



# Boosting Cross-lingual Knowledge Linking via Concept Annotation

Zhichun Wang<sup>1</sup>, Juanzi Li<sup>2</sup>, Jie Tang<sup>2</sup>

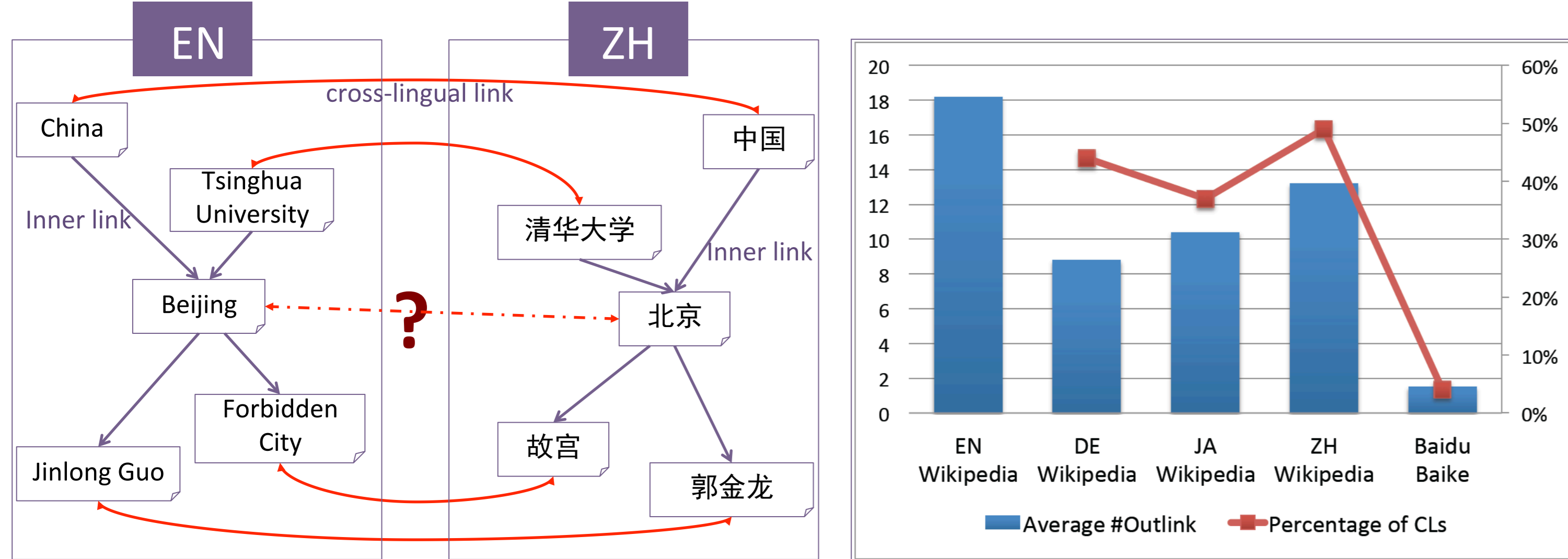
<sup>1</sup>Beijing Normal University, <sup>2</sup>Tsinghua University



