

**PAN 2015**

13th evaluation lab on uncovering plagiarism, authorship, and social software misuse

# JOINT TALK ON THREE DATA SUBMISSIONS TO TEXT ALIGNMENT AND ONE SOURCE RETRIEVAL ALGORITHM

Mostafa Dehghani

ICT Research Institute, ACECR, Iran

September, 10, 2015



# Outline of My Talk

# Outline of My Talk

2

## A. Data Submissions to Text Alignment:

- Developing Monolingual Persian Corpus for Extrinsic Plagiarism Detection Using Artificial Obfuscation
- Developing Monolingual English Corpus for Plagiarism Detection using Human Annotated Paraphrase Corpus
- Developing Bilingual Plagiarism Detection Corpus Using Sentence Aligned Parallel Corpus

# Outline of My Talk

2

## A. Data Submissions to Text Alignment:

- Developing Monolingual Persian Corpus for Extrinsic Plagiarism Detection Using Artificial Obfuscation
- Developing Monolingual English Corpus for Plagiarism Detection using Human Annotated Paraphrase Corpus
- Developing Bilingual Plagiarism Detection Corpus Using Sentence Aligned Parallel Corpus
- Evaluation of Text Reuse Corpora for Text Alignment Task of plagiarism Detection

# Outline of My Talk

2

## **A. Data Submissions to Text Alignment:**

- Developing Monolingual Persian Corpus for Extrinsic Plagiarism Detection Using Artificial Obfuscation
- Developing Monolingual English Corpus for Plagiarism Detection using Human Annotated Paraphrase Corpus
- Developing Bilingual Plagiarism Detection Corpus Using Sentence Aligned Parallel Corpus
- Evaluation of Text Reuse Corpora for Text Alignment Task of plagiarism Detection

## **B. Source Retrieval Plagiarism Detection based on Noun Phrase and Keyword Phrase Extraction**

A

## Data Submissions to Text Alignment

# Corpus Construction Steps

4

- 
- 
- 
- 
-

# Corpus Construction Steps

4

## ➤ Preprocessing

- 
- 
- 
-



# Corpus Construction Steps

4

- Preprocessing
- Clustering
- 
- 
-

# Corpus Construction Steps

4

- Preprocessing
- Clustering
- Fragment Extraction
- 
-

# Corpus Construction Steps

4

- Preprocessing
- Clustering
- Fragment Extraction
- Fragment Obfuscation

.

# Corpus Construction Steps

4

- Preprocessing
- Clustering
- Fragment Extraction
- Fragment Obfuscation
- Inserting Plagiarism Cases into Documents

# Developing Monolingual Persian Corpus for Extrinsic Plagiarism Detection Using Artificial Obfuscation

Data resource:

Wikipedia Articles

# Mono Lingual Persian Corpus

6

.

.

# Mono Lingual Persian Corpus

6

## ➤ Preprocessing

- Persian is one of the Indo-European languages which have borrowed its script from Arabic, a member of the Semitic language family

▪

# Mono Lingual Persian Corpus

6

## ➤ Preprocessing

- Persian is one of the Indo-European languages which have borrowed its script from Arabic, a member of the Semitic language family

## ➤ Clustering

- In this step, collection of Wikipedia documents clustered into different topically related groups
- A bipartite graph of documents-categories was created to cluster the documents
- In the next step, the Infomap community detection algorithm was applied to the graph and all communities were detected
- Finally, Documents within a community are considered as one cluster



# Mono Lingual Persian Corpus

7

➤ Fragment Extraction

▪

# Mono Lingual Persian Corpus

7

## ➤ Fragment Extraction

- Divided Documents into Two Categories:
  - 50% Source Documents
  - 50% Suspicious Documents : 25% with Plagiarism – 25% no Plagiarism

▪

# Mono Lingual Persian Corpus

7

## ➤ Fragment Extraction

- Divided Documents into Two Categories:
  - 50% Source Documents
  - 50% Suspicious Documents : 25% with Plagiarism – 25% no Plagiarism
- The task of the fragment extraction is to extract fragments from source documents.

