Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Searching the Web for Cross-lingual Parallel Data

Ahmed El-Kishky[†], Philipp Koehn[∗], Holger Schwenk[†]

Facebook AI[†], Johns Hopkins University[∗]
http://www.statmt.org/web-mining-tutorial/

July 25, 2020

Introduction

Corpora and WEB Crawling

Multilingual Represent.
LASER
Evaluation

Document Retrieval

Local Alignment

Global Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext Filtering

# Overview

**1** Background and Motivation

**2** Multilingual Corpora and Web Crawling

**3** Multilingual Representations
    LASER
    Evaluation

**4** Parallel Document Retrieval

**5** Local Sentence Alignment

**6** Global Sentence Alignment
    WikiMatrix
    CCMatrix
    WMT/TED

**7** Parallel Sentence Filtering
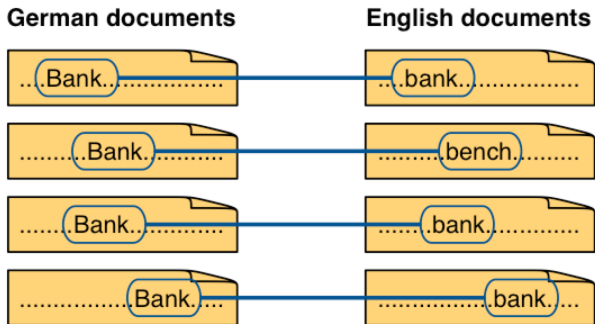
# Intro

For example, machine translation.



**José Salvador** Eu até já assinei a petição mas ainda a pouco tempo li que o presidente de junta que roubou e autorizou essa construção foi homenageado pelo povo das Cortes ...ESTRANHO

I have even signed the petition but I have only recently read that the president of the junta who stole and authorized this construction was honored by the people of the cortes... strange
Automatically Translated

# Learning from Data



**German documents**          **English documents**
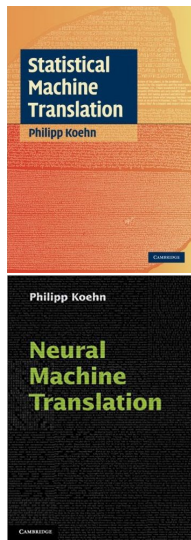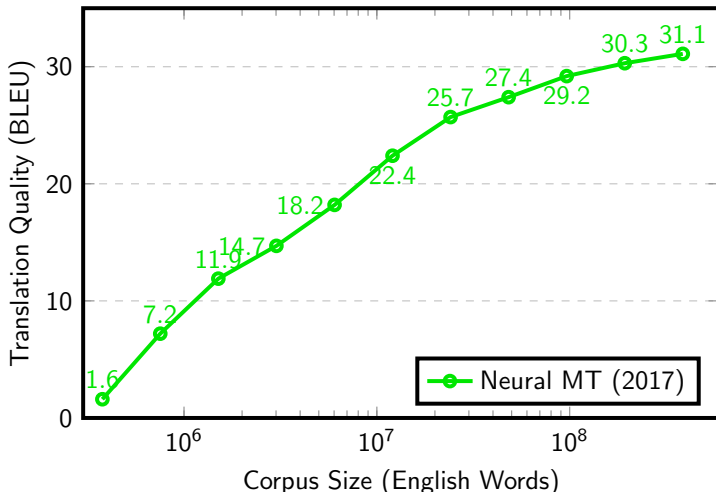
Needed: examples of translated sentences

# Data-Driven Machine Translation

- Given: parallel corpora
  (collections of translated sentences)

- Output: machine translation models

- Since ∼2000: statistical methods
- Since ∼2015: neural methods



Statistical Machine Translation
Philipp Koehn

Philipp Koehn
Neural Machine Translation

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# More Data is Better

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

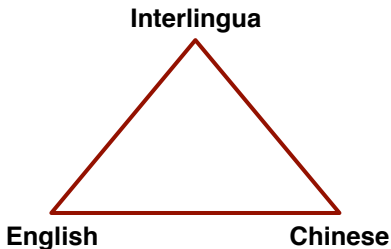Bitext
Filtering

# More Data is Better



(from Google)

# More Data is Better

Don't think about algorithms, get more data!

If you want to think, think about getting more data!

Eric Brill, 2001

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Towards Interlingua

Language-agnostic meaning representations.



**Interlingua**

**English**          **Chinese**

Parallel corpora give us two corners of this triangle

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Other Uses of Parallel Data

For example, multi-lingual hate speech detection.



- Annotate an English corpus
- Train a classifier
- But: use language-independent representations of input (trained on parallel data)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Naturally Occurring Data

- Translation is a common human activity
- Billion dollar industry that
    - localizes products and their documentation
    - makes information accessible in many languages
    - enables communication in multi-lingual organizations
    - translates books, TV shows, movies, ...
- We do not need to create this data.
- We just need to find it.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Large Pools of Data

- For instance, Europarl, 2005
    - well structured web site with clear mapping between translations
    - specialized scripts for crawling, text extraction, alignment
    - maintained structure: sessions, speakers, paragraphs

    *Europarl: A Parallel Corpus for Statistical Machine Translation*
    Koehn, MT Summit 2005

- Other efforts like this
    - Project Syndicate ("news commentary")
    - Global Voices
    - EU Bookstore
    - United Nations
    - Acquis Communitaire

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Commoncrawl

- The web on a hard drive

- Extraction pipeline

  *Dirt Cheap Web-Scale Parallel Text from the Common Crawl*,
  Smith, Saint Amand, Plamada, Koehn, Callison-Burch, Lopez,
  ACL 2013

  - detect document pairs based on URL
  - use HTML structure to check document matches
  - extract text (in chunks indicated by HTML)
  - sentence alignment
  - sentence filtering

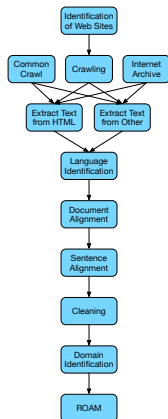- Decent amounts: French (120m words), German (80m),
  ..., Pashto (200k)

- Preview: we wil present methods to extract parallel sentences from CommonCrawl
- CCMatrix: largest collection of high quality mined bitexts
  - 4.5 billion parallel sentences in 39 languages
  - "matrix": aligned across all pairs, not just paired with English
- Extraction purely with retrieval over sentence embeddings

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Paracrawl

- Crawling the web for parallel data
  - funding from two Google grants (2014, 2016)
  - funding from the EU since 2017
- Currently in collaboration with Edinburgh, Alicante, Prompsit, TAUS, Omniscien Technology
- ⇒ Corpora with billions* of words for major languages

* amounts vary based on degree of filtering — for German–English, raw corpus has 4 billion sentence pairs, recommended corpus only 40 million deduplicated sentence pairs)

# Processing Pipeline

- Identifying multi-lingual web sites
- Crawling
- Text extraction from HTML and PDF
- Document alignment
- Sentence alignment
- Sentence pair repair (Bifixer)
- Sentence pair filtering

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Candidate Web Sites

- Extracted all text from CommonCrawl
  - Language ID on all of it

    *N-gram Counts and Language Models from the Common
    Crawl*, Buck, Heafield, Van Ooyen, LREC 2014

  $\Rightarrow$ List of web sites with content in multiple languages
  CommonCrawl has 1.6 million domains with de-en data,
  1.7 million for es-en, etc.

- Search for language name ("Chinese") or flags (en.gif)

- For low resource languages, crawl all web sites with
  language content

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Crawling

- Several off-the-shelf tools available
  - HTTrack: multi-platform tool for crawling
  - Heritrix: Internet Archive's web crawler
  - Creepy: Python library with basic resources for crawling
  - Wget: popular Unix tool

- Many practical problems
  - large sites
  - protected content
  - interference with web server operations
  - robots.txt

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Extract text

- Raw crawls: HTML, TXT, PDF, junk

- Converted into usable format, for each document
    - URL
    - language identification
    - raw HTML (base64)
    - extracted text (base64)

- Special challenges by formats such as PDF

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# A Web Page

HTML Source

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment

WikiMatrix

CCMatrix

WMT/TED

Bitext
Filtering

```
    load-average.html#comments" rel="nofollow">8 comments</a></span>&middot; <span>LAST UPDATED
    <abbr
51  class="updated" title="2013-08-08">August 8, 2013</abbr><span><p
52  class='headline_meta'> in <span><a
53  rel='tag' href="http://www.cyberciti.biz/tips/category/linux">Linux</a>, <a
54  rel='tag' href="http://www.cyberciti.biz/tips/category/monitoring">Monitoring</a>, <a
55  rel='tag' href="http://www.cyberciti.biz/tips/category/sys-admin">Sys admin</a></span>
    </p></div><div
56  class="format_text entry-content"><p><span
57  class="drop_cap">Y</span>es, I know we can use the <kbd>uptime</kbd> command to find out the
    system load average. The uptime command displays the current time, the length of time the
    system has been up, the number of users, and the load average of the system over the last 1,
    5, and 15 minutes. However, if you try to use the uptime command in script, you know how
    difficult it is to get correct load average. As the time since the last, reboot moves from
    minutes, to hours, and an even day after system rebooted. Just type the uptime command:<br
58  /> <span
59  id="more-631"></span><br
60  /> <code>$ uptime</code><br
61  /> Sample outputs:</p><pre>1:09:01  up 29 min,  1 user,  load average: 0.00, 0.00, 0.00</pre>
    <p>OR<br
62  /> <code>$ uptime</code><br
63  /> Sample outputs:</p><pre>2:13AM  up 34 days, 16:15, 36 users, load averages: 1.56, 1.89,
    2.06</pre><p>Traditionally, UNIX administrators used sed and other shell command in scripting
    to get correct value of load average. Here is my own modified hack to save the time<br
64  /> <code>$ uptime | awk -F'load averages:' '{ print $2 }'</code><br
65  /> OR better use the following code:<br
66  /> <code>$ uptime | awk -F'[a-z]:' '{ print $2}'</code><br
67  /> Output taken from my <strong>OS X desktop</strong>:</p><pre> 1.24 1.34 1.35</pre><p>Output
    taken from my <strong>Ubuntu</strong> Linux server:</p><pre> 0.00, 0.01, 0.05</pre><p>Output
    taken from my <strong>RHEL</strong> based server:</p><pre> 0.24, 0.27, 0.21</pre><p>Output
    taken from my <strong>FreeBSD</strong> based server:</p><pre> 0.71, 0.71, 0.58</pre><p>Please
    note that command works on all variant of UNIX operating systems.</p><h2>See also</h2>
    <ul><li>See <a
68  href="http://bash.cyberciti.biz/monitoring/chksysload.bash.php">chksysload.bash</a> script to
```

# Method 1: Strip Tags

LAST UPDATED August 8, 2013 in Linux , Monitoring , Sys admin Y es, I know we can use the uptime command to find out the system load average. The uptime command displays the current time, the length of time the system has been up, the number of users, and the load average of the system over the last 1, 5, and 15 minutes. However, if you try to use the uptime command in script, you know how difficult it is to get correct load average. As the time since the last, reboot moves from minutes, to hours, and an even day after system rebooted. Just type the uptime command: $ uptime Sample outputs: 1:09:01 up 29 min, 1 user, load average: 0.00, 0.00, 0.00

# Method 2: HTML Parser

LAST UPDATED August 8, 2013

in Linux, Monitoring, Sys admin

Y

es, I know we can use the uptime command to find out the system
load average. The uptime command displays the current time, the
length of time the system has been up, the number of users, and
the load average of the system over the last 1, 5, and 15
minutes. However, if you try to use the uptime command in
script, you know how difficult it is to get correct load
average. As the time since the last, reboot moves from minutes,
to hours, and an even day after system rebooted. Just type the
uptime command:

$ uptime

Sample outputs: 1:09:01 up 29 min, 1 user, load average: 0.00,
0.00, 0.00

# What Language?

Muitas intervenções alertaram para o facto de
a política dos sucessivos governos PS, PSD e
CDS, com cortes no financiamento das
instituições do Ensino Superior e com a
progressiva desresponsabilização do Estado
das suas funções, ter conduzido a uma
realidade de destruição da qualidade do Ensino
Superior público.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Clues: Letter N-Grams

Muitas intervenções alertaram para o facto de
a política dos sucessivos governos PS, PSD e
CDS, com cortes no financiamento das
instituições do Ensino Superior e com a
progressiva desresponsabilização do Estado
das suas funções, ter conduzido a uma
realidade de destruição da qualidade do Ensino
Superior público.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Align Documents

- Paracrawl method
    - translate foreign document into English
    - score based on n-gram matches
    - matching of URL
    - other features

    *Quick and Reliable Document Alignment with TF/IDF Cosine Distance*, Buck and Koehn, WMT 2016

- Shared task WMT 2016
    - n-gram matches on (translated) documents powerful
    - only very recently more research on topic

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Align Sentences

The man looks intently at the window.
The sees a shadow.
It was in the trees.
What was it?
He is alarmed and awake.

He has long lived in the woods.
He likes the isolation and solitude of his house.
It's small, but cozy.
The next village is miles away.
He only goes there once a week.

It just after dusk.
The hot sun finally set.
The forest was still abuzz in chatter.
Voices of birds and insects fill the air.
A comforting sound.
But the shadow was larger than those animals.
Only little creatures live here, not this.
It seemed almost as large as a man.
But why that?
Nobody comes ever here.
So the man's eyes keep looking.

As the minutes passed, nothing happens.
But then, cast against the bright moonlit, it returns.

Der Mann schaut aus dem Fenster.
Er sieht einen Schatten in den Bäumen.
Was war das?
Er war alarmiert und wach.

Er hat schon lange im Wald gelebt.
Er genießt die Einsamkeit des Hauses.
Es ist klein.
Aber es ist gemütlich.
Das nächste Dorf ist meilenweit entfernt.
Er geht dorthin nur einmal in Monat.

Es ist nach der Untergang der heißen Sonne.
Der Wald ist voller Geschwätz.
Stimmen von Vögeln und Insekten dringen herüber.

Aber der Schatten war größer als diese Tiere.
Nur Kleingetier lebt hier.
Nicht soetwas Großes.
Es erschien fast so groß wie ein Mensch.
Aber warum, wenn hier niemand jemals herkommt?
Der Mann schaut.
Sein Augen aus dem Fenster gerichtet.

Minuten vergehen, aber nichts passiert.
Dann plötzlich kehrt er im Mondschein zurück.

- Given: pair of documents
- Task: match sentence
- Allow 1-2 mappings etc.?
- Reordering of sentences?
- Several established tools (Hunalign, Bleualign, ...)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Corpus Cleaning

- Two objectives for clean corpus
- Fluency
    - well-formed language
- Adequacy
    - foreign and English sentence have same meaning, style, etc.
- Open question: what is harmful noise?
- Shared tasks at WMT 2018-2020

# Open-Source Code: Bitextor

`https://github.com/bitextor/bitextor`

- Bitextor: Integrated tool to execute entire pipeline

- Pipeline management for distributed computation

# ParaCrawl Release 6

# Multilingual Models

- 7 111 living languages

Native speakers

# Multilingual Models

- 7 111 living languages
- 40% are endangered

Native speakers

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Multilingual Models

- 7 111 living languages
- 40% are endangered
- 23 languages account for half the population

## Native speakers

# Multilingual Models

- 7 111 living languages
- 40% are endangered
- 23 languages account for half the population
- MT: $< 100$ languages

## Native speakers

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Multilingual Models

- 7 111 living languages
- 40% are endangered
- 23 languages account for half the population
- MT: $< 100$ languages
- Almost all NLP applications are mostly English (classification, sentiment analysis or NLI, Q&A, dialog, . . .)

Native speakers

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Multilingual Models

- 7 111 living languages
- 40% are endangered
- 23 languages account for half the population
- MT: $< 100$ languages
- Almost all NLP applications are mostly English (classification, sentiment analysis or NLI, Q&A, dialog, ...)

Native speakers



$\Rightarrow$ Input in foreign language is translated into English

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling
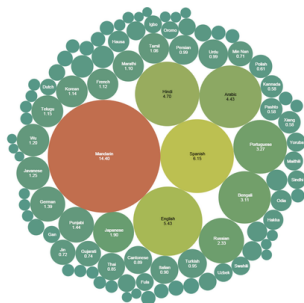
Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Multilingual Models



## Motivation

- Try to embed sentences written in many languages into one joint space
  - ⇒ cross-lingual transfer for various NLP applications
    - benefit of similarities among languages

# Multilingual Models



## Motivation

- Try to embed sentences written in many languages into one joint space
  - ⇒ cross-lingual transfer for various NLP applications
  - benefit of similarities among languages
- This gives us a highly semantic representation

# Multilingual Models



## Motivation

- Try to embed sentences written in many languages into one joint space
  - ⇒ cross-lingual transfer for various NLP applications
  - benefit of similarities among languages
- This gives us a highly semantic representation
- ⇒ Sentences with similar meaning are close
  (mono- or cross-lingual)

# Multilingual Models



Applications:

- zero-shot transfer
- bitext mining and filtering
- large-scale similarity search
- paraphrasing
- data augmentation
- ...

# Multilingual Models



Applications:

- zero-shot transfer
- bitext mining and filtering
- large-scale similarity search
- paraphrasing
- data augmentation
- …

# Multilingual Models

Some approaches:

- MUSE
  - unsupervised multilingual word embeddings
- LASER
  - supervised multilingual **sentence** embeddings
- XLM
  - unsupervised multilingual sentence embeddings
- Sentence BERT
  - fine-tuned for linguistic similarity

$\vdots$

# Multilingual Models

Some approaches:

- MUSE
  - unsupervised multilingual word embeddings
- LASER
  - supervised multilingual **sentence** embeddings
- XLM
  - unsupervised multilingual sentence embeddings
- Sentence BERT
  - fine-tuned for linguistic similarity

⋮

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# MUSE



## Principle

- Learn multilingual word embeddings
  without any aligned data
- fastText embeddings aligned in a common space
  - learn transformation of space $X$ to $Y$
- A. Conneau et al.,
  *Word Translation Without Parallel Data*, ICLR'18
- https://github.com/facebookresearch/MUSE

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# LASER: Architecture



Seq2seq approach with one joint encoder and decoder

- Based on `fairseq`

- Shared encoder and decoder for several languages

- No attention, but max-pooling

- Sentence representation is used at the input at each time step and to initialize decoder

- Also target language embedding

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# LASER: Architecture



## Training strategies

- *N*:1 translation is enough to learn a joint embedding
- No explicit criterion to enforce joint embedding
    - ranking loss
    - GAN to predict language
    - . . .

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# LASER: Architecture



## Training strategies

- $N$:1 translation is enough to learn a joint embedding
- No explicit criterion to enforce joint embedding
  - ranking loss
  - GAN to predict language
  - ...
- But $N$:1 doesn't cover target language (English)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# LASER: Architecture



### Training strategies

- *N*:1 translation is enough to learn a joint embedding
- No explicit criterion to enforce joint embedding
    - ranking loss
    - GAN to predict language
    - . . .
- But *N*:1 doesn't cover target language (English)
- Limited success with (noisy) autoencoder

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment

WikiMatrix

CCMatrix

WMT/TED

Bitext
Filtering

# LASER: Architecture



## Training strategies

- How to have a language at the input and output ?
  - in the past: $N \rightarrow$ (N-1)
- Two target languages are enough
  - English and Spanish
  - independently aligned
  - not all input languages need to be aligned to both
- Language pair is changed at each mini-batch
- Trained on 223M sentences of public bitexts

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment

WikiMatrix

CCMatrix

WMT/TED

Bitext
Filtering

# LASER: Architecture



## Encoder

- 5 layer BiLSTM (depth helps !)
- No information on input (or target) language
- Shared BPE tokens, 50k BPE operations
- No pretraining of BPE embeddings
- The training procedure makes no assumption on the encoder:
  $\Rightarrow$ transformers, convolutional, . . .