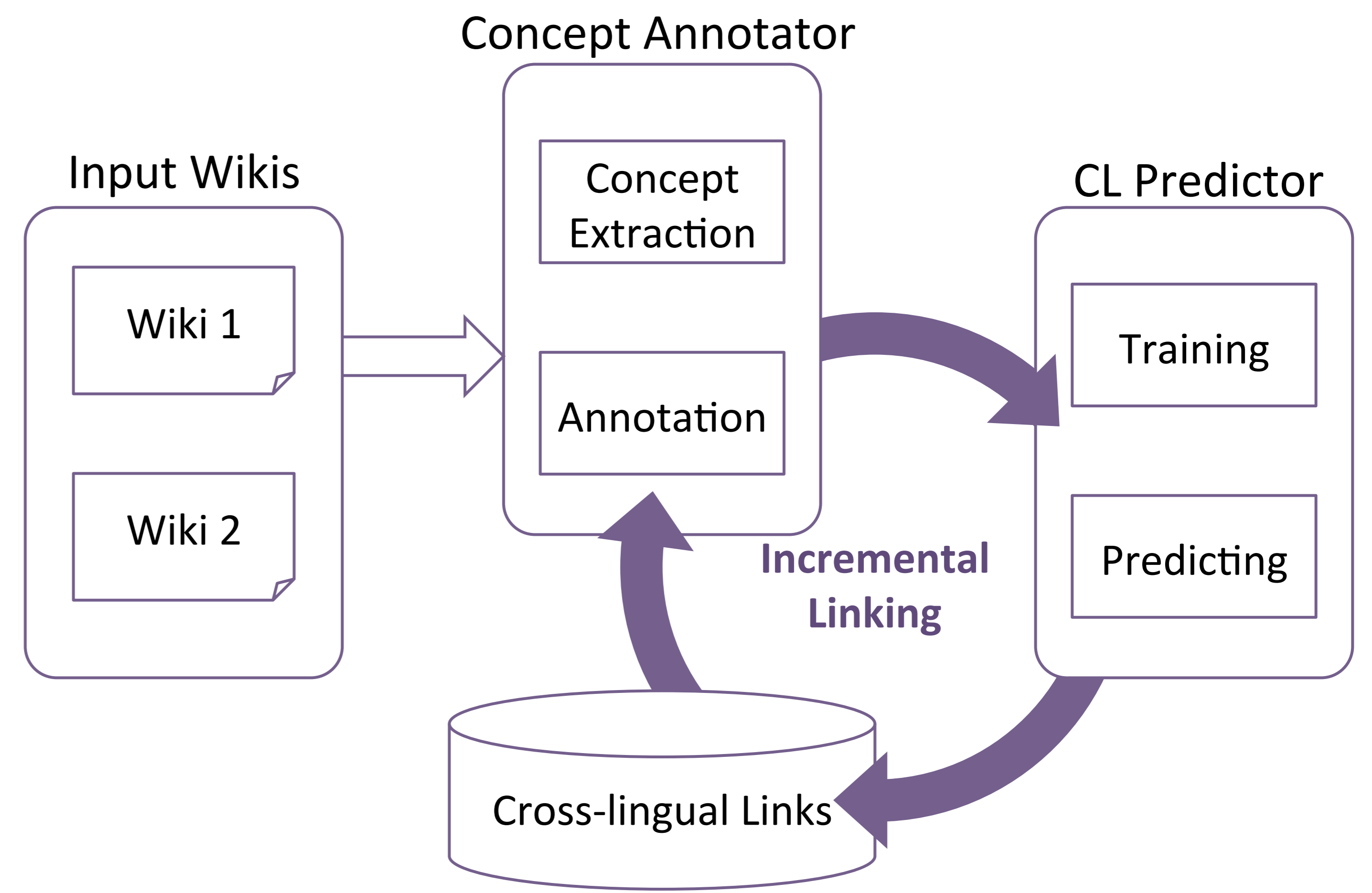
Poster my Pocket app

## Motivation

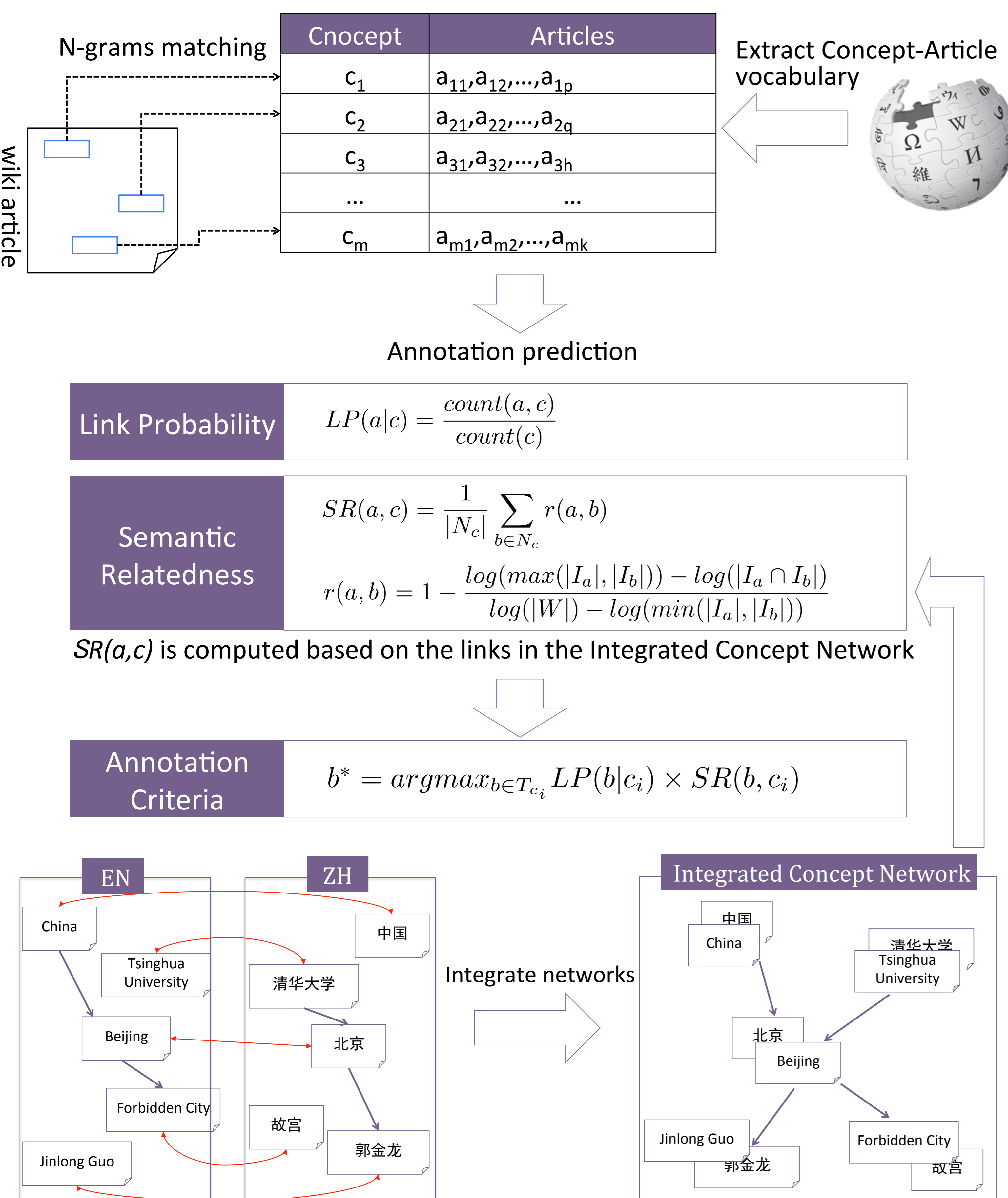
**Cross-lingual knowledge linking** is to automatically discover cross-lingual links (CLs) between wikis, which can largely enrich the cross-lingual knowledge and facilitate sharing knowledge across different languages. The seed CLs and the inner link structures are two important factors for finding new CLs. **A challenging problem is how to find large scale CLs between wikis (e.g. English Wikipedia and Baidu Baike) when there are insufficient seed CLs and inner links?**



