

# Finite-state morphological analysis for Marathi

---

Vinit Ravishankar <sup>1</sup> Francis M. Tyers <sup>2</sup>

<sup>1</sup>University of Malta, Malta

<sup>2</sup>Higher School of Economics, Moscow

- No decent, wide coverage, open-source analyser for Marathi (that we know of)

- No decent, wide coverage, open-source analyser for Marathi (that we know of)
- We wanted to enable machine translation for a language with traditionally rubbish MT

- No decent, wide coverage, open-source analyser for Marathi (that we know of)
- We wanted to enable machine translation for a language with traditionally rubbish MT
- Who doesn't like morphological analysers?

1. Apertium
2. Marathi
3. Approach
4. Evaluation

**Apertium**

---



**Apertium**

- Rule-based machine translation pipeline
- 46 translation pairs
- Chunking with XML rules that translate to inhumane `sed` chains
- Modular: you can use *hfst* instead of *Ittoolbox* for morph analysis



**Apertium**

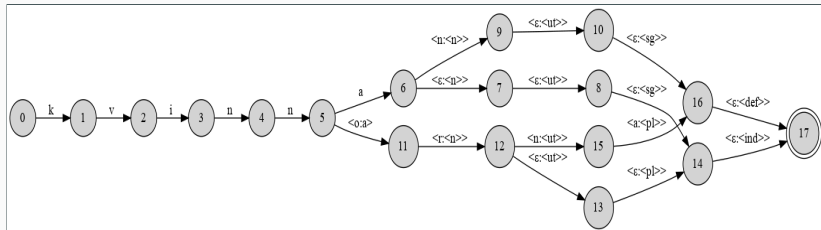
- Rule-based machine translation pipeline
- 46 translation pairs
- Chunking with XML rules that translate to inhumane `sed` chains
- Modular: you can use *hfst* instead of *Ittoolbox* for morph analysis (but we didn't)

---

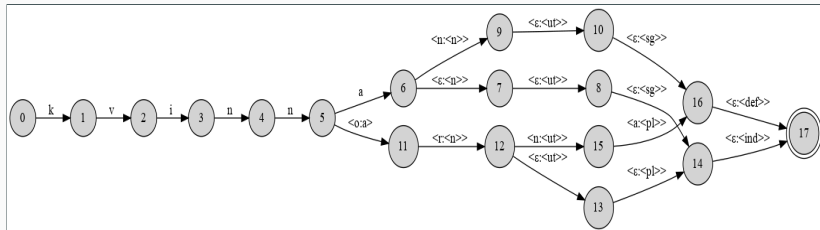
#apertium on irc.freenode.net



## Apertium's pet morphological analyser



## Apertium's pet morphological analyser



1. Define paradigms in XML
2. Add entries
3. Goto (1)

```
<pardef n="मुल/गण__n_m">
<e><p><l>गण</l><r>गण<s n="n"/><s n="m"/><s n="sg"/><s n="
  nom"/></r></p><par n="__emph"/></e>
<e><p><l>गण</l><r>गण<s n="n"/><s n="m"/><s n="sg"/><s n="
  obl"/><j/></r></p><par n="__clt"/></e>
<e><p><l>गण</l><r>गण<s n="n"/><s n="m"/><s n="sg"/></r></p>
  ><par n="__sg_case"/></e>
<e><p><l>गण</l><r>गण<s n="n"/><s n="m"/><s n="pl"/><s n="
  nom"/></r></p><par n="__emph"/></e>
<e><p><l>गण</l><r>गण<s n="n"/><s n="m"/><s n="pl"/><s n="
  obl"/><j/></r></p><par n="__clt"/></e>
<e><p><l>गण</l><r>गण<s n="n"/><s n="m"/><s n="pl"/><s n="
  obl"/><j/></r></p><par n="__pl_case"/></e>
</pardef>
```

```

<pardef n="kvinn/a__n_ut">
<e><p><l>a</l> <r>a<s n="n"/><s n="ut"/><s n="sg"/><s n="
    ind"/></r></p></e>
<e><p><l>an</l> <r>a<s n="n"/><s n="ut"/><s n="sg"/><s n="
    ind"/></r></p></e>
<e><p><l>or</l> <r>a<s n="n"/><s n="ut"/><s n="pl"/><s n="
    def"/></r></p></e>
<e><p><l>orna</l> <r>a<s n="n"/><s n="ut"/><s n="pl"/><s n
    ="def"/></r></p></e>
</pardef>

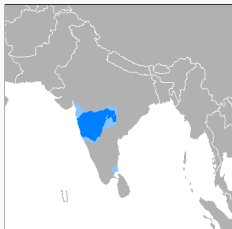
<e lm="kvinna"> <i>kvinn</i><par n="kvinn/a__n_ut"/></e>

