## **Xrphonetic: Akshara-based Phonetic String Similarity**

**Anonymous ACL submission** 

#### Abstract

001 Establishing String Similarity based on phonetics has been widely used in information re-002 trieval systems to identify differently spelled 004 but similar-sounding words. Another common 005 application often involves calculating a similarity score between two words coming from two 007 different sources which possibly can be two different spelling representations of the same word. A very interesting and common subset of this is estimating the phonetic similarity of two 011 words that are transliterated to Roman script from a different language. For such a use case, 012 it would be more effective if we can use the knowledge of the nature of the concerned writing system from which the words originated as people usually tend to carry over the nuances of the underlying writing system during transliteration. We propose Xrphonetic, a novel phonetic 019 similarity algorithm, for words transliterated to Roman script from languages using Abugidabased scripts by treating aksharas as the most fundamental atomic unit of words with consonant and vowel phonemes as its further sub-024 atomic units, and by having weighted phoneme mappings to get a more continuous spectrum of phonetic similarity.

## 1 Introduction

040

Phonetic string similarity is used to identify strings with different spellings but similar pronunciations.
A lot of proper nouns have multiple valid spellings but with similar pronunciations. For example, at various places, people might need to communicate their name verbally which may sometimes lead to inconsistencies of name spellings in different documents of the same person. In such a scenario, it is important to have a method to compare two names with similar pronunciations, irrespective of their spelling.

A lot of times words being compared may have been represented in the Roman script but have origin in a different language, and differences in spelling might have crept in during the transliteration process. Since different languages use scripts with different underlying writing systems, it can be effective to make use of the rules and structures of the underlying writing system of the originating language and script during comparison, as most of the time people tend to carry over these nuances during the transliteration process.

043

044

045

047

051

057

059

060

061

062

063

064

065

067

068

069

070

071

072

073

074

075

076

077

079

Broadly writing systems can be classified into Alphabets, Syllabaries, Logographies, Abjads, and Abugidas and have been covered extensively both from linguistic and computational points of view (Coulmas, 2003; Daniels and Bright, 1996; Sproat, 2000; Sproat, 2002; Sproat, 2003).

Emeneau, 1956, showed that most Indian languages use scripts derived from Brahmi, which can be classified as an Abugida-based system. The basic unit in these scripts is *Akshara* which more or less corresponds to a syllable, but the mapping between syllables and *Aksharas* is not exactly oneto-one (Singh, 2006). Usually, *Akshara* consists of a consonant followed by one or more vowels represented as diacritics.

Singh et al., 2007 highlighted the advantages of exploiting the characteristics of the Abugidabased writing system to enhance the performance of fuzzy text search for Indian languages in their corresponding scripts. We can exploit these same characteristics even while comparing words written in Roman script which is an Alphabetic writing system, if they have origin in a language with an Abugida-based writing system since people tend to preserve these characteristics regardless of the script of representation.

In this paper, we focus on the phonetic similarity of words transliterated to Roman script from languages that use an Abugida-based writing system.

## 2 Related Works

Soundex (Odell and Russell, 1918) is the most commonly used phonetic coding scheme. It converts the name into a four-character phonetic code with the aim to have the same code for similar-sounding names. It was mainly designed to match phonetically similar English surnames. There have been various other phonetic coding schemes developed ever since like Phonix (Gadd, 1988; Gadd, 1990), Double Metaphone (Philips, 2000), and Caverphone (Hood, 2002). There have been some works on modification of the Soundex to make it suitable for Indian languages (Chaware and Rao, 2011; Shah and Singh, 2014; Gautam et al., 2019).

083

087

091

100

101

102

103

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

Editex is a phonetic distance measure that combines the properties of edit distances with the lettergrouping strategy used by Soundex and Phonix (Zobel and Dart, 1996). Another method integrating approximate string matching with phonetic string similarity was presented by Ferri et al., 2018.

Phonetic codes derived from the above approaches usually do not have any resemblance to the actual phonetics of the words. While (Zobel and Dart, 1996) proposed the idea of leveraging string-to-pronunciation conversion algorithms to first convert the string into a phonetic representation and then do the matching on strings of phonemes, Kondrak (Kondrak, 2000; Kondrak, 2003) highlighted the idea of phonetic alignment and similarity scoring methods on the basis of multi-valued articulatory phonetic features.

