \times \begin{matrix} w_1 & w_2 & w_3 \\ \begin{pmatrix} S^S \end{pmatrix} \end{matrix} = \begin{matrix} w_1 & w_2 & w_3 \\ \begin{pmatrix} X_S(S^L, S^S) \end{pmatrix} \end{matrix}$$

- **Combining fusion operators**

- Applying one function to the result of another to produce a new fusion function

- **Second Degree**

- **Cross Feature Early Fusion:**  $X_F(S^T, E(M^S, M^L))$
- **Late Cross Feature Fusion:**  $L(M^T, X_F(S^T, M^T))$

Cross Feature Early Fusion

$$\begin{array}{c}
 \boxed{X_F(S^L, E(M^S, M^L))} \\
 \begin{array}{c}
 \begin{matrix} f_{S1} & f_{S2} & f_{S3} \\ w_1 & w_2 & w_3 \end{matrix} \begin{pmatrix} M^S \end{pmatrix} \parallel \begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ w_1 & w_2 & w_3 \end{matrix} \begin{pmatrix} M^L \end{pmatrix} = \begin{matrix} f_{S1} & f_{S2} & f_{S3} & f_{L1} & f_{L2} & f_{L3} \\ w_1 & w_2 & w_3 \end{matrix} \begin{pmatrix} E(M^S, M^L) \end{pmatrix} \\
 \begin{matrix} w_1 & w_2 & w_3 \end{matrix} \begin{pmatrix} S^L \end{pmatrix} \times \begin{matrix} f_{S1} & f_{S2} & f_{S3} & f_{L1} & f_{L2} & f_{L3} \\ w_1 & w_2 & w_3 \end{matrix} \begin{pmatrix} E(M^S, M^L) \end{pmatrix} = \begin{matrix} f_{S1} & f_{S2} & f_{S3} & f_{L1} & f_{L2} & f_{L3} \\ w_1 & w_2 & w_3 \end{matrix} \begin{pmatrix} X_F(S^L, E(M^S, M^L)) \end{pmatrix}
 \end{array}
 \end{array}$$

- **Combining fusion operators**

- Applying one function to the result of another to produce a new fusion function

- **Higher Degree**

- Triple Early Double Late Cross Feature Fusion:

$$E(M_L, E(E(M_T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$



## Higher Degree Operator

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

## Higher Degree Operator

The diagram illustrates a nested function call structure for a higher degree operator. It consists of several nested boxes and mathematical expressions:

- An outermost light blue box containing the expression  $E(M_L, \dots)$ .
- Inside the blue box, there is a light purple box containing  $E(E(M^T, \dots), L(M^L, X_F(S^S, M^L))))$ .
- Inside the purple box, there is a green box containing  $E(M^T, \dots)$ .
- Inside the green box, there is a red box containing  $L(M^T, X_F(S^T, M^T))$ .
- Inside the purple box, there is also a yellow box containing  $L(M^L, X_F(S^S, M^L))$ .

The overall expression is  $E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$ .

## Higher Degree Operator

$$E(M^L, E(E(M^T, L(M^T, X_F(S^T, M^T))))), L(M^L, X_F(S^S, M^L))))$$

$$L(M^L, X_F(S^S, M^L))$$

$$\begin{aligned} \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} w_1 w_2 w_3 \\ S^S \end{pmatrix} \times \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{L1} f_{L2} f_{L3} \\ M^L \end{pmatrix} &= \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{L1} f_{L2} f_{L3} \\ X_F(S^S, M^L) \end{pmatrix} \\ \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{L1} f_{L2} f_{L3} \\ M^L \end{pmatrix} + \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{L1} f_{L2} f_{L3} \\ X_F(S^S, M^L) \end{pmatrix} &= \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{L1} f_{L2} f_{L3} \\ L(M^L, X_F(S^S, M^L)) \end{pmatrix} \end{aligned}$$

## Higher Degree Operator

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$L(M^T, X_F(S^T, M^T))$$

$$\begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} w_1 w_2 w_3 \\ S^T \end{pmatrix} \times \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} f_{T2} f_{T3} \\ M^T \end{pmatrix} = \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} f_{T2} f_{T3} \\ X_F(S^T, M^T) \end{pmatrix}$$

$$\begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} f_{T2} f_{T3} \\ M^T \end{pmatrix} + \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} f_{T2} f_{T3} \\ X_F(S^T, M^T) \end{pmatrix} = \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} f_{T2} f_{T3} \\ L(M^T, X_F(S^T, M^T)) \end{pmatrix}$$