# LASER: Training languages

Afrikaans Albanian Amharic Arabic Armenian Aymara Azerbaijani Basque Belarusian Bengali Chavacano Chinese mandarin Coastal Swahili Croatian Burmese Bokma Berber languages Norwegian Breton Bulgarian Catalan Central Khmer Bosnian Czech Danish Dutch English Esperanto Estonian Finnish French Galician Georgian German Greek Hausa Hebrew Hindi Hungarian Icelandic Interlingua Interlingue Ido Indonesian Iranian Persian (Farsi) Italian Japanese Kabyle Kazakh Korean Kurdish Latavian Latin Lingua Franca Nova Lithuanian Low German / Saxon Macedonian Malay Malagasy Malayalam Marathi Maldivian (Divehi) Moldavian Russian Romanian Occitan (post 1500) Oriya Polish Portuguese Serbian Sindhi Sinhala Slovak Slovenian Somali Spanish Swedish Tagalog Tajik Tamil Tatar Telugu Thai Turkish Uighur Ukrainian Urdu Uzbek Vietnamese Wu Chinese Yue Chinese

Cross-lingual Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# LASER: Training languages

## 22 different writing scripts:

| Arabic | هناك العديد من اللغات في العالم. | Hebrew | יש נ שפות רבות בעולם. |
| Armenian | Աշխարհում շատ լեզուներ կան: | Kanjii | 世界にはたくさんの言語があります。 |
| Burmese | ကမ္ဘာပေါ်တွင်ဘာသာစကားများစွာရှိပါတယ်။ | Khmer | មានភាសាជាច្រើននៅលើពិភពលោក។ |
| Chinese | 世界上有很多种语言。 | Latin | There are many languages in the world. |
| Cyrillic | В мире много языков. | Malayalam | ലോകത്തിൽ അനേകം ഭാഷകൾ ഉണ്ട് |
| Devanagari | दुनिया में कई भाषाएं हैं। | Persian | زبان‌های بسیاری در جهان وجود دارد. |
| Eastern-Nagari | বিশ্বের অনেক ভাষা আছে। | Sinhala | ලෝකයේ බොහෝ භාෂාවන් පවතී. |
| Ge'ez | በዓለም ውስጥ ብዙ ቋንቋዎች አሉ. | Tamil | உலகில் பல மொழிகள் உள்ளன. |
| Georgian | მსოფლიოში ბევრი ენაა. | Telugu | ప్రపంచంలో అనేక భాషలు ఉన్నాయి. |
| Greek | Υπάρχουν πολλές γλώσσες στον κόσμο. | Thaana | *No free translation for Maldivian (Dhivehi)* |
| Hangul | 세계에는 많은 언어가 있습니다. | Thai | มีหลายภาษาในโลกนี้ |

- One single encoder can handle all these scripts
- All these sentences are close in the embedding space
- It is not necessary to specify the language or script
- Code-switching is also supported

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# LASER toolkit

Massively Multilingual Sentence Embeddings for Zero-Shot
Cross-Lingual Transfer and Beyond. Trans. Assoc. Comput.
Linguistics 7: 597-610 (2019)

## Well established in community, academia and industry

- https://github.com/facebookresearch/LASER/
- M. Artetxe and H. Schwenk, *Massively Multilingual
  Sentence Embeddings for Zero-Shot Cross-Lingual
  Transfer and Beyond*, TACL'19 and arXiv'18
- Fast and easy to use (2000 sentences/sec)
- One model for many applications
- Current SOTA for filtering and mining bitexts

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Applications of Multilingual Embeddings

- Zero-shot transfer in NLP
  - Use ML embeddings to train English NLP system
  - ⇒ apply it to other languages without any modification
  - classification, NLI, QA, . . .

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Applications of Multilingual Embeddings

- Zero-shot transfer in NLP
  - Use ML embeddings to train English NLP system
  - $\Rightarrow$ apply it to other languages without any modification
  - classification, NLI, QA, . . .
- Bitexts mining and filtering
  - sentence similarity $\sim$ distance in joint space

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Applications of Multilingual Embeddings

- Zero-shot transfer in NLP
  - Use ML embeddings to train English NLP system
  - $\Rightarrow$ apply it to other languages without any modification
  - classification, NLI, QA, ...
- Bitexts mining and filtering
  - sentence similarity $\sim$ distance in joint space
- Large-scale similarity search
  - index many sentences, search for closest ones
  - paraphrasing, data augmentation, ...

# Applications of Multilingual Embeddings

- Zero-shot transfer in NLP
  - Use ML embeddings to train English NLP system
  - $\Rightarrow$ apply it to other languages without any modification
  - classification, NLI, QA, . . .
- Bitexts mining and filtering
  - sentence similarity $\sim$ distance in joint space
- Large-scale similarity search
  - index many sentences, search for closest ones
  - paraphrasing, data augmentation, . . .

  **We always use the same LASER sentence
  embeddings, no task-specific fine-tuning**

# XNLI: Cross-Lingual NLI

- Fixed LASER embeddings
- NLI classifier trained on English only

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# XNLI: Cross-Lingual NLI

- Fixed LASER embeddings
- NLI classifier trained on English only
- Zero-shot transfer to any language supported by LASER

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# XNLI: Cross-Lingual NLI

- Fixed LASER embeddings
- NLI classifier trained on English only
- Zero-shot transfer to any language supported by LASER
- We can arbitrarily combine sentence in any language

| Premise | Hypothesis | Relation |
|---|---|---|
| **Bulgarian**<br>Никой не знаеше къде отидоха.<br>*Their destination was a secret.* | **Hindi**<br>उनका गंतव्य गुप्त था।<br>*Nobody knew where they went.* | Related<br>(line 210) |
| **Arabic**<br>مم ، ومذ ثمّ انتقلنا إلى منزلٍ جديد .<br>*Um, then we moved to a new house.* | **Swahili**<br>Tuliishi kwa nyumba moja maisha yetu yote.<br>*We stayed in the same house our whole lives.* | Opposite<br>(line 393) |
| **Thai**<br>สัปดาห์ต่อมา, หลานชายของฉันขอกีตาร์อะคูสติก<br>กินในวันเกิดของเขา<br>*The next week, my nephew asked for an acoustic guitar for his birthday.* | **Spanish**<br>Aprender a tocar la guitarra y comenzar una banda era todo lo que hablaba mi sobrino.<br>*Learning to play guitar and starting a band was all that my nephew talked about.* | Neutral<br>(line 4702) |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Bitext Mining Approach

## Margin criterion

Semantic similarity $\propto$ distance
$\Rightarrow$ mine parallel sentences

$$
\begin{aligned}
&\text{margin}(x, y) \\
&= \frac{\cos(x, y)}{\displaystyle\sum_{z \in \text{NN}_k(x)} \frac{\cos(x, z)}{2k} + \sum_{z \in \text{NN}_k(y)} \frac{\cos(y, z)}{2k}}
\end{aligned}
\tag{1}
$$

*(Artexe and Schwenk, arXiv Nov'18 and ACL'19)*

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Results for 93 Languages: BUCC

|  | TRAIN | | | | TEST | | | |
|---|---|---|---|---|---|---|---|---|
|  | de-en | fr-en | ru-en | zh-en | de-en | fr-en | ru-en | zh-en |
| *Azpeita et '17* | 83.3 | 78.83 | - | - | 83.7 | 79.5 | - | - |
| *Grégoire&Langlais '17* | - | 20.7 | - | - | - | 20 | - | - |
| *Zhang & Zweigenbaum '17* | - | - | - | 43.48 | - | - | - | 45.13 |
| *Azpeita et al. '18* | 84.3 | 80.6 | 80.9 | 76.5 | 85.5 | 81.5 | 81.3 | 77.5 |
| *Bouamor & Sajad '18* | - | 75.2 | - | - | - | 76.0 | - | - |
| *Leong & Chao '18* | - | - | - | 58.5 | - | - | - | 56 |
| *Schwenk ACL'18* | 76.1 | 74.9 | 73.3 | 71.6 | 76.9 | 75.8 | 73.8 | 71.6 |
| *Artetxe&Schwenk arXiv'18* | 94.8 | 91.9 | 90.9 | 91.0 | 95.6 | 92.9 | 92.0 | **92.6** |
| Proposed method | **95.4** | **92.4** | **92.3** | 91.2 | **96.2** | **93.9** | **93.3** | 92.3 |

- Significantly outperforms other systems of the BUCC eval

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment

WikiMatrix

CCMatrix

WMT/TED

Bitext
Filtering

# Results for 93 Languages: BUCC

|  | TRAIN | | | | TEST | | | |
|---|---|---|---|---|---|---|---|---|
|  | de-en | fr-en | ru-en | zh-en | de-en | fr-en | ru-en | zh-en |
| *Azpeita et '17* | 83.3 | 78.83 | - | - | 83.7 | 79.5 | - | - |
| *Grégoire&Langlais '17* | - | 20.7 | - | - | - | 20 | - | - |
| *Zhang & Zweigenbaum '17* | - | - | - | 43.48 | - | - | - | 45.13 |
| *Azpeita et al. '18* | 84.3 | 80.6 | 80.9 | 76.5 | 85.5 | 81.5 | 81.3 | 77.5 |
| *Bouamor & Sajad '18* | - | 75.2 | - | - | - | 76.0 | - | - |
| *Leong & Chao '18* | - | - | - | 58.5 | - | - | - | 56 |
| *Schwenk ACL'18* | 76.1 | 74.9 | 73.3 | 71.6 | 76.9 | 75.8 | 73.8 | 71.6 |
| *Artetxe&Schwenk arXiv'18* | 94.8 | 91.9 | 90.9 | 91.0 | 95.6 | 92.9 | 92.0 | **92.6** |
| Proposed method | **95.4** | **92.4** | **92.3** | 91.2 | **96.2** | **93.9** | **93.3** | 92.3 |

- Significantly outperforms other systems of the BUCC eval
- New system trained on 93 languages is better than dedicated system, limited to eval languages

# Generalization to New Languages

## System trained on the 21 languages of Europarl

|                   | De-En | Fr-En |
|-------------------|-------|-------|
| State-of-the-art  | 85.5  | 81.5  |
| Our approach      | 95.6  | 92.9  |

# Generalization to New Languages

System trained on the 21 languages of Europarl

|                  | De-En | Fr-En | Ru-En |
|------------------|-------|-------|-------|
| State-of-the-art | 85.5  | 81.5  | 81.3  |
| Our approach     | 95.6  | 92.9  | 62.0  |

- Good performance on Russian (precision=80%)
  **although Russian was not used during training**

# Generalization to New Languages

System trained on the 21 languages of Europarl

|  | De-En | Fr-En | Ru-En |
|---|---|---|---|
| State-of-the-art | 85.5 | 81.5 | 81.3 |
| Our approach | 95.6 | 92.9 | 62.0 |

- Good performance on Russian (precision=80%)
  **although Russian was not used during training**
- ⇒ Very promising to mine data for dialects and minority
  languages which are in the same family than a trained
  language, Gallician, Nepali, . . .

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Results for 93 Languages: Bitext Filtering

## WMT'19: Bitext filtering for low-resource conditions

- Filter very noisy Paracrawl crawled bitexts (40-60M)
- Evaluation by training SMT and NMT systems:
- Train: En/Ne (586k), En/Si (645k) + En/Hi (1.5M)

|       | Ne/En 1M | | Ne/En 5M | | Si/En 1M | | Si/En 5M | |
|       | SMT | NMT | SMT | NMT | SMT | NMT | SMT | NMT |
|-------|------|------|------|------|------|------|------|------|
| LASER | **4.21** | **6.88** | 4.63 | 2.84 | **4.27** | **6.39** | **4.94** | 4.02 |
| 2nd   | 4.10 | 5.48 | **4.74** | **3.43** | 4.19 | 4.97 | 4.62 | **4.44** |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Results for 93 Languages: Bitext Filtering

## WMT'19: Bitext filtering for low-resource conditions

- Filter very noisy Paracrawl crawled bitexts (40-60M)
- Evaluation by training SMT and NMT systems:
- Train: En/Ne (586k), En/Si (645k) + En/Hi (1.5M)

|  | Ne/En 1M | | Ne/En 5M | | Si/En 1M | | Si/En 5M | |
|---|---|---|---|---|---|---|---|---|
|  | SMT | NMT | SMT | NMT | SMT | NMT | SMT | NMT |
| LASER | **4.21** | 6.88 | 4.63 | 2.84 | **4.27** | 6.39 | **4.94** | 4.02 |
| 2nd | 4.10 | 5.48 | **4.74** | **3.43** | 4.19 | 4.97 | 4.62 | **4.44** |

- Overall best results by significant margin ($+25\%$)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment

WikiMatrix

CCMatrix

WMT/TED

Bitext
Filtering

# Results for 93 Languages: Bitext Filtering

## WMT'19: Bitext filtering for low-resource conditions

- Filter very noisy Paracrawl crawled bitexts (40-60M)
- Evaluation by training SMT and NMT systems:
- Train: En/Ne (586k), En/Si (645k) + En/Hi (1.5M)

|       | Ne/En 1M | | Ne/En 5M | | Si/En 1M | | Si/En 5M | |
|       | SMT | NMT | SMT | NMT | SMT | NMT | SMT | NMT |
|-------|------|------|------|------|------|------|------|------|
| LASER | **4.21** | **6.88** | 4.63 | 2.84 | **4.27** | **6.39** | **4.94** | 4.02 |
| 2nd   | 4.10 | 5.48 | **4.74** | **3.43** | 4.19 | 4.97 | 4.62 | **4.44** |

- Overall best results by significant margin (+25%)
- Good filtering is more important for NMT than SMT

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Results for 93 Languages: Bitext Filtering

## WMT'19: Bitext filtering for low-resource conditions

- Filter very noisy Paracrawl crawled bitexts (40-60M)
- Evaluation by training SMT and NMT systems:
- Train: En/Ne (586k), En/Si (645k) + En/Hi (1.5M)

|       | Ne/En 1M | | Ne/En 5M | | Si/En 1M | | Si/En 5M | |
|-------|------|------|------|------|------|------|------|------|
|       | SMT | NMT | SMT | NMT | SMT | NMT | SMT | NMT |
| LASER | **4.21** | **6.88** | 4.63 | 2.84 | **4.27** | **6.39** | **4.94** | 4.02 |
| 2nd   | 4.10 | 5.48 | **4.74** | **3.43** | 4.19 | 4.97 | 4.62 | **4.44** |

- Overall best results by significant margin ($+25\%$)
- Good filtering is more important for NMT than SMT
- Vishrav et al, *Low-Resource Corpus Filtering using Multilingual Sentence Embeddings*, WMT'19

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# XLM: Architecture



## Multilingual extension of BERT

- Unsupervised: Masked LM training (MLM)
  - Joint BPE or SentencePiece vocabulary
- Supervised: Cross-Lingual LM training (CLM)
  - attention can attend words in either language which encourages alignment

# XLM: Applications

## Results

- Trained unsupervised on 2.5 billon sentences of CC
- Supports 100 languages, very strong results on many English and cross-lingual tasks (GLUE, XNLI, . . . )
    - generally task specific fine-tuning
- G. Lample and A. Conneau,
  *Cross-lingual Language Model Pretraining*, NIPS'19
- A. Conneau et al.,
  *Unsupervised Cross-lingual Representation Learning at Scale*, ACL'20
- https://github.com/facebookresearch/XLM
- Application to similarity search requires some sort of fine-tuning

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Sentence BERT

### Recent research

- Several works aim in achieving transformer-based language agnostic sentence representations

- Feng et al., *Language-agnostic BERT Sentence Embedding*, arxiv Jul'20



- Very Interesting results, but no comparision with margin-based LASER mining

# parallel document retrieval

# Cross-lingual Document Retrieval

Finding pairs of documents that are
translations/near-translations of each other.

# Cross-Lingual Document Retrieval

Large, Heterogenous, Multilingual Web-Corpora

Language Identification

Cross-lingual Document Retrieval

Bi-Lingual Sentence Alignment

Sentence Pair Filtering

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Cross-lingual Document Retrieval



(a) English Webpage        (b) French Webpage

Figure: Two web documents that are translations of each other.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Cross-lingual Document Retrieval



(a) Arabic Webpage

(b) Spanish Webpage

Figure: Two web documents that are translations of each other.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Motivation

- Training data for information retrieval
  - Supervision for learning-to-rank
  - Supervision for retrieval
- Source of training data for learning multilingual representation
  - Cross-lingual word representations
  - Cross-lingual sentence representations
  - Cross-lingual document representation
- Source of training data for machine translation (BLEU goes up)
  - Mine parallel data for low-resource directions
  - Web parallel data covers a variety of domains

# Objective

Given a corpus of web-documents, automatically identify pairs
of documents that are translations of each other.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Objective



Figure: Documents are aligned 1-to-1 within each domain.

# Evaluation

- Recall only i.e. what percentage of the test-set pairs is found
- 1-1 rule; every document can only occur in one pair.

CC-Aligned

# CC-Aligned: A Massive Collection of Cross-Lingual Web Documents

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Motivation

Creating a large cross-lingual parallel document dataset can be valuable

- High-quality multilingual dataset can be used to benchmark document alignment algorithms

- Parallel dataset can be used for supervision for cross-lingual representation

- A large parallel dataset can be mined for parallel sentences for NMT training

# Insights

URL Signals for Parallel Web Documents

- URLs often contain language codes signifying the language a piece of web content is in
- URL structural information can be used as a signal for identifying parallel documents
  - `https://anonymizedURL.com`
  - `https://fr-fr.anonymizedURL.com`

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Insights

| Source URL | Target URL |
| --- | --- |
| **eng.**aaa.com | aaa.com |
| aaa.com/**en-gb**/b | aaa.com/**zh-cn**/b |
| aaa.com/**English**/b | aaa.com/**Yoruba**/b |
| aaa.com/b/**en** | aaa.com/b/**vi** |
| aaa.com/b/ | **thai.**aaa.com/b/ |
| aaa.com/b**&lang=english** | aaa.com/b**&lang=arabic** |
| aaa.com/b**?lang=en** | aaa.com/b**?lang=fr** |
| aaa.com/b | aaa.com/b**?lang=1** |

Table: URL matching via language identifiers.

# CCAligned Dataset

Multilingual Web Corpus → Language Identification → URL Web Alignment

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Common Crawl Corpus

CommonCrawl Corpus: An Open Repository of Web Data

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Common Crawl Corpus

CommonCrawl Corpus: An Open Repository of Web Data

- Text content
- Status information
- HTTP response code
- HTML title
- HTML meta tags
- RSS/Atom information
- All anchors/hyperlinks

# Corpus Statistics

Corpus Statistics

- 68 CommonCrawl Snapshots (every month 2013-2020)
- Each snapshot contains over 2 billion web-documents
- 169.4 billion web documents
- 107.8 million distinct web-domainis

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Preprocessing Steps

URL Normalization

- URL Normalization: removing the protocol and host name
  - https://www.aaa.com $\rightarrow$ aaa.com)
- Deduplicate based on normalized URL
  - URL that appears more than once, we select the instance that possesses the longest document content.
    1. document was deleted and gets shorter
    2. document is amended and gets longer.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# De-duplication Corpus Statistics

- 169.4 billion documents $\rightarrow$ 29.6 billion
- 83% reduction from raw corpus
- 107.8 million distinct web-domains

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Mining Parallel Documents

The next step is to mine parallel web documents.

1. De-duplicate CommonCrawl corpus
2. Perform language identification on each web-document.
3. Apply URL-Matching Heuristics

# Mined Parallel Documents

Parallel Cross-lingual Documents

1. 364 million aligned documents
   - 100M with English
   - 264M without English
2. 4598 language pairs
   - 98 with English
   - 4500 without English

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# CCAligned Dataset Quality

Human annotators evaluated quality of the mined documents

|  | Language | $P_{maj}$ | $K\alpha$ | $P_{adj}$ |
|---|---|---|---|---|
| **High** | German | 90.0 | 0.74 | 96.7 |
| | Chinese | 86.7 | 0.68 | 93.3 |
| **Mid** | Arabic | 83.3 | 0.72 | 90.0 |
| | Romanian | 76.7 | 0.50 | 96.7 |
| **Low** | Estonian | 83.3 | 0.68 | 90.0 |
| | Burmese | 86.7 | 0.88 | 100.0 |
| | **Avg** | 84.4 | 0.70 | 94.5 |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Dataset Analysis

Dataset Analysis

- High-precision collection of cross-lingual documents
- Dataset was constructed using **ONLY** URL-features
- Can one evaluate content-based alignment strategies on this dataset?

