

# Cross-Lingual Semantic Similarity Measure for Comparable Articles

{ Motaz.Saad,David.Langlois,Kamel.Smaili }@loria.fr

## 1. PROBLEM

### Retrieve comparable documents

What are comparable documents?

- Cross-lingual documents
- Aligned at document level.
- Related to the same topic but they are not necessarily translations of each other.
- Have same content in terms of topic, but not in terms of syntax and sentences.
- Our work on Arabic-English comparable docs.  
⇒ We use a cross-lingual similarity measure based on latent semantics

## 3. FORMULATION

- Use same approach as Littman et al.<sup>a</sup>.
- Monolingual term-document matrix  $X_a$

$$X_a = \begin{matrix} & d_1 & d_2 & \dots & d_n \\ t_1 & w_{11} & w_{12} & \dots & w_{1n} \\ t_2 & w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ t_m & w_{m1} & w_{m2} & \dots & w_{mn} \end{matrix} \quad (1)$$

- Cross-lingual term-document matrix  $X_{cl}$  (parallel or comparable corpus):

$$X_{cl} = \begin{matrix} & d_1^u & d_2^u & \dots & d_n^u \\ t_1^a & w_{11}^a & w_{12}^a & \dots & w_{1n}^a \\ t_2^a & w_{21}^a & w_{22}^a & \dots & w_{2n}^a \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ t_l^a & w_{l1}^a & w_{l2}^a & \dots & w_{ln}^a \\ t_1^e & w_{11}^e & w_{12}^e & \dots & w_{1n}^e \\ t_2^e & w_{21}^e & w_{22}^e & \dots & w_{2n}^e \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ t_m^e & w_{m1}^e & w_{m2}^e & \dots & w_{mn}^e \end{matrix} \quad (2)$$

- Each  $d_i^u$  in  $X_{cl}$  is the concatenation of the Arabic doc  $d_i^a$  and its corresponding English doc  $d_i^e$ .
- Weights  $w_{ij}$  in  $X_a$  &  $X_{cl}$  are the *tfidf*
- $X_{cl}$  enables LSI to learn the relationship between terms which are semantically related in the same language or between two languages.
- Apply SVD on  $X_a$  &  $X_{cl}$ :  $USV^T$

<sup>a</sup>Littman et al., "Automatic Cross-Language Information Retrieval Using Latent Semantic Indexing", 1998