Most of the works highlighted above while assessing the phonetic similarity don't take into account that many times these words might have origins in different writing systems. To bridge this gap, there have been works focusing on making use of the characteristics of the underlying writing systems for better estimation of phonetic similarity. Particularly, (Singh et al., 2007; Gupta et al., 2014) highlighted the advantages of using Akshara-based phonetic similarity for the native Indian scripts which fall under the Abugida-based writing system.

Furthermore, it is quite common for words having origin in an Abugida-based writing system to be represented in the Roman script which is an Alphabet-based writing system. Hindex (Prabhakar et al., 2021) has tried to exploit these ideas to provide a phonetic coding and similarity approach for words transliterated to Roman script from Indian languages utilizing the character-wise mapping of Soundex and Editex, and Levenstein edit distance (Levenshtein et al., 1966) as an approximate string matching algorithm. Hindex (Prabhakar et al., 2021) and other approaches have failed to fully leverage the characteristics of the Abugida-based writing systems in a unified manner. Even while they have implicitly highlighted Aksharas and phonemes as the primary building block of sound units for Abugida-based writing systems, the consistent and primal treatment across some of the most critical sub-tasks like segmentation, code mappings or even computing similarity scores is missing. For example, Hindex still leverages Alphabets to prepare the code mappings and scoring, which one way contradicts the whole notion of using *Aksharas* or phonemes as the most atomic building blocks for sound.

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

155

156

157

158

159

160

161

162

163

164

165

166

168

169

170

171

172

173

174

175

176

177

178

179

180

181

To fully utilize the highly phonetic nature of the Abugida-based writing systems for calculating the phonetic similarity of words with origin in such languages, we propose *Xrphonetic*, which treats *akshara* as the most fundamental atomic unit of words even if they are represented in a Roman script. Further, we contribute two novel improvements to the existing works:

- Historically, all existing code mapping-based approaches have taken an unweighted approach to mappings. Given that in the real world, different units of sound are not always equally similar or dissimilar to each other, we propose the introduction of code mappings which are grouped into three different groups each associated with an experimentally determined similarity score. Not only it provides a more natural way to arrive at the overall phonetic similarity score which is a continuous value, but it also helps in striking out a much better balance between the often contradictory precision and recall metrics.
- Building on top of *akshara*-based segmentation and phoneme mappings, Xrphonetic allows any edit-distance similarity algorithm to be used to arrive at the final similarity score. In addition to supporting the generally used edit-distance algorithms, we also introduce a new scoring function for calculating phonetic similarity for use cases requiring higher precision.

## 3 Xrphonetic

For languages that use an Abugida-based script *Akshara* forms the most fundamental and atomic unit of sound in a word which in turn consists of a

182consonant phoneme and a vowel phoneme (usually<br/>represented as diacritics unless at the start of the<br/>word). And these *aksharas* are usually written pho-<br/>netically consistent meaning two different sounds<br/>are rarely represented by the same *akshara*. So, for<br/>words originating in these languages, even if they<br/>are represented in Roman script, if we can prop-<br/>erly segment them into their constituent *aksharas*,<br/>we can utilize the highly phonetic nature of these<br/>scripts to make phonetic comparisons of the words.

Even though these languages when written in their native script are highly phonetic in nature, for words with origin in these languages when represented in Roman script, there can be multiple characters or groups of characters that can be used to represent the same or very similar sound. Hence, we prepare a list of consonant and vowel phoneme mappings with an assigned similarity score for each pair that depends on their perceived phonetic similarity and use it to calculate the similarity score of two words.

Once, we have the segmentation and the mappings, either simply a weighted edit distance algorithm (Levenshtein et al., 1966; Wagner and Fischer, 1974) can be used to compare the words treating *aksharas* as the fundamental unit instead of a character, or one can also use the specialized phonetic scoring function that we introduce in subsequent sections.

#### 3.1 Segmentation

192

193

195

196

197

199

201

204

207

210

211

212

213

214

215

216

217

218

221

230

We segment a string into its constituent *aksharas*, which is further subdivided into consonant and vowel phonemes. The motivation behind the twostep segmentation process is to have the scope of using different weights to consonant and vowel phoneme similarity scores while calculating the akshara similarity score since consonants usually carry a larger weight in determining the sound. Barring a few exceptions, consonant phonemes are always a single character whereas vowel phonemes can be composed of multiple characters. Each Akshara can be a single consonant phoneme, a single vowel phoneme, or a combination of a consonant and a vowel phoneme.

Alphabets "h", and "y" can act as either a consonant or a vowel depending on the context, and are being treated accordingly for segmentation.