▪

# Mono Lingual Persian Corpus

7

## ➤ Fragment Extraction

- Divided Documents into Two Categories:
  - 50% Source Documents
  - 50% Suspicious Documents : 25% with Plagiarism – 25% no Plagiarism
- The task of the fragment extraction is to extract fragments from source documents.

▪

Fragment Length	
<b>Short</b>	30 – 50 words
<b>Medium</b>	150 – 250 words
<b>Long</b>	300 – 500 words

# Mono Lingual Persian Corpus

7

## ➤ Fragment Extraction

- Divided Documents into Two Categories:
  - 50% Source Documents
  - 50% Suspicious Documents : 25% with Plagiarism – 25% no Plagiarism
- The task of the fragment extraction is to extract fragments from source documents.

▪

# Mono Lingual Persian Corpus

7

## ➤ Fragment Extraction

- Divided Documents into Two Categories:
  - 50% Source Documents
  - 50% Suspicious Documents : 25% with Plagiarism – 25% no Plagiarism
- The task of the fragment extraction is to extract fragments from source documents.

## ➤ Fragment Obfuscation

# Mono Lingual Persian Corpus

7

## ➤ Fragment Extraction

- Divided Documents into Two Categories:
  - 50% Source Documents
  - 50% Suspicious Documents : 25% with Plagiarism – 25% no Plagiarism
- The task of the fragment extraction is to extract fragments from source documents.

## ➤ Fragment Obfuscation

- Artificial Obfuscation
  - None (No Obfuscation)
  - Random Change of Order
  - POS-preserving Change of Order
  - Synonym Substitution
  - Addition / Deletion

# Mono Lingual Persian Corpus

8

- Inserting Plagiarism Cases into suspicious Documents



# Mono Lingual Persian Corpus

8

- Inserting Plagiarism Cases into suspicious Documents
  - In this step, according to suspicious document's length, one or more plagiarism cases are selected.

# Mono Lingual Persian Corpus

8

- Inserting Plagiarism Cases into suspicious Documents
  - In this step, according to suspicious document's length, one or more plagiarism cases are selected.

Plagiarism per Document	
<b>Little</b>	5% - 20%
<b>Medium</b>	20% - 50%
<b>Much</b>	50% - 80%
<b>Very Much</b>	80% - 100%

# Mono Lingual Persian Corpus

8

- Inserting Plagiarism Cases into suspicious Documents
  - In this step, according to suspicious document's length, one or more plagiarism cases are selected.

# Mono Lingual Persian Corpus

8

- Inserting Plagiarism Cases into suspicious Documents
  - In this step, according to suspicious document's length, one or more plagiarism cases are selected.
  - Each of selected cases inserted at random positions in suspicious document.

# Mono Lingual Persian Corpus

8

- Inserting Plagiarism Cases into suspicious Documents
  - In this step, according to suspicious document's length, one or more plagiarism cases are selected.
  - Each of selected cases inserted at random positions in suspicious document.
  - Each suspicious document and its corresponding source documents are selected from one cluster.

# Mono Lingual Persian Corpus

8

- Inserting Plagiarism Cases into suspicious Documents
  - In this step, according to suspicious document's length, one or more plagiarism cases are selected.
  - Each of selected cases inserted at random positions in suspicious document.
  - Each suspicious document and its corresponding source documents are selected from one cluster.

```
<?xml version="1.0" encoding="UTF-8" ?>
<document xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" ....
  <feature name="project-gutenberg" etext_number="7243" url="http://www.gutenberg.org/files/....
  <feature name="language" value="en" />
  <feature name="artificial-plagiarism" translation="false" obfuscation="none"
    this_offset="487" this_length="4218" source_reference="source-document03471.txt"
    source_offset="10866" source_length="4226" />

  <feature name="artificial-plagiarism" translation="false" obfuscation="low"
    this_offset="7507" this_length="1872" source_reference="source-document03471.txt"
    source_offset="4846" source_length="1792" />

  <feature name="artificial-plagiarism" translation="false" obfuscation="low"
    this_offset="10626" this_length="805" source_reference="source-document03471.txt"
    source_offset="2399" source_length="800" />
</document>
```