```

```
$ echo "kvinnorna" | lt-proc swe.automorf.bin  
^kvinnorna/kvinna<n><ut><pl><def>/kvinna<n><ut><pl><def  
  ><compound-R>$
```

**Marathi**

---



- /məɾɑtʰi:/ ; Indo-European language
- Spoken mainly in Maharashtra (cities == Mumbai, Pune)
- ~71m native speakers, ~82% literate
- Script - Devanagari / देवनागरी + extra chars

**Partially fusional:**



## Partially fusional:

- Paradigmatic stem alterations for non-nominative nouns:

*mulgā* → *mulā* 'boy' ; *maitriṅ* → *maitriṅī* 'girlfriend'

- Merged aspect and gender:

*basto*<vblex><impf><p3><m><sg> 'he sits'

*baslā*<vblex><perf><p3><m><sg> 'he sat'

**Partially agglutinative:**

## Partially agglutinative:

- Cases and postpositions agglutinate with the oblique:

*mulgā* → *mulā* → *mulā-lā*<sup>DAT</sup>

- Verbs can get longish:

bas le lyā m̄ samor ūn ac  
sit PERF OBL PL in front of ABL PART

‘From ahead of the people that were sitting (and from nowhere else)’

## 'Layered' morphology

Three potential layers for nouns (and gerunds)

1. Base noun/gerund:

*mulgā*<sup>NOM</sup> 'boy'

2. (Optional) case suffixes on the oblique:

*mulā*<sup>OBL</sup>-*lā*<sup>DAT</sup> 'to the boy'

3. Postpositions:

*mulā*<sup>OBL</sup>-*cyā*<sup>GEN</sup>-*samor*<sup>POST</sup> 'in front of the boy'

## Cases and postpositions

No clear distinction between the two, but several 'tests':

	<b>Cases</b>	<b>Postpositions</b>
<b>Attachment</b>	Oblique noun	Oblique noun/genitive
<b>Free morpheme</b>	Impossible	Possible
<b>Inflectable</b>	Yes	No

## Cases and postpositions

No clear distinction between the two, but several 'tests':

	<b>Cases</b>	<b>Postpositions</b>
<b>Attachment</b>	Oblique noun	Oblique noun/genitive
<b>Free morpheme</b>	Impossible	Possible
<b>Inflectable</b>	Yes	No
<b>Tradition</b>	Ugh	Ugh

## Quirks

- All tests have exceptions
- Pronouns are ugly/fusional

- Four layers can exist:

*mulā*<sup>OBL</sup>-*cyā*<sup>GEN</sup>-*var*<sup>POST</sup>-*ūn*<sup>ABL</sup> 'from above the boy'

- Cases can/have to be pushed to layer 3 (sometimes):

*mājhyā*<sup>OBL.GEN</sup>-*hūn*<sup>ABL</sup> 'from me'

- Multiple 'participles' with very specific uses (-at, -atā, -tānā, -tānṃā, -ṇe, -ū, -ūn, ...)
- Multiple negative verbs
- Nobody agrees on names:  
-ṇār(ā):  
<pros> (Wali) = <ptcp><fut> (Navalkar) = <agt> (Masica) = ☹ (us)
- Can't rip better-resourced languages like Hindi off
- Dire need for standards



## HOW STANDARDS PROLIFERATE: (SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC.)

SITUATION:  
THERE ARE  
14 COMPETING  
STANDARDS.