We use regular expressions to segment the string into *aksharas* and phonemes.

A few examples are shown in Table 1. In the sec-

String	Segmentation
vaishali	[('v', 'ai'), ('sh', 'a'), ('l', 'i')]
nayak	[('n', 'aya'), ('k', '')]
tamatar	[('t', 'a'), ('m', 'a'), ('t', 'a'), ('r', ")]
byaz	[('b', 'ya'), ('z', '')]
akanksha	[(``, `a`), (`k`, `a`), (`n`, ``), (`ksh`, `a`)]

Table 1: Segmentation Examples

ond column, each element in parentheses inside the list represents an Akshara, and the sub-elements represent consonant and vowel phonemes respectively. 232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

## 3.2 Mappings

We define three different groups of consonant and vowel mappings based on their perceived phonetic similarity with an empirically determined similarity score for each group. While comparing the phonemes, we check if the phoneme pair belongs to any of the phoneme mappings and use the corresponding similarity score.

Sometimes some characters maybe repeated or alphabet "h" maybe present without changing the sound, we also check if the phoneme pair belongs to any of the mappings after removing these and accordingly assign the score.

Details of consonant and vowel mappings are given in Appendix.

## 3.3 Specialized Phonetic Scoring Function

For comparing each akshara-pair, we assign a score of 1.0 for an exact match (exact match for both consonant and vowel pairs), else we score consonant and vowel phonemes separately. We then look if the consonant and vowel phoneme pairs can be matched according to some consonant and vowel mappings defined above. We assign some empirically derived scores for each level of mapping. If the pair doesn't match at any level of mapping, we assign a zero score in case of a consonant phoneme pair while we use Weighted Levenshtein similarity in case of vowel phoneme pair, using the same vowel mapping scores as above. Finally, we take a weighted average of consonant and vowel phoneme pairs' scores to get the Akshara pair score, giving an empirically determined higher weight to consonant pair.

We introduce a specialized phonetic similarity scoring approach which is more suitable in places where more precision is required. The idea is that

when talking about sound units even one highly dis-272 similar sound unit can change completely change 273 the sound of the word, and hence should bring 274 down the score in a non-linear way. To calculate the 275 final string similarity score, we divide the akshara similarity scores into matches and non-matches 277 keeping the threshold for a match at 0.85. Such a 278 split is to avoid scores averaging out in the case of 279 strings with just one akshara phonetic mismatch with lots of aksharas matching. We can combine 281 the scores of all matches and all non-matches separately, using either the mean or product of scores. The final similarity score is calculated as a mean of matches and non-matches scores. 285

## 4 Experiments and Results

287

290

291

294

295

297

299

310

311

313

315

316

317

318

319

We prepared two separate train and test datasets of twenty thousand word pairs each to evaluate our proposed algorithm and to determine optimal values of the hyperparameters empirically. We are using f1-score as the evaluation metric for all the experiments. The code for the Python implementation of the algorithm, synthetic dataset generation and evaluation, train and test datasets will be made available as open-source.

#### 4.1 Dataset Preparation

To prepare the datasets, we have used Dakshina Dataset (Roark et al., 2020) as the source dataset and used IndicXlit (Madhani et al., 2022) transliteration models to generate transliteration pairs.

We have used the lexicons dataset for the ten different Abugida languages available in the Dakshina dataset to generate one thousand positive and one thousand negative pairs for each language. For generating a positive sample, we take a random Indic script sample from lexicons, generate the top five Roman transliteration variants for it using IndicXlit model, and take any two variants randomly. For generating a negative sample, we take a random Roman script sample from lexicons, get another random Roman script sample from lexicons that is within one Levenshtein distance from it, generate the top five Indic script transliteration variants for words, and accept it as a negative sample if none of their top five transliteration variants match.

There were three different lexicon datasets for dev, train, and test in the Dakshina dataset. We have used the training dataset for generating the training dataset and used the dev and test set for generating the test set.

Algorithm	Precision	Recall	f1-score
Xrphonetic	0.7103	0.8507	0.7742
Double Metaphone	0.6455	0.7761	0.7048
Caverphone	0.6662	0.7101	0.6874
Soundex	0.5956	0.7990	0.6825
Phonix	0.5948	0.7961	0.6809
Editex	0.5138	0.7035	0.5938
Weighted Hindex	0.4450	0.5737	0.5012
Hindex	0.8422	0.2424	0.3764

Table 2: Evaluation Results

We have used the generated training dataset for finding the optimal hyperparameters and used the test dataset for evaluation. 321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

346

347

348

349

350

352

353

354

355

#### 4.2 Optimal Hyperparameters Search

For calculating the optimal hyperparamets, we took a range of values for each hyperparameter and used the training dataset to get the optimal value based on F1-score metric.