# Mono Lingual Persian Corpus

9

## ➤ Results

<b>Documents</b>	
The number of source documents:	1057
The number of suspicious documents:	
With plagiarism:	529
No plagiarism:	528
<b>Plagiarism Cases</b>	
The number of plagiarism cases:	
No obfuscation cases:	259
With obfuscation cases:	564
<b>Plagiarism per Document</b>	
The number of Little plagiarized documents:	301
The number of Medium plagiarized documents:	80
The number of Much plagiarized documents:	96
The number of Very much plagiarized documents:	52

# Developing Monolingual English Corpus for Plagiarism Detection using Human Annotated Paraphrase Corpus

## Data resources:

- Wikipedia Articles
- SemEval Dataset



# Mono Lingual English Corpus

11

➤ Clustering

.

# Mono Lingual English Corpus

11

- Clustering
- Fragment Extraction
  - Method 1: The fragments are extracted from source documents.
  - Method 2: The fragments are generated based on SemEval dataset sentences.

# Mono Lingual English Corpus

11

- Clustering
- Fragment Extraction
  - Method 1: The fragments are extracted from source documents.
  - Method 2: The fragments are generated based on SemEval dataset sentences.

Fragment Length	
<b>Short</b>	3 – 5 sentences
<b>Medium</b>	6 – 8 sentences
<b>Long</b>	9 – 12 sentences

# Mono Lingual English Corpus

11

- Clustering
- Fragment Extraction
  - Method 1: The fragments are extracted from source documents.
  - Method 2: The fragments are generated based on SemEval dataset sentences.

# Mono Lingual English Corpus

12

- Fragment Obfuscation

▪

# Mono Lingual English Corpus

12

- Fragment Obfuscation
  - Artificial Obfuscation

▪

# Mono Lingual English Corpus

12

## ➤ Fragment Obfuscation

- Artificial Obfuscation
- Simulated Obfuscation
  - The pairs of sentences from the SemEval dataset with their corresponding similarity score are used for constructing the simulated plagiarism cases.
  - To consider the degree of obfuscation in plagiarized fragments, a combination of sentences with a variety of similarity scores is used in a fragment.

▪

# Mono Lingual English Corpus

12

## ➤ Fragment Obfuscation

- Artificial Obfuscation
- Simulated Obfuscation
  - The pairs of sentences from the SemEval dataset with their corresponding similarity score are used for constructing the simulated plagiarism cases.
  - To consider the degree of obfuscation in plagiarized fragments, a combination of sentences with a variety of similarity scores is used in a fragment.

Degree	Similarity Scores of Sentences		
	3	4	5
Low	-	1% - 15%	85% - 100%
Medium	25% - 45%		55% - 75%
High	45% - 65%		35% - 55%



# Mono Lingual English Corpus

12

## ➤ Fragment Obfuscation

- Artificial Obfuscation
- Simulated Obfuscation
  - The pairs of sentences from the SemEval dataset with their corresponding similarity score are used for constructing the simulated plagiarism cases.
  - To consider the degree of obfuscation in plagiarized fragments, a combination of sentences with a variety of similarity scores is used in a fragment.

▪

# Mono Lingual English Corpus

12

## ➤ Fragment Obfuscation

- Artificial Obfuscation
- Simulated Obfuscation
  - The pairs of sentences from the SemEval dataset with their corresponding similarity score are used for constructing the simulated plagiarism cases.
  - To consider the degree of obfuscation in plagiarized fragments, a combination of sentences with a variety of similarity scores is used in a fragment.

## ➤ Inserting Plagiarism Cases into Documents

# Mono Lingual English Corpus

12

## ➤ Fragment Obfuscation

- Artificial Obfuscation
- Simulated Obfuscation
  - The pairs of sentences from the SemEval dataset with their corresponding similarity score are used for constructing the simulated plagiarism cases.
  - To consider the degree of obfuscation in plagiarized fragments, a combination of sentences with a variety of similarity scores is used in a fragment.