### Higher Degree Operator

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$E(M^T, L(M^T, X_F(S^T, M^T)))$$

$$\begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} & f_{T2} & f_{T3} \\ M^T \end{pmatrix} \parallel \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} & f_{T2} & f_{T3} \\ L(M^T, X_F(S^T, M^T)) \end{pmatrix} = \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{T1} & f_{T2} & f_{T3} & f_{T1} & f_{T2} & f_{T3} \\ E(M^T, L(M^T, X_F(S^T, M^T))) \end{pmatrix}$$

Higher Degree Operator

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L)))$$

$$\begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \left( \begin{matrix} f_{T1} & f_{T2} & f_{T3} \\ E(M^T, L(M^T, X_F(S^T, M^T))) \end{matrix} \right) \parallel \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \left( \begin{matrix} f_{L1} & f_{L2} & f_{L3} \\ L(M^L, X_F(S^S, M^L)) \end{matrix} \right) =$$

$$\begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \left( \begin{matrix} f_{L1} & f_{L2} & f_{L3} & f_{L1} & f_{L2} & f_{L3} \\ E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))) \end{matrix} \right)$$

## Higher Degree Operator

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L))))$$

$$\begin{aligned} \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{L1} f_{L2} f_{L3} \\ M^T \end{pmatrix} \parallel \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{L1} f_{L2} f_{L3} \\ E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L)))) \end{pmatrix} = \\ \begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix} \begin{pmatrix} f_{L1} f_{L2} f_{L3} f_{L1} f_{L2} f_{L3} \\ E(M_L, E(E(M^T, L(M^T, X_F(S^T, M^T))), L(M^L, X_F(S^S, M^L)))) \end{pmatrix} \end{aligned}$$

# **Contributions in Detail**

## **Finding Communities in the Network**

---



- **Language networks tend to be scale-free**
  - There are certain nodes (hubs) that are very well connected forming communities within the network

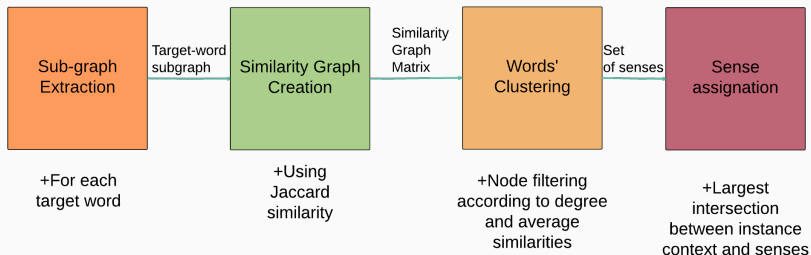
- **Language networks tend to be scale-free**
  - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
  - Hyperlex [VÓ4]
  - University of York (UoY) [KMo7]

- **Language networks tend to be scale-free**
  - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
  - Hyperlex [VÓ4]
  - University of York (UoY) [KM07]
- **Limitations of existing approaches**
  - Single typed networks
  - Large number of parameters

- **Language networks tend to be scale-free**
  - There are certain nodes (hubs) that are very well connected forming communities within the network
- **Seminal approaches**
  - Hyperlex [VÓ4]
  - University of York (UoY) [KMo7]
- **Limitations of existing approaches**
  - Single typed networks
  - Large number of parameters
- **Proposition**
  - Be able to exploit different types of linguistic information (lexical or syntactic co-occurrence)
  - Keep the number of parameters low and allow for their automatic adjusting according to the network's nature

# Finding Communities in the Network

## Proposed Method



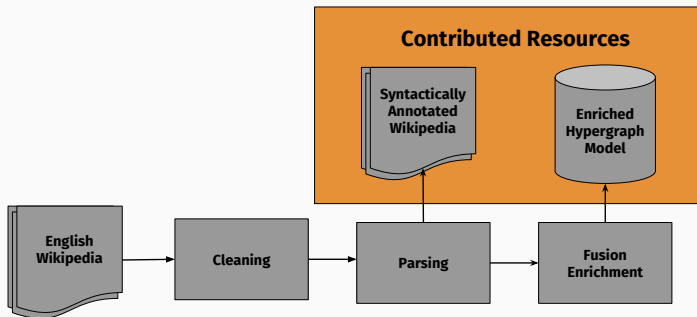
# **Applications to NLP**

## **Hypergraph Model Instantiation**

---

- **Apply our proposed linguistic model to a real world corpus**
  - Use the English Wikipedia as input and generate a textual structure following the proposed network model

- **Apply our proposed linguistic model to a real world corpus**
  - Use the English Wikipedia as input and generate a textual structure following the proposed network model
- **Steps performed**





Hypergraph Model Instantiation

Hypergraph Incidence Matrix

		CONSTITUENT			DEPENDENCY	SENTENCE
		NP <sub>1</sub> DT:NN	NP <sub>2</sub> NP:PP:PP	NP <sub>3</sub> NNS	nsubj contains      dobj contains	S <sub>1</sub>
N N	report	1			1	1
	copies		1	1	1	1
	minutes		1			1
	meetings		1			1
VB	contains					1