## 4.3 Results

We have evaluated Xrphonetic against Soundex, Double Metaphone, Phonix, Caverphone, Editex, Hindex and Weighted Hindex. For Soundex, Double Metaphone, Phonix and Caverphone, we are considering exact match of the phonetic code. For Hindex, we have considered exact match of phonetic codes and a threshold of 2 for the Edit distance. For Weighted Hindex, we have considered a threshold of 1 for the edit distance for phonetic codes and 2 for weighted edit distance for the original words. For Editex and Xrphonetic, we have considered a threshold of 0.85 on normalized similarity scores. We observe that Xrphonetic outperforms all the evaluated algorithms on F1-score metric as shown in Table 8.

## 5 Limitations

Xrphonetic provides a holistic way of providing similarity scores for Abugida-based word pairs using the underlying *aksharas* even if they are represented in Roman script. In order for Xrphonetic to work effectively, the correct identification and treatment of aksharas is paramount and requires a strong understanding of the underlying scripts. This creates a strong inherent dependency on language experts when one tries to provide similarity scores for words across different Abugida-based languages.

#### References

- Sandeep Chaware and Srikantha Rao. 2011. Rule-based phonetic matching approach for hindi and marathi. *International Journal of Research in Social Sciences*, 1(1):26–41.
- Florian Coulmas. 2003. Writing systems: An introduction to their linguistic analysis. Cambridge University Press.
- Peter T Daniels and William Bright. 1996. *The world's* writing systems. Oxford University Press on Demand.
- Murray B Emeneau. 1956. India as a lingustic area. *Language*, 32(1):3–16.
- Junior Ferri, Hegler Tissot, and Marcos Didonet Del Fabro. 2018. Integrating approximate string matching with phonetic string similarity. In Advances in Databases and Information Systems: 22nd European Conference, ADBIS 2018, Budapest, Hungary, September 2–5, 2018, Proceedings 22, pages 173–181. Springer.
- TN Gadd. 1988. 'fisching fore werds': phonetic retrieval of written text in information systems. *Program*, 22(3):222–237.
- TN Gadd. 1990. Phonix: The algorithm. *Program*, 24(4):363–366.
- Vishakha Gautam, Aayush Pipal, and Monika Arora. 2019. Soundex algorithm revisited for indian language. In International Conference on Innovative Computing and Communications: Proceedings of ICICC 2018, Volume 2, pages 47–55. Springer.
- Sandeep Gupta, Arun Pratap Srivastava, and Shashank Awasthi. 2014. Fast and effective searches of personal names in an international environment. *Int J Innov Res Eng Manag*, 1.
- David Hood. 2002. Caverphone: Phonetic matching algorithm. *Technical Paper CTP060902, University* of Otago, New Zealand.
- Grzegorz Kondrak. 2000. A new algorithm for the alignment of phonetic sequences. In 1st Meeting of the North American Chapter of the Association for Computational Linguistics.
- Grzegorz Kondrak. 2003. Phonetic alignment and similarity. *Computers and the Humanities*, 37:273–291.
- Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals.
  In *Soviet physics doklady*, volume 10, pages 707–710.
  Soviet Union.
- Yash Madhani, Sushane Parthan, Priyanka A. Bedekar, Ruchi Khapra, Vivek Seshadri, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh M. Khapra. 2022. Aksharantar: Towards building open transliteration tools for the next billion users. ArXiv, abs/2205.03018.

