



Sketch Engine pre-processing pipelines: towards on-the-fly tokenization of user queries

Matúš Kostka, Marek Medved'

`xkostka4@fi.muni.cz`, `marek.medved@sketchengine.eu`

Lexical Computing

June 30, 2023

Pipeline

SkE pipeline

Input

'eLex' conferences aim to explore innovative developments in the field of lexicography.

Input sentence containing single quotes ' and '.

Uninorm

'eLex' conferences aim to explore innovative developments in the field of lexicography.

Uninorm provides normalization of text and converts the content into NFKC normalization form. The ' and ' becomes ". Unifies search and improves stability.

SkE pipeline

Unitok

```
,  
<g/>  
eLex  
<g/>  
,  
conferences  
aim  
to  
explore  
innovative  
developments  
in  
the  
field  
of  
lexicography  
<g/>  
,
```

Unitok splits the input text into units suitable for further computational processing. It is an important data preparation step allowing us to perform more advanced tasks.

SkE pipeline

Tag Sentences

```
<s>  
,  
<g/>  
eLex  
<g/>  
,  
conferences  
aim  
to  
explore  
innovative  
developments  
in  
the  
field  
of  
lexicography  
<g/>  
,  
</s>
```

Tag Sentences identifies the boundaries of individual sentences inside the text.

SkE pipeline

Lemma

```
<s>
. .
<g/>
eLex eLex
<g/>
. .

conferences conference
aim aim
to to
explore explore
innovative innovative
developments development
in in
the the
field field
of of
lexicography lexicography
<g/>
. .
</s>
```

Lemmatizer assigns a lemma (base form of the word) to each word form in a corpus.

SkE pipeline

Tag

```

<s>
'   POS   '
<g/>
eLex NP   Elex
<g/>
'   POS   '
conferences      NNS conference
aim   VVP   aim
to    TO    to
explore  VV  explore
innovative JJ  innovative
developments NN  development
in      IN   in
the     DT   the
field  NN   field
of     IN   of
lexicography NN  lexicography
<g/>
.      SENT .
</s>

```

Tagger assigns special labels to each token in the corpus to indicate part of speech and other grammatical categories.

SkE pipeline

- **Post-processing:** simple tag and lemma modifications

https://sketchengine.eu NN <unknown>

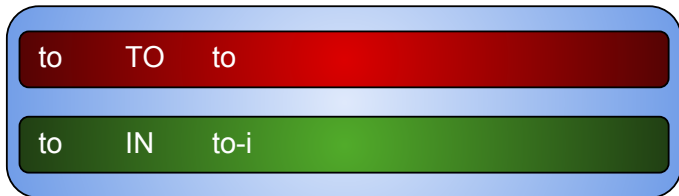
https://sketchengine.eu NP [url]-n

1st JJ 1st

1st JJ [number]-j

SkE pipeline

- **Post-processing:** assign tag 'TO' only if followed by an adverb or a verb, otherwise assign tag 'IN'



