# How can we capture multiword expressions?

*Seongmin Mun*[1], Guillaume Desagulier[2], Anne Lacheret[3] , Kyungwon Lee[4]

[1] Lifemedia Interdisciplinary Program, Ajou University, South Korea
[1,3] UMR 7114 MoDyCo - CNRS, University Paris Nanterre, France
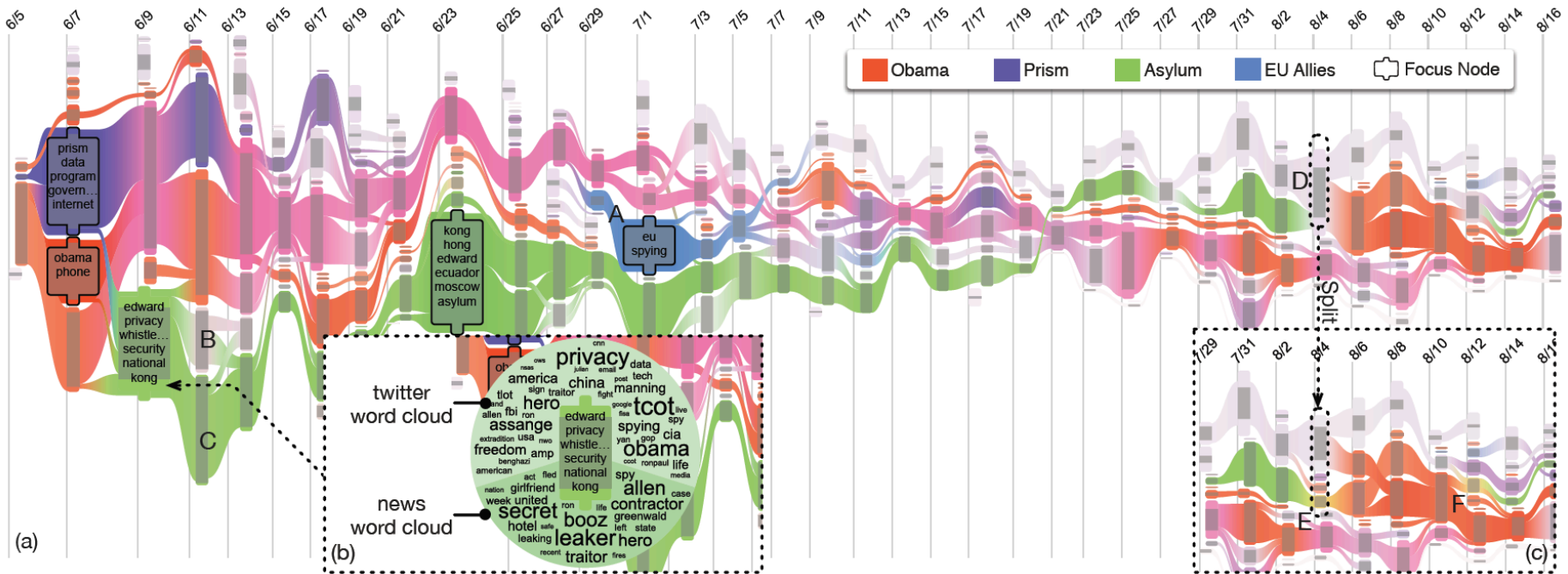[2] UMR 7114 MoDyCo - University Paris 8, CNRS, University Nanterre
[4] Department of Digital Media, Ajou University, South Korea

Université
Paris Nanterre  AJOU UNIVERSITY

# Introduction

Topics in a text corpus include features and information.

Analyzing these topics can improve a user's understanding of the corpus.

Université
Paris Nanterre  AJOU UNIVERSITY

# Previous studies



WEIWEI CUI SHIXIA LIU Z. W. H. W.: How hierarchical topics evolve in large text corpora. In IEEE Transactions on Visualization and Computer Graphics (2014), vol. 20, pp. 2281–2290.

# Research background and purpose

Topics can be broadly divided into two categories.

