

## ISMLA Project: Malayalam Glosser

Thora Daneyko

March 16, 2018

## Abstract

Miau

# 1 Introduction

Bla bla bla test മലയാളി bla bla.

## 2 About Malayalam

## 2.1 General information

## 2.2 NLP challenges

### 2.2.1 Tokenization

അരിപ്പെട്ടി *aripetti* ‘rice box’ (*ari+petti*), പാൽക്കപ്പി *pālkkuppi* ‘milk bottle’ (*pāl+kuppi*) (Asher and Kumari 1997, p. 397)

- (1) മേഘം പോലെ കറുപ്പുനിറഞ്ഞൊടുക്കിയവർ ആണ്.  
*Mēgham pōle karuppuniraññōtukūtiyavar ān*.

മേഘം പോലെ കറുപ്പ് നിറത്തോട് കൂടി അവർ ആണ്.  
mēgham pōle karupp̣ nirañ-ñ-ōṭ̣ kūṭiavar āṇ̣ .  
cloud like black be.full-PSTPART-SOC with they be .

‘They are black like clouds.’ (Vēnugōpālan 2009, p. 179)

- (2) അതിന് നിനക്കെന്താ?  
*Atin<sup>u</sup> ninakkentā?*

അതിന് നിനക്ക് എന്ത് ആണ്?  
 at-in<sup>u</sup> nin-akk<sup>u</sup> ent<sup>u</sup> āṇ<sup>u</sup> ?  
 that-DAT you-DAT what be ?

‘Why do you care?’ (Moag 1994, p. 165)

## 3 The Malayalam Glosser

Bla bla

### 3.1 Transliteration

#### 3.1.1 Supported Scripts

#### 3.1.2 Transliterators

MalayalamTranscriptor

### 3.2 Morphology Generation

MorphGen

### 3.3 Tokenization

MalayalamGlosser

### 3.4 Dictionary lookup

MalayalamDictionary

#### 3.4.1 Efficiency considerations

Considering that the dictionary may be very large and that the main function of the Glosser is to look words up in this dictionary, being able to load and query it very quickly is essential for the performance of the Glosser. Hence, I experimented with a few alternatives for storing the dictionary data and investigated their efficiency. The tests elaborated below are not very exact or well-designed and were only meant to quickly assess the usefulness of the considered methods.

#### **HashMap vs. ReverseTrie**

The straightforward way to represent a dictionary as a Java object is a **HashMap**. Apart from being readily available and easy to use, querying a **HashMap** is fast. However, this also means that all entries are stored as their complete `String` representation, which may consume quite a lot of space. Considering that the inflected forms of the words share most of their characters, a trie representation seemed quite suitable and might be able to save space compared to a simple **HashMap**. Since Malayalam is exclusively suffixing, I programmed a **ReverseTrie** which reads and retrieves the strings from last to first character, in order to save

as much space as possible. A useful side effect of this is that the tokenizer does not need to look up all suffixes of a compound word in the dictionary, but can simply do a suffix search of the **ReverseTrie** to get the longest contained suffix.

In order to compare the performance of a **HashMap** and **ReverseTrie** based dictionary, I measured the memory used by the program before loading the dictionary data and after creating the **HashMap** and **Trie** (calculated as `Runtime.totalMemory() - Runtime.freeMemory()` after a `System.gc()` call). Then I let the dictionary find the longest known suffix of the test String *aviteyullatariññu* (*avite ullat<sup>~</sup> ariññu* “knew (he) was there”) 1,000,000 times and measured the time needed by a **HashMap** and **ReverseTrie** based dictionary (calculated using `System.currentTimeMillis()`). Finally, I rewrote the tokenizer to also work with a **ReverseTrie** and tested how long tokenization of a short conversation from Moag (1994) took it with the two dictionary types.

Despite the many shared suffixes, the **HashMap** was smaller than the **ReverseTrie**, taking up 8,318,164.8 bytes on average during five test runs, while the **Trie** required 12,590,051.2 bytes. However, the memory used by the **HashMap** varied greatly, ranging from only 5,160,456 to 9,801,392 bytes, while the **Trie** always consumed almost exactly the same amount of memory. This indicates that the measurements might have been verfälscht by background processes such as the garbage collection. However, the **HashMap** still seems to be considerably smaller.