## ➤ Inserting Plagiarism Cases into Documents

Plagiarism per Document	
<b>Hardly</b>	5% - 20%
<b>Medium</b>	20% - 40%
<b>Much</b>	40% - 60%

# Mono Lingual English Corpus

13

## ➤ Results

<b>Statistics</b>	
<b>Documents</b>	
The number of source documents:	3309
The number of suspicious documents:	952
<b>Plagiarism per Document</b>	
Hardly (5% - 20%)	60%
Medium (20% - 40%)	25%
Much (40% - 60%)	15%
<b>Plagiarism Cases</b>	
The number of plagiarism cases:	
- No obfuscation cases:	10%
- Random obfuscation:	78%
- Simulated obfuscation:	12%
<b>Case Length Statistics</b>	
Short (3 – 5 sentences):	50%
Medium (6 – 8 sentences):	32%
Long (9 – 12 sentences):	18%

# Developing Bilingual Plagiarism Detection Corpus Using Sentence Aligned Parallel Corpus

## **Data resources:**

- Wikipedia Articles
- Persian-English Parallel Corpus

# Bilingual Persian-English Corpus

15

- Clustering

- 

-

# Bilingual Persian-English Corpus

15

## ➤ Clustering

### Parallel Sentences Clustering

1. Persian Wikipedia documents were indexed by the Apache Lucene library.
2. We built a query from each Persian sentence
3. The query was searched in the indexed documents and returns the top document.
4. A bipartite graph of return documents-categories was created. Then, the info-map community detection algorithm was applied to the graph and all communities were detected. Documents within a community are considered as one cluster.
5. Finally, parallel sentences were assigned to the documents in the same cluster.

▪

# Bilingual Persian-English Corpus

15

## ➤ Clustering

### Parallel Sentences Clustering

1. Persian Wikipedia documents were indexed by the Apache Lucene library.
2. We built a query from each Persian sentence
3. The query was searched in the indexed documents and returns the top document.
4. A bipartite graph of return documents-categories was created. Then, the info-map community detection algorithm was applied to the graph and all communities were detected. Documents within a community are considered as one cluster.
5. Finally, parallel sentences were assigned to the documents in the same cluster.

### Documents Clustering

- For each cluster of return documents in the previous stage, the categories of documents have been extracted and considered as label of that cluster.
- The basic documents collected into different topically related clusters based on their categories. The documents are assigned to the cluster with maximum common categories.



# Bilingual Persian-English Corpus

16

- Fragment Extraction

▪

# Bilingual Persian-English Corpus

16

## ➤ Fragment Extraction

- Plagiarism cases are constructed from parallel sentences.
- Source fragments were generated from sentences in the English language and plagiarized fragments were constructed by Persian sentences paired with English sentences.

▪

# Bilingual Persian-English Corpus

16

## ➤ Fragment Extraction

- Plagiarism cases are constructed from parallel sentences.
- Source fragments were generated from sentences in the English language and plagiarized fragments were constructed by Persian sentences paired with English sentences.

Fragment Length	
<b>Short</b>	3 – 5 sentences
<b>Medium</b>	5 – 10 sentences
<b>Long</b>	10 – 15 sentences

# Bilingual Persian-English Corpus

16

## ➤ Fragment Extraction

- Plagiarism cases are constructed from parallel sentences.
- Source fragments were generated from sentences in the English language and plagiarized fragments were constructed by Persian sentences paired with English sentences.

▪

# Bilingual Persian-English Corpus

16

## ➤ Fragment Extraction

- Plagiarism cases are constructed from parallel sentences.
- Source fragments were generated from sentences in the English language and plagiarized fragments were constructed by Persian sentences paired with English sentences.

## ➤ Fragment Obfuscation

- To consider the degree of obfuscation in plagiarized fragments, a combination of sentences with different similarity score were chosen.

# Bilingual Persian-English Corpus

16

## ➤ Fragment Extraction

- Plagiarism cases are constructed from parallel sentences.
- Source fragments were generated from sentences in the English language and plagiarized fragments were constructed by Persian sentences paired with English sentences.

## ➤ Fragment Obfuscation

- To consider the degree of obfuscation in plagiarized fragments, a combination of sentences with different similarity score were chosen.

