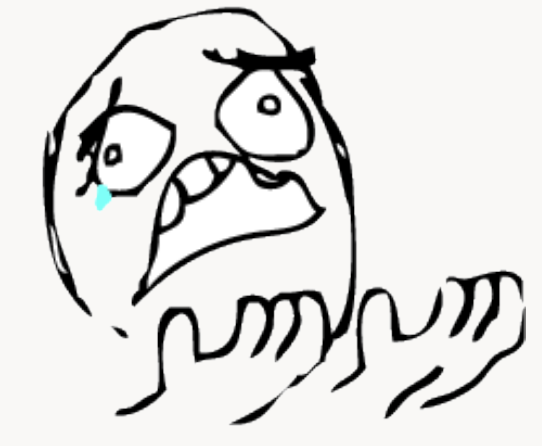


# Context-Dependent Multilingual Lexical Lookup for Under-Resourced Languages

**Input:** Text  $Q$  = sequence of lexical items (LIs)  $\{w_1, w_2, \dots, w_n\}$

**Output:** Ranked list of translations sets for each  $w_i$

## Motivation: Under-Resourced Languages

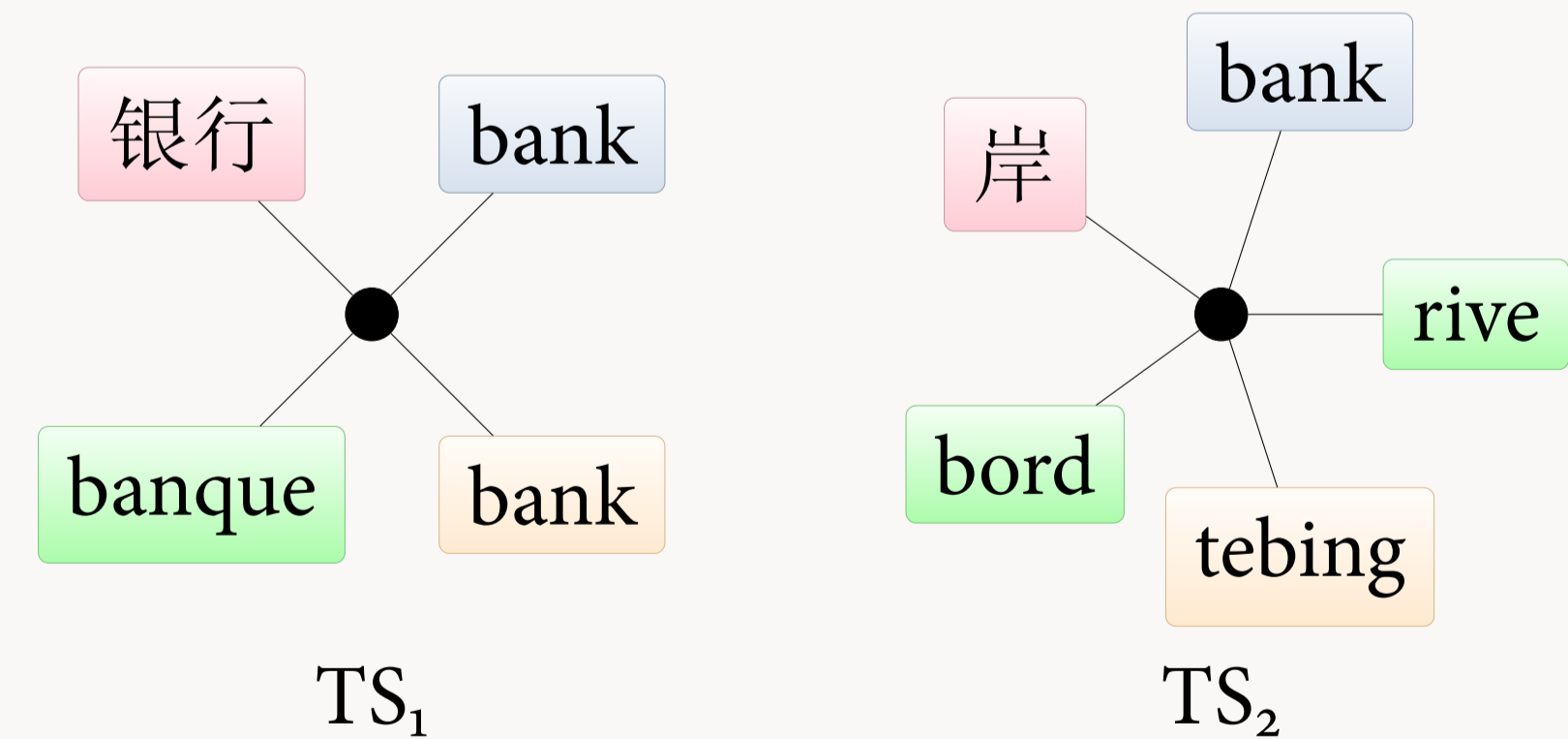
- WSD and translation selection typically require rich lexical knowledge and/or corpus resources
  - Gloss texts, subject codes, semantic networks
  - Corpora: (un)tagged; mono-/bilingual; aligned/comparable
- Whither under-resourced languages? 
  - Lack of rich lexical resources
  - Small corpora