## Proposed Framework



## Annotating Concepts Within Each Wiki



## Predicting New Cross-lingual Links Between Wikis

### Language-independent Features

Outlink similarity	$f_1(a,b) = \frac{2 \cdot  \phi_{1 \rightarrow 2}(O(a)) \cap O(b) }{ \phi_{1 \rightarrow 2}(O(a))  +  O(b) }$
	$f_2(a,b) = \frac{2 \cdot  \phi_{1 \rightarrow 2}(O^+(a)) \cap O^+(b) }{ \phi_{1 \rightarrow 2}(O^+(a))  +  O^+(b) }$
Inlink similarity	$f_3(a,b) = \frac{2 \cdot  \phi_{1 \rightarrow 2}(I(a)) \cap I(b) }{ \phi_{1 \rightarrow 2}(I(a))  +  I(b) }$
	$f_4(a,b) = \frac{2 \cdot  \phi_{1 \rightarrow 2}(I^+(a)) \cap I^+(b) }{ \phi_{1 \rightarrow 2}(I^+(a))  +  I^+(b) }$
Category similarity	$f_5(a,b) = \frac{2 \cdot  \phi_{1 \rightarrow 2}(C(a)) \cap C(b) }{ \phi_{1 \rightarrow 2}(C(a))  +  C(b) }$
	$f_6(a,b) = \frac{1}{nm} \sum_{i=1}^n \sum_{j=1}^m \lambda(c_i^a, c_j^b)$

### Regression Model for Predicting New CLs

For each article  $a$  in  $W_1$ , the article  $b^*$  in  $W_2$  that maximizes the following score function  $s(a, b^*)$  and satisfies  $s(a, b^*) > 0$  is predicted as the corresponding article of  $a$ .

$$s(a, b) = \omega_0 + \vec{\omega} \cdot \vec{f}_{a,b} = \omega_0 + \omega_1 \times f_1(a, b) + \dots + \omega_6 \times f_6(a, b)$$

A regression model is used to learn the weights of features based on a set of known CLs. The optimal weights should satisfies:

$$\forall a_i \in A, \forall b' \in (B - \{b_i\}), s(a_i, b_i) - s(a_i, b') > 0$$

$$\vec{\omega} \cdot (\vec{f}_{a_i, b_i} - \vec{f}_{a_i, b'}) > 0$$

Therefore, a training data set is generated to feed the linear regression algorithm to learn the weights:

$$D = \{(x_i, y_i)\}_{i=1}^n \text{ where } x_i = (\vec{f}_{a_i, b_i} - \vec{f}_{a_i, b_{j \neq i}}) \quad y_i = 1$$