Degree	Similarity scores of sentences in fragments		
	1 - 0.85	0.85 - 0.65	0.65 - 0.85
<b>Low</b>	100%	-	-
<b>Medium</b>	55% - 75%	25% - 45%	-
<b>High</b>	35% - 55%	-	45% - 65%

# Bilingual Persian-English Corpus

17

- Inserting Plagiarism Cases into Documents

# Bilingual Persian-English Corpus

17

- Inserting Plagiarism Cases into Documents
  - In this step, according to suspicious document's length, one or more plagiarism cases are selected.
  - Persian documents considering as suspicious documents and source documents are English documents.



# Bilingual Persian-English Corpus

17

## ➤ Inserting Plagiarism Cases into Documents

- In this step, according to suspicious document's length, one or more plagiarism cases are selected.
- Persian documents considering as suspicious documents and source documents are English documents.

Plagiarism per Document	
<b>Low</b>	5% - 20%
<b>Medium</b>	20% - 40%
<b>High</b>	40% - 60%

# Bilingual Persian-English Corpus

17

## ➤ Inserting Plagiarism Cases into Documents

- In this step, according to suspicious document's length, one or more plagiarism cases are selected.
- Persian documents considering as suspicious documents and source documents are English documents.
- English fragment inserted at random positions in source documents and its corresponding Persian fragments has been inserted into suspicious documents.
- Each suspicious document and its corresponding source documents are selected from one cluster.

Plagiarism per Document	
<b>Low</b>	5% - 20%
<b>Medium</b>	20% - 40%
<b>High</b>	40% - 60%

# Bilingual Persian-English Corpus

18

## ➤ Results

<b>Documents</b>	
The number of source documents (English):	19973
The number of suspicious documents (Persian):	
• With plagiarism:	3571
• No plagiarism:	3571
<b>Plagiarism cases</b>	
The number of plagiarism cases:	11200
<b>Plagiarism per Document</b>	
The number of Little plagiarized documents	2035
The number of Medium plagiarized documents	536
The number of Much plagiarized documents	642
The number of Very much plagiarized documents	58

# Evaluation of Text Reuse Corpora for Text Alignment Task of plagiarism Detection

Evaluation of Corpus Submissions to PAN 2015

# Corpora Statistical Information

# Corpora Statistical Information

20

	<b>cheema15</b>	<b>hanif15</b>	<b>Kong15</b>	<b>Alvi15</b>	<b>Palkovskii15</b>
<b>Type of Corpus</b>	Mono-Lingual	Bi-Lingual	Mono-Lingual	Mono-Lingual	Mono-Lingual
<b>Source-Suspicious Language</b>	English-English	Urdu-English	Chinese-Chinese	English-English	English- English
<b>Resource Documents</b>	Gutenberg books and Wikipedia	Wikipedia pages	Chinese thesis and <a href="http://wenku.baidu.com/">http://wenku.baidu.com/</a> website	“The Complete Grimm’s Fairy Tales” book	Internet web pages crawling

# Corpora Statistical Information

# Corpora Statistical Information

	Cheema15	Hanif15	Kong15	Alvi15	Palkovskii15
<b>Number of Docs</b>					
• <b>Suspicious Docs</b>	248	250	4	90	1175
• <b>Source Docs</b>	248	250	78	70	1950
<b>Length of Docs (in chars)</b>					
• <b>Min Length</b>	2263	361	394	514	519
• <b>Max Length</b>	22471	74083	121829	45222	517925
• <b>Average Length</b>	7239	4382	42839	7718	6512
<b>Length of Plagiarisms Cases (in chars)</b>					
• <b>Min Length</b>	134	78	62	259	157
• <b>Max Length</b>	2439	849	2748	1160	14336
• <b>Average Length</b>	503	361	423	464	782



# Corpora Statistical Information

# Corpora Statistical Information

20

Obfuscation Strategies	Cheema15	Hanif15	Kong15	Alvi15	Palkovskii15
<b>Simulated</b>	123	135	-	-	-
<b>Real</b>	-	-	109	-	-
<b>Automatic</b>	-	-	-	25	-
<b>Retelling-Human</b>	-	-	-	25	-
<b>Character-Substitution</b>	-	-	-	25	-
<b>Translation</b>	-	-	-	-	618
<b>Summary</b>	-	-	-	-	1292
<b>Random</b>	-	-	-	-	626
<b>None</b>	-	-	-	-	624
<b>Sum</b>	123	135	109	75	3160