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Content-based Alignment: Direct Embedding

Direct Embedding (DE) with LASER

- Embed the entire document using LASER embedding
- Each document $d$ has its dense vector representation $\mathbf{v}_d$

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Content-based Alignment: Sentence Average Embedding

Sentence Averaging (SA) with LASER

1. Decompose each document into sentences
2. Embed each sentence using LASER
3. document embedding by averaging the sentence vectors into a document vector $\mathbf{v}_d$

$$\mathbf{v}_d = \frac{1}{n} \sum_{i=1}^{n} \mathbf{v}_{s_i} \tag{2}$$

# Content-based Alignment: Weighted Sentence Average Embedding

Weighted Sentence Averaging (WSA) with LASER. Try common information retrieval tricks

1. **Sentence Length (SL)**: Longer sentences more important than shorter
2. **Inverse Document Frequency (IDF)**: More frequent sentences may be unimportant
3. **SL-IDF**: Combine both

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Content-based Alignment: Weighted Sentence Average Embedding

Weighted Sentence Averaging (WSA) with LASER

$$\mathbf{v}_d = \frac{1}{n} \sum_{i=1}^{n} \mathbf{w}_{s_i} \times \mathbf{v}_{s_i} \tag{3}$$

$$SL_{s_i} = \frac{|s_i|}{\sum_{s \in d} count(s) \times |s|} \tag{4}$$

$$IDF_{s_i} = \log \frac{N+1}{1 + |\{d \in D : s \in d\}|} \tag{5}$$

$$SLIDF_{s_i} = SL_{s_i} \times IDF_{s_i} \tag{6}$$

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Scoring Function

Cross-lingual Document Similarity

- dense document representations for each document from the source and target sets

- score pairs to evaluate how semantically similar documents are

- given two documents $a$ and $b$, compute their semantic similarity using a cosine similarity

$$sim(a, b) = \frac{\mathbf{v}_a \cdot \mathbf{v}_b}{||\mathbf{v}_a|| \; ||\mathbf{v}_b||} \qquad (7)$$

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Competitive Matching

Ensuring the pairs are 1-to-1 (each aligned document is in at most one pair)

- each document in the source document set, $D_s$ is paired with each document in the target set, $D_t$
- $D_s \times D_t$ scored pairs – a fully connected bipartite graph
- expected output assumes each page in the non-dominant language has a translated or comparable counterpart
- $min(|D_s|, |D_t|)$ expected number of aligned pairs
- Hungarian algorithm $\mathcal{O}(max(|D_s||D_t|)^3)$ ... intractable

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Competitive Matching

**Algorithm 1:** Competitive Matching

---

1 **Input:** $P = \{(d_s, d_t) | d_s \in D_s, d_t \in D_t\}$
2 **Output:** $P' = \{(d_{s,i}, d_{t,i}), ...\} \subset P$

3 $scored \leftarrow \{(p, score(p)) \text{ for } p \in P\}$
4 $sorted \leftarrow sort(scored)$ in descending order
5 $aligned \leftarrow \varnothing$
6 $S_s \leftarrow \varnothing$
7 $S_t \leftarrow \varnothing$
8 **for** $d_s, d_t \in sorted$ **do**
9      if $d_s \notin S_s \wedge d_t \notin S_t$
10      $aligned \leftarrow aligned \cup \{(d_s, d_t)\}$
11      $S_s \leftarrow S_s \cup d_s$
12      $S_t \leftarrow S_t \cup d_t$
13 **end**
14 **return** $aligned$

---

Cross-lingual Mining

A. El-Kishky, P. Koehn, H. Schwenk

Introduction

Corpora and WEB Crawling

Multilingual Represent.
LASER
Evaluation

Document Retrieval

Local Alignment

Global Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext Filtering

# Alignment Results High-Resource

| | **Recall** | | | | |
|---|---|---|---|---|---|
| **Language** | **DE** | **SA** | **SL** | **IDF** | **SLIDF** |
| French | 0.39 | **0.84** | 0.83 | 0.82 | **0.84** |
| Spanish | 0.34 | 0.53 | 0.55 | **0.58** | 0.57 |
| Russian | 0.06 | 0.48 | 0.50 | **0.61** | 0.60 |
| German | 0.52 | 0.74 | **0.76** | 0.74 | **0.76** |
| Italian | 0.22 | 0.54 | 0.55 | 0.55 | **0.57** |
| Portuguese | 0.17 | 0.36 | 0.39 | 0.33 | **0.40** |
| Dutch | 0.28 | 0.51 | 0.54 | 0.52 | **0.56** |
| Indonesian | 0.11 | 0.36 | **0.48** | 0.43 | **0.48** |
| Polish | 0.17 | 0.38 | 0.41 | **0.44** | 0.42 |
| Turkish | 0.12 | 0.30 | 0.34 | **0.45** | 0.41 |
| Swedish | 0.19 | 0.37 | 0.37 | 0.38 | **0.39** |
| Danish | 0.27 | 0.46 | 0.65 | 0.60 | **0.67** |
| Czech | 0.15 | 0.36 | **0.41** | 0.32 | **0.41** |
| Bulgarian | 0.07 | 0.34 | 0.37 | 0.40 | **0.44** |
| Finnish | 0.06 | 0.24 | 0.32 | 0.43 | **0.44** |
| Norwegian | 0.13 | 0.26 | 0.33 | 0.33 | **0.38** |
| **Macro-AVG** | 0.20 | 0.41 | 0.45 | 0.47 | **0.49** |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Alignment Results Mid-Resource

| | Recall | | | | |
|---|---|---|---|---|---|
| **Language** | DE | SA | SL | IDF | SLIDF |
| Romanian | 0.15 | 0.39 | 0.40 | 0.40 | **0.41** |
| Vietnamese | 0.06 | 0.13 | 0.18 | 0.15 | **0.23** |
| Ukrainian | 0.05 | 0.49 | 0.70 | 0.70 | **0.74** |
| Greek | 0.05 | 0.22 | 0.24 | **0.34** | 0.30 |
| Korean | 0.06 | 0.49 | 0.47 | 0.49 | **0.51** |
| Arabic | 0.04 | 0.26 | 0.46 | 0.42 | **0.51** |
| Croatian | 0.16 | 0.32 | 0.36 | 0.34 | **0.36** |
| Slovak | 0.20 | 0.37 | **0.44** | 0.41 | 0.42 |
| Thai | 0.02 | 0.15 | 0.28 | 0.19 | **0.35** |
| Hebrew | 0.05 | 0.19 | 0.30 | 0.27 | **0.33** |
| Hindi | 0.04 | 0.03 | 0.33 | 0.28 | **0.43** |
| Hungarian | 0.15 | 0.41 | 0.39 | 0.39 | **0.46** |
| Lithuanian | 0.11 | 0.61 | 0.72 | 0.74 | **0.80** |
| Slovenian | 0.13 | 0.20 | 0.26 | 0.31 | **0.33** |
| Farsi | 0.06 | 0.22 | 0.37 | 0.40 | **0.49** |
| | | | | | |
| **Macro-AVG** | 0.09 | 0.28 | 0.39 | 0.39 | **0.44** |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Alignment Results Low-Resource

| | Recall | | | | |
|---|---|---|---|---|---|
| **Language** | DE | SA | SL | IDF | SLIDF |
| Estonian | 0.28 | 0.57 | 0.62 | 0.58 | **0.64** |
| Bengali | 0.05 | 0.47 | **0.59** | 0.51 | 0.58 |
| Albanian | 0.23 | 0.56 | 0.60 | 0.57 | **0.61** |
| Macedonian | 0.02 | 0.16 | **0.22** | 0.19 | 0.08 |
| Urdu | 0.06 | **0.29** | 0.23 | 0.27 | 0.24 |
| Serbian | 0.06 | 0.46 | **0.58** | 0.47 | 0.56 |
| Azerbaijani | 0.08 | 0.27 | 0.28 | **0.34** | 0.27 |
| Armenian | 0.02 | 0.08 | 0.13 | 0.12 | **0.17** |
| Belarusian | 0.07 | 0.26 | 0.44 | 0.36 | **0.51** |
| Georgian | 0.06 | 0.18 | 0.23 | **0.25** | **0.25** |
| Tamil | 0.02 | 0.13 | 0.19 | 0.23 | **0.34** |
| Marathi | 0.02 | 0.13 | **0.20** | 0.10 | 0.16 |
| Kazakh | 0.05 | 0.16 | 0.24 | 0.25 | **0.33** |
| Mongolian | 0.03 | 0.01 | 0.05 | 0.10 | **0.22** |
| Burmese | 0.01 | **0.35** | 0.18 | 0.08 | 0.26 |
| Bosnian | 0.18 | 0.49 | 0.64 | 0.50 | **0.65** |
| **Macro-AVG** | 0.08 | 0.29 | 0.34 | 0.31 | **0.37** |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Downstream Mining

From the aligned documents, can do further sentence-level
mining.

- From this dataset, mined **2.25** billion parallel sentences
  covering 4598 language pairs
- 950 million pairs are sentences paired with English
  sentences
- 1.3 billion pairs are non-English sentence pairs

# Follow-up Research

From CCAligned aligned documents, there are many open
research problems that can leverage this data

- Mine more, higher quality parallel sentences from the
  CCAligned documents
- Use CCAligned documents as supervision for supervised
  document alignment (mine parallel documents using
  high-recall method)
- Leverage parallel documents to learn cross-linigual
  document representations and cross-lingual document
  retrieval

Cross-lingual Sentence Mover's
Distance

# Massively Multilingual Document Alignment with Cross-lingual Sentence Mover's Distance

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Motivation & Insight

- Motivation
    - Cross-lingual retrieval based on content is more general than using metadata (URL, timestamp, etc)
    - CCAligned is high-precision. For more training data (especially for low-resource direction, need a high-recall approach)
- Insight
    - Creating document level fixed representations may be destructive for variable-length documents.
    - Averaging sentence embedding places equal importance to all sentences
    - How well sentences match up between document pairs is a good signal for parallel documents

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Earth Mover's Distance

- measure of the distance between two probability distributions over a region $D$
- For example: if the distributions are interpreted as two different ways of piling up a certain amount of dirt over the region D
  - the EMD is the minimum cost of turning one pile into the other
  - the cost is assumed to be amount of dirt moved times the distance by which it is moved

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Earth Mover's Distance



- red distribution: "dirt"
- blue distribution: "holes"

The distance between points (ground distance) can be Euclidean distance, Manhattan...

Cross-lingual Mining

A. El-Kishky, P. Koehn, H. Schwenk

Introduction

Corpora and WEB Crawling

Multilingual Represent.
LASER
Evaluation

Document Retrieval

Local Alignment

Global Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext Filtering

# Cross-lingual Sentence Mover's Distance

- Each document has a distribution over sentences
  - multinomial distribution - normalize bag of sentences (nBOS)
- Euclidean distance between source document sents and target document sents
  - Leverage LASER embeddings to compute Euclidean distances

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Weighting Sentences Based on Importance

- XL-SMD requires a distribution over sentences for each document
- Each sentence has probability mass allocated to it.
- 4weighting schemes for each sentence investigated
  - Uniform weighting (each sentence equally weighted)
  - Sentence length (Longer sentences = more mass)
  - Inverse document frequency (IDF)
  - SL-IDF

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Document Mass Normalization

Normalizing the mass to unit measure in both the source and target documents each each document has a legitimate distribution and the induced distance metric is valid.

$$d'_{A,i} = \frac{d_{A,i}}{\sum\limits_{s \in A} d_{A,s}} \qquad (8)$$

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Optimal Transport

- $\Delta(i,j)$ is distance between the $i_{th}$ and $j_{th}$ sentences
- $V$ denote vocab size (sentences within a document pair)
- $\Delta(i,j) = ||v_i - v_j||$

$$XLSMD(A, B) = \min_{T \geq 0} \sum_{i=1}^{V} \sum_{j=1}^{V} T_{i,j} \times \Delta(i,j) \qquad (9)$$

subject to:

$$\forall i \sum_{j=1}^{V} T_{i,j} = d_{A,i}$$

$$\forall j \sum_{i=1}^{V} T_{i,j} = d_{B,j}$$

Where $T \in R^{V \times V}$ is a nonnegative matrix, where each $T_{i,j}$ denotes how much of sentence $i$ in document $A$ is assigned to sentences $j$ in document $B$, and constraints ensure the flow of a given sentence cannot exceed its allocated mass.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Greedy Mover's Distance

Solving the optimal transport problem is of cubic complexity and slow. Can it be approximated?

- find the two closest sentences and moves as much mass between the two sentences as possible
- the algorithm moves to the next two closest pairs
- terminates when all mass has been moved between the source and target document
- maintains mass constraints

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Approximation Performance

How does this approximation compare to the exact?

| Method | Kendall-Tau | Recall | MAE | Runtime (s) |
|--------|-------------|--------|-------|-------------|
| Exact-XLSMD | 1.00 | 0.69 | 0.000 | 0.402 |
| Relaxed-XLSMD | 0.70 | 0.58 | 0.084 | 0.031 |
| Greedy-XLSMD | 0.98 | 0.69 | 0.010 | 0.107 |

Table: Comparing exact XLSMD computation to approximation schemes for computing XLSMD on 10 webdomains.

# Approximate Distances



Distance Computations

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Approximate Performance

Which approximate computation works better?

| Approximation | Low | Mid | High | All |
|---------------|------|------|------|------|
| Relaxed-XLSMD | 0.44 | 0.43 | 0.50 | 0.46 |
| Greedy-XLSMD  | 0.54 | 0.50 | 0.56 | 0.54 |

Table: Document alignment performance of fast methods for approximating the same variant of XLSMD.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Alignment Results High-Resource

| Language | Recall | | | | | |
|---|---|---|---|---|---|---|
| | DE | SA | SMD | SL | IDF | SLIDF |
| French | 0.39 | 0.84 | 0.81 | 0.84 | 0.83 | **0.85** |
| Spanish | 0.34 | 0.53 | 0.59 | 0.63 | 0.62 | **0.64** |
| Russian | 0.06 | 0.64 | 0.69 | 0.69 | 0.70 | **0.71** |
| German | 0.52 | 0.74 | **0.78** | 0.76 | 0.77 | 0.77 |
| Italian | 0.22 | 0.47 | 0.55 | 0.56 | 0.56 | **0.59** |
| Portuguese | 0.17 | 0.36 | 0.39 | **0.41** | 0.38 | 0.40 |
| Dutch | 0.28 | 0.49 | 0.54 | 0.54 | 0.54 | **0.56** |
| Indonesian | 0.11 | 0.47 | 0.49 | 0.52 | 0.51 | **0.53** |
| Polish | 0.17 | 0.38 | 0.45 | 0.45 | **0.46** | 0.46 |
| Turkish | 0.12 | 0.38 | 0.52 | 0.56 | 0.57 | **0.59** |
| Swedish | 0.19 | 0.40 | 0.44 | 0.44 | **0.46** | 0.45 |
| Danish | 0.27 | 0.62 | 0.63 | **0.69** | 0.65 | **0.69** |
| Czech | 0.15 | 0.40 | 0.43 | **0.44** | **0.44** | 0.43 |
| Bulgarian | 0.07 | 0.43 | 0.52 | 0.54 | **0.55** | 0.52 |
| Finnish | 0.06 | 0.47 | 0.51 | 0.51 | **0.54** | 0.52 |
| Norwegian | 0.13 | 0.33 | 0.37 | 0.39 | **0.42** | 0.41 |
| **AVG** | 0.20 | 0.50 | 0.54 | 0.56 | 0.56 | **0.57** |

Cross-lingual Mining

A. El-Kishky, P. Koehn, H. Schwenk

Introduction

Corpora and WEB Crawling

Multilingual Represent.
LASER
Evaluation

Document Retrieval

Local Alignment

Global Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext Filtering

# Alignment Results Mid-Resource

| Language | DE | SA | SMD | SL | IDF | SLIDF |
|---|---|---|---|---|---|---|
| | | | **Recall** | | | |
| Romanian | 0.15 | 0.40 | 0.44 | 0.43 | **0.45** | 0.43 |
| Vietnamese | 0.06 | 0.28 | 0.29 | 0.29 | 0.29 | **0.32** |
| Ukrainian | 0.05 | 0.68 | 0.67 | 0.78 | 0.78 | **0.82** |
| Greek | 0.05 | 0.31 | 0.47 | 0.48 | **0.49** | **0.49** |
| Korean | 0.06 | 0.34 | 0.60 | 0.54 | **0.61** | 0.60 |
| Arabic | 0.04 | 0.32 | 0.63 | 0.59 | **0.65** | 0.61 |
| Croatian | 0.16 | 0.37 | 0.40 | 0.40 | **0.41** | 0.40 |
| Slovak | 0.20 | 0.41 | 0.46 | **0.46** | **0.46** | 0.44 |
| Thai | 0.02 | 0.19 | 0.41 | 0.33 | **0.47** | 0.41 |
| Hebrew | 0.05 | 0.18 | 0.39 | **0.43** | 0.41 | 0.41 |
| Hindi | 0.04 | 0.27 | 0.34 | **0.54** | 0.52 | 0.53 |
| Hungarian | 0.15 | 0.49 | 0.50 | **0.54** | 0.51 | **0.54** |
| Lithuanian | 0.11 | 0.73 | 0.79 | 0.79 | **0.80** | **0.80** |
| Slovenian | 0.13 | 0.33 | 0.34 | 0.35 | **0.36** | **0.36** |
| Persian | 0.06 | 0.32 | 0.56 | 0.57 | 0.53 | **0.59** |
| | | | | | | |
| **AVG** | 0.09 | 0.37 | 0.49 | 0.50 | **0.52** | 0.52 |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
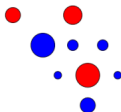Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment
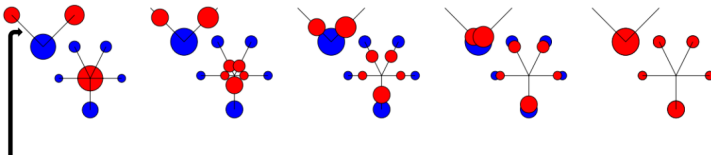
Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Alignment Results Low-Resource

| | Recall | | | | | |
|---|---|---|---|---|---|---|
| **Language** | DE | SA | SMD | SL | IDF | SLIDF |
| Estonian | 0.28 | 0.52 | 0.69 | 0.66 | **0.74** | 0.72 |
| Bengali | 0.05 | 0.32 | 0.78 | 0.72 | 0.77 | **0.79** |
| Albanian | 0.23 | 0.56 | **0.66** | 0.65 | 0.65 | **0.66** |
| Macedonian | 0.02 | 0.33 | 0.32 | 0.36 | **0.38** | 0.33 |
| Urdu | 0.06 | 0.22 | **0.60** | **0.60** | 0.49 | 0.56 |
| Serbian | 0.06 | 0.59 | **0.75** | 0.74 | 0.74 | 0.71 |
| Azerbaijani | 0.08 | 0.34 | 0.74 | 0.74 | **0.75** | 0.74 |
| Armenian | 0.02 | 0.18 | 0.32 | 0.35 | 0.34 | **0.38** |
| Belarusian | 0.07 | 0.47 | 0.67 | 0.69 | **0.73** | 0.71 |
| Georgian | 0.06 | 0.24 | 0.46 | **0.48** | 0.45 | 0.45 |
| Tamil | 0.02 | 0.20 | 0.51 | 0.45 | 0.51 | **0.53** |
| Marathi | 0.02 | 0.11 | 0.43 | **0.46** | 0.33 | 0.39 |
| Kazakh | 0.05 | 0.31 | 0.44 | **0.46** | 0.45 | 0.45 |
| Mongolian | 0.03 | 0.13 | 0.18 | 0.22 | 0.21 | **0.23** |
| Burmese | 0.01 | 0.10 | 0.26 | 0.33 | **0.46** | **0.46** |
| Bosnian | 0.18 | 0.64 | 0.61 | 0.69 | 0.65 | **0.72** |
| **AVG** | 0.08 | 0.33 | 0.53 | 0.54 | 0.54 | **0.55** |