- **Characteristics of the enriched space**

- Sparsity is reduced
- Semantic relatedness differs according to the representation space

	<b>Lexical Features (5.49%)</b> $M^L$	<b>Syntactic Features (4.97%)</b> $M^S$	<b>Early Fusion (5.23%)</b> $E(M^L, M^S)$	$X_F$ <b>Fusion (16.75%)</b> $X_F(S^S, M^L)$	$X_F$ <b>Fusion (13.45%)</b> $X_F(S^L, M^S)$
<b>priest</b>	priests nun canton sailor burial	monk regent aedile seer meek	sailor regent nuclei nun relic	vassal regent nun sailor monk	sailor fluent dean nuclei chorus

# **Applications to NLP**

## **Solving Named Entity Recognition**

---

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- Location (LOC)
- Organization (ORG)
- Person (PER)
- Miscellaneous (MISC)
- None (O)

- **NER Objective**

- The goal is to automatically discover mentions that belong to a well-defined semantic category.

- **Classic entities types**

- Location (LOC)
- Organization (ORG)
- Person (PER)
- Miscellaneous (MISC)
- None (O)

- **Our goal**

- We assess the effectiveness of the classic fusion methods and propose new hybrid combinations

### Example Phrase

*Australian scientist discovers star with telescope*

### Three different types of features

Word	Features	Feature Type
Australian	word:Australian, word+1:scientist, ...	<b>Lexical (L)</b>
scientist	Australian/JJ/amod, discovers/VBZ/nsubj_inv	<b>Syntactic (S)</b>
discover	discover, no-capital-letter, prf:dis, suf:ver, VBZ	<b>Standard (T)</b>

- **Preprocessing**
  - Normalize numbers



- **Preprocessing**

- Normalize numbers

- **Test Corpora**

- CoNLL-2003 (CONLL): Train: 219,554 lines. Test: 50,350 lines
- Wikiner (WNER): 3.5 million words.
- Wikigold (WGLD): 41,011 words.

- **Preprocessing**
  - Normalize numbers
- **Test Corpora**
  - CoNLL-2003 (CONLL): Train: 219,554 lines. Test: 50,350 lines
  - Wikiner (WNER): 3.5 million words.
  - Wikigold (WGLD): 41,011 words.
- **Learning Algorithm**
  - Structured Perceptron

- **Preprocessing**
  - Normalize numbers
- **Test Corpora**
  - CoNLL-2003 (CONLL): Train: 219,554 lines. Test: 50,350 lines
  - Wikiner (WNER): 3.5 million words.
  - Wikigold (WGLD): 41,011 words.
- **Learning Algorithm**
  - Structured Perceptron
- **Evaluation Metric**
  - F-measure
  - Evaluated with a 5-fold CV (WNER and WGLD)

## Solving Named Entity Recognition

### Evaluation Baselines (F-measure)

A	Single Features		
	CONLL	WNER	WGLD
$M^T$	77.41	77.50	59.66
$M^L$	69.40	69.17	52.34
$M^S$	32.95	28.47	25.49

A	B	Early Fusion (EF)		
		CONLL	WNER	WGLD
$M^L$	$M^S$	72.01	70.59	59.38
$M^L$	$M^T$	78.13	79.78	61.96
$M^S$	$M^T$	77.70	78.10	60.93
$M^L$	$E(M^S, M^T)$	78.90	80.04	63.20

## Solving Named Entity Recognition Evaluation (F-measure)

A	B	Baseline (EF)		
		CONLL	WNER	WGLD
$M^L$	$E(M^S, M^T)$	78.90	80.04	63.20

### First Degree Fusion

		Cross Feature Fusion ( $X_F F$ )		
		CONLL	WNER	WGLD
$S^L$	$M^T$	49.90	<b>70.27</b>	<b>62.69</b>
$S^S$	$M^T$	47.27	51.38	48.53
$S^T$	$M^L$	<b>52.89</b>	62.21	50.15

A	B	Baseline (EF)		
		CONLL	WNER	WGLD
$M^L$	$E(M^S, M^T)$	78.90	80.04	63.20

Second Degree Fusion

A	B	Early Cross Feature Fusion (EX <sub>F</sub> F)		
		CONLL	WNER	WGLD
$M^T$	$X_F(S^S, M^L)$	49.58	77.32	61.69

## Solving Named Entity Recognition Evaluation (F-measure)

A	B	Baseline (EF)		
		CONLL	WNER	WGLD
$M^L$	$E(M^S, M^T)$	78.90	80.04	63.20

## Second Degree Fusion

A	B	Late Cross Feature Fusion ( $LX_F F$ )		
		CONLL	WNER	WGLD
$M^T$	$X_F(S^S, M^T)$	<b>56.53</b>	62.27	52.39