Université
Paris Nanterre  AJOU UNIVERSITY

# Research background and purpose

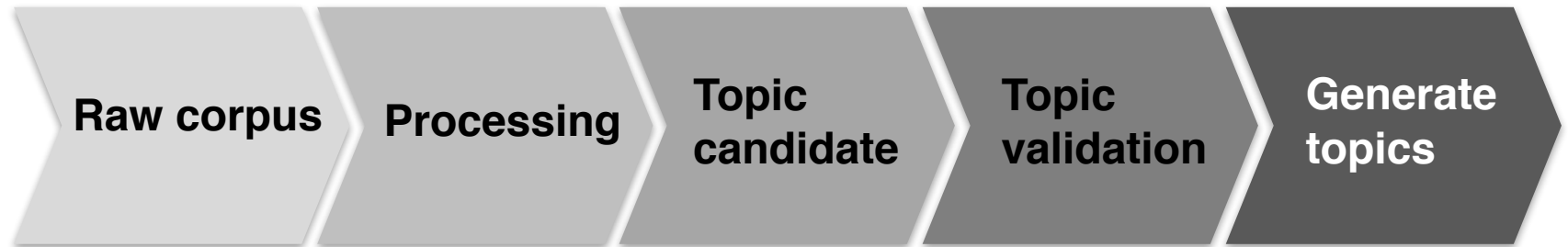"With profound gratitude and great humility, I accept your nomination for the presidency of the United States."

Université Paris Nanterre

AJOU UNIVERSITY

# Research background and purpose

"With profound ***gratitude*** and great humility, I accept your nomination for the presidency of the United States."

***Gratitude*** ➡ meaning that can be expressed in one word

Université Paris Nanterre  AJOU UNIVERSITY

# Research background and purpose

"With profound gratitude and great humility, I accept your nomination for the presidency of the ***United States***."

***United States*** ➡ meaning must be described using a combination of words.

# Research background and purpose

How can we capture multiword expressions?

To this aim, we designed an algorithm.

Université
Paris Nanterre  AJOU UNIVERSITY

# Data processing

**Raw corpus** → **Processing** → **Topic candidate** → **Topic validation** → **Generate topics**

Université
Paris Nanterre  AJOU UNIVERSITY

# Data processing

Raw corpus → Processing → Topic candidate → Topic validation → Generate topics

Université Paris Nanterre  AJOU UNIVERSITY

# Data processing

**Raw corpus**
(U.S. president speeches)



https://millercenter.org/the-presidency/presidential-speeches

# Data processing


Raw corpus | Processing | Topic candidate | Topic validation | Generate topics

**Raw corpus**
(U.S. president speeches)

Université Paris Nanterre    AJOU UNIVERSITY

# Data processing

**Raw corpus**
(U.S. president speeches)

Université Paris Nanterre  AJOU UNIVERSITY

# Data processing

Raw corpus → Processing → Topic candidate → Topic validation → Generate topics

Université Paris Nanterre   AJOU UNIVERSITY

# Data processing

**Processing**

- N-grams
- POS tagging

**Pre-processing**

- Cleaning with RegExp
- Lemmatization
- Tokenization
- Lowercasing

N-gram method is a contiguous sequence of *N* items from a given sequence of text.

Université Paris Nanterre   AJOU UNIVERSITY

# Data processing

## Processing

- N-grams
- POS tagging

## Pre-processing

- Cleaning with RegExp
- Lemmatization
- Tokenization
- Lowercasing

"Time flies like an arrow."

Université Paris Nanterre    AJOU UNIVERSITY

# Data processing

**Processing**

- N-grams
- POS tagging

**Pre-processing**

- Cleaning with RegExp
- Lemmatization
- Tokenization
- Lowercasing

"Time flies like an arrow."

Unigram : Time, flies, like, an, arrow.
Bigram   : Time flies, flies like, like an, an arrow.
Trigram  : Time flies like, flies like an, like an arrow.

