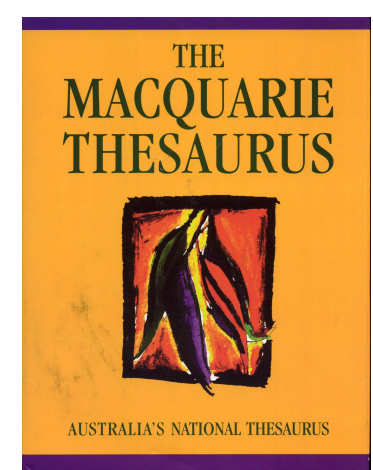# TOR, TORMD: Distributional Profiles of Concepts for Unsupervised Word Sense Disambiguation

Saif Mohammad, Graeme Hirst, and Philip Resnik

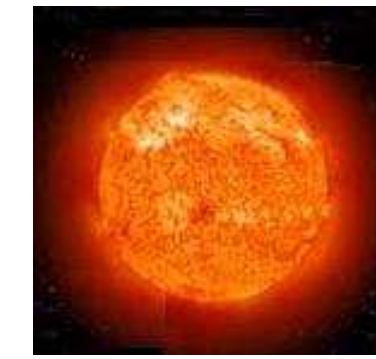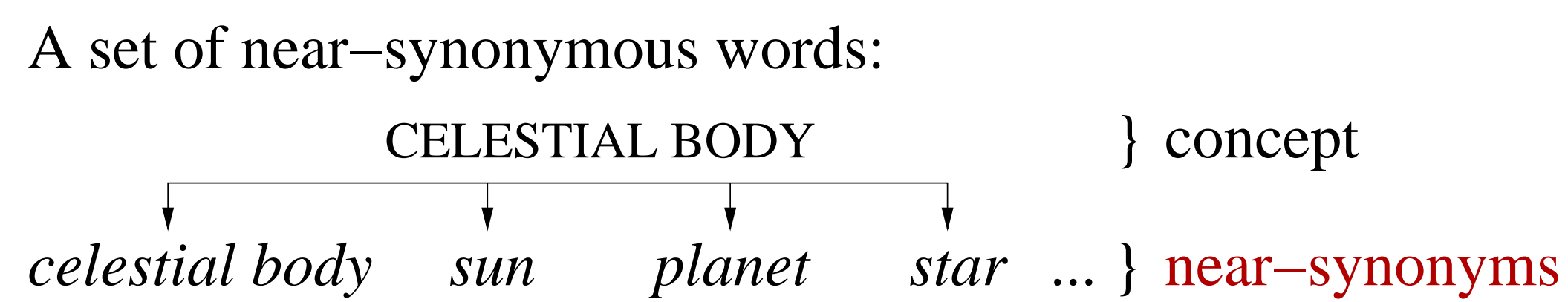University of Toronto and University of Maryland

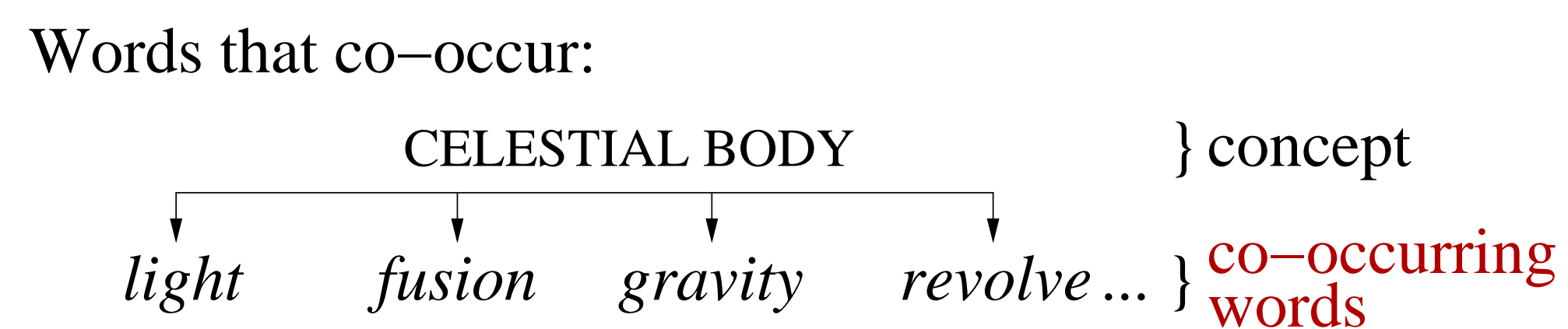{smm,gh}@cs.toronto.edu, resnik@umiacs.umd.edu