## Solving Named Entity Recognition Evaluation (F-measure)

A	B	Baseline (EF)		
		CONLL	WNER	WGLD
$M^L$	$E(M^S, M^T)$	78.90	80.04	63.20

## High Degree Fusion

Triple Early Double Late Cross Feature Fusion (EEELX <sub>F</sub> LX <sub>F</sub> )				
		CONLL	WNER	WGLD
$M^L_{\alpha=0.95}$	$\hat{b}_{EEELX_F LX_F}$	79.67	81.79	67.05

$$\hat{b}_{EEELX_F LX_F} = E(E(M^T, L(M^T, X_F(S^S, M^T))), L(M^L, X_F(S^S, M^L)))$$



- Split the operator in four different models

$$\begin{array}{c}
 \overbrace{\hspace{15em}}^{M_4} \\
 \overbrace{\hspace{10em}}^{M_2} \\
 E_{\alpha=0.95}(\underbrace{M^L, M^T}_{M_1}, L(M^T, X_F(S^S, M^T)), L(M^L, X_F(S^S, M^L))) \\
 \underbrace{\hspace{15em}}_{M_3}
 \end{array}$$

- Split the operator in four different models

$$\begin{array}{c}
 \overbrace{\hspace{15em}}^{M_4} \\
 \overbrace{\hspace{10em}}^{M_2} \\
 E_{\alpha=0.95}(\underbrace{M^L, M^T}_{M_1}, L(M^T, X_F(S^S, M^T)), L(M^L, X_F(S^S, M^L))) \\
 \underbrace{\hspace{15em}}_{M_3}
 \end{array}$$

$$M_1 \quad M^L$$

$$M_2 \quad E_{\alpha}(M^L, M^T)$$

$$M_3 \quad E_{\alpha}(M^L, M^T, L(M^T, X_F(S^S, M^T)))$$

$$M_4 \quad E_{\alpha}(M^L, M^T, L(M^T, X_F(S^S, M^T)), L(M^L, X_F(S^S, M^L)))$$

- **Error Analysis Model**

- To facilitate the interpretation, we change the prediction model to a logistic regression with  $L_1$  normalization, which also benefits from the enriched spaces

- **Error Analysis Model**

- To facilitate the interpretation, we change the prediction model to a logistic regression with  $L_1$  normalization, which also benefits from the enriched spaces

- **Procedure**

- We find an error on a model and then see if this error was fixed in the next evolved model
- We study the weights assigned to each feature and see if those added by the fusion make the model change its decision

- **Error Analysis Model**

- To facilitate the interpretation, we change the prediction model to a logistic regression with  $L_1$  normalization, which also benefits from the enriched spaces

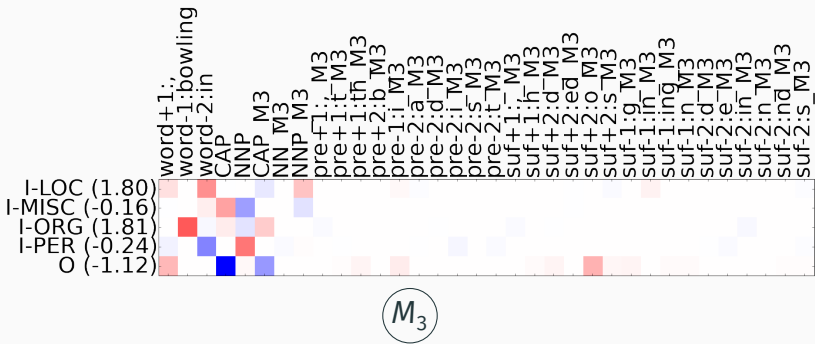
- **Procedure**

- We find an error on a model and then see if this error was fixed in the next evolved model
- We study the weights assigned to each feature and see if those added by the fusion make the model change its decision

- **Experiment**

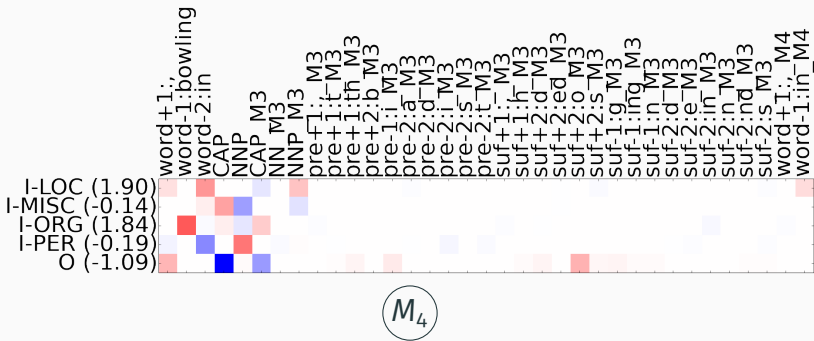
- We follow the location name *Green* from  $M_3$  (incorrectly classified as ORG) to  $M_4$  (correctly classified as LOC)