Université
Paris Nanterre   AJOU UNIVERSITY

# Data processing

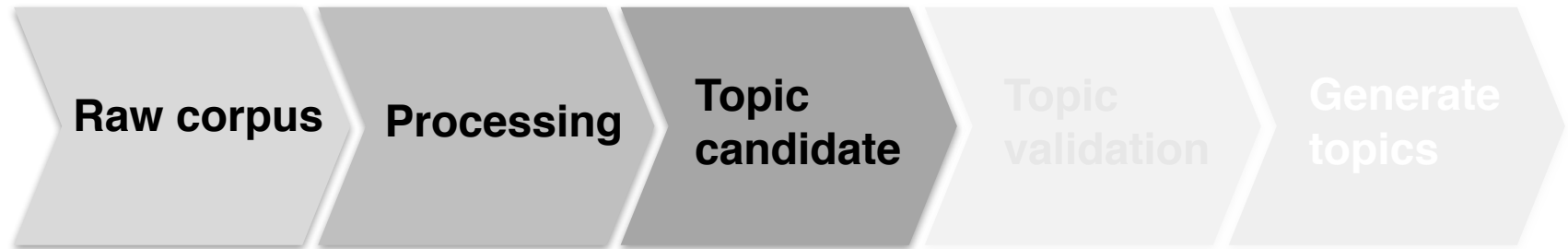## Processing

- N-grams
- POS tagging

### Pre-processing

- Cleaning with RegExp
- Lemmatization
- Tokenization
- Lowercasing

**2.George_W._Bush_3gram_result.txt**
```
1   number,text,count
2   1,the united states,128
3   2,men and women,73
4   3,the middle east,73
5   4,the american people,70
6   5,and we will,65
7   6,of the world,43
8   7,in the middle,41
9   8,one of the,39
10  9,in the world,39
11  10,weapons of mass,38
12  11,members of congress,38
13  12,and that is,38
14  13,i want to,38
15  14,it is the,36
16  15,the united nations,36
17  16,of our country,34
18  17,of the united,34
19  18,a lot of,34
20  19,thank you for,33
21  20,ask you to,33
22  21,is going to,32
23  22,of mass destruction,32
24  23,i ask you,32
25  24,want to thank,31
26  25,around the world,30
```

**1.Barack_Obama_2gram_result.txt**
```
1   number,text,count
2   1,of the,802
3   2,in the,787
4   3,to the,429
5   4,that is,421
6   5,of our,402
7   6,it is,379
8   7,and the,378
9   8,we have,375
10  9,for the,336
11  10,we can,319
12  11,that we,316
13  12,we will,305
14  13,to be,299
15  14,on the,289
16  15,the world,280
17  16,going to,256
18  17,we are,251
19  18,and i,240
20  19,is not,229
21  20,that the,221
22  21,want to,213
23  22,will be,212
24  23,the united,212
25  24,and we,212
26  25,is the,207
```

Université Paris Nanterre  AJOU UNIVERSITY

# Data processing



Raw corpus → Processing → Topic candidate → Topic validation → Generate topics

Université Paris Nanterre  AJOU UNIVERSITY

# Data processing

Raw corpus → Processing → **Topic candidate** → Topic validation → Generate topics
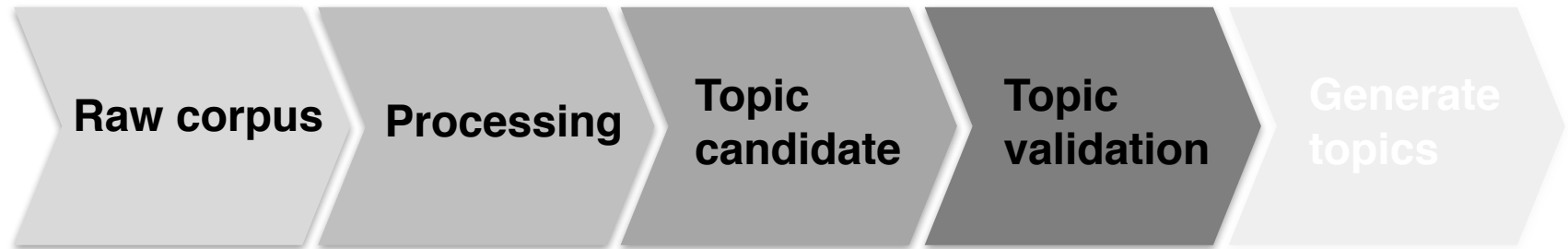
**Topic candidate extraction & filtering**

