

## Objective of this Konkani Phonetic Transcription System is to built a resource that would phonetically transcribe integers in Konkani language and generate their corresponding written form in the Devanagari script

Algorithms

### Abstract

In this work, we have built a resource that would phonetically transcribe Konkani Integers and generate their written form in the Devanagari script. The algorithm developed in this work made robust enough to automatically give a written form of any Konkani numeral in the Devanagari script, along with its equivalent IPA transcription. We have tried to keep the phonetic transcriptions as closer to their natural pronunciations. This is done for the purpose of capturing the general tendency of the language. So, for example, while the number '8' आठ [ath] is written with an aspirated retroflex consonant  $\sigma$  [t<sup>h</sup>], the final consonant [t<sup>h</sup>] is heard without aspiration in the actual speech. This loss of aspiration at word final position generally happens across all the consonants of the language, in the Konkani varieties spoken in the state of Goa in India.

## Introduction

India is a multilingual country having various languages and dialects. An Indian language -Konkani which belongs to Indo-Aryan language family, is an official language of the state of Goa in India. Anyone who understands a certain script can read any language written in that script. However, when it comes to the reading (pronunciation) of numerals, one needs to follow certain rules. Following are some representative examples in Konkani where we can see a combination of text with numerals. Such data instances in a text corpus could pose a big challenge for any system that aims to transcribe text accurately.

- . "ता. १ फेब्रुवारी २०२२" [ta. ek februvari don həd͡zar baviːs] ('dtd. 1<sup>st</sup> February 2022'). The above character and numeral combination refers to a specific date and a year. While 9 'one' in १ फेब्रुवारी '1st February' might be easily transcribed as ek (februvari) by any system, the numeral 2022 which is a year, has to be read and transcribed as दोन हजार बावीस [don hədzar baviːs] "two thousand twenty-two" and not as दोन शुन्य दोन दोन [don ʃunj<sup>ə</sup> don don] "two zero two two".
- 2. "सकाळी 08:00 ते 10:00 वरांमेरेन" [səkalī aːth tɛ dha vərāmeren] '(from) morning 8 to 10 a.m.' (Lit. morning 8 to 10 hours till). This phrase specifies a certain time of the day. The system needs to acknowledge this context of time and generate a string that reads the numerals as hours (and minutes in some other temporalcontext).
- "गोंयचें क्षेत्रफळ 3701 चौखण किलोमिटर आसा." [gõjtJɛ̃ kʃetrə'fəl tiːn həd͡zar satʃɛ ek t͡soʊk<sup>h</sup>əŋ kilomitar asa ] ('The (total) area of Goa is 3701 sq. kms.'). The numeral in the above sentence specifies the area of the region of Goa. The numbers should be read as one unit, i.e., as तीन हजार सातशे एक [tiːn həd͡zar satʃɛ ek] (Lit. "Three thousand Seven Hundred One") and not as individual numbers
- 4. "माशेलाचो पिन कोड 403 107." [maʃɛlat͡sɔɔ pin'kod t͡ʃar ʃunjə tiːn ek ʃunjə saːt] ('The pin code of Marcela is 403 107'). Postal Index Number (PIN or simply PIN Code) refers to the six-digit number used by India Post in its postal code system. More commonly, the numbers indicating such a code are read by spelling out the numerals as discrete units.
- 5. "ताचो फोन नंबर 9850 403 107" [ tat͡ʃɔ fon nəmbər nəv aːtʰ pãːt͡s ʃunjə t͡ʃaːr ʃunjə tin eːk ʃunjə saːt / ta͡tʃɔ fon nəmbər nain eṭ faiʊ d͡ziro foːr d͡ziro t<sup>h</sup>ri vən d͡ziro sɛvən]('His phone number is 9850 403 107'). Phone numbers can be read differently by different speakers. However, reading the numbers as discrete units is a good way to spell out the long number string.

The table 1 represents the position mapping rules that were used for phonetic transcription.

Sr No	Integer length	IPA	Devanagari Transcription
1	3	[3]	शे
2	4	[hədzar]	हजार
3	6	[lak <sup>h</sup> ]	लाख
4	8	[koti]	कोटी
5	10	[ərəb]	अरब
6	12	[k <sup>h</sup> ərəb]	खरब

Table 1. Position mapping rules.

### **Scope of the Work**

We have made an effort to develop an automatic system for Konkani language that gives the phonetic as well as Devanagari transcription of a given integer. This is the first kind of work that aims to automatically transcribe Konkani numerals appearing in different contexts into the officially recognised Devanagari script along with the pronunciation of the numerals (in IPA).

Konkani Integer Phonetic Transcription System

<sup>1</sup>Discipline of Computer Science and Technology, Goa Business School, Goa University, India <sup>2</sup>Department of Konkani, Govt. College of Arts Science and Commerce Quepem, Goa India

**Data:** *integer* **Result:** *transcription\_text*  $y \leftarrow$  "";  $X \leftarrow input;$  $N \leftarrow len(X);$ if N > 12 then  $R \leftarrow assign\ last\ 11\ didgits;$  $L \leftarrow X/10^{11}$ ;  $y \leftarrow left\_trans(L) + pos\_mapping(12) + right\_trans(R)$  $y \leftarrow right\_trans(X);$ Algorithm 1: Integer transcription. **Data:** integer **Result:** transcription\_text  $y \leftarrow "";$  $X \leftarrow input;$  $N \leftarrow len(X);$ if N < 2 then if N == 0 then  $y \leftarrow$  "": else  $y \leftarrow int\_mapping();$ end if N < 3 then if N == 100 then  $y \leftarrow int\_mapping()$ else  $y \leftarrow int\_mapping() + pos\_mappings(3) + right\_trans(R)$ end else if  $N \leq 5$  then  $y \leftarrow int\_mapping() + pos\_mapping(4) + right\_trans(R)$ else if  $N \leq 7$  then  $y \leftarrow int\_mapping() + pos\_mapping(6) + right\_trans(R)$ 

else if  $N \leq 9$  then  $y \leftarrow int\_mapping() + pos\_mapping(8) + right\_trans(R)$ else  $y \leftarrow int\_mapping() + pos\_mapping(10) + right\_trans(R)$ end end end ena

end

Algorithm 2: right\_transcription.

**Data:** *integer* **Result:** transcription\_text  $y \leftarrow "$  $X \leftarrow input;$  $N \leftarrow len(X);$ if N > 11 then  $R \leftarrow assign \ last \ 9 \ digits;$  $L \leftarrow X/10^9$ ;  $y \leftarrow left\_trans(L) + pos\_mapping(10) + right\_trans(R)$ else  $y \leftarrow right\_trans(X)$ end

**Algorithm 3:** left\_transcription.

SIGUL 2023 : 2nd Annual Meeting of the Special Interest Group on Under-resourced Languages : a Satellite Workshop of Interspeech 2023, Dublin, Ireland, 18-20 August 2023

# Edna Vaz Fernandes<sup>23</sup> Hanumant Redkar<sup>1</sup> Teja Kundaikar<sup>1</sup> Ramdas Karmali <sup>1</sup> Jyoti D. Pawar<sup>1</sup>

## The Transcription System

This work presents an automatic phonetic transcription system for Konkani Integers. The system takes an integer as an input and generates its representation in word (the written form) along with its phonetic transcription. Figure 1 diagrammatically presents the transcription system.



languages, 1950.

[5] The Constitution of India, Eighth schedule, Article(s): 344(1) and 351, Description: Official