The location *Green* is classified as ORG by  $M_3$ . It is fixed by  $M_4$ , classifying it as LOC



Solving Named Entity Recognition  
Analyzing the Best Fusion Operator

The location *Green* is classified as ORG by  $M_3$ . It is fixed by  $M_4$ , classifying it as LOC



# **Applications to NLP**

## **Solving Word Sense Induction and Disambiguation**

---



- **WSI/WSD Objective**

- The goal is to determine a set of possible senses to a given word according to its possible contexts (WSI). Then, assigning a correct sense to a particular instance of said word (WSD)

- **WSI/WSD Objective**

- The goal is to determine a set of possible senses to a given word according to its possible contexts (WSI). Then, assigning a correct sense to a particular instance of said word (WSD)

- **Our goals**

- Assess the effectiveness of the fusion enriched spaces
- Evaluate the pertinence of our community discovering algorithm

## Experimental Protocol

- **Feature Space**
  - Lexical (L) and Syntactic (S) Features
- **Preprocessing**
  - Remove very frequent and very infrequent words

### Experimental Protocol

- **Feature Space**

- Lexical (L) and Syntactic (S) Features

- **Preprocessing**

- Remove very frequent and very infrequent words

- **Test Corpora**

- Semeval Competition 2007: Train 219,554 lines. Test 50,350 lines

- **Feature Space**
  - Lexical (L) and Syntactic (S) Features
- **Preprocessing**
  - Remove very frequent and very infrequent words
- **Test Corpora**
  - Semeval Competition 2007: Train 219,554 lines. Test 50,350 lines
- **Clustering Algorithm**
  - Spectral Clustering
  - Proposed Community Algorithm

- **Feature Space**
  - Lexical (L) and Syntactic (S) Features
- **Preprocessing**
  - Remove very frequent and very infrequent words
- **Test Corpora**
  - Semeval Competition 2007: Train 219,554 lines. Test 50,350 lines
- **Clustering Algorithm**
  - Spectral Clustering
  - Proposed Community Algorithm
- **Evaluation Metrics**
  - Supervised Recall (SR)
  - Unsupervised F-measure (UF)

- **Feature Space**
  - Lexical (L) and Syntactic (S) Features
- **Preprocessing**
  - Remove very frequent and very infrequent words
- **Test Corpora**
  - Semeval Competition 2007: Train 219,554 lines. Test 50,350 lines
- **Clustering Algorithm**
  - Spectral Clustering
  - Proposed Community Algorithm
- **Evaluation Metrics**
  - Supervised Recall (SR)
  - Unsupervised F-measure (UF)
  - **Proposed H-measure**

$$\text{H-measure} = \frac{1}{2} \left( 2 * \frac{SR * UF}{SR + UF} + \frac{\delta}{\delta + |\#cl - \delta|} \right)$$

$\delta$  is the average true number of senses of the words in a test corpus

## Supervised Recall

Fusion Operation / System	Recall (%)			#cl
	all	nouns	verbs	
Single Features				
$M^L$	79.20	82.10	75.80	4.13
$M^S$	79.10	81.60	76.20	4.47
Early Fusion (EF)				
$E(M^L, M^S)$	78.70	81.11	76.10	4.46
Late Cross Feature Fusion (LX <sub>F</sub> F)				
$L(M^S, X_F(S^L, M^S))$	78.60	81.10	75.80	4.22
$L(M^L, X_F(S^L, M^L))$	79.50	82.80	75.70	3.96

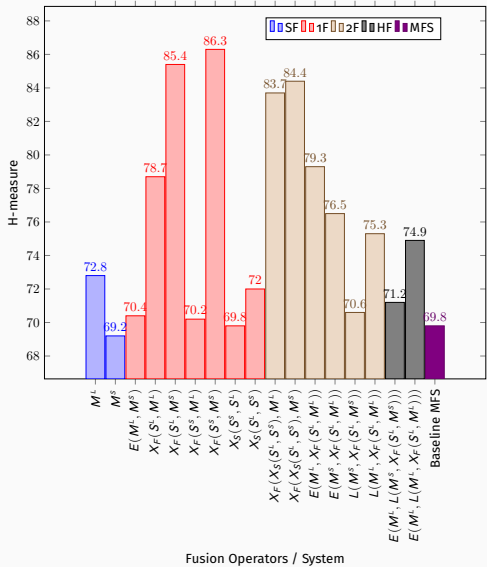


## Unsupervised F-measure

Fusion Operation / System	F-measure (%)			#cl
	all	nouns	verbs	
Single Features				
$M^L$	72.70	76.90	67.90	4.13
$M^S$	69.30	69.40	69.20	4.47
Early Fusion (EF)				
$E(M^L, M^S)$	74.00	76.66	71.11	4.46
Cross Feature Fusion ( $X_F F$ )				
$X_F(S^S, M^L)$	78.90	80.70	76.90	1.08