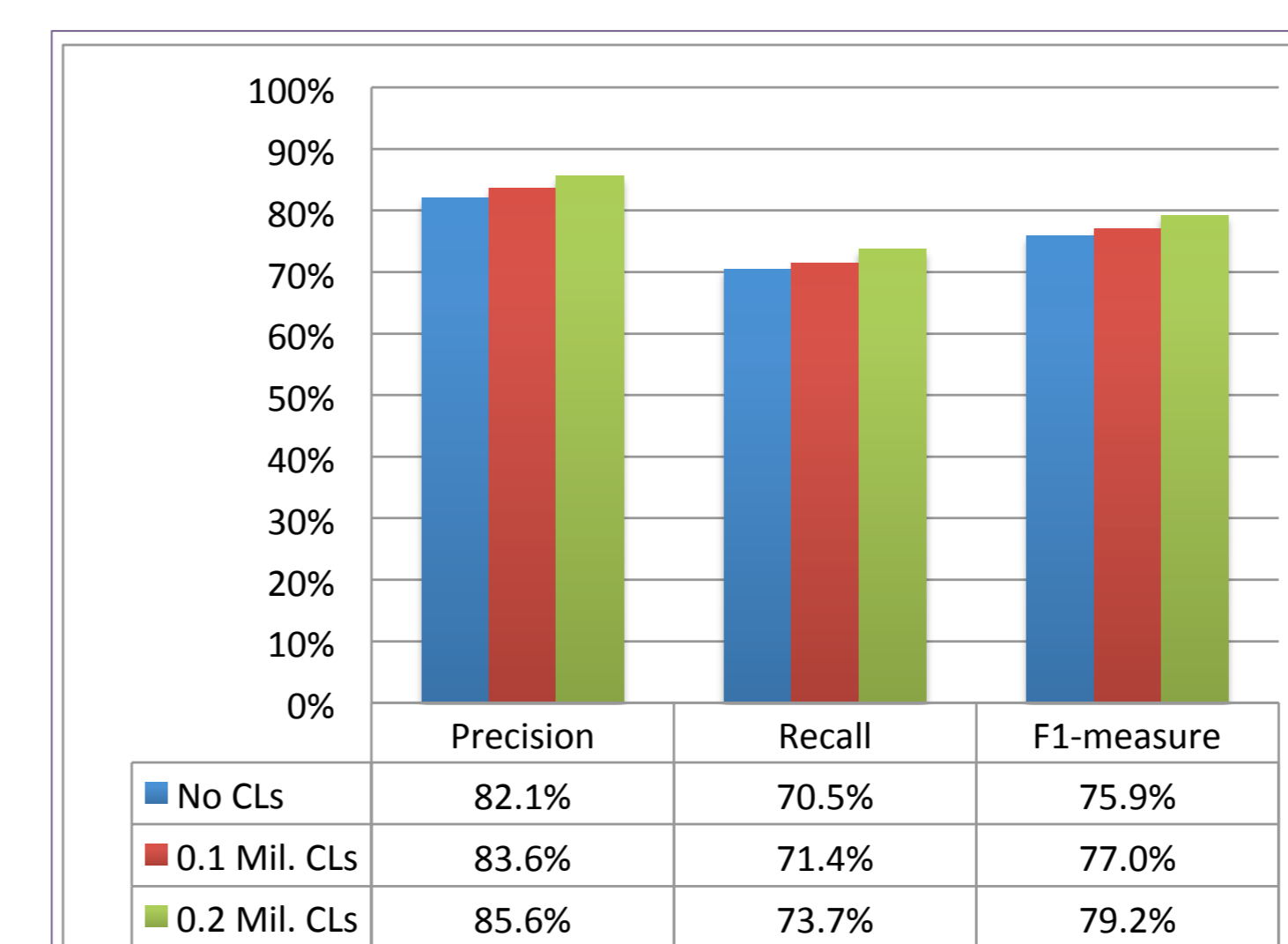
## Experiments

The datasets of English Wikipedia (4 million articles) and Chinese Wikipedia (499 thousand articles) that are archived in August 2012 has been used to evaluate the proposed approach. There are 239,309 cross-lingual links between two wikis.