WMT 2016 Shared Task

# WMT 2016 Shared Task

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# WMT 2016 Shared Task:
# Challenges

Big-ish websites

- E.g. cinedoc.org: 50k English, 50k French pages
- Makes 2.5B possible pairs
- Only allowed to pick 50k

Language detection unreliable

- Made sure test set can be found
- Some participants ran their own pipelines

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# WMT 2016 Shared Task:
# Challenges II

Near duplicates

- Removed pages when text was exactly the same
- www.taize.fr/fr article10921.html
- www.taize.fr/fr article10921.html?chooselang=1
- Almost identical

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Submissions

- 11 participating groups
- 19 submissions
- Up to 95% recall (NovaLincs-URL-Coverage)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# NovaLincs

- use a phrase table from a phrase-based statistical machine translation system to compute coverage scores
- based on the ratio of phrase pairs covered by a document pair.
- NOVALINCS-COVERAGE (88.6%)
- NOVALINCSCOVERAGE-URL (85.8%) coverage first then URL
- NOVALINCS-URL-COVERAGE (95.0%) URL first then coverage

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# YODA

- uses the machine translation of the French document, and finds the English corresponding document based on bigram and 5-gram matches, assisted by a heuristics based on document length ratio

- YODA: (93.9%)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# UEdin1

- uses cosine similarity between tf/idf weighted vectors, extracted by collecting n-grams from the English and machine translated French text. compare many hyperparameters such as weighting schemes and two pair selection algorithms.
- UEdin1: (89.1%)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Submission Results

| Name | Predicted pairs | Pairs after 1-1 rule | Found pairs | Recall % |
|------|-----------------|----------------------|-------------|----------|
| ADAPT | 61 094 | 61 094 | 644 | 26.8 |
| ADAPT–v2 | 69 518 | 69 518 | 651 | 27.1 |
| BadLuc | 681 610 | 263 133 | 1 905 | 79.3 |
| DOCAL | 191 993 | 191 993 | 2 128 | 88.6 |
| ILSP–ARC–pv42 | 291 749 | 287 860 | 2 040 | 84.9 |
| JIS | 323 929 | 28 903 | 48 | 2.0 |
| Medved | 155 891 | 155 891 | 1 907 | 79.4 |
| NovaLincs-coverage-url | 207 022 | 207 022 | 2 060 | 85.8 |
| NovaLincs-coverage | 235 763 | 235 763 | 2 129 | 88.6 |
| **NovaLincs-url-coverage** | 235 812 | 235 812 | 2 281 | 95.0 |
| UA PROMPSIT bitextor 4.1 | 95 760 | 95 760 | 748 | 31.1 |
| UA PROMPSIT bitextor 5.0 | 157 682 | 157 682 | 2 001 | 83.3 |
| UEdin1 cosine | 368 260 | 368 260 | 2 140 | 89.1 |
| UEdin2 LSI | 681 744 | 271 626 | 2 062 | 85.8 |
| UEdin2 LSI–v2 | 367 948 | 367 948 | 2 105 | 87.6 |
| UFAL–1 | 592 337 | 248 344 | 1 953 | 81.3 |
| UFAL–2 | 574 433 | 178 038 | 1 901 | 79.1 |
| UFAL–3 | 574 434 | 207 358 | 1 938 | 80.7 |
| UFAL–4 | 1 080 962 | 268 105 | 2 023 | 84.2 |
| YSDA | 277 896 | 277 896 | 2 021 | 84.1 |
| YODA | 318 568 | 318 568 | 2 256 | 93.9 |
| Baseline | 148 537 | 148 537 | 1 436 | 59.8 |

Figure: Documents are aligned 1-to-1 within each domain.

# Shared Task Insights

- Machine translated text helpful
- Finding matching n-grams works well
- Big boost by combination with URL-matching baseline

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# sentence alignment

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Sentence Alignment



"Local" alignment: limited to document pairs

- given: document pair
- output: matching sentence pairs

We also respect the order of sentences.

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Church and Gale (1993)



- Consider only the lengths of the sentences

$$\text{abs}(\log \frac{\text{length}_e}{\text{length}_f})$$

- Find the Viterbi path that with the best length ratios
- Additional cost factors for alignments other than 1-1

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Church and Gale (1993)



- Consider only the lengths of the sentences

$$\text{abs}(\log \frac{\text{length}_e}{\text{length}_f})$$

- Find the Viterbi path that with the best length ratios
- Additional cost factors for alignments other than 1-1

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Use of Dictionaries



The man looks intently at the window.
The sees a shadow.
It was in the trees.
What was it?
He is alarmed and awake.

Der Mann schaut aus dem Fenster.
Er sieht einen Schatten in den Bäumen.
Was war das?
Er war alamiert und wach.

- Given a word translation dictionary
  man = Mann; window = Fenster; shaddow = Schatten;
  trees = Bäumen; alarmed = alamiert
- Find matching word pairs
- Score sentence pairs based on number of matches
- **Hunalign:** (Varga et al., 2005) tool using this feature
- **Gargantua:** unsupervised induction of translation
  dictionary

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Translate and Match



- Use machine translation to translate foreign sentences
- Match translation with English
- Use of standard machine translation metric to assess match: BLEU score
- **Bleualign** (Sennrich and Volk, 2010)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Use of Sentence Embeddings

- Multilingual sentence embeddings, e.g., LASER
- Sentences with similar meaning have similar embedding
  — **independent** of language
- Comparison based on Cosine distance

$$
c(x, y) = \frac{(1 - \cos(x, y)) \ N(x) \ N(y)}{\sum\limits_{s=1}^{S} 1 - \cos(x, y_s) + \sum\limits_{s=1}^{S} 1 - \cos(x_s, y)}
$$

- **Vecalign** (Thompson and Koehn, 2019)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Efficient Algorithm



- Complexity of alignment via dynamic programming: $O(n^2)$
- Coarse to fine algorithm: $O(n)$
  embedding for block = average of sentence embeddings

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

## Evaluation: Text + Bild

| Algorithm | O( ) | P | R | $F_1$ |
|---|---|---|---|---|
| Gargantua | $N^2$ | 0.48 | 0.54 | 0.51 |
| Hunalign w/o lexicon | **N** | 0.59 | 0.70 | 0.64 |
| Hunalign w/ lexicon | **N** | 0.61 | 0.73 | 0.66 |
| Church and Gale | $N^2$ | 0.71 | 0.72 | 0.72 |
| Moore | ‡ | 0.86 | 0.71 | 0.78 |
| Bleualign | $N^2$ | 0.83 | 0.78 | 0.81 |
| Bleualign-NMT | $N^2$ | 0.85 | 0.83 | 0.84 |
| Coverage-Based | $N^2$ | 0.85 | 0.84 | 0.85 |
| Vecalign | **N** | **0.89** | **0.90** | **0.90** |

‡O( ) is data dependent

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Evaluation: Bible

| Languages | Verse-level $F_1$ | |
|---|---|---|
| | Vecalign | Hunalign |
| Afrikaans–Arabic | **0.863** | 0.339 |
| Afrikaans–Tagalog | **0.922** | 0.775 |
| Arabic–Norwegian | **0.787** | 0.406 |
| Arabic–Somali | **0.634** | 0.067 |
| Turkish–Somali | **0.533** | 0.331 |
| Norwegian–Somali | **0.697** | 0.687 |
| Somali–Afrikaans | **0.782** | 0.738 |
| Tagalog–Norwegian | **0.874** | 0.764 |
| Turkish–Afrikaans | **0.703** | 0.401 |
| Turkish–Tagalog | **0.647** | 0.247 |

# Evaluation: CommonCrawl

| Language Pair | LASER-only | Vecalign (best setup) |
|---|---|---|
| English–Portuguese | 31.5 | 32.9 (+1.4) |
| Portuguese–English | 36.0 | 38.8 (+2.8) |
| English–Bulgarian | 29.6 | 32.6 (+3.0) |
| Bulgarian–English | 20.6 | 22.3 (+1.7) |
| English–Estonian | 14.0 | 15.0 (+1.0) |
| English–Georgian | 8.6 | 9.1 (+0.5) |
| English–Urdu | 10.9 | 12.5 (+1.6) |
| English–Marathi | 10.0 | 10.3 (+0.3) |
| English–Burmese | 8.0 | 9.0 (+1.0) |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Evaluation: Paracrawl

- Task: align sentences in document pairs (subset of ParaCrawl data)

| Language | Web Domains | Document Pairs | English Tokens |
|----------|-------------|----------------|----------------|
| German | 21,806 | 17,109,018 | 10,788,923,009 |
| Czech | 12,179 | 6,661,650 | 4,089,806,440 |
| Hungarian | 5,560 | 2,770,432 | 1,504,698,348 |
| Estonian | 5,129 | 2,301,309 | 1,427,328,440 |
| Maltese | 933 | 303,198 | 134,232,546 |

Cross-lingual Mining

A. El-Kishky, P. Koehn, H. Schwenk

Introduction

Corpora and WEB Crawling

Multilingual Represent.
LASER
Evaluation

Document Retrieval

Local Alignment

Global Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext Filtering

# Evaluation: Paracrawl

- Results: BLEU scores for best subset (English token count)

| Language | Hunalign | Vecalign | Bleualign |
|----------|----------|----------|-----------|
| German | 35.1 (100m) | **35.8 (150m)** | 35.0 (100m) |
| Czech | 21.0 (50m) | **21.2 (50m)** | 21.0 (50m) |
| Hungarian | 16.5 (30m) | **16.8 (30m)** | 16.6 (15m) |
| Estonian | **21.8 (20m)** | 21.6 (20m) | 21.4 (20m) |
| Maltese | 33.5 (5m) | **34.1 (7m)** | 30.3 (2m) |

- Best results with Vecalign, except for Estonian

# global sentence alignment

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Bitext Mining in Wikipedia

## Some statistics

- 300 different languages
- Huge differences in size:
    - 1M+ articles: 15 languages
      (major European languages, ru, vi, ja, zh)
    - 100k+ articles: 47 languages
    - 10k+ articles: 81 languages
    - long tail . . .
- English by far the biggest (5.8M articles, 208M sentences)
- Cebuano has many articles produced by a bot
  (5.4M articles, 67M sentences)

# Bitext Mining in Wikipedia

## Local mining



- Only articles with link
- + Seems logical
- + Very fast
- – Ignored articles
- – Many simple sentences

A. El-Kishky,
P. Koehn,
H. Schwenk

# Bitext Mining in Wikipedia

## Local mining



## Global mining



- Only articles with link
+ Seems logical
+ Very fast
– Ignored articles
– Many simple sentences

- Always consider all sent.
– Increased complexity
± Lower recall ?
+ Generic: any corpus

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Bitext Mining in Wikipedia

## Global mining

- Compare **all** sentences of two Wikipedia
- Computationally more challenging: $134M \times 51M$ distances
- + Ability to handle two languages even though there are only few articles in common
- + Margin criterion:
  excludes short sentences which differ in NE only
- Potentially increased risk of misalignment and a lower recall

We chose global mining for this study (more generic)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Bitext Mining in Wikipedia

Processing pipeline

- Sentence splitting (very difficult for Thai)

- Deduplication

- Language identification with fasttext

$\rightarrow \approx 600M$ sentences for $> 180$ languages
  (each with more than 50k sentences)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Bitext Mining in Wikipedia

## Processing pipeline

- Sentence splitting (very difficult for Thai)

- Deduplication

- Language identification with fasttext

$\rightarrow \approx 600M$ sentences for $> 180$ languages
(each with more than 50k sentences)

## Complexity issues

- English/German Wikipedia:
    - 134M $\times$ 51M sentences
    - 513 + 204GB memory to store LASER embeddings
    - $6.8 \times 10^{15}$ distance calculations

$\Rightarrow$ Optimization and compression are needed !!

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Efficient Mining with FAISS

## FAISS library

- Library for efficient similarity search and clustering of dense vectors
- https://github.com/facebookresearch/faiss
- Mainly used for indexing images
  but can operate on any arbitrary vectors
- ⇒ Used here for efficient large-scale bitext mining
- Can be scaled to search in billions of sentences

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Efficient Mining with FAISS

FAISS index types

- Define $N$ Voronoi cells



- Quantizers:
  - PCA, not enough compression
  - Product: `OPQ64,IVF32768,PQ64`, 55x compression
  - Scalar: `PCAR128,IVF32768,SQ8`, 28x compression

# Efficient Mining with FAISS

FAISS index types

- Define $N$ Voronoi cells

- Quantizers:
  - PCA, not enough compression
  - Product: `OPQ64,IVF32768,PQ64`, 55x compression
  - Scalar: `PCAR128,IVF32768,SQ8`, 28x compression

- English FAISS index: 9.2GB

# Efficient Mining with FAISS

FAISS index types

- Define *N* Voronoi cells

- Quantizers:
  - PCA, not enough compression
  - Product: `OPQ64,IVF32768,PQ64`, 55x compression
  - Scalar: `PCAR128,IVF32768,SQ8`, 28x compression

- English FAISS index: 9.2GB

- English/German mining: 3h30 on 8 GPUS

# Efficient Mining with FAISS

Overall complexity

# Efficient Mining with FAISS

### Overall complexity

- Deduplication and LID
  - a couple of hours, run on parallel on standard server

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Efficient Mining with FAISS

## Overall complexity

- Deduplication and LID
  - a couple of hours, run on parallel on standard server
- Sentence embeddings with LASER (>7M sents/h)
  - total of 100h, can be run in parallel on cluster

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Efficient Mining with FAISS

## Overall complexity

- Deduplication and LID
    - a couple of hours, run on parallel on standard server
- Sentence embeddings with LASER ($>$7M sents/h)
    - total of 100h, can be run in parallel on cluster
- Train and create FAISS index for each language (CPU)
    - English $\approx$ 4h, total 21h

# Efficient Mining with FAISS

## Overall complexity

- Deduplication and LID
  - a couple of hours, run on parallel on standard server
- Sentence embeddings with LASER ($>$7M sents/h)
  - total of 100h, can be run in parallel on cluster
- Train and create FAISS index for each language (CPU)
  - English $\approx$ 4h, total 21h
- Mine bitext for each language pair
  - total of $\approx$ 1000h

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Efficient Mining with FAISS

## Overall complexity

- Deduplication and LID
  - a couple of hours, run on parallel on standard server
- Sentence embeddings with LASER ($>$7M sents/h)
  - total of 100h, can be run in parallel on cluster
- Train and create FAISS index for each language (CPU)
  - English $\approx$ 4h, total 21h
- Mine bitext for each language pair
  - total of $\approx$ 1000h
- $\Rightarrow$ **Total of 43 days on one GPU**

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Efficient Mining with FAISS

## Overall complexity

- Deduplication and LID
  - a couple of hours, run on parallel on standard server
- Sentence embeddings with LASER (>7M sents/h)
  - total of 100h, can be run in parallel on cluster
- Train and create FAISS index for each language (CPU)
  - English $\approx$ 4h, total 21h
- Mine bitext for each language pair
  - total of $\approx$ 1000h
- $\Rightarrow$ **Total of 43 days on one GPU**

  or much less on many GPUs . . .

Cross-lingual Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Threshold Optimization

## Precision/recall trade-off

- Margin-based mining has only one parameter:
  the margin between the closest and the average distance
    - large margin: high precision, low recall
    - small margin: lower precision, higher recall
- We have no gold-alignments to optimize this parameter

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Threshold Optimization

## Precision/recall trade-off

- Margin-based mining has only one parameter:
  the margin between the closest and the average distance
    - large margin: high precision, low recall
    - small margin: lower precision, higher recall
- We have no gold-alignments to optimize this parameter

## Task oriented threshold optimization

- Mine bitexts for thresholds in range [1.01–1.06]
- Train NMT systems for increasing amounts of data
- Evaluate each one and keep best one

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Threshold Optimization on Europarl

## BLEU score for NMT trained on Wikipedia only



Precision/recall trade-off

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Threshold Optimization on Europarl

## BLEU score for NMT trained on Wikipedia only



## Precision/recall trade-off

- Threshold on margin of 1.04 best for most conditions

# Threshold Optimization on Europarl

## BLEU score for NMT trained on Wikipedia only



## Precision/recall trade-off

- Threshold on margin of 1.04 best for most conditions

| Bitexts | de-en | de-fr | cs-de | cs-fr |
|---|---|---|---|---|
| Mined | 1.0M | 372k | 201k | 219k |
| Wikipedia | 24.4 | 22.7 | 13.1 | 16.3 |
| Europarl | 1.0M | 370k | 200k | 220k |
| | 21.2 | 21.1 | 12.6 | 19.2 |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Threshold Optimization on Europarl

## BLEU score for NMT trained on Wikipedia only



## Precision/recall trade-off

- Threshold on margin of
  1.04 best for most
  conditions

- WikiMatrix bitexts
  outperform Europarl

| Bitexts | de-en | de-fr | cs-de | cs-fr |
|---|---|---|---|---|
| Mined | 1.0M | 372k | 201k | 219k |
| Wikipedia | 24.4 | 22.7 | 13.1 | 16.3 |
| Europarl | 1.0M | 370k | 200k | 220k |
| | 21.2 | 21.1 | 12.6 | 19.2 |

# WikiMatrix: 85 Languages, 1620 Pairs

Introduction

Corpora and
WEB Crawl

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment

WikiMatrix

CCMatrix

WMT/TED

Bitext
Filtering

- With English: Indonesian 1M,
  Hebrew 545k, Farsi 303k or Marathi 124k

# WikiMatrix: 85 Languages, 1620 Pairs

Introduction

Corpora and
WEB Crawl

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment

WikiMatrix

CCMatrix

WMT/TED

Bitext
Filtering

| Code | Name | Language Family |
|---|---|---|
| ar | Arabic | Semitic |
| az | Azerbaijani | Turkic |
| ba | Bashkir | Turkic |
| be | Belarusian | Slavic |
| bg | Bulgarian | Slavic |
| bn | Bengali | Indo-Aryan |
| bs | Bosnian | Slavic |
| ca | Catalan | Romance |
| cs | Czech | Slavic |
| da | Danish | Germanic |
| de | German | Germanic |
| el | Greek | Hellenic |
| en | English | Germanic |
| eo | Esperanto | constructed |
| es | Spanish | Romance |
| et | Estonian | Uralic |
| eu | Basque | Isolate |
| fa | Farsi | Iranian |
| fi | Finnish | Uralic |
| fr | French | Romance |
| gl | Galician | Romance |
| he | Hebrew | Semitic |
| hi | Hindi | Indo-Aryan |
| hr | Croatian | Slavic |
| hu | Hungarian | Uralic |
| id | Indonesian | Malayo-Polynesian |
| is | Icelandic | Germanic |
| it | Italian | Romance |
| ja | Japanese | Japonic |
| kk | Kazakh | Turkic |
| ko | Korean | Korean |
| lt | Lithuanian | Baltic |
| mk | Macedonian | Slavic |
| ml | Malayalam | Dravidian |
| mr | Marathi | Indo-Aryan |
| ne | Nepali | Indo-Aryan |
| nl | Dutch | Germanic |
| no | Norwegian | Germanic |
| oc | Occitan | Romance |
| pl | Polish | Slavic |
| pt | Portuguese | Romance |
| ro | Romanian | Romance |
| ru | Russian | Slavic |
| sh | Serbo-Croatian | South Slavic |
| si | Sinhala | Indo-Aryan |
| sk | Slovak | Slavic |
| sl | Slovenian | Slavic |
| sq | Albanian | Albanian |
| sr | Serbian | Slavic |
| sv | Swedish | Germanic |
| sw | Swahili | Niger-Congo |
| ta | Tamil | Dravidian |
| te | Telugu | Dravidian |
| tg | Tagalog | Malayo-Polynesian |
| tr | Turkish | Turkic |
| tt | Tatar | Turkic |
| uk | Ukrainian | Slavic |
| vi | Vietnamese | Vietic |
| zh | Chinese | Chinese |

- Russian/Ukrainian 2.5M, Catalan/Spanish 1.6M

# WikiMatrix: 85 Languages, 1620 Pairs

Introduction

Corpora and
WEB Crawl

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment

WikiMatrix

CCMatrix

WMT/TED

Bitext
Filtering

- Between Romance languages fr, es, it and pt 480k–923k

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawl

Multilingual
Represent.