# Solving Word Sense Induction and Disambiguation

## Spectral Clustering Evaluation: H-measure



## Supervised Recall

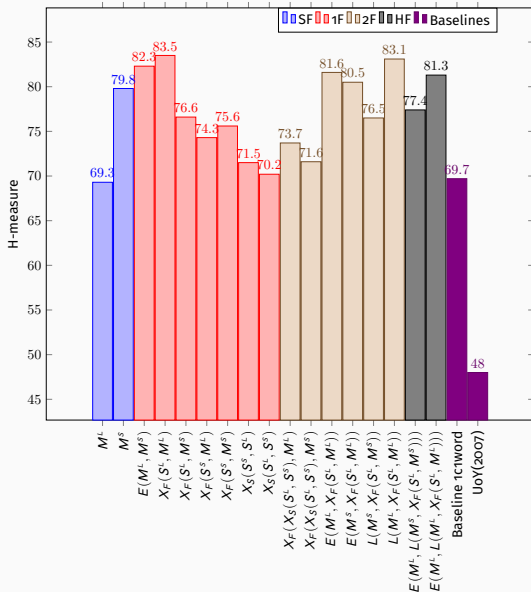
Fusion Operation / System	Recall (%)			#cl
	all	nouns	verbs	
Single Features				
$M^L$	78.70	81.00	76.00	4.21
$M^S$	78.41	80.30	76.10	2.26
Early Fusion (EF)				
$E(M^L, M^S)$	78.80	81.00	76.40	2.43
Cross Feature Fusion ( $X_F$ )				
$X_F(S^L, M^L)$	79.10	81.60	76.40	1.73

Unsupervised F-measure

Fusion Operation / System	F-measure (%)			#cl
	all	nouns	verbs	
Single Features				
$M^L$	63.80	61.30.90	66.50	4.21
$M^S$	75.90	78.80	72.60	2.26
Early Fusion (EF)				
$E(M^L, M^S)$	76.90	80.20	73.10	2.43
Cross Feature Fusion ( $X_S F$ )				
$X_F(S^S, M^L)$	78.30	79.70	76.80	1.10

# Solving Word Sense Induction and Disambiguation

## Proposed Algorithm Evaluation: H-measure



# Conclusions

---

- **Hypergraph linguistic model to hold heterogeneous information**
  - Hypergraphs allow a multi-layered representation of text within a single resource.
  - The Wikipedia-based instantiation serves as a NLP system starting point

- **Hypergraph linguistic model to hold heterogeneous information**
  - Hypergraphs allow a multi-layered representation of text within a single resource.
  - The Wikipedia-based instantiation serves as a NLP system starting point
- **Multimedia fusion techniques to combine and densify representations**
  - High-degree combinations of linguistic representations reduce sparsity
  - These fusion spaces achieve improvements on NER and WSI/WSD compared to single features and trivial fusion



- **Hypergraph linguistic model to hold heterogeneous information**
  - Hypergraphs allow a multi-layered representation of text within a single resource.
  - The Wikipedia-based instantiation serves as a NLP system starting point
- **Multimedia fusion techniques to combine and densify representations**
  - High-degree combinations of linguistic representations reduce sparsity
  - These fusion spaces achieve improvements on NER and WSI/WSD compared to single features and trivial fusion
- **Finding semantically-related communities on linguistic networks**
  - The proposed community finding method improves over similar algorithms while being simpler and allowing for heterogeneous features

- **Hypergraph Linguistic Model**
  - A dataframe-like structure specialized on linguistic information based on the proposed model
  - Defining inter-features similarities measures within the network

- **Hypergraph Linguistic Model**
  - A dataframe-like structure specialized on linguistic information based on the proposed model
  - Defining inter-features similarities measures within the network
- **Combining Features and Dealing with Sparsity**
  - Finding a more principled way to determine what type of context with what type of fusion operation according to the task at hand
  - Exploring fusion with other types of features (other modalities)

- **Hypergraph Linguistic Model**

- A dataframe-like structure specialized on linguistic information based on the proposed model
- Defining inter-features similarities measures within the network

- **Combining Features and Dealing with Sparsity**

- Finding a more principled way to determine what type of context with what type of fusion operation according to the task at hand
- Exploring fusion with other types of features (other modalities)

- **Applications to NLP**

- Comparison with other distributional representations (word embeddings)
- Using the large Wikipedia-based network as a background corpus to further enrich domain-specific corpora
- Test more feature weighting schemes, validate findings on more datasets