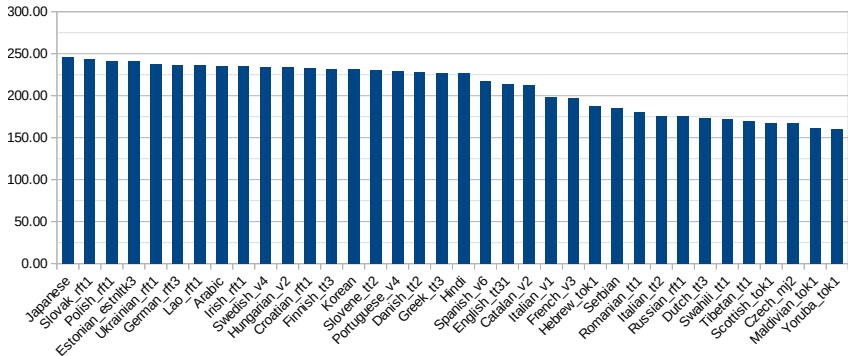
Results

Hardware setup

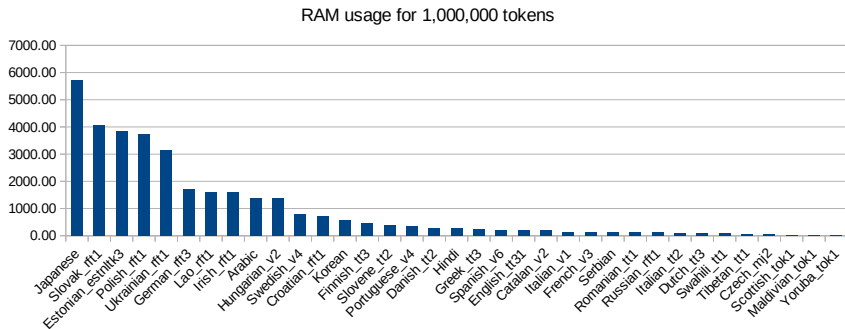
- RAM: 252 GB
- CPU: 32 threads

CPU usage

CPU usage for 1,000,000 tokens

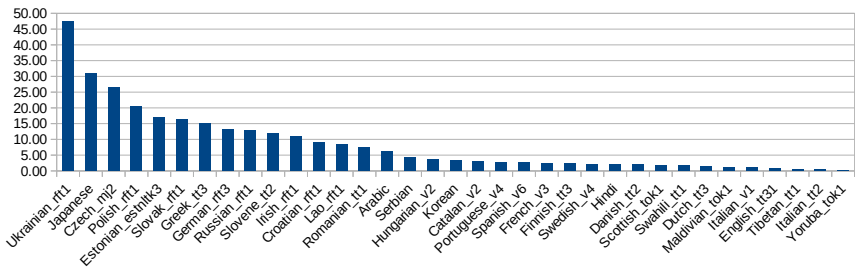


RAM usage



Execution time

Execution time for 1,000,000 tokens



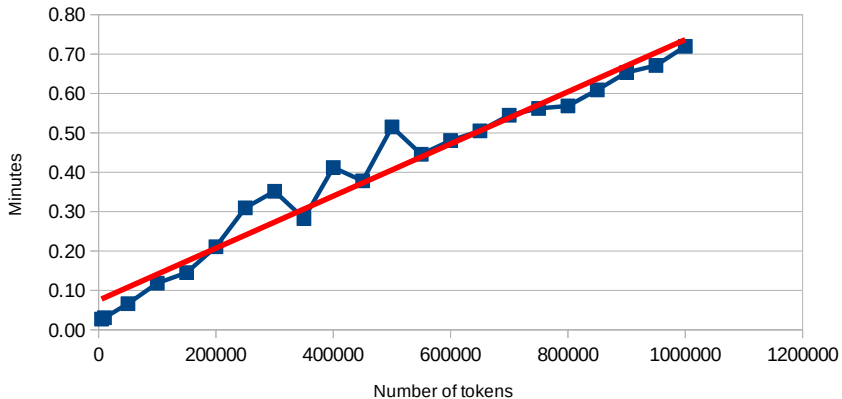
Statistics

Table: SM: Overall stats 1,000,000 tokens

	Min value	Max value	Average	Median
Execution time (sec)	11.77	2851.95	475.74	171.24
CPU usage (%)	57	224	111	105
RAM usage (GB)	0.007	5.600	0.866	0.243