LASER

Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment

WikiMatrix

CCMatrix

WMT/TED

Bitext
Filtering

# WikiMatrix: 85 Languages, 1620Pairs

| | Name | Language Family | size |
|---|---|---|---|
| ar | Arabic | Semitic | 6516 |
| az | Azerbaijani | Turkic | 1873 |
| ba | Bashkir | Turkic | 536 |
| be | Belarusian | Slavic | 1690 |
| bg | Bulgarian | Slavic | 3327 |
| bn | Bengali | Indo-Aryan | 1412 |
| bs | Bosnian | Slavic | 1060 |
| ca | Catalan | Romance | 8332 |
| cs | Czech | Slavic | 7434 |
| da | Danish | Germanic | 2681 |
| de | German | Germanic | 50944 |
| el | Greek | Hellenic | 3211 |
| en | English | Germanic | 134431 |
| eo | Esperanto | constructed | 2371 |
| es | Spanish | Romance | 25202 |
| et | Estonian | Uralic | 2303 |
| eu | Basque | Isolate | 2259 |
| fa | Farsi | Iranian | 3954 |
| fi | Finnish | Uralic | 7428 |
| fr | French | Romance | 34494 |
| gl | Galician | Romance | 2221 |
| he | Hebrew | Semitic | 6962 |
| hi | Hindi | Indo-Aryan | 1353 |
| hr | Croatian | Slavic | 2229 |
| hu | Hungarian | Uralic | 7702 |
| id | Indonesian | Malayo-Polynesian | 3899 |
| is | Icelandic | Germanic | 488 |
| it | Italian | Romance | 21025 |
| ja | Japanese | Japonic | 13614 |
| kk | Kazakh | Turkic | 1684 |
| ko | Korean | Koreanic | 5400 |
| lt | Lithuanian | Baltic | |
| mk | Macedonian | Slavic | 1487 |
| ml | Malayalam | Dravidian | 272 |
| mr | Marathi | Indo-Aryan | 503 |
| ne | Nepali | Indo-Aryan | 235 |
| nl | Dutch | Germanic | 14091 |
| no | Norwegian | Germanic | 1364 |
| oc | Occitan | Romance | 472 |
| pl | Polish | Slavic | 16270 |
| pt | Portuguese | Romance | 13354 |
| ro | Romanian | Romance | 6004 |
| ru | Russian | Slavic | 32537 |
| sh | Serbo-Croatian | South Slavic | 2069 |
| si | Sinhala | Indo-Aryan | 320 |
| sk | Slovak | Slavic | 1904 |
| sl | Slovenian | Slavic | 3838 |
| sq | Albanian | Albanian | 740 |
| sr | Serbian | Slavic | 4226 |
| sv | Swedish | Germanic | 17906 |
| sw | Swahili | Niger-Congo | 235 |
| ta | Tamil | Dravidian | 1629 |
| te | Telugu | Dravidian | 1509 |
| tg | Tagalog | Malayo-Polynesian | 312 |
| tr | Turkish | Turkic | 4067 |
| tt | Tatar | Turkic | 449 |
| uk | Ukrainian | Slavic | 11585 |
| vi | Vietnamese | Vietic | 6727 |
| zh | Chinese | Chinese | 12308 |

- Japanese/Korean 222k, Japanese/Russian 196k,
  Indonesian/Vietnamese 146k, Hebrew/fr,es,it,ru 120–150k

# Large-Scale Bitext Mining

Scaling up !

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Large-Scale Bitext Mining

## Scaling up !

- Can we apply the same global mining approach to a much bigger corpus ?

# Large-Scale Bitext Mining

## Scaling up !

- Can we apply the same global mining approach to a much bigger corpus ?

- 10 snapshot of curated common crawl corpus (Wenzek et al, arxiv'19)

- 36 billion unique sentences (50× bigger than Wikipedia)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Large-Scale Bitext Mining

Scaling up !

- Can we apply the same global mining approach to a much bigger corpus ?

- 10 snapshot of curated common crawl corpus (Wenzek et al, arxiv'19)

- 36 billion unique sentences (50× bigger than Wikipedia)

⇒ Substantial computational and storage challenges
  - Mining Russian against Japanese: $3 \times 2.9$ billion sentences
  - $\approx 8.7 \cdot 10^{18}$ distances (6 months on 8 GPUs)
  - optimized and highly parallelized processing

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
**CCMatrix**
WMT/TED

Bitext
Filtering

# Mining in the Whole Internet

## CCMatrix

- 36 billion sentences collected on the Internet in 39 languages
- $\Rightarrow$ More than 4.5 billion parallel sentences in 39 languages

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Mining in the Whole Internet

## CCMatrix

- 36 billion sentences collected on the Internet in 39 languages

$\Rightarrow$ More than 4.5 billion parallel sentences in 39 languages



$\Rightarrow$ By far the largest collection of high quality mined bitexts

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Mining in the Whole Internet

## CCMatrix

- 36 billion sentences collected on the Internet in 39 languages
- ⇒ More than 4.5 billion parallel sentences in 39 languages



⇒ By far the largest collection of high quality mined bitexts

- Expected to cover many topics: politics, sports, tourism, daily life, . . .

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Complexity Optimization

## Example: mining French/English



- Split monolingual texts into many parts
- Calculate forward and backward distances in parallel
- ⇒ Extract bitexts when all distances are available

Cross-lingual Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and WEB Crawling

Multilingual Represent.
    LASER
    Evaluation

Document Retrieval

Local Alignment

Global Alignment
    WikiMatrix
    CCMatrix
    WMT/TED

Bitext Filtering

# CCMatrix

| ISO | Name | Family | Size | bg | cs | da | de | el | en | es | fa | fi | fr | he | hi | hu | id | it | ja | ko | ms | nl | no | pl | pt | ru | sv | tr | uk | vi | zh | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ar | Arabic | Arabic | 196 | 3.0 | 3.9 | 2.7 | 7.5 | 3.3 | 6.5 | 10.0 | 3.1 | 2.7 | 23.8 | 2.2 | 1.4 | 2.7 | 4.1 | 5.8 | 5.0 | 2.5 | 1.5 | 5.1 | 2.5 | 4.5 | 6.7 | 9.2 | 5.6 | 5.5 | 1.5 | 4.2 | 5.4 | 141.7 |
| bg | Bulgarian | Slavic | 68 | - | 6.1 | 3.7 | 9.9 | 4.3 | 3.7 | 10.7 | 2.3 | 3.6 | 11.4 | 2.1 | 1.5 | 3.8 | 3.8 | 7.4 | 5.7 | 2.8 | 1.3 | 6.9 | 3.0 | 7.2 | 7.5 | 17.4 | 7.6 | 5.8 | 2.3 | 4.4 | 5.0 | 154.1 |
| cs | Czech | Slavic | 303 | - | - | 5.9 | 18.3 | 5.4 | 9.8 | 15.5 | 2.9 | 6.1 | 17.3 | 3.1 | 2.0 | 6.1 | 5.3 | 11.2 | 8.0 | 4.0 | 2.0 | 11.6 | 4.9 | 13.2 | 10.7 | 18.1 | 12.9 | 8.6 | 2.6 | 6.0 | 7.0 | 228.7 |
| da | Danish | Germanic | 109 | - | - | - | 12.6 | 3.8 | 4.5 | 10.2 | 2.0 | 4.8 | 12.0 | 2.3 | 1.5 | 3.7 | 3.9 | 7.3 | 5.6 | 2.9 | 1.4 | 9.5 | 9.6 | 6.5 | 7.4 | 9.2 | 15.2 | 5.7 | 1.5 | 4.2 | 4.9 | 164.6 |
| de | German | Germanic | 1728 | - | - | - | - | 9.8 | 67.3 | 38.8 | 4.8 | 11.3 | 50.0 | 5.6 | 3.2 | 11.0 | 9.6 | 29.5 | 11.6 | 6.2 | 3.5 | 33.2 | 10.4 | 20.5 | 23.4 | 29.3 | 29.3 | 15.5 | 3.4 | 9.7 | 11.8 | 497.5 |
| el | Greek | Hellenic | 144 | - | - | - | - | - | 5.6 | 12.2 | 2.2 | 3.6 | 12.9 | 2.3 | 1.4 | 3.7 | 3.7 | 8.5 | 5.2 | 2.6 | 1.4 | 6.9 | 3.0 | 6.2 | 8.4 | 9.9 | 7.3 | 5.6 | 1.7 | 4.2 | 4.7 | 150.1 |
| en | English | Germanic | 8677 | - | - | - | - | - | - | 86.3 | 2.5 | 4.1 | 94.1 | 1.5 | 0.7 | 3.6 | 13.4 | 31.3 | 33.7 | 7.2 | 0.8 | 23.8 | 3.8 | 16.0 | 33.1 | 72.4 | 43.8 | 26.8 | 1.6 | 18.5 | 17.6 | 634.2 |
| es | Spanish | Romance | 1534 | - | - | - | - | - | - | - | 5.5 | 9.7 | 70.9 | 5.9 | 3.2 | 9.5 | 12.4 | 44.3 | 11.6 | 6.2 | - | 23.3 | 8.8 | 19.6 | 59.4 | 32.4 | 22.3 | 15.2 | 4.0 | 11.9 | 13.2 | 573.1 |
| fa | Farsi | Iranian | 192 | - | - | - | - | - | - | - | - | 2.0 | 5.5 | 1.7 | 1.2 | 1.9 | 3.1 | 3.6 | 3.5 | 2.0 | 1.3 | 3.6 | 1.9 | 3.2 | 4.1 | 5.6 | 4.0 | 4.9 | 1.1 | 3.3 | 3.4 | 86.3 |
| fi | Finnish | Uralic | 132 | - | - | - | - | - | - | - | - | - | 11.1 | 2.2 | 1.4 | 4.2 | 3.8 | 7.1 | 6.2 | 3.0 | 1.4 | 8.1 | 4.1 | 6.8 | 7.1 | 9.9 | 13.8 | 6.2 | 1.7 | 4.4 | 5.2 | 155.8 |
| fr | French | Romance | 1869 | - | - | - | - | - | - | - | - | - | - | 6.8 | 3.5 | 10.3 | 11.9 | 46.2 | 12.6 | 6.9 | 4.2 | 32.1 | 9.9 | 21.1 | 37.9 | 31.9 | 27.6 | 17.4 | 4.2 | 12.5 | 14.0 | 619.8 |
| he | Hebrew | Semitic | 70 | - | - | - | - | - | - | - | - | - | - | - | 1.2 | 1.9 | 2.8 | 4.0 | 5.3 | 2.5 | 1.1 | 4.2 | 2.0 | 3.6 | 4.3 | 6.4 | 5.1 | 4.4 | 1.2 | 3.6 | 3.6 | 92.9 |
| hi | Hindi | Indo-Aryan | 48 | - | - | - | - | - | - | - | - | - | - | - | - | 1.3 | 1.9 | 2.3 | 2.7 | 1.6 | 0.9 | 2.4 | 1.4 | 2.1 | 2.6 | 3.4 | 3.0 | 3.2 | 0.8 | 1.9 | 2.4 | 56.0 |
| hu | Hungarian | Uralic | 148 | - | - | - | - | - | - | - | - | - | - | - | - | - | 3.2 | 7.0 | 5.2 | 2.6 | 1.3 | 7.1 | 3.0 | 7.1 | 6.8 | 9.6 | 7.4 | 5.6 | 1.7 | 3.7 | 4.6 | 139.6 |
| id | Indonesian | Malayo-Polynesian | 366 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 7.4 | 5.9 | 3.5 | 4.4 | 7.6 | 3.7 | 6.0 | 9.1 | 9.9 | 8.6 | 8.1 | 1.7 | 7.9 | 6.3 | 172.9 |
| it | Italian | Romance | 686 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 8.9 | 4.7 | 2.5 | 16.6 | 6.1 | 14.7 | 25.4 | 20.5 | 16.0 | 10.5 | 2.8 | 8.0 | 8.6 | 368.4 |
| ja | Japanese | Japonic | 2944 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 3.3 | | 8.9 | 5.1 | 7.7 | 9.1 | 11.6 | 11.3 | 12.1 | 2.8 | 6.5 | 13.5 | 228.7 |
| ko | Korean | Koreanic | 778 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 1.9 | 4.8 | 2.6 | 4.0 | 4.9 | 6.0 | 7.1 | 8.4 | 1.4 | 5.2 | 6.3 | 113.7 |
| ms | Malay | Malayo-Polynesian | 25 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.6 | 1.3 | 2.3 | 2.8 | 3.7 | 3.6 | 3.4 | 0.8 | 3.2 | 2.8 | 60.8 |
| nl | Dutch | Germanic | 510 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 7.8 | 12.9 | 15.5 | 17.7 | 20.8 | 11.0 | 2.7 | 7.2 | 8.4 | 322.2 |
| no | Norwegian | Germanic | 109 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 5.5 | 6.4 | 8.1 | 13.8 | 5.2 | 1.4 | 3.9 | 4.3 | 143.8 |
| pl | Polish | Slavic | 505 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 13.5 | 22.9 | 13.8 | 9.1 | 3.4 | 6.5 | 7.1 | 267.1 |
| pt | Portuguese | Romance | 729 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 20.9 | 15.7 | 11.0 | 3.0 | 8.8 | 9.5 | 375.2 |
| ru | Russian | Slavic | 3047 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 18.9 | 15.3 | 31.2 | 10.4 | 13.0 | 475.0 |
| sv | Swedish | Germanic | 1200 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.8 | 10.6 | 10.4 | | 358.5 |
| tr | Turkish | Turkic | 1382 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.5 | 10.4 | 10.0 | 247.4 |
| uk | Ukrainian | Slavic | 110 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 0.2 | 2.2 | 88.6 |
| vi | Vietnamese | Vietic | 1172 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 9.1 | 179.2 |
| zh | Chinese | Chinese | 2512 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 214.3 |

Table 1: CCMatrix: size of mined sentences **(in millions)** for each language pair.

- French/Spanish: 71M

Cross-lingual Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# CCMatrix

| ISO | Name | Family | Size | bg | cs | da | de | el | en | es | fa | fi | fr | he | hi | hu | id | it | ja | ko | ms | nl | no | pl | pt | ru | sv | tr | uk | vi | zh | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ar | Arabic | Arabic | 196 | 3.0 | 3.9 | 2.7 | 7.5 | 3.3 | 6.5 | 10.0 | 3.1 | 2.7 | 23.8 | 2.2 | 1.4 | 2.7 | 4.1 | 5.8 | 5.0 | 2.5 | 1.5 | 5.1 | 2.5 | 4.5 | 6.7 | 9.2 | 5.6 | 5.5 | 1.5 | 4.2 | 5.4 | 141.7 |
| bg | Bulgarian | Slavic | 68 | - | 6.1 | 3.7 | 9.9 | 4.3 | 3.7 | 10.7 | 2.3 | 3.6 | 11.4 | 2.1 | 1.5 | 3.8 | 3.8 | 7.4 | 5.7 | 2.8 | 1.3 | 6.9 | 3.0 | 7.2 | 7.5 | 17.4 | 7.6 | 5.8 | 2.3 | 4.4 | 5.0 | 154.1 |
| cs | Czech | Slavic | 303 | - | - | 5.9 | 18.3 | 5.4 | 9.8 | 15.5 | 2.9 | 6.1 | 17.3 | 3.1 | 2.0 | 6.1 | 5.3 | 11.2 | 8.0 | 4.0 | 2.0 | 11.6 | 4.9 | 13.2 | 10.7 | 18.1 | 12.9 | 8.6 | 2.6 | 6.0 | 7.0 | 228.7 |
| da | Danish | Germanic | 109 | - | - | - | 12.6 | 3.8 | 4.5 | 10.2 | 2.0 | 4.8 | 12.0 | 2.3 | 1.5 | 3.7 | 3.9 | 7.3 | 5.6 | 2.9 | 1.4 | 9.5 | 9.6 | 6.5 | 7.4 | 9.2 | 15.2 | 5.7 | 1.5 | 4.2 | 4.9 | 164.6 |
| de | German | Germanic | 1728 | - | - | - | - | 9.8 | 67.3 | 38.8 | 4.8 | 11.3 | 50.0 | 5.6 | 3.2 | 11.0 | 9.6 | 29.5 | 11.6 | 6.2 | 3.5 | 33.2 | 10.4 | 20.5 | 23.4 | 29.3 | 29.3 | 15.5 | 3.8 | 9.7 | 11.8 | 497.5 |
| el | Greek | Hellenic | 144 | - | - | - | - | - | 5.6 | 12.2 | 2.2 | 3.6 | 12.9 | 2.3 | 1.4 | 3.7 | 3.7 | 8.5 | 5.2 | 2.6 | 1.4 | 6.9 | 3.0 | 6.2 | 8.4 | 9.9 | 7.3 | 5.6 | 1.7 | 4.2 | 4.7 | 150.1 |
| en | English | Germanic | 8677 | - | - | - | - | - | - | 86.3 | 2.5 | 4.1 | 94.1 | 1.5 | 0.7 | 3.6 | 13.4 | 31.3 | 33.7 | 7.2 | 0.8 | 23.8 | 3.8 | 16.0 | 33.1 | 72.4 | 43.8 | 26.8 | 1.6 | 18.5 | 17.6 | 634.2 |
| es | Spanish | Romance | 1534 | - | - | - | - | - | - | - | 5.5 | 9.7 | 70.9 | 5.9 | 3.2 | 9.5 | 12.4 | 44.3 | 11.6 | 6.2 | - | 23.3 | 8.8 | 19.6 | 59.4 | 32.4 | 22.3 | 15.2 | 4.0 | 11.9 | 13.2 | 573.1 |
| fa | Farsi | Iranian | 192 | - | - | - | - | - | - | - | - | 2.0 | 5.5 | 1.7 | 1.2 | 1.9 | 3.1 | 3.6 | 3.5 | 2.0 | 1.3 | 3.6 | 1.9 | 3.2 | 4.1 | 5.6 | 4.0 | 4.9 | 1.1 | 3.3 | 3.4 | 86.3 |
| fi | Finnish | Uralic | 132 | - | - | - | - | - | - | - | - | - | 11.1 | 2.2 | 1.4 | 4.2 | 3.8 | 7.1 | 6.2 | 3.0 | 1.4 | 8.1 | 4.1 | 6.8 | 7.1 | 9.9 | 13.8 | 6.2 | 1.7 | 4.4 | 5.2 | 155.8 |
| fr | French | Romance | 1869 | - | - | - | - | - | - | - | - | - | - | 6.8 | 3.5 | 10.3 | 11.9 | 46.2 | 12.6 | 6.9 | 4.2 | 32.1 | 9.9 | 21.1 | 37.9 | 31.9 | 27.6 | 17.4 | 4.2 | 12.5 | 14.0 | 619.8 |
| he | Hebrew | Semitic | 70 | - | - | - | - | - | - | - | - | - | - | - | 1.2 | 1.9 | 2.8 | 4.0 | 5.3 | 2.5 | 1.1 | 4.2 | 2.0 | 3.6 | 4.3 | 6.4 | 5.1 | 4.4 | 1.2 | 3.6 | 3.6 | 92.9 |
| hi | Hindi | Indo-Aryan | 48 | - | - | - | - | - | - | - | - | - | - | - | - | 1.3 | 1.9 | 2.3 | 2.7 | 1.6 | 0.9 | 2.4 | 1.4 | 2.1 | 2.6 | 3.4 | 3.0 | 3.2 | 0.8 | 1.9 | 2.4 | 56.0 |
| hu | Hungarian | Uralic | 148 | - | - | - | - | - | - | - | - | - | - | - | - | - | 3.2 | 7.0 | 5.2 | 2.6 | 1.3 | 7.1 | 3.0 | 7.1 | 6.8 | 9.6 | 7.4 | 5.6 | 1.7 | 3.7 | 4.6 | 139.6 |
| id | Indonesian | Malayo-Polynesian | 366 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 7.4 | 5.9 | 3.5 | 4.4 | 7.6 | 3.7 | 6.0 | 9.1 | 9.9 | 8.6 | 8.1 | 1.7 | 7.9 | 6.3 | 172.9 |
| it | Italian | Romance | 686 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 8.9 | 4.7 | 2.5 | 16.6 | 6.1 | 14.7 | 25.4 | 20.5 | 16.0 | 10.5 | 2.8 | 8.0 | 8.6 | 368.4 |
| ja | Japanese | Japonic | 2944 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 3.3 | 8.9 | 5.1 | 7.7 | 9.1 | 11.6 | 11.3 | 12.1 | 2.8 | 6.5 | 13.5 | 228.7 |
| ko | Korean | Koreanic | 778 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 1.9 | 2.6 | 4.0 | 4.9 | 6.0 | 7.1 | 8.4 | 1.4 | 5.2 | 6.3 | 113.7 |
| ms | Malay | Malayo-Polynesian | 25 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.6 | 1.3 | 2.3 | 3.7 | 3.6 | 3.4 | 0.8 | 3.2 | 2.8 | 60.8 |
| nl | Dutch | Germanic | 510 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 7.8 | 12.9 | 15.5 | 17.7 | 20.8 | 11.0 | 2.7 | 7.2 | 8.4 | 322.2 |
| no | Norwegian | Germanic | 109 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 5.5 | 6.4 | 8.1 | 13.8 | 5.2 | 1.4 | 3.9 | 4.3 | 143.8 |
| pl | Polish | Slavic | 505 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 13.5 | 22.9 | 13.8 | 9.1 | 3.4 | 6.5 | 7.1 | 267.1 |
| pt | Portuguese | Romance | 729 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 20.9 | 15.7 | 11.0 | 3.0 | 8.8 | 9.5 | 375.2 |
| ru | Russian | Slavic | 3047 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 18.9 | 15.3 | 31.2 | 10.4 | 13.0 | 475.0 |
| sv | Swedish | Germanic | 1200 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.8 | 10.6 | 10.4 | 358.5 |
| tr | Turkish | Turkic | 1382 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.5 | 10.4 | 10.0 | 247.4 |
| uk | Ukrainian | Slavic | 110 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 0.2 | 2.2 | 88.6 |
| vi | Vietnamese | Vietic | 1172 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 9.1 | 171.9 |
| zh | Chinese | Chinese | 2512 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 214.3 |

Table 1: CCMatrix: size of mined sentences **(in millions)** for each language pair.