## How to represent a concept?

**As a category from a thesaurus**

A set of near–synonymous words:

CELESTIAL BODY } concept

*celestial body   sun   planet   star* ... } near–synonyms

CELESTIAL BODY

**By its usage in text**

Words that co–occur:

CELESTIAL BODY } concept

*light   fusion   gravity   revolve* ... } co–occurring words

↓ **Combining the two** ↓

## Central Idea: DISTRIBUTIONAL PROFILES OF CONCEPTS

**Capturing co-occurrence associations between words and concepts**

CELESTIAL BODY   CELEBRITY } concepts

strong association   weak association

· *space* · · *star* · · · } text

**Distributional Profile of a concept (DPC)**

CELESTIAL BODY: *space 0.36, light 0.27, revolve 0.14,...*
(planet, sun, star,...)

concept   near–synonyms   co–occurring words with strength of association

**How to create these DPCs?**

↓ ↓

*You know someone by the company they keep.*

**Create Word–Category Co-occurrence Matrix (WCCM)**

categories/concepts →

$$\begin{array}{c|ccccc} & c_1 & c_2 & \ldots & c_j & \ldots \\ \hline w_1 & m_{11} & m_{12} & \ldots & m_{1j} & \ldots \\ w_2 & m_{21} & m_{22} & \ldots & m_{2j} & \ldots \\ \vdots & \vdots & \vdots & \ddots & \vdots & \\ w_i & m_{i1} & m_{i2} & \ldots & m_{ij} & \ldots \\ \vdots & \vdots & \vdots & & \ddots & \end{array}$$

words ↓

**Unsupervised naïve Bayes word sense classifier**

desired concept $= \underset{c_j \in C}{\arg\max} \; P(c_j) \prod_{w_i \in W} P(w_i | c_j)$

The WCCM can be used to estimate probabilities (in an unsupervised manner), that are traditionally calculated using sense–annotated data.

$\dfrac{\sum_i m_{ij}}{\sum_{i,j} m_{ij}}$   $\dfrac{m_{ij}}{\sum_i m_{ij}}$

Base WCCM: $m_{ij}$ is the number of times word $w_i$ co-occurs with any word *that has $c_j$ as a sense*.

Bootstrapped WCCM: $m_{ij}$ is the number of times word $w_i$ co-occurs with any word *used in sense $c_j$*.

↓ **Apply this monolingually** ↓

## English Lexical Sample Task

Accuracies were markedly better than the random baseline—an increase of more than twenty percentage points.

**Results**



### Conclusions

- Placed first among unsupervised systems in the Chinese–English Task.
- Only about 1 percentage point behind the best in the English Lexical Task.
- Cross-lingual DPCs can help automatic machine translation.
- DPCs create simple yet powerful baselines for WSD.