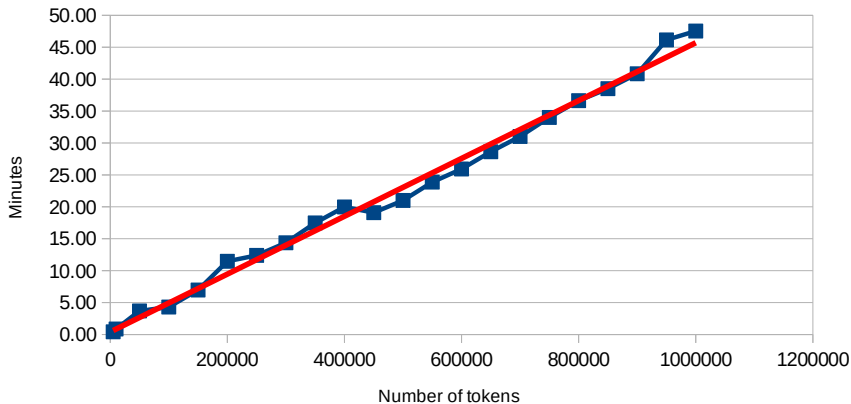
English

Linear regression English

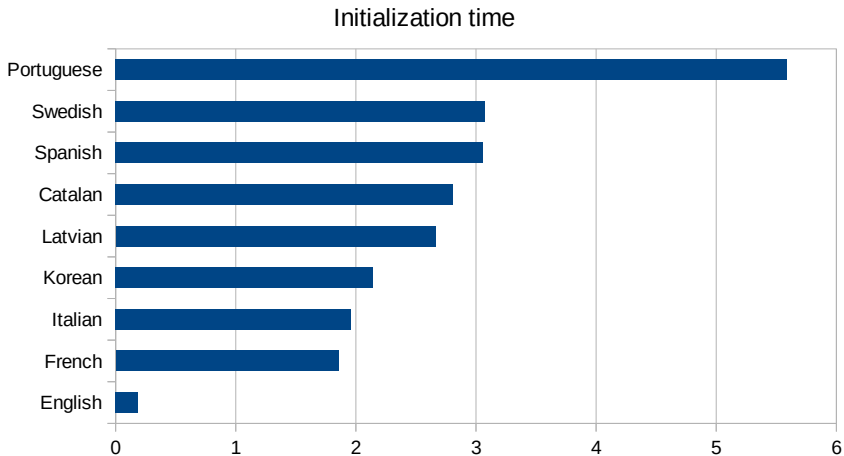


Ukrainian

Linear regression Ukrainian



Initialization time

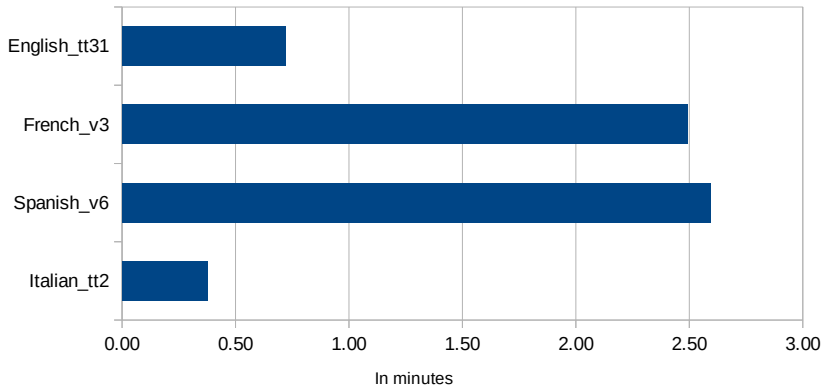


Conclusions

- the less performing pipeline inside SkE is the **Tagalog**
- RAM and CPU heavy usage pipelines are **Thai** (thai_sw1) and **Hebrew** (yap_he_v1)
- languages with **different alphabets** as Latin are usually slower
- **12%** of pipelines can be used in the multi-threaded setup
- **positive fact** is that all pipelines are in a linear relationship with the number of tokens
- **state-of-the-art** pipeline performance for the most used languages inside SkE

Top 4 used languages in SkE

Execution time for the top four languages, 1,000,000 tokens



Future work

- improve the time complexity
- implement parallel processing for more pipelines
- improve pipelines stability



**SKETCH
ENGiNE**