14?! RIDICULOUS!  
WE NEED TO DEVELOP  
ONE UNIVERSAL STANDARD  
THAT COVERS EVERYONE'S  
USE CASES.



SOON:

SITUATION:  
THERE ARE  
15 COMPETING  
STANDARDS.

## Approach

---

- (Partially) scraped from an analyser created at the International Institute of Information Technology, Hyderabad
- Modified to fit a consistent paradigm system
- Verbal paradigms created, nominals significantly extended - hard-coded postpositions replaced with joins
- Manually adding entries to improve coverage

- (Partially) scraped from an analyser created at the International Institute of Information Technology, Hyderabad
- Modified to fit a consistent paradigm system
- Verbal paradigms created, nominals significantly extended - hard-coded postpositions replaced with joins
- Manually adding entries to improve coverage
- Caffeine

### Grey areas!

- Non-genitive case information is on the root morpheme:

^मुलाला/मुलगा<n><m><sg><dat>\$ 'to the boy'

- Genitive cases and postpositions are treated as separate morphemes:

^मुलासमोर/मुलगा<n><m><sg><obl>+समोर<post><adv>\$ 'in front of the boy'

- We pretend pronouns are more agglutinative than they are:

^माझा/मी<prn><p1><mf><sg>+चा<gen><m><sg>\$ 'my (single mularculine thing)'

- Treat causativity as non-productive (uh oh!)
- Ignore a tonne of derivational morphology; upset people
- Rely on intuition to separate loanwords from foreign words
- Light verbs - fake words are now real words, eg. *māhit asṇe* 'to know'

## Evaluation

---

## Gold standard:

```
^wound/wound<n><sg>/wind<vblex><pp>/wind<vblex><past>  
/wound<vblex><inf>/wound<vblex><pres>$
```

## Analyser:

```
^wound/wound<n><sg>/wound<vblex><inf>/wound<vblex><pres>  
/wound<adj>$
```

$$P = \frac{3}{3 + 1} = 0.75$$



## Gold standard:

```
^wound/wound<n><sg>/wind<vblex><pp>/wind<vblex><past>  
/wound<vblex><inf>/wound<vblex><pres>$
```

## Analyser:

```
^wound/wound<n><sg>/wound<vblex><inf>/wound<vblex><pres>  
/wound<adj>$
```

$$R = \frac{3}{3 + 2} = 0.6$$

	<b>Tokens</b>	<b>Coverage</b>	<b>Mean ambig.</b>
<b>Wikipedia</b>	4.0M	80.2	1.7
<b>Bible</b>	751k	80.7	1.9

## Results

	<b>Tokens</b>	<b>Coverage</b>	<b>Mean ambig.</b>
<b>Wikipedia</b>	4.0M	80.2	1.7
<b>Bible</b>	751k	80.7	1.9

	<b>Precision</b>	<b>Recall</b>
<b>Known tokens</b>	0.97	0.97
<b>All tokens</b>	0.97	0.71

- Bring coverage up to the 90s
- Write extensive documentation and guidelines that can apply to other Indic languages
- Use the analyser on a Universal Dependencies treebank; evaluate
- Create a proof-of-concept translation pair with moderate lexical coverage

- Bring coverage up to the 90s
- Write extensive documentation and guidelines that can apply to other Indic languages
- Use the analyser on a Universal Dependencies treebank; evaluate
- Create a proof-of-concept translation pair with moderate lexical coverage (Marathi - Crimean Tatar???)

- Bring coverage up to the 90s
- Write extensive documentation and guidelines that can apply to other Indic languages
- Use the analyser on a Universal Dependencies treebank; evaluate
- Create a proof-of-concept translation pair with moderate lexical coverage (Marathi–Crimean Tatar???) (Marathi - Gujarati)

**Thanks!**