- Norwegian/Swedish: 14M

Cross-lingual Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# CCMatrix

| ISO | Name | Family | Size | bg | cs | da | de | el | en | es | fa | fi | fr | he | hi | hu | id | it | ja | ko | ms | nl | no | pl | pt | ru | sv | tr | uk | vi | zh | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ar | Arabic | Arabic | 196 | 3.0 | 3.9 | 2.7 | 7.5 | 3.3 | 6.5 | 10.0 | 3.1 | 2.7 | 23.8 | 2.2 | 1.4 | 2.7 | 4.1 | 5.8 | 5.0 | 2.5 | 1.5 | 5.1 | 2.5 | 4.5 | 6.7 | 9.2 | 5.6 | 5.5 | 1.5 | 4.2 | 5.4 | 141.7 |
| bg | Bulgarian | Slavic | 68 | - | 6.1 | 3.7 | 9.9 | 4.3 | 3.7 | 10.7 | 2.3 | 3.6 | 11.4 | 2.1 | 1.5 | 3.8 | 3.8 | 7.4 | 5.7 | 2.8 | 1.3 | 6.9 | 3.0 | 7.2 | 7.5 | 17.4 | 7.6 | 5.8 | 2.3 | 4.4 | 5.0 | 154.1 |
| cs | Czech | Slavic | 303 | - | - | 5.9 | 18.3 | 5.4 | 9.8 | 15.5 | 2.9 | 6.1 | 17.3 | 3.1 | 2.0 | 6.1 | 5.3 | 11.2 | 8.0 | 4.0 | 2.0 | 11.6 | 4.9 | 13.2 | 10.7 | 18.1 | 12.9 | 8.6 | 2.6 | 6.0 | 7.0 | 228.7 |
| da | Danish | Germanic | 109 | - | - | - | 12.6 | 3.8 | 4.5 | 10.2 | 2.0 | 4.8 | 12.0 | 2.3 | 1.5 | 3.7 | 3.9 | 7.3 | 5.6 | 2.9 | 1.4 | 9.5 | 9.6 | 6.5 | 7.4 | 9.2 | 15.2 | 5.7 | 1.5 | 4.2 | 4.9 | 164.6 |
| de | German | Germanic | 1728 | - | - | - | - | 9.8 | 67.3 | 38.8 | 4.8 | 11.3 | 50.0 | 5.6 | 3.2 | 11.0 | 9.6 | 29.5 | 11.6 | 6.2 | 3.5 | 33.2 | 10.4 | 20.5 | 23.4 | 29.3 | 29.3 | 15.5 | 3.8 | 9.7 | 11.8 | 497.5 |
| el | Greek | Hellenic | 144 | - | - | - | - | - | 5.6 | 12.2 | 2.2 | 3.6 | 12.9 | 2.3 | 1.4 | 3.7 | 3.7 | 8.5 | 5.2 | 2.6 | 1.4 | 6.9 | 3.0 | 6.2 | 8.4 | 9.9 | 7.3 | 5.6 | 1.7 | 4.2 | 4.7 | 150.1 |
| en | English | Germanic | 8677 | - | - | - | - | - | - | 86.3 | 2.5 | 4.1 | 94.1 | 1.5 | 0.7 | 3.6 | 13.4 | 31.3 | 33.7 | 7.2 | 0.8 | 23.8 | 3.8 | 16.0 | 33.1 | 72.4 | 43.8 | 26.8 | 1.6 | 18.5 | 17.6 | 634.2 |
| es | Spanish | Romance | 1534 | - | - | - | - | - | - | - | 5.5 | 9.7 | 70.9 | 5.9 | 3.2 | 9.5 | 12.4 | 44.3 | 11.6 | 6.2 | - | 23.3 | 8.8 | 19.6 | 59.4 | 32.4 | 22.3 | 15.2 | 4.0 | 11.9 | 13.2 | 573.1 |
| fa | Farsi | Iranian | 192 | - | - | - | - | - | - | - | - | 2.0 | 5.5 | 1.7 | 1.2 | 1.9 | 3.1 | 3.6 | 3.5 | 2.0 | 1.3 | 3.6 | 1.9 | 3.2 | 4.1 | 5.6 | 4.0 | 4.9 | 1.1 | 3.3 | 3.4 | 86.3 |
| fi | Finnish | Uralic | 132 | - | - | - | - | - | - | - | - | - | 11.1 | 2.2 | 1.4 | 4.2 | 3.8 | 7.1 | 6.2 | 3.0 | 1.4 | 8.1 | 4.1 | 6.8 | 7.1 | 9.9 | 13.8 | 6.2 | 1.7 | 4.4 | 5.2 | 155.8 |
| fr | French | Romance | 1869 | - | - | - | - | - | - | - | - | - | - | 6.8 | 3.5 | 10.3 | 11.9 | 46.2 | 12.6 | 6.9 | 4.2 | 32.1 | 9.9 | 21.1 | 37.9 | 31.9 | 27.6 | 17.4 | 4.2 | 12.5 | 14.0 | 619.8 |
| he | Hebrew | Semitic | 70 | - | - | - | - | - | - | - | - | - | - | - | 1.2 | 1.9 | 2.8 | 4.0 | 5.3 | 2.5 | 1.1 | 4.2 | 2.0 | 3.6 | 4.3 | 6.4 | 5.1 | 4.4 | 1.2 | 3.6 | 3.6 | 92.9 |
| hi | Hindi | Indo-Aryan | 48 | - | - | - | - | - | - | - | - | - | - | - | - | 1.3 | 1.9 | 2.3 | 2.7 | 1.6 | 0.9 | 2.4 | 1.4 | 2.1 | 2.6 | 3.4 | 3.0 | 3.2 | 0.8 | 1.9 | 2.4 | 56.0 |
| hu | Hungarian | Uralic | 148 | - | - | - | - | - | - | - | - | - | - | - | - | - | 3.2 | 7.0 | 5.2 | 2.6 | 1.3 | 7.1 | 3.0 | 7.1 | 6.8 | 9.6 | 7.4 | 5.6 | 1.7 | 3.7 | 4.6 | 139.6 |
| id | Indonesian | Malayo-Polynesian | 366 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 7.4 | 5.9 | 3.5 | 4.4 | 7.6 | 3.7 | 6.0 | 9.1 | 9.9 | 8.6 | 8.1 | 1.7 | 7.9 | 6.3 | 172.9 |
| it | Italian | Romance | 686 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 8.9 | 4.7 | 2.5 | 16.6 | 6.1 | 14.7 | 25.4 | 20.5 | 16.0 | 10.5 | 2.8 | 8.0 | 8.6 | 368.4 |
| ja | Japanese | Japonic | 2944 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 3.3 | 8.9 | 5.1 | 7.7 | 9.1 | 11.6 | 11.3 | 12.1 | 2.8 | 6.5 | 13.5 | 228.7 |
| ko | Korean | Koreanic | 778 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 1.9 | 4.8 | 2.6 | 4.0 | 4.9 | 6.0 | 7.1 | 8.4 | 1.4 | 5.2 | 6.3 | 113.7 |
| ms | Malay | Malayo-Polynesian | 25 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.6 | 1.3 | 2.3 | 3.7 | 3.6 | 3.4 | 0.8 | 3.2 | 2.8 | 60.8 |
| nl | Dutch | Germanic | 510 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 7.8 | 12.9 | 15.5 | 17.7 | 20.8 | 11.0 | 2.7 | 7.2 | 8.4 | 322.2 |
| no | Norwegian | Germanic | 109 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 5.5 | 6.4 | 8.1 | 13.8 | 5.2 | 1.4 | 3.9 | 4.3 | 143.8 |
| pl | Polish | Slavic | 505 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 13.5 | 22.9 | 13.8 | 9.1 | 3.4 | 6.5 | 7.1 | 267.1 |
| pt | Portuguese | Romance | 729 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 20.9 | 15.7 | 11.0 | 3.0 | 8.8 | 9.5 | 375.2 |
| ru | Russian | Slavic | 3047 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 18.9 | 15.3 | 31.2 | 10.4 | 13.7 | 475.0 |
| sv | Swedish | Germanic | 1200 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.8 | | 10.6 | 10.4 | 358.5 |
| tr | Turkish | Turkic | 1382 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.5 | 10.4 | 10.0 | 247.4 |
| uk | Ukrainian | Slavic | 110 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 0.2 | 2.2 | 88.6 |
| vi | Vietnamese | Vietic | 1172 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 9.1 | 190.2 |
| zh | Chinese | Chinese | 2512 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 214.3 |

Table 1: CCMatrix: size of mined sentences **(in millions)** for each language pair.

- Chinese/Japanese: 13.5M

# CCMatrix

| ISO | Name | Family | Size | bg | cs | da | de | el | en | es | fa | fi | fr | he | hi | hu | id | it | ja | ko | ms | nl | no | pl | pt | ru | sv | tr | uk | vi | zh | Total |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ar | Arabic | Arabic | 196 | 3.0 | 3.9 | 2.7 | 7.5 | 3.3 | 6.5 | 10.0 | 3.1 | 2.7 | 23.8 | 2.2 | 1.4 | 2.7 | 4.1 | 5.8 | 5.0 | 2.5 | 1.5 | 5.1 | 2.5 | 4.5 | 6.7 | 9.2 | 5.6 | 5.5 | 1.5 | 4.2 | 5.4 | 141.7 |
| bg | Bulgarian | Slavic | 68 | - | 6.1 | 3.7 | 9.9 | 4.3 | 3.7 | 10.7 | 2.3 | 3.6 | 11.4 | 2.1 | 1.5 | 3.8 | 3.8 | 7.4 | 5.7 | 2.8 | 1.3 | 6.9 | 3.0 | 7.2 | 7.5 | 17.4 | 7.6 | 5.8 | 2.3 | 4.4 | 5.0 | 154.1 |
| cs | Czech | Slavic | 303 | - | - | 5.9 | 18.3 | 5.4 | 9.8 | 15.5 | 2.9 | 6.1 | 17.3 | 3.1 | 2.0 | 6.1 | 5.3 | 11.2 | 8.0 | 4.0 | 2.0 | 11.6 | 4.9 | 13.2 | 10.7 | 18.1 | 12.9 | 8.6 | 2.6 | 6.0 | 7.0 | 228.7 |
| da | Danish | Germanic | 109 | - | - | - | 12.6 | 3.8 | 4.5 | 10.2 | 2.0 | 4.8 | 12.0 | 2.3 | 1.5 | 3.7 | 3.9 | 7.3 | 5.6 | 2.9 | 1.4 | 9.5 | 9.6 | 6.5 | 7.4 | 9.2 | 15.2 | 5.7 | 1.5 | 4.2 | 4.9 | 164.6 |
| de | German | Germanic | 1728 | - | - | - | - | 9.8 | 67.3 | 38.8 | 4.8 | 11.3 | 50.0 | 5.6 | 3.2 | 11.0 | 9.6 | 29.5 | 11.6 | 6.2 | 3.5 | 33.2 | 10.4 | 20.5 | 23.4 | 29.3 | 29.3 | 15.5 | 3.8 | 9.7 | 13.4 | 497.5 |
| el | Greek | Hellenic | 144 | - | - | - | - | - | 5.6 | 12.2 | 2.2 | 3.6 | 12.9 | 2.3 | 1.4 | 3.7 | 3.7 | 8.5 | 5.2 | 2.6 | 1.4 | 6.9 | 3.0 | 6.2 | 8.4 | 9.9 | 7.3 | 5.6 | 1.7 | 4.2 | 4.7 | 150.1 |
| en | English | Germanic | 8677 | - | - | - | - | - | - | 86.3 | 2.5 | 4.1 | 94.1 | 1.5 | 0.7 | 3.6 | 13.4 | 31.3 | 33.7 | 7.2 | 0.8 | 23.8 | 3.8 | 16.0 | 33.1 | 72.4 | 43.8 | 26.8 | 1.6 | 18.5 | 17.6 | 634.2 |
| es | Spanish | Romance | 1534 | - | - | - | - | - | - | - | 5.5 | 9.7 | 70.9 | 5.9 | 3.2 | 9.5 | 12.4 | 44.3 | 11.6 | 6.2 | - | 23.3 | 8.8 | 19.6 | 59.4 | 32.4 | 22.3 | 15.2 | 4.0 | 11.9 | 13.2 | 573.1 |
| fa | Farsi | Iranian | 192 | - | - | - | - | - | - | - | - | 2.0 | 5.5 | 1.7 | 1.2 | 1.9 | 3.1 | 3.6 | 3.5 | 2.0 | 1.3 | 3.6 | 1.9 | 3.2 | 4.1 | 5.6 | 4.0 | 4.9 | 1.1 | 3.3 | 3.4 | 86.3 |
| fi | Finnish | Uralic | 132 | - | - | - | - | - | - | - | - | - | 11.1 | 2.2 | 1.4 | 4.2 | 3.8 | 7.1 | 6.2 | 3.0 | 1.4 | 8.1 | 4.1 | 6.8 | 7.1 | 9.9 | 13.8 | 6.2 | 1.7 | 4.4 | 5.2 | 155.8 |
| fr | French | Romance | 1869 | - | - | - | - | - | - | - | - | - | - | 6.8 | 3.5 | 10.3 | 11.9 | 46.2 | 12.6 | 6.9 | 4.2 | 32.1 | 9.9 | 21.1 | 37.9 | 31.9 | 27.6 | 17.4 | 4.2 | 12.5 | 14.0 | 619.8 |
| he | Hebrew | Semitic | 70 | - | - | - | - | - | - | - | - | - | - | - | 1.2 | 1.9 | 2.8 | 4.0 | 5.3 | 2.5 | 1.1 | 4.2 | 2.0 | 3.6 | 4.3 | 6.4 | 5.1 | 4.4 | 1.2 | 3.6 | 3.6 | 92.9 |
| hi | Hindi | Indo-Aryan | 48 | - | - | - | - | - | - | - | - | - | - | - | - | 1.3 | 1.9 | 2.3 | 2.7 | 1.6 | 0.9 | 2.4 | 1.4 | 2.1 | 2.6 | 3.4 | 3.0 | 3.2 | 0.8 | 1.9 | 2.4 | 56.0 |
| hu | Hungarian | Uralic | 148 | - | - | - | - | - | - | - | - | - | - | - | - | - | 3.2 | 7.0 | 5.2 | 2.6 | 1.3 | 7.1 | 3.0 | 7.1 | 6.8 | 9.6 | 7.4 | 5.6 | 1.7 | 3.7 | 4.6 | 139.6 |
| id | Indonesian | Malayo-Polynesian | 366 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 7.4 | 5.9 | 3.5 | 4.4 | 7.6 | 3.7 | 6.0 | 9.1 | 9.9 | 8.6 | 8.1 | 1.7 | 7.9 | 6.3 | 172.9 |
| it | Italian | Romance | 686 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 8.9 | 4.7 | 2.5 | 16.6 | 6.1 | 14.7 | 25.4 | 20.5 | 16.0 | 10.5 | 2.8 | 8.0 | 8.6 | 368.4 |
| ja | Japanese | Japonic | 2944 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 3.3 | 8.9 | 5.1 | 7.7 | 9.1 | 11.6 | 11.3 | 12.1 | 2.8 | | 6.5 | 13.5 | 228.7 |
| ko | Korean | Koreanic | 778 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 1.9 | 4.8 | 2.6 | 4.0 | 4.9 | 6.0 | 7.1 | 8.4 | 1.4 | 5.2 | 6.3 | 113.7 |
| ms | Malay | Malayo-Polynesian | 25 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.6 | 1.3 | 2.3 | 2.8 | 3.7 | 3.6 | 3.4 | 0.8 | 3.2 | 2.8 | 60.8 |
| nl | Dutch | Germanic | 510 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 7.8 | 12.9 | 15.5 | 17.7 | 20.8 | 11.0 | 2.7 | 7.2 | 8.4 | 322.2 |
| no | Norwegian | Germanic | 109 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 5.5 | 6.4 | 8.1 | 13.8 | 5.2 | 1.4 | 3.9 | 4.3 | 143.8 |
| pl | Polish | Slavic | 505 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 13.5 | 22.9 | 13.8 | 9.1 | 3.4 | 6.5 | 7.1 | 267.1 |
| pt | Portuguese | Romance | 729 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 20.9 | 15.7 | 11.0 | 3.0 | 8.8 | 9.5 | 375.2 |
| ru | Russian | Slavic | 3047 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 18.9 | 15.3 | 31.2 | 10.4 | 13.0 | 475.0 |
| sv | Swedish | Germanic | 1200 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.8 | | 10.6 | 10.4 | 358.5 |
| tr | Turkish | Turkic | 1382 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 2.5 | 10.4 | 10.0 | 247.4 |
| uk | Ukrainian | Slavic | 110 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 0.2 | 2.2 | 88.6 |
| vi | Vietnamese | Vietic | 1172 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 9.1 | |
| zh | Chinese | Chinese | 2512 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 214.3 |

Table 1: CCMatrix: size of mined sentences **(in millions)** for each language pair.