- Frequency counting
- Filters :
  - ✓ Stopwords
  - ✓ Thresholds

```
public class Remove {

    public static String[] stopwords_three = {
        "it",
        "is",
        "we",
        "are",
        "this",
        "be",
        "may",
        "would",
        "am",
        "more",
        "than",
        "do",
        "that",
        "can",
        "not",
        "could",
        "sould",
        "shall",
        "will",
        "were",
        "was",
        "might",
        "all",
        "so",
        "you",
        "he",
        "him",
        "his",
        "she",
        "her",
        "your",
        "me",
```

```
public static String[] stopwords = {
    "t",
    "stv",
    "de",
    "ss",
    "el",
    "ho",
    "em",
    "men",
    "ere",
    "ad",
    "la",
    "pro",
    "fe",
    "wit",
    "vi",
    "ted",
    "eve",
    "iv",
    "era",
    "ear",
    "va",
    "ive",
    "led",
    "owe",
    "tho",
    "gi",
    "a",
    "will",
    "able",
    "about",
    "above",
    "abst",
```

Université Paris Nanterre  AJOU UNIVERSITY

# Data processing



Raw corpus → Processing → Topic candidate → Topic validation → Generate topics

Université Paris Nanterre  AJOU UNIVERSITY

# Data processing

Raw corpus → Processing → Topic candidate → **Topic validation** → Generate topics

## Topic validation

- Human annotation

- Matching with **Dictionaries**

## English dictionaries

- THE DEVIL'S DICTIONARY ((C)1911 Released April 15 1993)
- Easton's 1897 Bible Dictionary
- Elements database 20001107
- The Free On-line Dictionary of Computing (27 SEP 03)
- U.S. Gazetteer (1990)
- The Collaborative International Dictionary of English v.0.44
- Hitchcock's Bible Names Dictionary (late 1800's)
- Jargon File (4.3.1, 29 June 2001)
- Virtual Entity of Relevant Acronyms (Version 1.9, June 2002)
- WordNet (r) 2.0
- CIA World Factbook 2002
- User Dictionary

Université Paris Nanterre  AJOU UNIVERSITY

# Data processing

Raw corpus → Processing → Topic candidate → Topic validation → **Generate topics**

Université Paris Nanterre    AJOU UNIVERSITY

# Visual system



http://ressources.modyco.fr/sm/MultiwordVis/

Université Paris Nanterre  AJOU UNIVERSITY

# Ambiguous sentence

"Shall I wake him up?"

# Ambiguous sentence

"Shall I **wake** him **up**?"

We can't extract wake up if we only use N-gram algorithm.

Université
Paris Nanterre  AJOU UNIVERSITY

# Dependency tag



Dependency tag can provide a simple description of the grammatical relationships in a sentence.

# Improving algorithm

```
Result of dependency graph below

dependency graph:
-> wake/VBP (root)
  -> Shall/NNP (nsubj)
    -> I/PRP (dep)
  -> him/PRP (dobj)
  -> up/RP (compound:prt)
  -> ?/. (punct)
```

```
Result of multiword candidates

wake Shall
Shall I
wake Shall I
wake him
wake up
wake ?
```

