# Automating Dictionary Production: a Tagalog-English-Korean Dictionary from Scratch

Vít Baisa, Marek Blahuš, Michal Cukr, Ondřej Herman, **Miloš Jakubíček**, Vojtěch Kovář, Marek Medveď, Michal Měchura, Pavel Rychlý, Vít Suchomel

LEXiCAL COMPUTING    MASARYK UNIVERSITY

Brno, Czech Republic

eLex 2019, Sintra, Portugal

Automating dictionary production: How far can we get?

data + tools + people = dictionary

# Tagalog → English + Korean

# Tagalog → English + Korean

from scratch, automatically generated draft, post-editing

# Tagalog → English + Korean

from scratch, automatically generated draft, post-editing
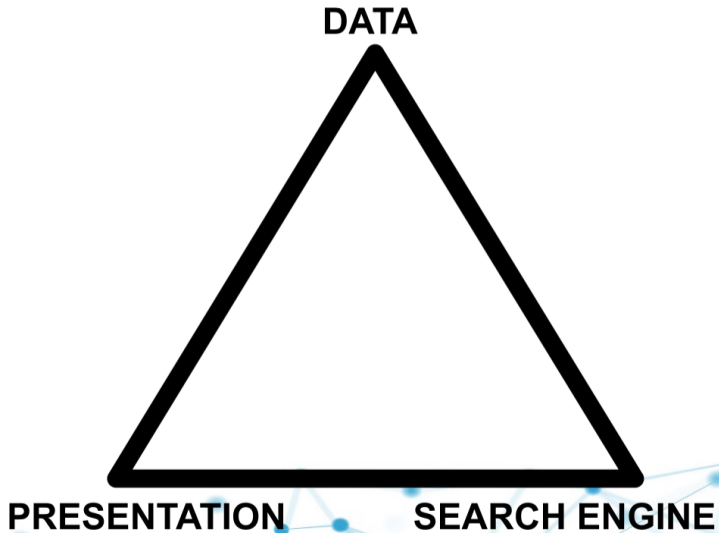
## 15 + 35 = 50

# Tagalog → English + Korean

from scratch, automatically generated draft, post-editing

## 15 + 35 = 50

### < 9

Naver corporation, the biggest search engine in South Korea

freely available at https://dict.naver.com/, soon[TM]

Urdu and Lao in the pipeline

# Entry structure

- headword
- inflected forms
- audio pronunciation
- word sense
  - identified with a disambiguating gloss
  - up to 10 collocations
  - up to 10 synonyms/antonyms
  - a picture, if appropriate
  - 3 example sentences
  - English translation of the gloss and one example
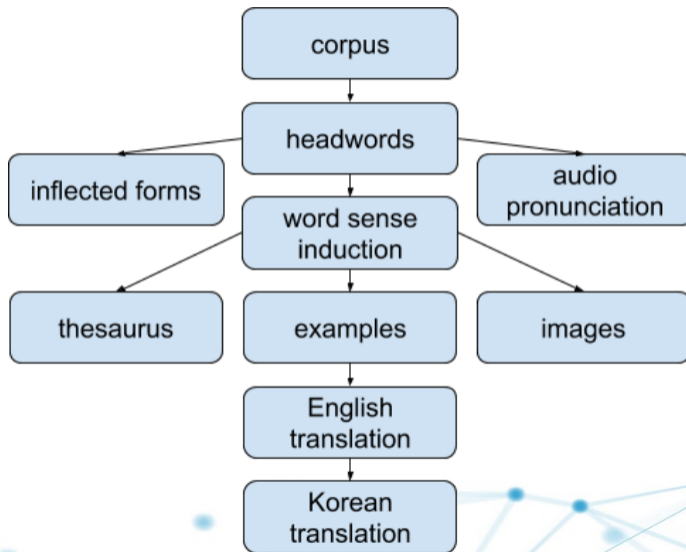  - Korean translation of the gloss and one example

# The recipe

- a big web corpus: 650 - 420 = 230
- PoS tagger
- lemmatizer
- sketch grammar
- Sketch Engine
- Lexonomy + custom editing widgets
- 1−8 developers
- 5−15 editors
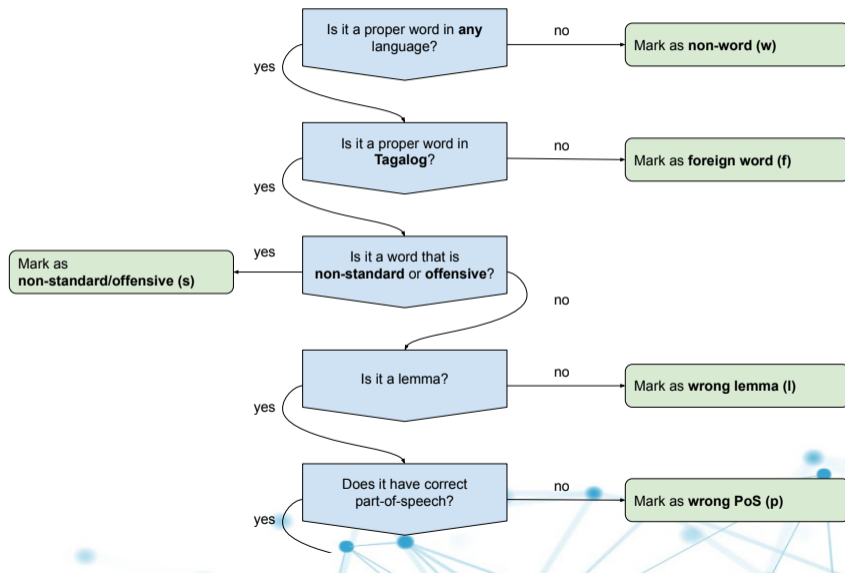
## cook but not boil