## Comparable Bilingual Corpora

- Idea:** leverage from *richer-resourced language pairs*
- Mine translation context data from 'rich' language pair
- Comparable corpora is easier to obtain
- E.g. concatenated English & Malay Wikipedia articles
- 62 993 documents, 67 499 terms**

## Multilingual Lexcion

- Translation sets**, like multilingual synsets
- Corresponding to coarse-grained concepts
- Member lexical items (LIs) from different languages



- 24 371 English, 13 226 Chinese, 35 640 Malay, 17 063 French, 14 687 Thai, 5629 Iban LIs**
- Uses translation context data from **rich language pair** for *entire translation set*

## Extracting Translation Context Knowledge

- Document = bilingual article pair as bag-of-words
- Preprocess corpus (segment, stop word removal, lemmatise)
- Run latent semantic indexing (LSI) on corpus
  - 1000 factors, 45 minutes (MacBook Pro 2.3 GHz, 4 GB RAM)
  - One term vector for each LI in both languages
- $V(TS) = \sum(\text{all available term vectors of member LIs})$
- $\therefore V(TS_1) = V(\langle\langle\text{bank}\rangle\rangle_{\text{eng}}) + V(\langle\langle\text{bank}\rangle\rangle_{\text{msa}})$
- $V(TS_2) = V(\langle\langle\text{bank}\rangle\rangle_{\text{eng}}) + V(\langle\langle\text{tebing}\rangle\rangle_{\text{msa}})$

## Context-Dependent Lexical Lookup

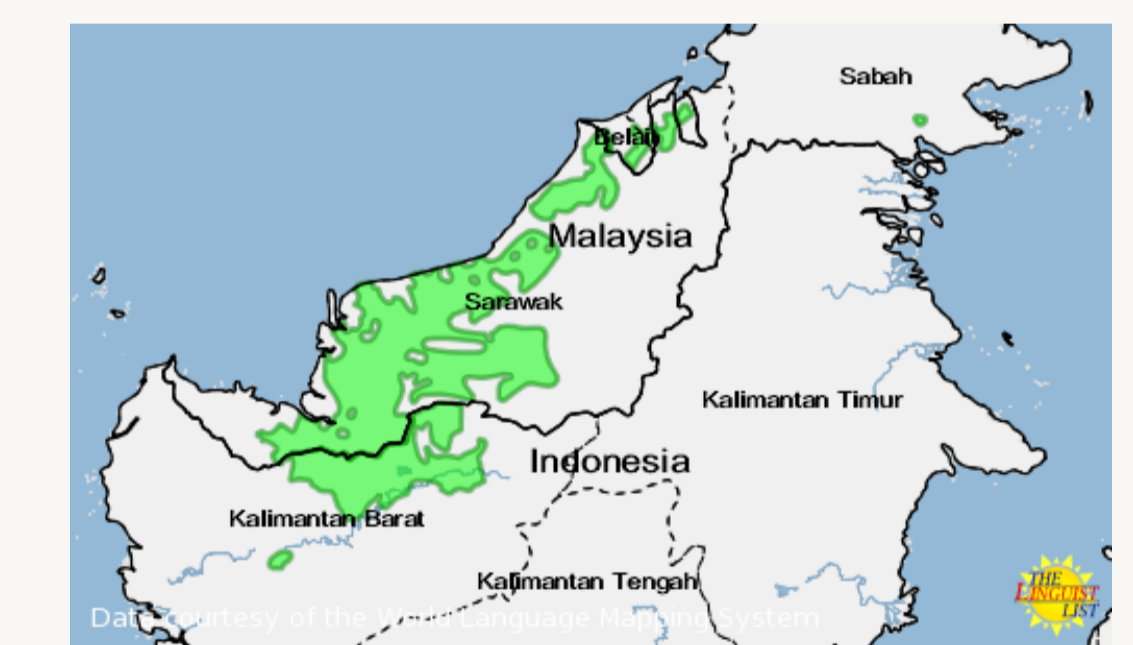
- For input text  $Q$ ,  $V_Q = \sum V(w_i)$
- If  $Q$  language not in training corpus,  $V_Q = \sum V(\text{all TS containing } w_i)$
- For each LI  $w_i$  in  $Q$ ,  $TS_{w_i} = \{t_1, \dots, t_n\}$  set of all TS containing  $w_i$
- $TS_{w_i}$  sorted by  $CSim(V_t, V_Q) = \frac{V_t \cdot V_Q}{|V_t| \times |V_Q|}$