As expected, the **ReverseTrie** outperformed the **HashMap** on the looped suffix search of *aviteyullatariññu*. The **Map** took an average of 999 milliseconds during five test runs, while the **Trie** only needed 312.4 ms. However, the performance of the **Trie** was very unstable, ranging from 140 to 518 ms between runs, while the **HashMap** always needed between 908 and 1049 ms, which is still much slower than the slowest suffix search of the **Trie**.

On a real Malayalam text, where only few words are long compounds such as *aviteyullatariññu*, both methods were equally fast. During 10 glossings of the Moag conversation, the **Map** based tokenization took 156.3 ms on average and the **Trie** based tokenization 161.7 ms. Both ran very stable.

All in all, the **HashMap** seems to be the better choice, since it is smaller than the **Trie** and equally fast on normal Malayalam texts. The **Trie** is faster when tokenizing long compound words, which however are not frequent enough to justify preferring it over the **HashMap**.

## File storage vs. Serialization

Loading the dictionary data into the underlying **HashMap** (or **ReverseTrie**) takes a considerable amount of time at launch. Hence, I considered serializing the **Map** or **Trie** object to be able to load it quicker. Since the Java serialization is known to be rather slow, I used the FST Fast Serialization library for my tests. I first read the dictionary data from the text file and created the **HashMap** and **ReverseTrie** from it, measuring the time needed. Then I serialized the two objects and took the time required to deserialize them.

During five test runs, parsing the text file into an object took 278.2 ms on average for the `HashMap` and 310.6 ms for the `Trie`. Deserializing the same objects required 563.4 ms on average for the `HashMap` and 339.8 ms for the `Trie`. Loading the data from a text file is thus faster than deserializing a previously created object.

The file storing the serialized `ReverseTrie` was twice as large as the file with the `HashMap`. This confirms my assertions from the previous section that the `Trie` takes more space than the `HashMap`.

### 3.5 UI Design

## 4 Conclusion

## References

- Asher, Ronald E. and T. C. Kumari (1997). *Malayalam*. Psychology Press.
- Bindu, MS and Sumam Mary Idicula (2011). “High Order Conditional Random Field Based Part of Speech Tagger and Ambiguity Resolver for Malayalam-a Highly Agglutinative Language.” In: *International Journal of Advanced Research in Computer Science* 2.5.
- Devadath, VV et al. (2014). “A Sandhi Splitter for Malayalam.” In: *Proceedings of the 11th International Conference on Natural Language Processing*, pp. 156–161.
- Gamliel, Ophira, ed. (forthcoming). *God’s Own Language. Malayalam Grammar Text Book*. Draft from July 2016.
- Jayan, Jisha P, RR Rajeev, S Rajendran, et al. (2011). “Morphological analyser and morphological generator for malayalam-tamil machine translation.” In: *International Journal of Computer Applications* 13.8, pp. 0975–8887.
- Kuncham, Prathyusha et al. (2015). “Statistical sandhi splitter for agglutinative languages.” In: *International Conference on Intelligent Text Processing and Computational Linguistics*. Springer, pp. 164–172.
- Manju, K, S Soumya, and Sumam Mary Idicula (2009). “Development of a POS tagger for Malayalam-an experience.” In: *Advances in Recent Technologies in Communication and Computing, 2009. ARTCom’09. International Conference on*. IEEE, pp. 709–713.
- Moag, Rodney F. (1994). *Malayalam: A University Course and Reference Grammar*. Austin: University of Texas, Center for Asian Studies.
- Nisha, M and PC Reghu Raj (2016). “Sandhi Splitter for Malayalam Using MBLP Approach.” In: *Procedia Technology* 24, pp. 1522–1527.
- Rajeev, RR and Elizabeth Sherly (2007). “Morph analyser for malayalam language: A suffix stripping approach.” In: *Proceedings of 20th Kerala Science Congress*.
- Sebastian, Mary Priya and G Santhosh Kumar (n.d.). “Machine Learning Approach to Suffix Separation on a Sandhi Rule Annotated Malayalam Data Set.” In:
- Suneera, CM and MT Kala (n.d.). “A Rule Based Approach For Malayalam-English Translation.” In:

Vēṇugōpālan, P., ed. (2009). *Āścaryacūḍāmaṇi. Sampūrṇamāya āṭṭaparakā-  
vum kramadīpikayum*. Tiruvanantapuram: Mārgi.