# Manual Evaluation of Corpora

- 

- 

- 

-

# Manual Evaluation of Corpora

21

- Manually investigate twenty pairs of corresponding source and suspicious fragments in each corpus
  - 
  - 
  -

# Manual Evaluation of Corpora

21

- Manually investigate twenty pairs of corresponding source and suspicious fragments in each corpus
  - Changes in syntactic structure between source and plagiarized passage
  - 
  -

# Manual Evaluation of Corpora

21

- Manually investigate twenty pairs of corresponding source and suspicious fragments in each corpus
  - Changes in syntactic structure between source and plagiarized passage
  - Concept preserving from source passage to plagiarized passage

.

# Manual Evaluation of Corpora

21

- Manually investigate twenty pairs of corresponding source and suspicious fragments in each corpus
  - Changes in syntactic structure between source and plagiarized passage
  - Concept preserving from source passage to plagiarized passage
  - Distribution of obfuscation types in suspicious documents

# Automatic Evaluation of Corpora

22

- 

- 

- 

- 

- 

- 

-



# Automatic Evaluation of Corpora

22

- Evaluating two remained obfuscation scenarios:
  - Real obfuscation from Kong15 corpus
  - Summary obfuscation from Palkovskii15 corpus

.

.

.

.

# Automatic Evaluation of Corpora

22

- Evaluating two remained obfuscation scenarios:
  - Real obfuscation from Kong15 corpus
  - Summary obfuscation from Palkovskii15 corpus
- For Kong15 corpus

.

.

.

# Automatic Evaluation of Corpora

22

- Evaluating two remained obfuscation scenarios:
  - Real obfuscation from Kong15 corpus
  - Summary obfuscation from Palkovskii15 corpus
- For Kong15 corpus
  - All source and correspond suspicious fragments are extracted, and the total number of similar “characters n-grams” between source and suspicious plagiarized passages are calculated for n in range of one to four .

# Automatic Evaluation of Corpora

22

- Evaluating two remained obfuscation scenarios:
  - Real obfuscation from Kong15 corpus
  - Summary obfuscation from Palkovskii15 corpus
- For Kong15 corpus
  - All source and correspond suspicious fragments are extracted, and the total number of similar “characters n-grams” between source and suspicious plagiarized passages are calculated for n in range of one to four .
- For evaluation of summary obfuscation

# Automatic Evaluation of Corpora

22

- Evaluating two remained obfuscation scenarios:
  - Real obfuscation from Kong15 corpus
  - Summary obfuscation from Palkovskii15 corpus
- For Kong15 corpus
  - All source and correspond suspicious fragments are extracted, and the total number of similar “characters n-grams” between source and suspicious plagiarized passages are calculated for n in range of one to four .
- For evaluation of summary obfuscation
  - From the point of “concept preserving” measure, we have extracted 10% of top words from source fragments based on tf.idf weight.

# B

## Source Retrieval based on Noun and Keyword Phrase Extraction

Data resources:

External PD Corpus of PAN 2011

# Approach in Use: Five Steps

24

- 
- 
- 
- 
-

# Approach in Use: Five Steps

24

- Suspicious Document Chunking

- 
- 
- 
-



# Approach in Use: Five Steps

24

- Suspicious Document Chunking
- Noun Phrase and Keyword Phrase Extraction
- 
- 
-

# Approach in Use: Five Steps

24

- Suspicious Document Chunking
- Noun Phrase and Keyword Phrase Extraction
- Query Formulation
- 
-

# Approach in Use: Five Steps

24

- Suspicious Document Chunking
- Noun Phrase and Keyword Phrase Extraction
- Query Formulation
- Search Control

.

# Approach in Use: Five Steps

24

- Suspicious Document Chunking
- Noun Phrase and Keyword Phrase Extraction
- Query Formulation
- Search Control
- Document Filtering and Downloading

# Suspicious Document Chunking

25

.