- Hindi with Chinese, Japanese, Korean, Indonesian: ≈2M

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Example Parallel Sentences

| | |
|---|---|
| Japanese | 国内レベルで進捗状況を監視するためには、良質かつアクセス可能な適時のデータ収集や、地域的なフォローアップと検証が必要となります<br><br>*(Monitoring progress at the national level requires quality, accessible and timely data collection and regional follow-up and verification.)* |
| Russian | Для мероприятий по отслеживанию прогресса на национальном уровне необходимо обеспечить сбор качественных, доступных и актуальных данных, а также проведение последующей деятельности и обзора на региональном уровне.<br><br>*(For activities to track progress at the national level, it is necessary to ensure the collection of quality, accessible and relevant data, as well as follow-up and review at the regional level.)* |

Cross-lingual Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Example Parallel Sentence

| | |
|---|---|
| Malay | Tahun ketiga pengajian biasanya dibelanjakan ke luar negara di institusi rakan kongsi di Timur Tengah atau Afrika Utara. *(The third year of study is usually spent abroad at partner institutions in the Middle East or North Africa.)* |
| Chinese | 研究的第三年通常是在中东或北非伙伴机构在国外度过。 *(The third year of the study is usually spent abroad in partner institutions in the Middle East or North Africa.)* |

# Example Parallel Sentence

|  | |
|---|---|
| Arabic | ذهب الحراس إلى الأبواب البلوط الكبيرة وهيلين وجاك نزلا إلى الاسطبلات لجعل لأجل اثنين من الخيول. |
| | *(The guards went to the large oak doors, Helen and Jack came down to stables to make for two horses.)* |
| Hebrew | השומרים חזרו אל דלתות עץ האלון הגדולות והלנה וג'ק ירדו לאורוות להכין שני סוסים. |
| | *(The guards returned to the large oak doors and Helena and Jack went down to the stables to make two horses.)* |

# Example Parallel Sentence

| Finish | Tarkista aina väliin, että olet oikealla tiellä. |
| | *(Always check in between that you are on the right track.)* |
| Tamil | எப்போதும் உறுதிப்படுத்திக் கொள்ளுங்கள் நீங்கள் சரியான வழியில் செல்வீர்கள்.|
| | *(Always make sure you go the right way.)* |

Cross-lingual
Mining

A. El-Khshky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# 11-way Parallel Sentences

| En | You should clean the refrigerator once a month. | Visiting a sick friend. |
|---|---|---|
| Ar | وأخيراً تذكري أنه يجب عليكي تنظيف الثلاجة مرة واحدة في الشهر. | زرت صديقاً مريضا. |
| De | Den Kühlschrank sollten Sie einmal im Monat saubermachen. | Ein Besuch in einem kranken Freund |
| Fr | Il est recommandé de nettoyer le réfrigérateur une fois par mois. | visite à un ami malade. |
| Id | Sebulan sekali kulkas harus dibersihkan. | Kunjungi teman yang sakit |
| Ja | １ヶ月に１回くらいは冷蔵庫の蔵ざらえをしなきゃ。 | 病の友達を訪ねる |
| Ko | 한 달에 한 번 정도는 냉장고 청소를 해주는 게 좋다. | 아픈 친구를 보는 심정으로 |
| Ru | Холодильник следует размораживать раз в месяц. | Посещение больного друга. |
| Tr | Buzdolabını boşaltarak ayda bir kez temizleyin. | Hasta bir dostu ziyaret etmek. |
| Vi | Vì vậy, mỗi tháng bạn nên vệ sinh tủ lạnh một lần. | Thăm người bạn THÂN bệnh |
| Zh | 如果有必要，你可以一个月清理一次冰箱。 | 探望一个生病的朋友。 |

| En | When we breathe quickly we also build up oxygen in our blood. |
|---|---|
| Ar | عندما نتنفس بسرعة نقوم ببناء الأكسجين في دمائنا. |
| De | Wenn wir schnell atmen, bauen wir auch Sauerstoff in unserem Blut auf. |
| Fr | Lorsque nous respirons rapidement, nous créons également de l'oxygène dans notre sang. |
| Id | Ketika kita bernapas dengan cepat, kita juga membangun oksigen dalam darah kita. |
| Ja | 私たちが素早く呼吸すると、血液中に酸素も蓄積します。 |
| Ko | 우리가 빨리 숨을 쉬면 우리도 피 속에 산소를 축적합니다. |
| Ru | Когда мы дышим быстро, мы также накапливаем кислород в нашей крови. |
| Tr | Khi chúng ta thở nhanh, chúng ta cũng tích tụ oxy trong máu. |
| Vi | Çabucak nefes aldığımızda, kanımızda da oksijen biriktiririz. |
| Zh | 当我们快速呼吸时，我们的血液中也会积聚氧气。 |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# 11-way Parallel Sentences

| En | With the growing importance of world trade and the global community, business executives and legal professionals are expected to look beyond national jurisdictions and understand issues of international law and international commercial law. |
|---|---|
| Ar | مع تزايد أهمية التجارة العالمية والمجتمع العالمي، ومن المتوقع أن تنظر إلى أبعد السلطات القضائية الوطنية وفهم قضايا القانون الأوروبي والدولي المستشارين القانونيين. |
| De | Da Handel und Unternehmen immer globaler werden, wird erwartet, dass Rechtsberater über nationale Zuständigkeiten hinausblicken und Fragen des europäischen und internationalen Rechts verstehen. |
| Fr | Avec l'importance croissante du commerce mondial et la communauté mondiale, consultants juridiques devraient regarder au-delà des juridictions nationales et de comprendre les questions de droit européen et international. |
| Id | Dengan semakin pentingnya perdagangan dunia dan masyarakat global, konsultan hukum diharapkan untuk melihat melampaui yurisdiksi nasional dan memahami masalah hukum Eropa dan internasional. |
| Ja | 法律コンサルタントは、貿易とビジネスがますますグローバル化するにつれて、国の管轄権を超えて、欧州および国際法の問題を理解することが期待されています。 |
| Ko | 무역 및 비즈니스가 전 세계적으로 증가함에 따라 법률 컨설턴트는 국가 관할권을 넘어서서 유럽 및 국제법 문제를 이해할 것으로 예상됩니다. |
| Ru | С ростом важности мировой торговли и мирового сообщества, юридические консультанты, как ожидается, искать за пределами национальной юрисдикции и понимания вопросов европейского и международного права. |
| Tr | Ticaret ve iş dünyası gittikçe küreselleştikçe, hukuk müşavirlerinin ulusal yargıların ötesine geçmesi ve Avrupa ve uluslararası hukuk konularını anlamaları beklenmektedir. |
| Vi | Với tầm quan trọng ngày càng tăng của thương mại thế giới và cộng đồng quốc tế, tư vấn pháp luật được dự kiến để nhìn xa hơn khu vực pháp lý quốc gia và hiểu các vấn đề của pháp luật châu Âu và quốc tế. |
| Zh | 随着世界贸易和全球社会的重要性日益增加，法律顾问有望超越国家管辖和了解欧洲和国际法律的问题。 |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Evaluating CCMatrix

## WMT'19

- De-facto standard for NMT progress, strong competition
- Train NMT systems **on mined data only**,
  no human bitexts

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Evaluating CCMatrix

## WMT'19

- De-facto standard for NMT progress, strong competition
- Train NMT systems **on mined data only**,
  no human bitexts
- Newstest 2018:

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Evaluating CCMatrix

## WMT'19

- De-facto standard for NMT progress, strong competition
- Train NMT systems **on mined data only**,
  no human bitexts
- Newstest 2018:



- We outperform all best single systems, +3.8 BLEU en-de

# Evaluating CCMatrix

## WMT'19

- De-facto standard for NMT progress, strong competition
- Train NMT systems **on mined data only**,
  no human bitexts
- Newstest 2019:

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
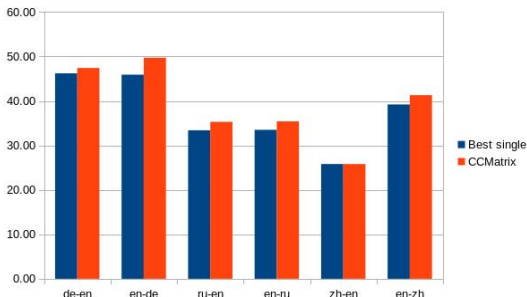WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Evaluating CCMatrix

## WMT'19
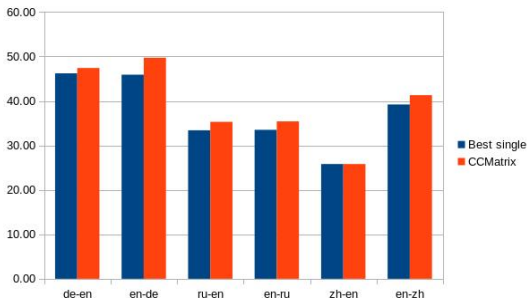
- De-facto standard for NMT progress, strong competition
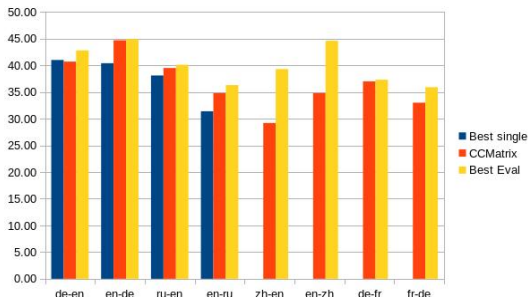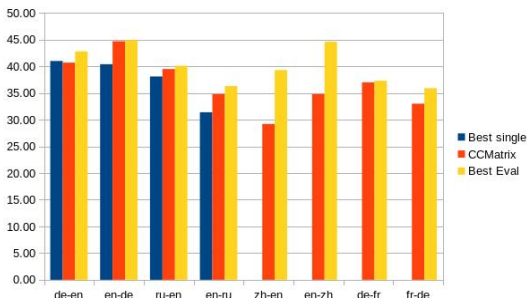- Train NMT systems **on mined data only**,
  no human bitexts
- Newstest 2019:



- en-de/de-fr: on-pair with eval system (BT, sys.comb)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# BLEU Scores on TED Test Sets

| | ar | bg | cs | da | de | el | en | es | fa | fi | fr | he | hi | id | it | ja | ko | ms | nl | no | pl | pt | ru | sv | tr | uk | vi | zh |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ar | - | 16.6 | 11.5 | 14.9 | 15.5 | 17.5 | 28.7 | 22.4 | 8.7 | 7.3 | 19.6 | 10.9 | 12.6 | 17.5 | 18.6 | 8.5 | 2.8 | 10.2 | 15.8 | 16.2 | 10.2 | 20.7 | 13.8 | 18.8 | 7.6 | 7.1 | 19.0 | 9.9 |
| bg | 9.5 | - | 19.3 | 24.9 | 22.9 | 24.4 | 36.3 | 27.8 | 9.8 | 11.9 | 26.1 | 13.7 | 14.4 | 22.0 | 23.3 | 9.8 | 3.4 | 12.9 | 21.8 | 22.2 | 15.0 | 26.1 | 19.3 | 24.1 | 9.9 | 13.9 | 21.8 | 10.5 |
| cs | 7.0 | 21.9 | - | 21.8 | 21.0 | 20.0 | 29.2 | 24.0 | 7.9 | 13.8 | 24.1 | 10.8 | 13.9 | 18.8 | 20.5 | 9.9 | 3.5 | 8.5 | 21.3 | 22.6 | 16.1 | 22.5 | 18.1 | 22.6 | 9.2 | 11.4 | 20.2 | 10.3 |
| da | 6.5 | 25.3 | 18.1 | - | 26.8 | 23.5 | 44.1 | 28.9 | 8.4 | 14.3 | 27.6 | 12.6 | 15.8 | 22.2 | 24.3 | 10.3 | 3.5 | 14.7 | 28.2 | 32.1 | 16.3 | 27.0 | 18.9 | 33.4 | 10.3 | 9.9 | 19.9 | 9.3 |
| de | 9.2 | 24.0 | 19.8 | 30.4 | - | 22.3 | 35.8 | 28.1 | 9.7 | 13.1 | 27.6 | 13.7 | 18.4 | 22.5 | 24.8 | 11.7 | 3.4 | 15.5 | 26.9 | 20.5 | 15.5 | 26.9 | 18.6 | 26.6 | 11.4 | 12.3 | 22.6 | 11.3 |
| el | 9.8 | 24.4 | 17.2 | 25.0 | 21.3 | - | 35.8 | 28.4 | 9.3 | 11.8 | 26.7 | 13.8 | 16.0 | 21.5 | 23.8 | 9.7 | 3.2 | 13.4 | 22.1 | 24.1 | 13.9 | 26.4 | 17.5 | 24.1 | 10.1 | 10.3 | 22.2 | 10.9 |
| en | 16.6 | 35.7 | 24.5 | 42.2 | 32.6 | 34.6 | - | 42.4 | 15.7 | 17.6 | 36.6 | 24.8 | 25.2 | 33.4 | 34.1 | 12.3 | 5.8 | 24.8 | 33.3 | 43.2 | 18.4 | 41.2 | 21.9 | 38.2 | 16.1 | 19.2 | 29.5 | 15.1 |
| es | 11.6 | 26.4 | 19.4 | 28.7 | 25.1 | 26.2 | 41.1 | - | 10.9 | 14.4 | 30.7 | 15.7 | 17.9 | 24.2 | 29.9 | 10.9 | 4.5 | 15.5 | 26.2 | 25.0 | 15.7 | 32.8 | 19.4 | 26.9 | 11.4 | 13.9 | 24.7 | 12.3 |
| fa | 7.0 | 14.7 | 10.7 | 14.4 | 14.2 | 14.6 | 31.6 | 21.6 | - | 5.7 | 17.2 | 8.0 | 9.1 | 16.9 | 15.2 | 7.6 | 2.3 | 9.3 | 14.2 | 11.8 | 8.5 | 17.5 | 12.3 | 15.6 | 8.1 | 6.4 | 17.7 | 8.8 |
| fi | 4.7 | 14.2 | 13.2 | 17.7 | 14.6 | 13.8 | 21.6 | 18.1 | 4.0 | - | 16.1 | 7.8 | 11.9 | 13.3 | 14.1 | 9.5 | 2.8 | 2.2 | 16.1 | 13.7 | 10.0 | 14.8 | 11.6 | 16.6 | 7.2 | 6.3 | 16.1 | 8.0 |
| fr | 10.4 | 25.9 | 19.7 | 29.5 | 25.4 | 26.6 | 40.2 | 32.6 | 10.6 | 13.8 | - | 15.6 | 18.8 | 25.5 | 29.6 | 10.8 | 4.6 | 13.7 | 26.7 | 30.8 | 14.8 | 31.2 | 20.0 | 27.2 | 12.0 | 13.4 | 24.6 | 10.3 |
| he | 9.6 | 19.0 | 13.8 | 19.8 | 18.2 | 18.9 | 33.4 | 24.6 | 7.7 | 9.7 | 23.2 | - | 12.8 | 18.8 | 20.2 | 7.9 | 3.1 | 8.4 | 18.3 | 17.2 | 11.7 | 22.7 | 15.9 | 21.1 | 8.1 | 7.5 | 18.8 | 9.2 |
| hi | 4.1 | 10.3 | 7.9 | 11.8 | 14.2 | 11.4 | 24.3 | 16.5 | 3.6 | 6.5 | 17.2 | 6.9 | - | 13.4 | 12.9 | 6.5 | 1.9 | 6.7 | 12.9 | 9.6 | 7.1 | 15.4 | 12.4 | 13.9 | 6.8 | 3.8 | 15.3 | 6.5 |
| id | 8.3 | 19.9 | 14.3 | 21.4 | 19.6 | 19.1 | 31.9 | 24.7 | 9.9 | 9.8 | 23.7 | 11.7 | 16.4 | - | 20.2 | 10.1 | 4.6 | 19.0 | 20.8 | 21.2 | 12.8 | 23.5 | 16.0 | 21.2 | 10.2 | 9.6 | 23.7 | 11.5 |
| it | 11.0 | 24.2 | 18.4 | 26.4 | 24.1 | 25.4 | 30.4 | 32.5 | 10.2 | 13.0 | 30.4 | 14.1 | 16.8 | 23.1 | - | 10.7 | 3.8 | 13.2 | 24.8 | 25.3 | 15.4 | 30.8 | 18.7 | 26.9 | 11.2 | 12.1 | 23.7 | 11.4 |
| ja | 3.8 | 7.4 | 5.9 | 8.2 | 7.8 | 7.7 | 11.8 | 11.6 | 4.3 | 4.6 | 10.7 | 4.0 | 9.1 | 9.6 | 9.2 | - | 3.1 | 5.5 | 8.4 | 7.4 | 6.1 | 9.9 | 7.3 | 8.1 | 4.6 | 3.3 | 12.0 | 6.9 |
| ko | 4.2 | 8.6 | 7.4 | 9.8 | 10.0 | 8.8 | 15.4 | 13.7 | 5.1 | 5.6 | 13.2 | 5.2 | 10.8 | 12.5 | 10.5 | 11.0 | - | 6.7 | 10.2 | 9.9 | 6.7 | 12.1 | 8.8 | 11.3 | 6.1 | 3.8 | 14.3 | 8.1 |
| ms | 7.4 | 11.9 | 8.8 | 14.2 | 13.3 | 13.6 | 29.1 | | 9.0 | 5.2 | 16.9 | 6.4 | 13.0 | 20.9 | 18.0 | 8.6 | 2.5 | - | 14.0 | 14.8 | 8.9 | 19.4 | 12.9 | 20.1 | 7.8 | 4.8 | 23.8 | 10.1 |
| nl | 8.7 | 21.8 | 18.0 | 28.3 | 25.6 | 21.9 | 35.8 | 28.3 | 9.3 | 12.8 | 27.9 | 13.4 | 16.9 | 22.7 | 23.7 | 10.4 | 3.7 | 12.2 | - | 18.1 | 15.5 | 26.6 | 17.2 | 25.4 | 10.6 | 9.8 | 22.0 | 10.7 |
| no | 9.7 | 23.3 | 19.8 | 34.0 | 21.9 | 24.2 | 45.2 | 27.0 | 5.3 | 12.4 | 25.6 | 12.7 | 13.5 | 22.9 | 26.8 | 10.8 | 3.5 | 16.8 | 19.1 | - | 13.1 | 27.6 | 18.8 | 32.5 | 9.7 | 11.0 | 16.4 | 9.7 |
| pl | 6.5 | 17.3 | 16.2 | 19.5 | 16.8 | 15.9 | 22.3 | 20.0 | 6.5 | 10.2 | 20.0 | 9.3 | 12.9 | 16.3 | 17.4 | 9.4 | 3.0 | 9.8 | 17.5 | 14.4 | - | 18.8 | 15.4 | 17.4 | 7.6 | 10.6 | 17.7 | 8.5 |
| pt | 11.5 | 26.6 | 19.3 | 29.1 | 25.0 | 27.3 | 43.0 | 35.7 | 11.0 | 13.0 | 31.9 | 16.0 | 18.5 | 25.9 | 30.5 | 10.6 | 4.2 | 15.1 | 26.3 | 24.3 | 16.2 | - | 19.7 | 26.6 | 11.5 | 13.2 | 25.1 | 10.0 |
| ru | 8.0 | 19.6 | 15.9 | 18.8 | 18.1 | 18.3 | 24.2 | 21.8 | 8.6 | 10.6 | 22.1 | 11.2 | 16.0 | 17.8 | 19.0 | 10.1 | 3.5 | 11.7 | 18.2 | 19.3 | 14.1 | 20.7 | - | 19.7 | 8.4 | 21.1 | 19.3 | 10.3 |
| sv | 11.1 | 24.4 | 20.2 | 34.3 | 26.3 | 24.4 | 40.8 | 28.9 | 9.4 | 14.9 | 27.7 | 15.4 | 18.7 | 24.3 | 26.3 | 10.8 | 4.6 | 15.7 | 25.7 | 30.3 | 14.7 | 28.3 | 18.8 | - | 12.0 | 11.2 | 24.9 | 12.1 |
| tr | 7.5 | 16.4 | 12.8 | 16.2 | 15.4 | 16.6 | 25.0 | 20.3 | 8.7 | 9.6 | 19.9 | 9.2 | 15.8 | 17.5 | 16.9 | 9.6 | 4.1 | 10.8 | 16.6 | 14.7 | 10.8 | 18.3 | 13.3 | 18.2 | - | 7.3 | 19.2 | 9.9 |
| uk | 5.0 | 17.2 | 12.0 | 13.8 | 14.7 | 13.2 | 23.1 | 18.6 | 5.4 | 8.0 | 17.7 | 6.9 | 8.7 | 12.8 | 15.3 | 7.0 | 1.7 | 6.2 | 13.1 | 12.9 | 12.3 | 17.2 | 23.0 | 13.6 | 5.6 | - | 14.5 | 7.3 |
| vi | 8.0 | 16.8 | 12.9 | 17.1 | 16.5 | 17.0 | 25.8 | 21.7 | 8.7 | 9.3 | 21.0 | 9.9 | 15.7 | 21.3 | 18.3 | 9.6 | 4.3 | 16.5 | 17.6 | 16.1 | 11.0 | 20.6 | 14.1 | 19.0 | 9.4 | 9.1 | - | 10.8 |
| zh | 6.3 | 11.8 | 9.3 | 11.2 | 12.2 | 12.3 | 18.3 | 16.0 | 6.9 | 7.4 | 15.2 | 7.6 | 12.3 | 14.8 | 13.4 | 9.6 | 3.5 | 9.7 | 12.8 | 12.6 | 8.4 | 14.0 | 11.2 | 13.8 | 6.8 | 6.1 | 18.1 | - |

- Same NMT system for all language pairs
  (despite huge difference in bitext size)
- Best: BLEU 45.2 for Norwegian/English

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# CCMatrix: What's Next ?