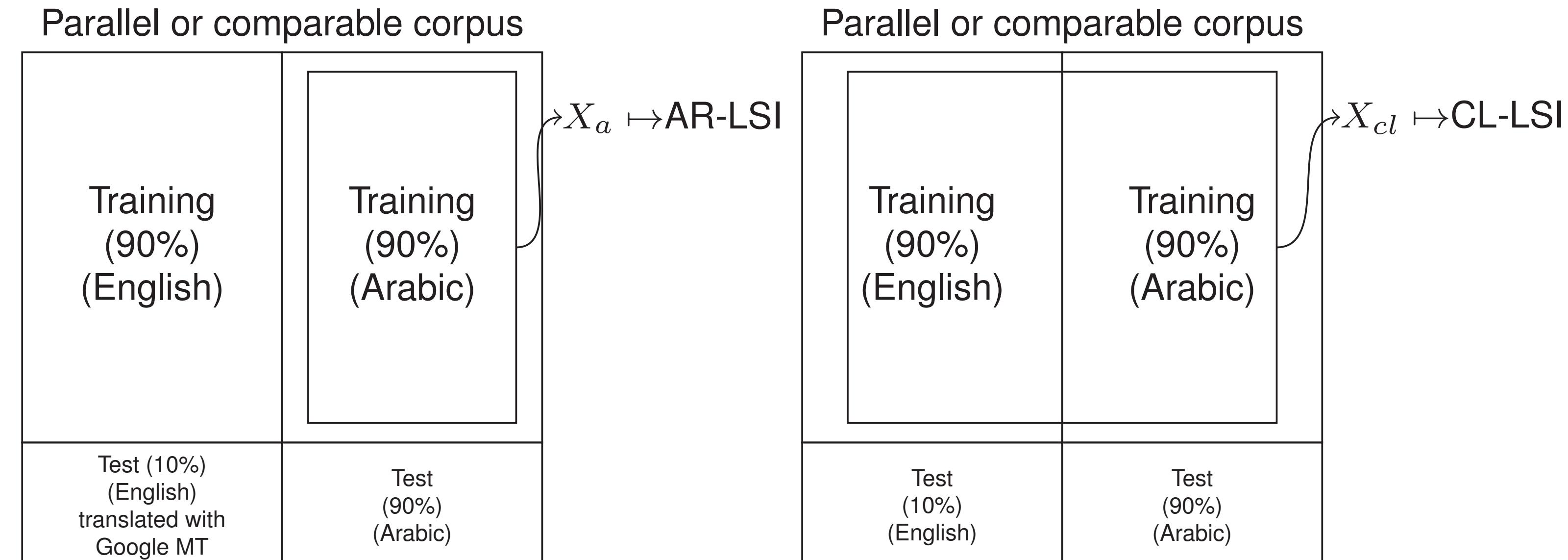
## 5. RESULTS ON PARALLEL CORPUS

Retrieving parallel document mates using AR-LSI & CL-LSI

Corpus	Method	R@1	R@5	R@10
<b>Newspapers</b>				
AFP	AR-LSI	0.94	0.96	0.99
	CL-LSI	<b>0.97</b>	<b>0.99</b>	<b>0.99</b>
ANN	AR-LSI	0.80	0.91	0.94
	CL-LSI	<b>0.82</b>	<b>0.92</b>	<b>0.94</b>
ASB	AR-LSI	0.79	0.90	0.92
	CL-LSI	<b>0.85</b>	<b>0.92</b>	<b>0.97</b>
Medar	AR-LSI	0.56	0.76	0.81
	CL-LSI	<b>0.61</b>	<b>0.78</b>	<b>0.85</b>
NIST	AR-LSI	0.78	0.87	0.92
	CL-LSI	0.71	0.82	0.84
<b>United Nations Resolutions</b>				
UN	AR-LSI	0.97	1.00	1.00
	CL-LSI	<b>0.98</b>	0.99	<b>1.00</b>
<b>Talks</b>				
TED	AR-LSI	0.52	0.73	0.82
	CL-LSI	<b>0.60</b>	<b>0.83</b>	<b>0.92</b>
<b>Movie Subtitles</b>				
OST	AR-LSI	0.39	0.61	0.72
	CL-LSI	0.33	0.76	0.85

## 2. CROSS-LINGUAL SIMILARITY MEASURE

- The cross-lingual similarity measure is based on Latent Semantic Indexing (LSI).
- Two ways:



## 3. PROCEDURES

### Algorithm 1: Retrieving Arabic documents using AR-LSI

```

Input:  $C_e$ : English corpus,  $C_a$ : Arabic corpus,  $n$ : number of docs to retrieve
1  $C'_e \leftarrow \emptyset; C'_a \leftarrow \emptyset;$ 
2 foreach doc  $a_j$  in  $C_a$  do
3    $a'_j \leftarrow a_j^t US^{-1}$ ; put  $a'_j$  in  $C'_a$ ;
4 foreach doc  $e_i$  in  $C_e$  do
5    $a'_{ei} \leftarrow \text{translate}(e_i); a'_{ei} \leftarrow a_{ei}^t US^{-1};$ 
6    $R \leftarrow \text{retrieve}(a'_{ei}, C'_a, n)$  // retrieve top-n similar docs to  $e'_i$  from  $C'_a$ 
7   evaluate( $R$ ) // check if  $a'_i$  in  $R$ 

```

### Algorithm 2: Retrieving Arabic documents using CL-LSI

```

Input:  $C_e$ : English corpus,  $C_a$ : Arabic corpus,  $n$ : number of docs to retrieve
1  $C'_e \leftarrow \emptyset; C'_a \leftarrow \emptyset;$ 
2 foreach doc  $a_j$  in  $C_a$  do
3    $a'_j \leftarrow a_j^t US^{-1}$ ; put  $a'_j$  in  $C'_a$ ;
4 foreach doc  $e_i$  in  $C_e$  do
5    $e'_{ei} \leftarrow e_i^t US^{-1};$ 
6    $R \leftarrow \text{retrieve}(e'_{ei}, C'_a, n)$  // retrieve top-n similar docs to  $e'_i$  from  $C'_a$ 
7   evaluate( $R$ ) // check if  $a'_i$  in  $R$ 

```

### Procedure retrieve( $d_{si}$ , $C_t$ , $n$ )

```

Input:  $d_{si}$ : source doc,  $C_t$ : target corpus,  $n$ : number of docs to retrieve
1  $R \leftarrow \emptyset$ ; // a list of retrieved docs
2 foreach doc  $d_{tj}$  in  $C_t$  do
3    $sim \leftarrow \cos(d_{si}, d_{tj})$ ; put  $(j, sim)$  in  $R$ ;
4   Sort  $R$  in descending order according to  $sim$  values;
5 return top  $n$  elements of  $R$ ;

```

## 6. RESULTS ON COMPARABLE CORPUS

Retrieving comparable documents using CL-LSI

Corpus	R@1	R@5	R@10
Wikipedia	0.42	0.84	0.94
Euro-news	0.84	0.99	1.0

Comparing corpora using CL-LSI

Corpus	parallel	comparable
avg(cos)	0.53	0.46

## 7. CONCLUSION

- CL-LSI is competitive to monolingual LSI
- The advantage of CL-LSI is that it does not need machine translation !
- CL-LSI is language independent approach
- Future work: use CL-LSI to align English-Arabic comparable documents.