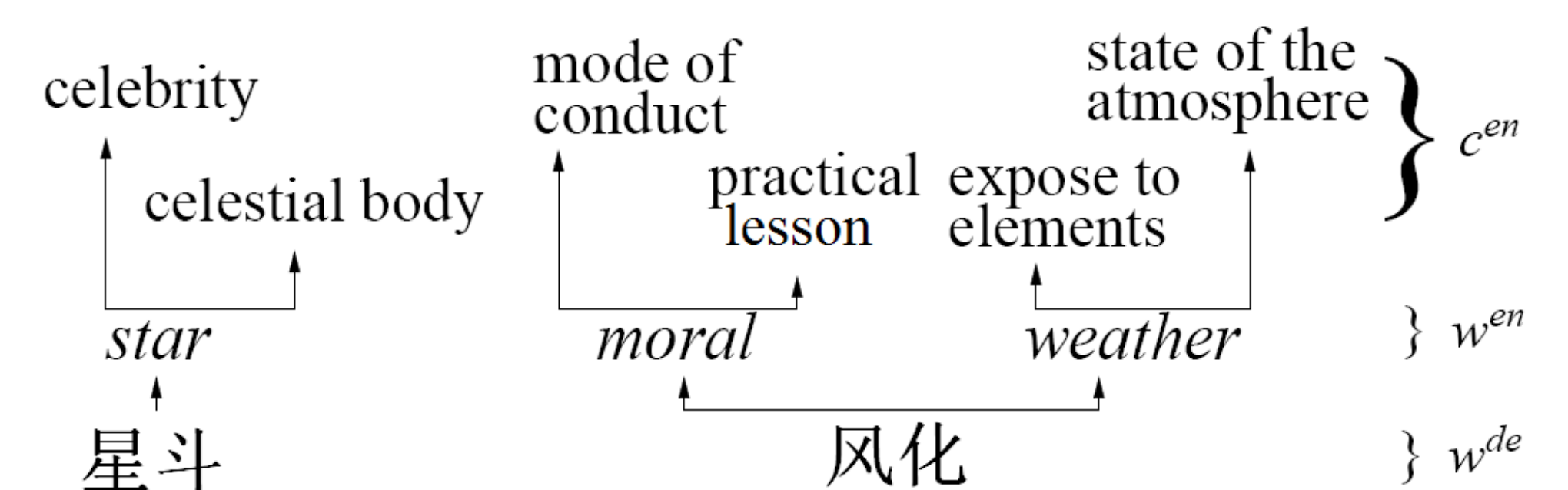↑ *Apply this cross-lingually* ↑

## Multilingual Chinese–English Lexical Sample Task

**Words having 'celestial body' as cross-lingual candidate sense**

celestial body } $c^{en}$

*celestial body   sun   star* ... } $w^{en}$

天体   日   太阳   星 星斗 } $w^{ch}$

**Cross-lingual candidate senses of Chinese words 星斗 and 风化**

celebrity   mode of conduct   state of the atmosphere

celestial body   practical lesson   expose to elements } $c^{en}$

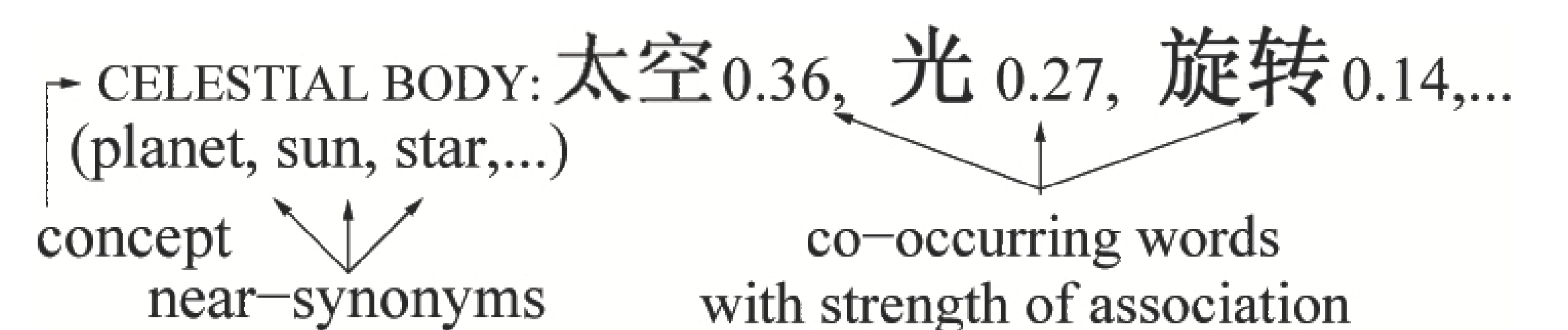*star*   *moral*   *weather* } $w^{en}$

星斗   风化 } $w^{de}$

**So effectively we have Chinese words with English senses** ⇓⇓

**Capturing co-occurrence associations between Chinese words and English concepts**

CELESTIAL BODY   CELEBRITY } concepts

strong association   weak association

· 太空 · · 星 · · · } text

**Cross-lingual Distributional Profiles of Concepts**

CELESTIAL BODY: 太空 0.36, 光 0.27, 旋转 0.14,...
(planet, sun, star,...)

concept   near–synonyms   co–occurring words with strength of association

**Chinese word–English category co-occurrence matrix**

$$\begin{array}{c|ccccc} & c_1^{en} & c_2^{en} & \ldots & c_j^{en} & \ldots \\ \hline w_1^{ch} & m_{11} & m_{12} & \ldots & m_{1j} & \ldots \\ w_2^{ch} & m_{21} & m_{22} & \ldots & m_{2j} & \ldots \\ \vdots & \vdots & \vdots & \ddots & \vdots & \\ w_i^{ch} & m_{i1} & m_{i2} & \ldots & m_{ij} & \ldots \\ \vdots & \vdots & \vdots & & \ddots & \end{array}$$