- Edmundo-Pavel Soriano-Morales, Julien Ah-Pine, Sabine Loudcher: **Fusion Techniques for Named Entity Recognition and Word Sense Induction and Disambiguation**. DS 2017
- Edmundo-Pavel Soriano-Morales, Julien Ah-Pine, Sabine Loudcher: **Using a Heterogeneous Linguistic Network for Word Sense Induction and Disambiguation**. CICLING 2016
- Edmundo-Pavel Soriano-Morales, Julien Ah-Pine, Sabine Loudcher: **Hypergraph Modelization of a Syntactically Annotated English Wikipedia Dump**. LREC 2016
- Adrien Guille, Edmundo-Pavel Soriano-Morales, Ciprian-Octavian Truica: **Topic modeling and hypergraph mining to analyze the EGC conference history**. EGC 2016
- Adrien Guille, Edmundo-Pavel Soriano-Morales: **TOM: A library for topic modeling and browsing**. EGC 2016
- Julien Ah-Pine, Edmundo-Pavel Soriano-Morales: **A Study of Synthetic Oversampling for Twitter Imbalanced Sentiment Analysis**. DMNLP@PKDD/ECML 2016
- Sabine Loudcher, Wararat Jakawat, Edmundo-Pavel Soriano-Morales, Cécile Favre: **Combining OLAP and information networks for bibliographic data analysis: a survey**. Scientometrics 103(2) 2015

# Thank you for your attention

- Edmundo-Pavel Soriano-Morales, Julien Ah-Pine, Sabine Loudcher: **Fusion Techniques for Named Entity Recognition and Word Sense Induction and Disambiguation**. DS 2017
- Edmundo-Pavel Soriano-Morales, Julien Ah-Pine, Sabine Loudcher: **Using a Heterogeneous Linguistic Network for Word Sense Induction and Disambiguation**. CICLING 2016
- Edmundo-Pavel Soriano-Morales, Julien Ah-Pine, Sabine Loudcher: **Hypergraph Modelization of a Syntactically Annotated English Wikipedia Dump**. LREC 2016
- Adrien Guille, Edmundo-Pavel Soriano-Morales, Ciprian-Octavian Truica: **Topic modeling and hypergraph mining to analyze the EGC conference history**. EGC 2016
- Adrien Guille, Edmundo-Pavel Soriano-Morales: **TOM: A library for topic modeling and browsing**. EGC 2016
- Julien Ah-Pine, Edmundo-Pavel Soriano-Morales: **A Study of Synthetic Oversampling for Twitter Imbalanced Sentiment Analysis**. DMNLP@PKDD/ECML 2016
- Sabine Loudcher, Wararat Jakawat, Edmundo-Pavel Soriano-Morales, Cécile Favre: **Combining OLAP and information networks for bibliographic data analysis: a survey**. Scientometrics 103(2) 2015

# Appendix

---

- **Creation of the linguistic network**
  - After preprocessing, we build a HLM  $G_{tw}$  that contains the co-occurent (lexically and syntactically) words for a target word  $tw$ .



## Proposed Method: Step Two

- **Computing the similarity between nodes**

- $G_{tw}$  is represented as a bipartite graph  $B_{tw}$ . Left nodes  $U$  represent words and right nodes  $W$  correspond to the hyperedges. An edge from a node  $u$  to a node  $w$  depicts the incidence of node  $u$  in hyperedge  $w$ .
- A similarity matrix  $S_{tw}$  of dimension  $|U| \times |U|$  is calculated using the Jaccard similarity: given  $n_{i,j} \in U$ , then  $Jaccard(i, j) = \frac{|N(i) \cap N(j)|}{|N(i) \cup N(j)|}$ .
- Induce a new incidence matrix  $F_{tw}$  from  $S_{tw}$  containing only the closest neighbours to each word  $n_i \in U$ . Each of these hyperedges represent a set of words that are deemed similar between them according to their Jaccard index value, which must be equal or higher than an assigned threshold  $th_1$ .

## Proposed Method: Step Three

- **Clustering words together**

- We select the top  $c$ -nodes in  $F_{tw}$  according to their degree. These nodes are candidate hubs, which must surpass a second threshold  $th_2$  to be considered as proper hubs. We use the average Jaccard measure defined for each node  $n$  as:

$$AvgJaccard(n) = \frac{1}{|hedges(n)|} \sum_{h \in hedges(n)} \frac{\sum_{\substack{i \in h \\ j \in h; i \neq j}} Jaccard(i, j)}{|h|}$$

where  $hedges(n)$  is the set of hyperedges  $n$  is incident in and its cardinality is defined as  $|hedges(n)|$ .  $|h|$  is the number of nodes in hyperedge  $h$ .

- Accepted hubs represent senses alongside with their co-occurrent words. The final set of senses is called  $SoS_{tw}$ .