Scaling even further

- Scaling to 32 crawls, 100 languages
- $\Rightarrow$ $\approx$ 10 billion bitexts
- Further improvements on WMT evaluation

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# CCMatrix: What's Next ?

## Scaling even further

- Scaling to 32 crawls, 100 languages
- $\Rightarrow$ $\approx$ 10 billion bitexts
- Further improvements on WMT evaluation

## Sharing our results

- Looking for means to share these CCMatrix bitexts

# parallel sentence filtering

# Filtering for What?

- We have intuitive notions of useful training data
    - source and target match in meaning
    - both are well-formed text

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Filtering for What?

- We have intuitive notions of useful training data
    - source and target match in meaning
    - both are well-formed text
- But: the right question is: does it help to build a better MT system
- We do not know how to answer that

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Types of Noise

- Misaligned sentences
- Disfluent language (from MT, bad translations)
- Wrong language data (e.g., French in German–English corpus)
- Untranslated sentences
- Short segments (e.g., dictionaries)
- Mismatched domain

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Mismatched Sentences

- Artificial created by randomly shuffling sentence order
- Added to existing parallel corpus in different amounts

| 5% | 10% | 20% | 50% | 100% |
|---|---|---|---|---|
| 24.0 | 24.0 | 23.9 | 26.1 23.9 | 25.3 23.4 |
| -0.0 | -0.0 | -0.1 | -1.1 -0.1 | -1.9 -0.6 |

- Bigger impact on NMT (green, left) than SMT (blue, right)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Misordered Words

- Artificial created by randomly shuffling words in each sentence

| | 5% | 10% | 20% | 50% | | 100% | |
|---|---|---|---|---|---|---|---|
| **Source** | 24.0<br>-0.0 | 23.6<br>-0.4 | 23.9<br>-0.1 | 26.6<br>-0.6 | 23.6<br>-0.4 | 25.5<br>-1.7 | 23.7<br>-0.3 |
| **Target** | 24.0<br>-0.0 | 24.0<br>-0.0 | 23.4<br>-0.6 | 26.7<br>-0.5 | 23.2<br>-0.8 | 26.1<br>-1.1 | 22.9<br>-1.1 |

- Similar impact on NMT than SMT, worse for source reshuffle

# Untranslated Sentences



|  | 5% | 10% | 20% | 50% | 100% |
|---|---|---|---|---|---|
| Source | 17.6  23.8 / −0.2 / −9.8 | 11.2  23.9 / −0.1 / −16.0 | 5.6  23.8 / −0.2 / −21.6 | 3.2  23.4 / −0.6 / −24.0 | 3.2  21.1 / −2.9 / −24.0 |
| Target | 27.2 / −0.0 | 27.0 / −0.2 | 26.7 / −0.5 | 26.8 / −0.4 | 26.9 / −0.3 |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Copy Noise

- Harmfulness of copy noise also discovered by Ott, Auli, Granger, Ranzato (Facebook FAIR)
  - noticed link to beam search decoding

  - proposed remedies at inference time

- Motivated overlap penalty as feature in data filtering

# Zipporah: A Fast and Scalable Data Cleaning System for Noisy Web-Crawled Parallel Corpora

# Zipporah: Motivation & Objective

Motivation

- Often have large pool of noisy parallel data
- Need to perform fast data-selection to select a higher-quality subset of this data

Objective

- Design a function to rank the sentence pairs
- Select the best sentences under some size constraint

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Features

**Features**

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Features

**Features**

- Adequacy: how good the translation is

| French | English | adequacy |
|---|---|---|
| Je suis Hainan. | I am Hainan. | ✓ |
| Je suis Hainan. | The weather is quite good today. | ✗ |
| - - - - - - | - - - - - - | ✓ |

- Fluency: measures how good a sentence is

| French | English | fluency |
|---|---|---|
| - - - - - - | - - - - - - | ✗✗ |
| Je suis Hainan. | The weather is not quite good today. | ✓✓ |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Adequacy Features

- Dictionary
  - P(I—Je) = 1
  - P(am—suis) = 3/4
  - P(follow—suis) = 1/4

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Adequacy Features

$$XEnt(p, q) = \sum_i p(i) log \frac{1}{q(j)} \qquad (10)$$

- A(en—fr) = 1/3 * log 3 + 1/3 * log 4 + 1/3 log 3 = 1.1945
- Also compute A(fr—en), given e2f dictionary
- For each sentence pair, define A(fr—en) + A(en—fr) as the adequacy feature
- Small when the translation is good (and literal)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Fluency Features

$$F(s) = -\frac{log(p_{LM}(s))}{lengh(s)} \tag{11}$$

- Tarn ngram LMs for both languages
- For each sentence, we compute the $F(s)$
- For each sentence pair, define $F(en) + F(fr)$ as the fluency feature
- Small when the sentence pair is fluent

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Scoring Function

- Goal: train a classifier to distiinguish between good and bad data
- Have good data (true parallel sentences)
- Need bad data. Preferably one that covers all types f bad data in the feature space
- Auto generate bad data from good data

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Generating Bad Data

Starting from a good dev corpus

- Shuffly individual words within sentences (bad fluency)
- shuffle sentences (bad adequacy)
- Shuffle both (bad both)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Bad Data vs Good Data

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Logistic Regression Classifier

- Task: To separate the good parallel data from (synthetic) bad parallel data
- Method: Logistic regression classifier with polynomials of features.
- Use the trained weights to compute a signed-distance to the decision boundary as score

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
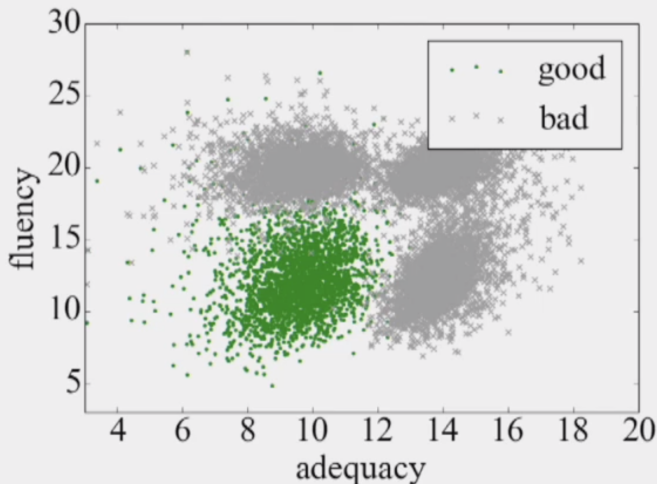CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Baselines

- Random selection
- QE Clean: Uses LM scores and word-alignment scores to perform data selection

# Zipporah: Results

French-English: Ted Talks Dataset

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Zipporah: Results

German-English: Newstest 11 Datasett

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

WMT Shared Task on Sentence
Pair Filtering

# WMT Shared Task on Sentence Pair Filtering

# WMT Shared Task on Sentence Pair Filtering

- Shared Task in 2018: High Resouce
  - German–English
  - 1 billion words of noisy parallel data
  - 100+ million words of clean parallel data

- Shared Task in 2019: Low Resource
  - Sinhala–English and Nepali–English
  - 50-60 million words of noisy parallel data
  - 3-4 million words of *relatively* clean parallel data

# Task Definition

- Given
  - very noise web crawled corpus
  - sentence-aligned
  - 50-60 billion English words

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Task Definition

- Given
  - very noise web crawled corpus
  - sentence-aligned
  - 50-60 billion English words
- Submission: sentence-level quality score for each sentence pair

# Task Definition

- Given
    - very noise web crawled corpus
    - sentence-aligned
    - 50-60 billion English words
- Submission: sentence-level quality score for each sentence pair
- Evaluation
    - subselection of training corpus based on quality threshold
        - 1 million English words
        - 5 million English words
    - machine translation performance on undisclosed test sets
        - statistical machine translation (Moses)
        - neural machine translation (fairseq)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Provided Resources

- Noisy parallel corpus
    - English sentence
    - foreign sentence
    - Hunalign score
- Training data

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Provided Resources

- Noisy parallel corpus
  - English sentence
  - foreign sentence
  - Hunalign score
- Training data
- Development pack
  - script to subsample corpora
  - Moses configuration file to build and test SMT system
  - Fairseq scripts to build and test NMT system
  - Development and test sets: Wikipedia translations

Cross-lingual Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Clean Parallel Corpora

| Nepali | Sentence Pairs | English Words |
|---|---:|---:|
| Bible (two translations) | 61,645 | 1,507,905 |
| Global Voices | 2,892 | 75,197 |
| Penn Tree Bank | 4,199 | 88,758 |
| GNOME/KDE/Ubuntu | 494,994 | 2,018,631 |
| Nepali Dictionary | 9,916 | 25,058 |

| Sinhala | Sentence Pairs | English Words |
|---|---:|---:|
| Open Subtitles | 601,164 | 3,594,769 |
| GNOME/KDE/Ubuntu | 45,617 | 150,513 |

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Development and Test Sets

- Evaluation on translations of Wikipedia content

  **Two New Evaluation Datasets for Low-Resource Machine Translation: Nepali-English and Sinhala-English**, Francisco Guzmán, Peng-Jen Chen, Myle Ott, Juan Pino, Guillaume Lample, Philipp Koehn, Vishrav Chaudhary, Marc'Aurelio Ranzato, *arXiv:1902.01382*

|  | Nepali | | Sinhala | |
|---|---|---|---|---|
|  | **Sentence Pairs** | **English Words** | **Sentence Pairs** | **English Words** |
| dev | 2,559 | 46,274 | 2,898 | 53,479 |
| dev test | 2,835 | 51,458 | 2,766 | 50,985 |
| test | 2,924 | 54,062 | 2,905 | 52,851 |

# Participants

| Acronym | Participant and System Description Citation |
|---|---|
| AFRL | Air Force Research Lab, USA |
| DiDi | DiDi, USA |
| Facebook | Facebook, USA |
| Helsinki | University of Helsinki, Finland |
| IITP | Indian Institute of Technology Patna, India |
| Webinterpret | WebInterpret Inc., USA |
| NRC | National Research Council, Canada |
| Stockholm | Stockholm University, Sweden |
| SUNY Buffalo | State University of New York, USA |
| Sciling | Sciling S.L., Spain |
| TALP-UPC | TALP, Universitat Politècnica de Catalunya, Spain |

# methods

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# AFRL

### AFRL

- Uses coverage metric and quality metric.

- Coverage metric discourages addition of sentence pairs that have vocab already included in selected set

- Quality metric based on comparing machine translation of foreign sentences with English sent using METEOR MT metric

DiDi

### DiDi

- Dual cross-entropy based on monolingual language models to find pairs where each sentence has similar probability

- Cynical data selection that prefers to select representative subset

- Length-ratio and using character-set based language identification

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Facebook

Facebook

- Ensemble
- Matching of cross-lingual sentence embeddings feature
- Dual cross entropy based on neural translation model scores
- Open source Ziporah and Bicleaner

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# NRC

### NRC

- Filtering rules based on lang ID, length ratio, mismatched numbers, near duplicates
- Cross-lingual semantic evaluation metric (Yisi-2) that uses:
    - cross-lingual word embeddings
    - transformer model language model pretrainined based on XLM
    - optimized to distinguish between clean parallel data and synthetic noisy parallel data
- Reranking to increase coverage

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Sciling

### Sciling

- Build translation models on clean data

- translate non-english to English in noisy data

- Similarity between machine translation and given English sentence

- Filtering rules for sent length, source-target overlap, and lang identification

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Stockholm

Stockholm

- Filtering (excessive numbers, too few words, sentence length, too long, etc)

- Mono-lingual word embeddings with FastText

- learn projection between emebdding spaces based on word alignment from parallel data

- Cosine similarity between English word to best matching projection of the word

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# TALP-UPC

TALP-UPC

- Monolingual word embeddings with FastText
- Unsupervised word ealignment
- Word mover's distance between sentences
- Filtering rules (sent length, lang identification, num mismatches)

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Helsinki

### TALP-UPC

- Clean the clean parallel data using filter rules (sent length, sents with long words, XML, HTML, tags, wrong script)
- Obtain word alignments from this clean data
- Noisy parallel data is scored using word alignments
- Filtered with language models, lang identifiication, ratio of chars in correct script, punctuation, number matching, length mismatch
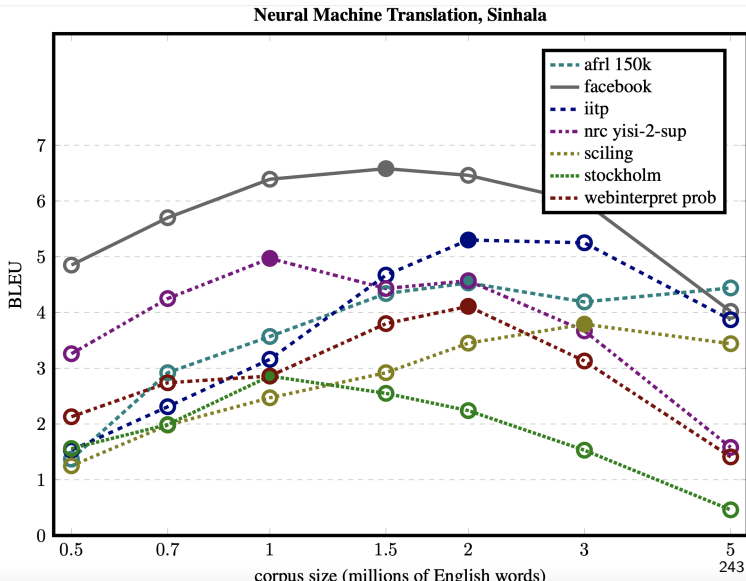
# Webinterpret

Webinterpret

- Filtering rules based on language identification and sent length

- Coverage ranking incrementally adds sentence pairs to increase vocan and ngram coverage

- Adequacy ranking considers IBM Model 1 word translation scores

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Different Subset Sizes

## Neural Machine Translation, Sinhala



**Neural Machine Translation, Sinhala**

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Different Subset Sizes

## Statistical Machine Translation, Sinhala



Statistical Machine Translation, Sinhala

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Different Subset Sizes

## Neural Machine Translation, Nepali



**Neural Machine Translation, Nepali**

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Different Subset Sizes

## Statistical Machine Translation, Nepali



Statistical Machine Translation, Nepali

corpus size (millions of English words)

# Things Learned

## Commonalities Learned from Submissions

# Pre-Filtering Rules

- Discard some data based on deterministic filtering rules
    - too short or too long
    - too many non-words
    - average token length is too short or too long
    - mismatched lengths
    - names, numbers, email addresses, URLs do not match between both sides
    - too similar, indicating simple copying
    - language identification

# Embeddings

- Cross-lingual sentence embeddings
    - central to best performing system
    - LASER (Artexte and Schwenk, 2018)
- Word embeddings
    - monolingual spaces, mapped unsupervised or using dictionaries
    - bilingually trained word embeddings

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Use of Machine Translation Models

- Quality scores on translations
    - translate foreign into English
    - score with METEOR, BLEU, Levenshtein distance
- Cross-entropy filtering
    - force-translate foreign into given English
    - consider translation model score

# Scoring Functions

- *N*-gram or neural language models on clean data
- Language models trained on the provided raw data as contrast
- Neural translation models
- Bag-of-words lexical translation probabilities
- Off-the-shelf tools: Zipporah, Bicleaner

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Learning Weights for Scoring Functions

- Large number of scoring functions $\rightarrow$ averaging scores inadequate
- Learning weights to optimize MT quality computationally intractable
- Solution: train classifier to distinguish between good and bad sentence pairs
  - good sentence pairs from clean corpus
  - bad sentence pairs from provided data, or synthetic noise

A. El-Kishky,
P. Koehn,
H. Schwenk

Low-Resource Corpus Filtering
using Multilingual Sentence
Embeddings

# Low-Resource Corpus Filtering using Multilingual Sentence Embeddings

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Approach

- Leverage LASER multilingual embeddings as a tool to measure parallel sentence quality
- Margin-based scoring function to score sentence pairs

255 / 258

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Scoring Function

$$\frac{2k\ cos(x,y)}{\sum_{y' \in NN_k(x)} cos(x,y') + \sum_{x' \in NN_k(y)} cos(x'y))}$$

where

- $NN_k(x)$ denotes the k nearest neighbors of x in the other language and analogously for $NN_k(y)$.

- pool of sentences are deduplicated

- Global: pool of neighbors can be from global (all = clean + noisy data)

- Local: pool of neighbors can be from local (only from noisy data)

255 / 258

# Results

## Dev Test Results

| Method | ne-en | | si-en | |
|---|---|---|---|---|
| | **1M** | **5M** | **1M** | **5M** |
| **Zipporah** | | | | |
| base | 5.03 | 2.09 | 4.86 | 4.53 |
| + LID | 5.30 | 1.53 | 5.53 | 3.16 |
| + Overlap | 5.35 | 1.34 | 5.18 | 3.14 |
| **Dual X-Ent.** | | | | |
| base | 2.83 | 1.88 | 0.33 | $4.63^+$ |
| + LID | 2.19 | 0.82 | 6.42 | 3.68 |
| + Overlap | 2.23 | 0.91 | 6.65 | 4.31 |
| **Bicleaner** | | | | |
| base | 5.91 | $2.54^+$ | 6.20 | 4.25 |
| + LID | 5.88 | 2.09 | 6.36 | 3.95 |
| + Overlap | $6.12^+$ | 2.14 | $6.66^+$ | 3.26 |
| **LASER** | | | | |
| *local* | *7.37\** | **3.15** | *7.49\** | 5.01 |
| *global* | 6.98 | *2.98\** | 7.27 | 4.76 |
| **Ensemble** | | | | |
| ALL | 6.17 | 2.53 | **7.64** | **5.12** |
| LASER *glob. + loc.* | **7.49** | 2.76 | 7.27 | *5.08\** |

Bold=best scores, Italics\*= runner up

Cross-lingual
Mining

A. El-Kishky,
P. Koehn,
H. Schwenk

Introduction

Corpora and
WEB Crawling

Multilingual
Represent.
LASER
Evaluation

Document
Retrieval

Local
Alignment

Global
Alignment
WikiMatrix
CCMatrix
WMT/TED

Bitext
Filtering

# Results

Test Results

| **Method** | **ne-en** | | **si-en** | |
|---|---|---|---|---|
| | **1M** | **5M** | **1M** | **5M** |
| Main - Ensemble | 6.8 | 2.8 | **6.4** | 4.0 |
| Constr. - LASER *local* | **6.9** | 2.5 | 6.2 | 3.8 |
| Best (other) | 5.5 | **3.4** | 5.0 | **4.4** |