### Name Translation in SMT: Learning When to Transliterate

Authors: U. Hermjakob, K. Knight and H. Daume III

Presenter: Waleed Ammar

### Contributions

1. Method for transliterating words

2. Identification of words to be transliterated

3. Integration in SMT system

4. Metric to evaluate NE translation

### Contributions

1. Method for transliterating words

2. Identification of words to be transliterated

3. Integration in SMT system

4. Metric to evaluate NE translation

### How to Transliterate a word? The Problem

- 1. Phonetic ambiguity in source language
  - سلم transliterates to "silm"
  - or salam, sallam, sallim or sollam?
  - سلم .8٧ سَلَّمَ –

- 2. Proper spelling in target language
  - جون کلارك transliterates to "Jon Clark"
  - ... or John Clark?

### How to Transliterate a word? Search-Oriented Approach

1. Define translit-cost(ar-word, en-word)

- 2. Define score(en-word|ar-word)
  - = log-freq(en-word) / 20 translit-cost(e,a)

- 3. Transliteration
  - = argmax<sub>en-word</sub> score(en-word|ar-word)

## How to Transliterate a word? translit-cost(ar-word, en-word)

- Think Levenshtein edit distance
- 732 rules:
  - Substituting "q" for "ق" costs 0
  - Substituting "k" for "ق" costs 0.2
  - Substituting "k\$" for "\$ق" costs 0
  - Substituting "sh" for "ش" costs 0.1
  - **–** ...
- No rule applies => match fails
- Accumulated cost > threshold => match fails

### How to Transliterate a word? translit-cost caveats

- 700+ rules per LP
- Longest applicable rule
- Context
- Multi-token transliterations
- Style flags

#### How to Transliterate a word?

### argmax<sub>en-word</sub>

- Transliteration
  - = argmax<sub>en-word</sub> score(en-word|ar-word)
- Too many candidate en-words (3.5M)
- Solution:
  - Map English n-grams to consonant skeletons
     Rachmaninoff rkmnnf, rmnnf, rsmnnf, rtsmnnf
  - Build a reverse index(skeleton) → English n-grams
  - Map an Arabic word to a consonant skeleton
  - Use the reverse-index to find candidate English n-grams
  - Compute score() for each candidate

# How to Transliterate a word? Evaluation

### Contributions

1. Method for transliterating words

2. Identification of words to be transliterated

3. Integration in SMT system

4. Metric to evaluate NE translation

وبينما حمل الجيش السوداني عبر الناطق الرسمي باسمه الصوارمي خالد سعد الحركة الشعبية مسؤولية أمن وسلامة الرهائن الصينيين، طالب السفير الصيني بالخرطوم ليو شياو فوانغ الحكومة السودانية ببذل قصارى جهدها لإنقاذ رعيا بلاده.

#### 1. Named Entities

ر عبا بلاده.

Translit: Al-Gaysh Sudanese

Translit: Al-Haraka Al-Shabeyya armi Khaled

Translit: Liu Xiao Fuang

Correct: Liu Xiao Fuang

الصوارمي خالد سعد الحركة الصيني بالخرطوم ليو
الرهائن الصينيين، طالب السفير الصيني بالخرطوم ليو
شياو فوانغ الحكومة السودانية ببذل قصارى جهدها لإنقاذ

- 1. Named Entities
- 2. Out of Vocab

وبينما حمل الجيش السوداني عبر الناطق الرسمي باسمه الصوارمي خالد سعد الحركة الشعبية مسؤولية أمن وسلامة الرهائن الصينيين، طالب السفير الصيني بالخرطوم ليو شياو فوانغ الحكومة السودانية ببذل قصارى جهدها لإنقاذ رعيا بلاده.

- 1. Named Entities
- 2. Out of Vocab

وبينما حمل الجيش السوداني عبر الناطق الرسمي باسمه الصوارمي خالد سعد الحركة الشعبية مسؤولية أمن وسلامة الرهائن الصينيين، طالب السفير الصيني بالخرطوم ليو شياو فوانغ الحكومة السودانية ببذل قصارى جهدها لإنقاذ رعيا بلاده.