## Example with Iban Input (Top Ranked Translation Sets Shown)

*Iban is a Bornean language with 600 000 speakers.*

$Q$  = 'Lelaki nya tikah enggau emperaja iya, siko dayang ke ligung'

<b>lelaki</b> zho: 男性 tha: ตัวผู้ fra: homme, mâle, masculin msa: lelaki, jantan, eng: male	<b>tikah</b> zho: 结婚 tha: สมรส, ออกเรือน fra: épouser, se marier msa: menikahi, mengahwini eng: marry, wed	<b>emperaja</b> zho: 情人 tha: คู่ควง, ยอดรัก msa: kekasih eng: sweetheart	<b>dayang</b> zho: 女孩子, 姑娘 tha: กัญญา, สาวน้อย msa: pemudi, perawan, dara eng: girl	<b>ligung</b> zho: 可爱 tha: น่าเกลียดน่าชัง, น่ารักน่าชัง fra: joli, mignon msa: comel eng: cute, pretty
--	---	--	---	--



Courtesy of LL-MAP  
llmap.org/languages/iba/static\_map.html

$Q$  = 'Udah ujan nya ngetu terbubuh, matahari enggau emperaja lalu ayan ba langit'

<b>ujan</b> zho: 雨 tha: ฝนตก, พายุ, สายฝน fra: flotte, pluie msa: hujan eng: rain	<b>ngetu</b> zho: 停止 tha: ยั้ง, ขาดช่วง, แวะ fra: arrêter, cesser, finir msa: berhenti eng: cease, quit, stop	<b>matahari</b> zho: 太阳, 日 tha: ดวงตะวัน, พระอาทิตย์, ภากร fra: solaire, soleil msa: matahari eng: solar, sun	<b>emperaja</b> zho: 彩虹 tha: รุ้งกินน้ำ, สายรุ้ง, อินทธรณู fra: arc-en-ciel msa: pelangi eng: rainbow	<b>ayan</b> zho: 看得见 msa: terlihat eng: perceptible, visible	<b>langit</b> zho: 天空 tha: คณางค์, ทิฆัมพร, เวทะ fra: ciel msa: langit eng: sky, welkin
--	--	--	--	---	--

## Some Quick Results

- 80 text sentences:  $\langle\langle\text{bank}\rangle\rangle_{\text{eng}}$ ,  $\langle\langle\text{plant}\rangle\rangle_{\text{eng}}$ ,  $\langle\langle\text{kabinet}\rangle\rangle_{\text{msa}}$ ,  $\langle\langle\text{mangga}\rangle\rangle_{\text{msa}}$ ,  $\langle\langle\text{谷}\rangle\rangle_{\text{zho}}$ ,  $\langle\langle\text{emperaja}\rangle\rangle_{\text{iba}}$
- Strategies
  - wiki-lsi Proposed strategy
  - base-freq Baseline: most frequent translation
  - goog-tr Google Translate
- Metrics: Precision, Mean Reciprocal Rank

Strategy	Incl. Eng. & Iban		W/o Eng. & Iban	
	Precision	MRR	Precision	MRR
wiki-lsi	0.650	0.810	0.690	0.845
base-freq	0.550	0.771	0.524	0.762
goog-tr	0.797	0.812	0.690	0.708

## Conclusions

- Trained on *bilingual* comparable corpus
- But can be used for *multilingual* inputs
- May not be highly accurate, but *fast, cheap* for under-resourced languages