Université
Paris Nanterre   AJOU UNIVERSITY

# Improving algorithm

```
Final result below

0. wake is meaningful : wake
1. shall is meaningful : shall
2. i is meaningful : i
3. up is meaningful : up
4. shall i is meaningful : shall i
5. him is meaningful : him
```

```
Final result below

0. wake is meaningful : wake
1. shall i is meaningful : shall i
2. i is meaningful : i
3. wake up is meaningful : wake up
4. up is meaningful : up
5. him is meaningful : him
6. shall is meaningful : shall
```
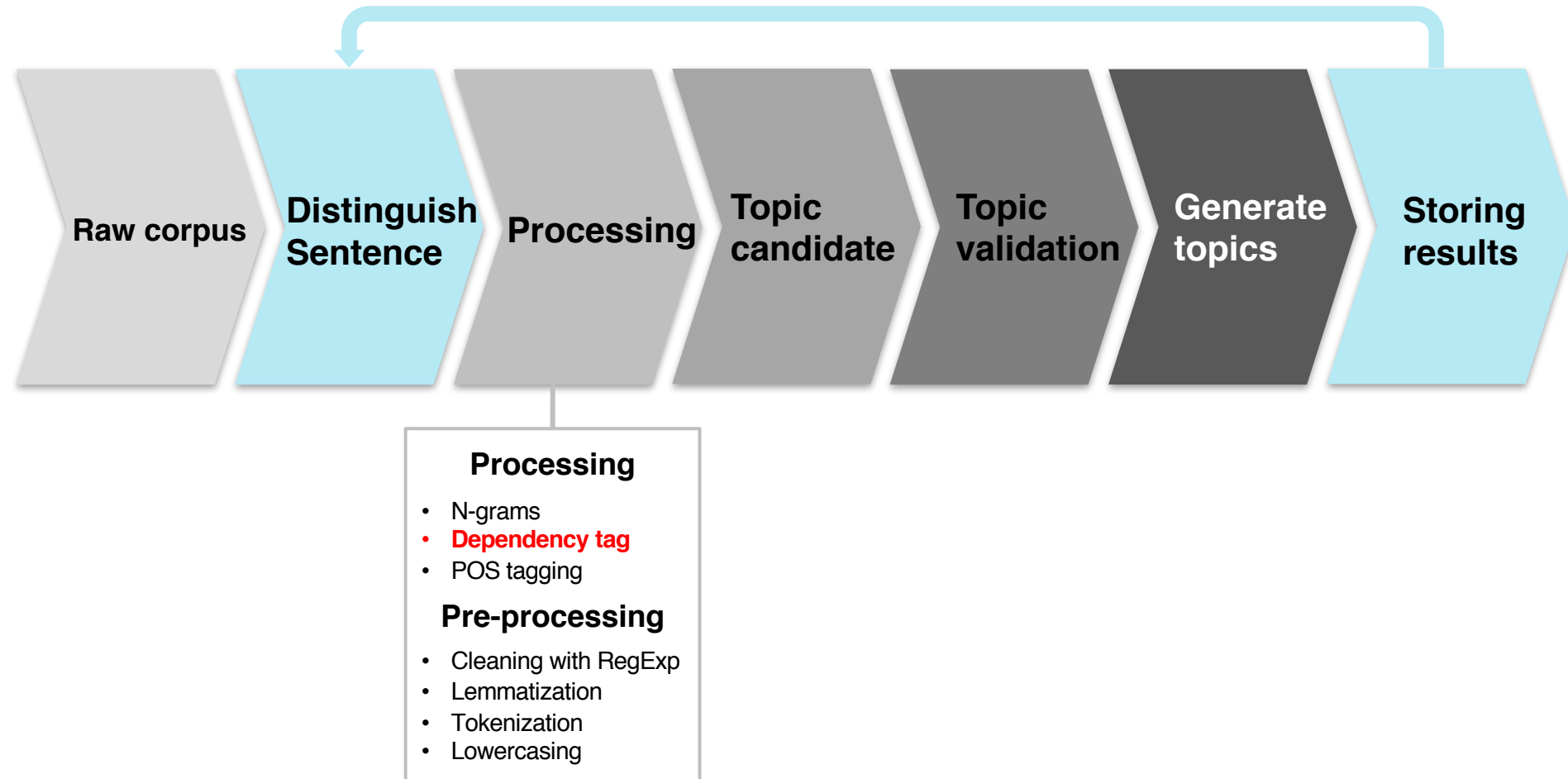
N-gram                    Dependency tag

# Data processing



Raw corpus → Distinguish Sentence → Processing → Topic candidate → Topic validation → Generate topics → Storing results

**Processing**

- N-grams
- **Dependency tag**
- POS tagging

**Pre-processing**

- Cleaning with RegExp
- Lemmatization
- Tokenization
- Lowercasing

Université Paris Nanterre    AJOU UNIVERSITY

# Q&A

# Thank you for listening.
## stat34@ajou.ac.kr