small, single, focused tasks

training, multiple annotations, IAA
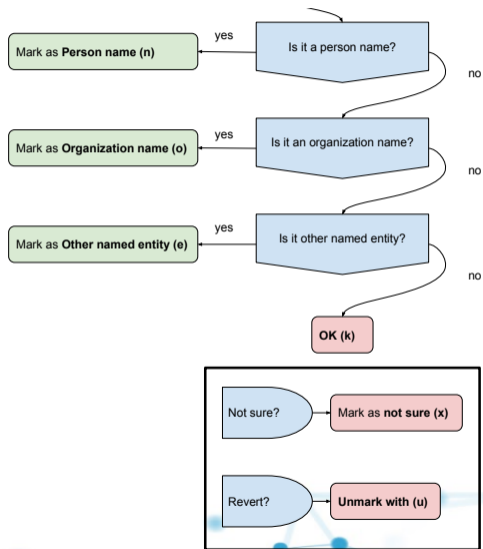
tasks $\Rightarrow$ batches in Lexonomy

central database in NVH (name-value hierarchy), Měchura (2020)

# Corpus

- 650 - 420 = 230
- modified Stanford PoS tagger
- improved stemmer
- sketch grammar

# Headwords



Is it a proper word in **any** language? → no → Mark as **non-word (w)**

yes ↓

Is it a proper word in **Tagalog**? → no → Mark as **foreign word (f)**

yes ↓

Is it a word that is **non-standard** or **offensive**? → yes → Mark as **non-standard/offensive (s)**

no ↓

Is it a lemma? → no → Mark as **wrong lemma (l)**

yes ↓

Does it have correct part-of-speech? → no → Mark as **wrong PoS (p)**

yes ↓

reassigning lemmas

# Pronunciation

# Word senses

# Word senses

**bago**ADJECTIVE

Senses:
- ▸ sense 1 named: [moderno] ✗
- ▸ sense 2 named: [dati] ✗

[ADD SENSE]

Translations:

| | | | |
|---|---|---|---|
| new | ✗ | **1** | 2 |
| modified | ✗ | **1** | 2 |
| prior | ✗ | 1 | **2** |

[ADD TRANSLATION]

---

## cluster 1

Mark all: [ 1 ] [ 2 ] [ NEW ] [ MIXED ] [ ERROR ]

| example usage | actions | | | | | collocate | relation to headword | concordance |
|---|---|---|---|---|---|---|---|---|
| *mga bagong bayani* | **1** | 2 | NEW | MIXED | ERROR | bayani_NOUN | nouns modified by "bago" | ⸗ |
| *ang bagong prinsesa* | **1** | 2 | NEW | MIXED | ERROR | prrinsesa_NOUN | nouns modified by "bago" | ⸗ |
| *bagong superhero* | **1** | 2 | NEW | MIXED | ERROR | superhero_NOUN | nouns modified by "bago" | ⸗ |

## cluster 2

Mark all: [ 1 ] [ 2 ] [ NEW ] [ MIXED ] [ ERROR ]

| example usage | actions | | | | | collocate | relation to headword | concordance |
|---|---|---|---|---|---|---|---|---|
| *mga bagong sibol na* | **1** | 2 | NEW | MIXED | ERROR | sibol_NOUN | nouns modified by "bago" | ⸗ |
| *mga bagong halaman* | **1** | 2 | NEW | MIXED | ERROR | halaman_NOUN | nouns modified by "bago" | ⸗ |
| *bagong puno ng* | **1** | 2 | NEW | MIXED | ERROR | puno_NOUN | nouns modified by "bago" | ⸗ |

- Wikidata + Wiktionary + Wikipedia ("Wikimedia projects")
- Pixabay
- Google Custom Search

**baboy**   PoS: **noun**  |  sense: **hayop**  |

translations:   pig

Choose best image for the above headword and sense:

○ None of the following images is good

clustering and re-clustering

from Sketch Engine, on sense level

GDEX

GT + MB

# Lessons learned

- biggest issues not in technology but in human resource (crowd) and data management
- technologies mature enough to be beneficial
- good and big corpus is the best start for everything
- the better and bigger, the easiest the rest

# Conclusions

- How far can we get?
- almost all steps can be efficiently automated and then post-edited
  - if it pays off
- Tagalog almost finished and soon to be published
- Urdu, Lao in the pipeline
- more detailed information on post-editing statistics to be published