# Structured Perceptron

---

**Algorithm 1:** Training phase of the Structured Perceptron

---

**Input:** Data  $x \in \mathcal{X}$

**Input:** True labels  $y \in \mathcal{Y}$

**Input:** Max number of iterations  $\text{MaxIteration}$

**Output:** A vector of lerned weights  $w$

```
1 for Iteration = 1 ... MaxIterations do
2   foreach  $(x, y) \in \mathcal{X}, \mathcal{Y}$  do
3      $\hat{y} = \arg \max_{\hat{y} \in \mathcal{Y}} w \cdot \Phi(x, \hat{y})$ 
4     if  $\hat{y} \neq y$  then
5        $w \leftarrow w + \Phi(x, y) - \Phi(x, \hat{y})$ 
6     end
7   end
8 end
9 return  $w$ 
```

---

The normalized Laplacian of an affinity (symmetric and positive) matrix  $W \in \mathbb{R}^{n \times n}$ , with  $w_{ij} = w_{ji} \geq 0$ , is defined as:

$$\mathcal{L} = I - D^{-1/2} W D^{-1/2} \quad (1)$$

where  $I$  is the identity matrix and  $D$  is the degree matrix of  $W$ .  $D$  is defined as the diagonal matrix with the degrees  $d_1, \dots, d_n$  on the diagonal. As  $W$  may not be an adjacency matrix, we define the degrees of each row in the matrix as:  $d_i = \sum_{j=1}^n w_{ij}$ .

## Spectral Clustering 2

Given a symmetric and positive similarity matrix  $W \in \mathbb{R}^{n \times n}$ , and a number of desired clusters  $k$ , the steps required to perform spectral clustering are:

1. Obtain the normalized Laplacian  $\mathcal{L}$  as indicated in Equation 1.
2. Obtain the first  $k$  eigenvectors  $u_{1 \dots k}$  of  $\mathcal{L}$ .
3. Store said eigenvectors as columns in a matrix  $V \in \mathbb{R}^{n \times k}$ . This matrix is akin to a lower-dimension projection of the original similarity matrix  $W$ .
4. Cluster the points in  $V_i$  with  $k$ -means. The clusters found and their members correspond to the cluster of the spectral algorithm.

## Precision and Recall for NER

### CONLL

Method	ALL_F	ALL_P	ALL_R
$M^T$	77.41	77.39	77.42
$M^L$	69.4	80.73	60.86
$M^S$	32.95	53.79	23.75
$E(M^T, M^L, M^S$	78.9	78.82	78.99
$E(M_L, E(E(M_T, L(M^T, X_F(S^T, M^T)))) ,$ $L(M^L, X_F(S^S, M^L))))$	79.67	80.45	78.9

## Precision and Recall for NER

### WNER

Method	ALL_F	ALL_P	ALL_R
$M^T$	77.5	77.83	77.18
$M^L$	69.17	79.07	61.47
$M^S$	28.47	38.36	22.45
$E(M^T, M^L, M^S$	80.04	80.26	79.83
$E(M_L, E(E(M_T, L(M^T, X_F(S^T, M^T)))),$ $L(M^L, X_F(S^S, M^L))))$	81.79	82.28	81.32

## Precision and Recall for NER

---

### WGLD

---

Method	ALL_F	ALL_P	ALL_R
$M^T$	59.66	60.37	58.75
$M^L$	52.34	68.42	42.38
$M^S$	25.49	36.55	19.56
$E(M^T, M^L, M^S$	63.2	63.88	62.54
$E(M_L, E(E(M_T, L(M^T, X_F(S^T, M^T)))),$ $L(M^L, X_F(S^S, M^L))))$	67.05	69.63	64.64

---



## Word: authority

- Dependencies
  - **process**: [neighborhood, lawyer, idea, seizure, council, subsidiary, need, collector, court, office]
  - **cabinet**: [create, trade, stability, manager, swine, department, misconduct, settlement, economist, math]
- Lexical
  - **shop**: [shop, sketch, young, month, pareo, woman, moscow, opposite, tahitian, handler, verso]
  - **supply**: [supply, justice, money, hugo, telephone, authority, initiative, alberta, bundesbank, utility, impact]
  - **evidence**: [council, machine, court, august, district, instance, fulham, auditor, hammersmith, plant]

### Word: authority (4 senses in Gold Standard)

- Best fusion operator:  $X_F(S^L, M^L)$ 
  - **block**: [allow, including, study, told, seek, make, support, claim, provide, lawyers]
  - **veto**: [says, court, companies, years, does, law, loans, congress, trading, days]
  - **federal**: [federal, president, new, u.s., line-item, banks, local, company, airline, transportation]
  - **government**: [legislation, million, bush, year, people, billion, secretary, department, officials, house]