.

.

.

# Suspicious Document Chunking

25

- Segmentation of suspicious documents into parts called chunks
  - 
  - 
  -

# Suspicious Document Chunking

25

- Segmentation of suspicious documents into parts called chunks
- No fixed pattern to put one plagiarism fragment per chunk

.

.

# Suspicious Document Chunking

25

- Segmentation of suspicious documents into parts called chunks
- No fixed pattern to put one plagiarism fragment per chunk
- Sufficient length of chunks, In order to comprise:
  1. At least one plagiarism fragment per chunk,
  2. And Maximum numbers of extracted queries from the chunks.

.



# Suspicious Document Chunking

25

- Segmentation of suspicious documents into parts called chunks
- No fixed pattern to put one plagiarism fragment per chunk
- Sufficient length of chunks, In order to comprise:
  1. At least one plagiarism fragment per chunk,
  2. And Maximum numbers of extracted queries from the chunks.
- Individual sentences sets of 500 words Chunks as results.

# Noun phrase and keyword phrase Extraction

- 
- 
-

# Noun phrase and keyword phrase Extraction

26

Operation number	Operation Description
1	Selection of top 80% long sentences (based on length in chars)
2	Selection of top 80% sentences (based on number of nouns)
3	Selection of top three sentences (based on average tf.idf1 values)
4	Selection of top three sentences (based on number of words with highest values)

- 
- 
-

# Noun phrase and keyword phrase Extraction

26

Operation number	Operation Description
1	Selection of top 80% long sentences (based on length in chars)
2	Selection of top 80% sentences (based on number of nouns)
3	Selection of top three sentences (based on average tf.idf1 values)
4	Selection of top three sentences (based on number of words with highest values)

➤ Scenario1: Operation 1 → Operation 2 → Operation 3 for noun phrase extraction

·  
·

# Noun phrase and keyword phrase Extraction

26

Operation number	Operation Description
1	Selection of top 80% long sentences (based on length in chars)
2	Selection of top 80% sentences (based on number of nouns)
3	Selection of top three sentences (based on average tf.idf1 values)
4	Selection of top three sentences (based on number of words with highest values)

- Scenario1: Operation 1 → Operation 2 → Operation 3 for noun phrase extraction
- Scenario2: Operation 1 → Operation 2 → Operation 4 for keyword phrase extraction

# Noun phrase and keyword phrase Extraction

26

Operation number	Operation Description
1	Selection of top 80% long sentences (based on length in chars)
2	Selection of top 80% sentences (based on number of nouns)
3	Selection of top three sentences (based on average tf.idf1 values)
4	Selection of top three sentences (based on number of words with highest values)

- Scenario1: Operation 1 → Operation 2 → Operation 3 for noun phrase extraction
- Scenario2: Operation 1 → Operation 2 → Operation 4 for keyword phrase extraction
- Three sentences from each scenario1 and scenario2 selected to query formulation

# Query Formulation

27

.

.

.

.

# Query Formulation

27

➤ From each selected sentence, one query is extracted.

.

.

.



# Query Formulation

27

- From each selected sentence, one query is extracted.
- The threshold for the number of words in each query is limited to ten.

.

.

# Query Formulation

27

- From each selected sentence, one query is extracted.
- The threshold for the number of words in each query is limited to ten.
- Selection of high weighted terms to reach the ChatNoir limitation.

.

# Query Formulation

27

- From each selected sentence, one query is extracted.
- The threshold for the number of words in each query is limited to ten.
- Selection of high weighted terms to reach the ChatNoir limitation.
- The terms are placed next to each other based on the order in sentence.

# Download Filtering and Search Control

- 

- 

- 

- 

- 

- 

- 

- 

- 

-

# Download Filtering and Search Control

28

## ➤ Download Filtering

- 

- 

- 

- 

- 

- 

- 

- 

-

# Download Filtering and Search Control

28

- **Download Filtering**

- 14 top results are selected for each query

- 

- 

- 

- 

- 

- 

- 

-

# Download Filtering and Search Control

28

## ➤ **Download Filtering**

- 14 top results are selected for each query
- The query is divided into two sub-queries:
  - Snippet with the length of 500 characters are extracted as a sub-query.
  - Snippets are combined with each other and make a passage.