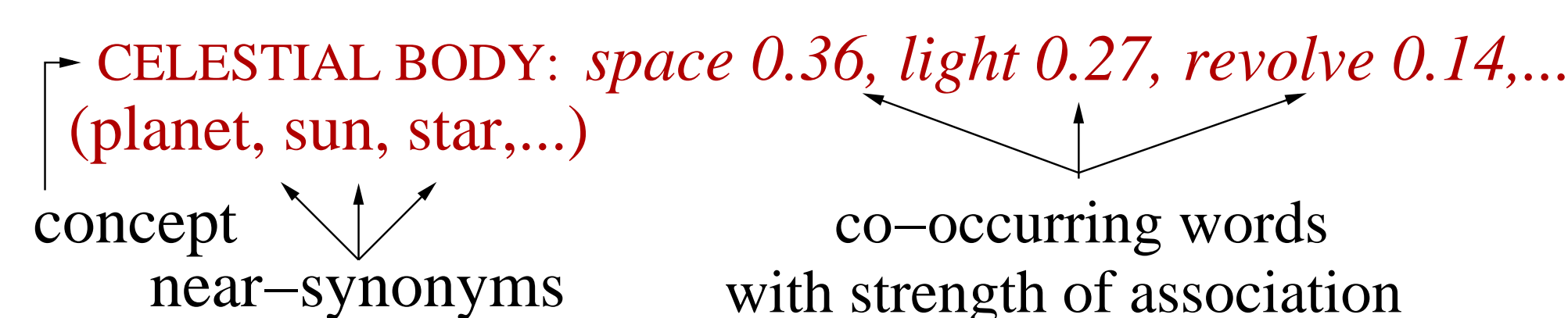Base WCCM: $m_{ij}$ is the number of times Chinese word $w_i^{ch}$ co-occurs with a word *that has $c_j$ as English sense*.