- 1. Named Entities
- 2. Out of Vocab
- 3. Transliterate-me tagger

وبيما حمل الجيش السوداني عير الناطق الرسمي باسمه الصوارمي خالد سعد الحركة الشعبية مسؤولية أمن وسلامة الرهائن الصينيين، طالب السفير الصيني بالخرطوم ليو شياو فوانغ الحكومة السودانية ببذل قصارى جهدها لإنقاذ رعيا بلاده.

1. Take a bitext

صِقِلِّية هي أكبر جزيرة في البحر الأبيض المتوسط

Sicily is the largest island in the Mediterranean Sea

- 1. Take a bitext
- 2. Find transliteration pairs (P:99.5%, R:95%)

صِقِلِّية هي أكبر جزيرة في البحر الأبيض المتوسط

Sicily is the largest island in the Mediterranean Sea

- 1. Take a bitext
- 2. Find transliteration pairs
- 3. Tag Arabic words that can be transliterated

صِقِلِّية هي أكبر جزيرة في البحر الأبيض المتوسط

Sicily is the largest island in the Mediterranean Sea

- 1. Take a bitext
- 2. Find transliteration pairs
- 3. Tag Arabic words that can be transliterated
- 4. Get rid of English

صِقِلِّية هي أكبر جزيرة في البحر الأبيض المتوسط

- 1. Take a bitext
- 2. Find transliteration pairs
- 3. Tag Arabic words that can be transliterated
- 4. Get rid of English
- 5. Split tagged Arabic sentences into train/test

- 1. Take a bitext
- 2. Find transliteration pairs
- 3. Tag Arabic words that can be transliterated
- 4. Get rid of English
- 5. Split tagged Arabic sentences into train/test
- 6. Train the "transliterate-me" tagger

صِقِلِّية هي أكبر جزيرة في البحر الأبيض المتوسط

- Surface form, f(surface\_form)
- Previous two words, g(prev\_two\_words)
- Next two words, h(prev\_two\_words)
- Prefixes
- Suffixes

Total = 250 features

#### Which words to transliterate?

#### Transliterate-me Caveats

- Transliteration matching:
  - Multi-token names
  - Many-to-one links
  - Capital English words
  - Stopwords removed
  - Country/nationality name equivalence
  - Split prefixes
- Tagger:
  - "Stat section"
  - Averaged perceptron

#### Which words to transliterate?

#### Transliterate-me Evaluation

Reference	Precision	Recall	F-meas.
Raw test corpus	87.4%	95.7%	91.4%

• Test set = 10K sents

### Contributions

1. Identification of words to be transliterated

2. Method for transliterating words

3. Integration in SMT system

4. Metric to evaluate NE translation

# Integration in End-to-End SMT Approach

- Tag Arabic sentence to be translated using "transliterate-me" tagger
- 2. Transliterate tagged items with unreliable counts in bitext
- 3. Add transliterations to the phrase table with the "transliterated" feature set to 1
- 4. At runtime, transliteration entries in the phrase table compete with original entries

### Contributions

1. Identification of words to be transliterated

2. Method for transliterating words

3. Integration in SMT system

4. Metric to evaluate NE translation

#### How to Know We Did a Good Job?

### NEWA: Named Entity Weak Accuracy

 Traditional metrics make no distinction btw dropping a comma and dropping a NE

• NEWA = # of NEs correctly translated

Total # of NEs

 A named entity is correctly translated if one of its proper translations appears in the translated sentence

### **End-to-End Evaluation**

Metric: NEWA

NE type: all

• NE count: 1730

• Sentence count: 637

Gold Standard	BBN GS
Human 1	87.0%
Human 2	85.3%
Human 3	90.4%
Human 4	86.5%
SMT System	80.4%

### **End-to-End Evaluation**

- NE types: PERson, GeoPoliticalEntity, ORGanization, FACility, PER.Nominal, LOCation
- BLEU: 50.70 → 50.96

NE Type	Count	Baseline	SMT with
		SMT	Transliteration

30

#### Criticism

- Transliteration rules and their costs are handcrafted
- Can't generate new words
- Unsupported claim about applicability to other language pairs
- Re-annotation partially simulate the transliteration algorithm by using transliterator output and ngram frequencies