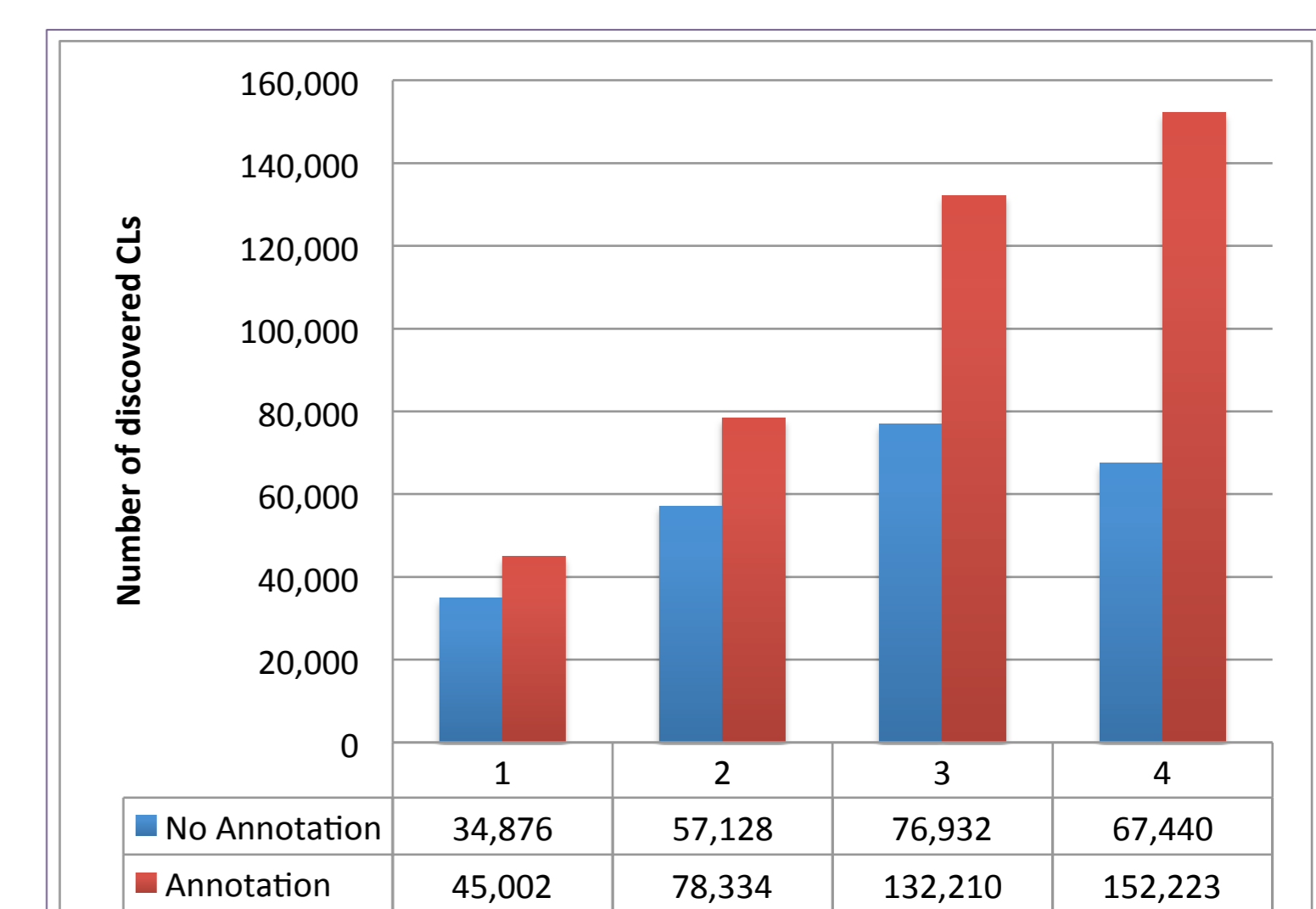
Results of cross-lingual link prediction

#Seed CLs	Model	Before Annotation			After Annotation		
		Precision	Recall	F1-measure	Precision	Recall	F1-measure
0.05 Mil. CLs	SVM	92.1	35.0	50.7	78.5	37.2	50.5
	RM	93.3	36.0	<b>52.0</b>	92.4	38.6	<b>54.5</b>
0.10 Mil. CLs	SVM	79.7	35.0	48.6	86.9	50.4	63.8
	RM	84.6	37.4	<b>51.9</b>	96.6	49.3	<b>65.3</b>
0.15 Mil. CLs	SVM	80.9	35.9	49.7	88.1	57.3	69.5
	RM	93.5	38.2	<b>54.2</b>	93.7	56.2	<b>70.2</b>
0.20 Mil. CLs	SVM	84.7	37.3	51.8	88.8	68.1	77.1
	RM	94.5	37.9	<b>54.1</b>	95.9	67.2	<b>79.0</b>

Our regression model (RM) is compared with SVM classification model.



Results of concept annotation



Results of incremental linking