Bootstrapped WCCM: $m_{ij}$ is the number of times Chinese word $w_i^{ch}$ co-occurs with a word *used in English sense $c_j$*.
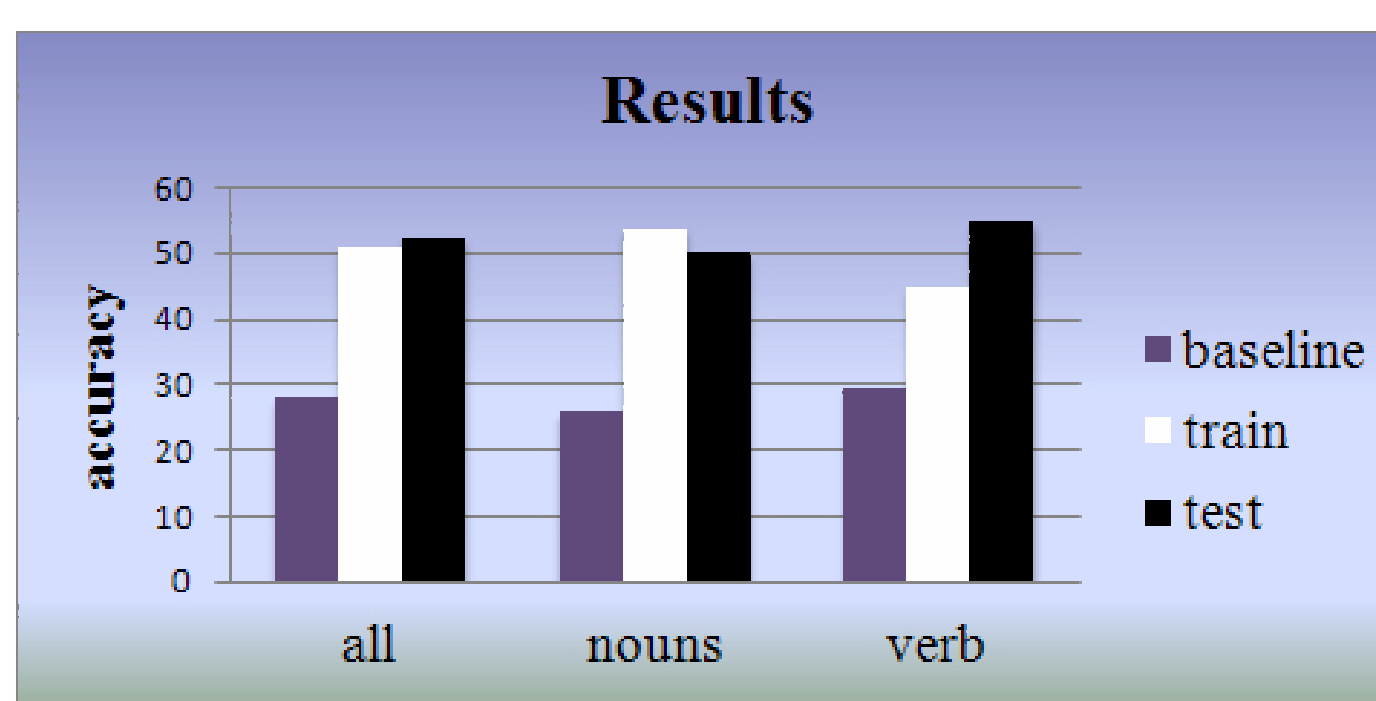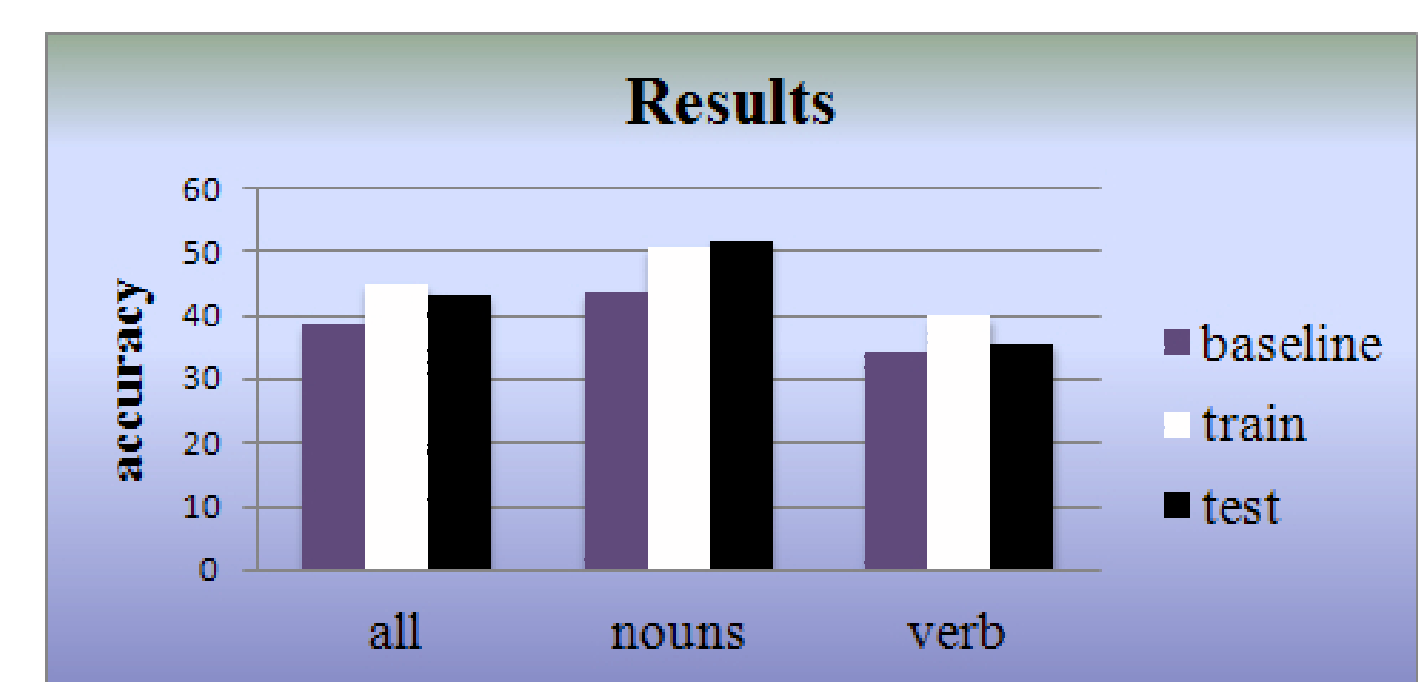
**Results**



See how cross-lingual DPCs can be used to obtain state-of-the-art semantic distance accuracies in a resource-poor language using a knowledge source from a resource-rich one.

**Come to EMNLP's Friday morning session**