# Download Filtering and Search Control

28

## ➤ **Download Filtering**

- 14 top results are selected for each query
- The query is divided into two sub-queries:
  - Snippet with the length of 500 characters are extracted as a sub-query.
  - Snippets are combined with each other and make a passage.
- If the resulted passage contains at least 50% words of the query
  - The related document is downloaded
  - The document is maintained for search control operation



# Download Filtering and Search Control

28

## ➤ **Download Filtering**

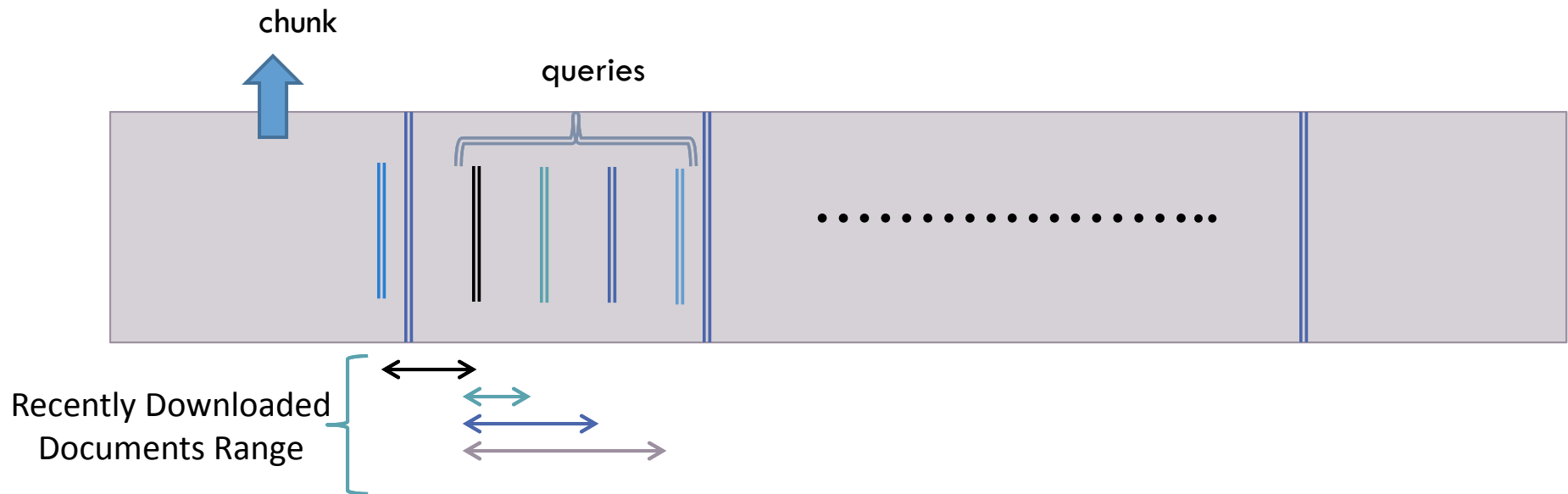
- 14 top results are selected for each query
- The query is divided into two sub-queries:
  - Snippet with the length of 500 characters are extracted as a sub-query.
  - Snippets are combined with each other and make a passage.
- If the resulted passage contains at least 50% words of the query
  - The related document is downloaded
  - The document is maintained for search control operation

## ➤ **Search Control**

- Drop a query when at least 60% of its terms are contained in recently downloaded documents set

# Search Control

- Drop a query when at least 60% of its terms are contained in recently downloaded documents set



# Evaluation

- 
-

# Evaluation

30

Downloads	F1	No Detection	Precision	Queries	Recall	Runtime
183.3	0.115	1	0.07539	43.5	0.41381	8:32:37

•

•

# Evaluation

30

Downloads	F1	No Detection	Precision	Queries	Recall	Runtime
183.3	0.115	1	0.07539	43.5	0.41381	8:32:37

- Highest rank in “No Detection” measure.
- Highest rank in “Runtime” measure.



**Thank you for Your Attention**