Margaret Odell and Robert Russell. 1918. The soundex coding system. US Patents, 1261167:9.	409 410
Lawrence Philips. 2000. The double metaphone search algorithm. <i>C/C++ users journal</i> , 18(6):38–43.	411 412
Dinesh Kumar Prabhakar, Sukomal Pal, and Chiranjeev	413
Kumar. 2021. Query expansion for transliterated text	414
retrieval. <i>Transactions on Asian and Low-Resource</i>	415
<i>Language Information Processing</i> , 20(4):1–34.	416
Brian Roark, Lawrence Wolf-Sonkin, Christo Kirov,	417
Sabrina J. Mielke, Cibu Johny, Isin Demirsahin, and	418
Keith Hall. 2020. Processing South Asian languages	419
written in the Latin script: the Dakshina dataset. In	420
<i>Proceedings of the 12th Language Resources and</i>	421
<i>Evaluation Conference</i> , pages 2413–2423, Marseille,	422
France. European Language Resources Association.	423
Rima Shah and Dheeraj Kumar Singh. 2014. Improve-	424
ment of soundex algorithm for indian language based	425
on phonetic matching. <i>International Journal of Com-</i>	426
<i>puter Science, Engineering and Applications (IJC-</i>	427
<i>SEA) Vol</i> , 4.	428
Anil Kumar Singh. 2006. A computational phonetic	429
model for indian language scripts. In <i>Constraints on</i>	430
<i>spelling changes: Fifth international workshop on</i>	431
<i>writing systems</i> , pages 1–19. Nijmegen, The Nether-	432
lands.	433
Anil Kumar Singh, Harshit Surana, and Karthik Gali.	434
2007. More accurate fuzzy text search for languages	435
using abugida scripts. <i>Improving Non English Web</i>	436
<i>Searching (iNEWS'07)</i> , page 71.	437
Richard Sproat. 2000. A computational theory of writ-	438
ing systems. Cambridge University Press.	439
<ul> <li>Richard Sproat. 2002. Brahmi scripts. In <i>Constraints</i> on Spelling Changes: Fifth International Workshop on Writing Systems, Nijmegen, The Netherlands.</li> <li>Richard Sproat. 2003. A formal computational analysis</li> </ul>	440 441 442 443
of indic scripts. In International symposium on indic	444
scripts: past and future, Tokyo. Citeseer.	445
Robert A Wagner and Michael J Fischer. 1974. The string-to-string correction problem. <i>Journal of the ACM (JACM)</i> , 21(1):168–173.	446 447 448
Justin Zobel and Philip Dart. 1996. Phonetic string	449
matching: Lessons from information retrieval. In	450
<i>Proceedings of the 19th annual international ACM</i>	451
<i>SIGIR conference on Research and development in</i>	452
<i>information retrieval</i> , pages 166–172.	453

## **A** Phoneme Mappings

454

The details of consonant and vowel phoneme map-<br/>ping are provided in following tables.455456

# 360 361 362 363 364

357

36

371

373

374

400

401

402

403

404

405

406

407

408

phoneme1	phoneme2
q	k
W	v
Z	j
ph	f
hr	r
wr	r
ksh	х
ks	х
ck	k
ck	c
ng	n
sc	S

Table 3:	High	Similarity	Consonant	Pairs
	$\mathcal{O}$	2		

phoneme1	phoneme2
с	k
ch	k
Z	S
W	b
V	b
v	bh
р	f
ksh	ch
ksh	chh
ksh	sh
х	ch
х	chh
Z	g
g	j
zh	1
k	g
th	dh
d	t
ch	s

Table 4: Medium Similarity Consonant Pairs

phoneme1	phoneme2
r	d
с	S
р	b
b	f

Table 5: Low Similarity Consonant Pairs

phoneme1	phoneme2
00	u
ee	i
ea	ii
ea	ee
ea	i
ea	e
ea	ie
e	i
i	У
e	У
ae	у
ai	у
ai	ei
ai	aya
ai	ay
ey	ay
ey	i
ey	ie
e	ay
e	ey
ae	e
ae	ai
ae	а
ai	ae
ai	e
ou	au
ou	0
au	0
au	ao
u	0
а	e
oe	oya
oe	oy
oe	oai
ei	i
ie	i
ei	e
ei	ey
ai	ey
ie	e
ie	У
ia	iya
ow	au
ow	ou
ow	0
oa	0
yoo	yu
yoo	eu
yu	eu
eu	u
yu	u
yoo	u

phoneme1	phoneme2
а	u
ia	aya
ia	ya
aye	aya
а	0
а	ah
i	ai

Table 7: Medium Similarity Vowel Pairs

phoneme1	phoneme2
у	,,
а	"
h	ha
e	"
u	,,

Table 8: Low Similarity Vowel Pairs

**B** Hyperparameters Values

457

458

459

460

461

462

463

464

465

466

467

468

469

Optimal values of empirically determined hyperparameters are as follows

- consonant weight = 0.65
- default scorer = "specialized xrphonetic"
  - score combination mode = "mean"
    - consonant diff "h" penalty = 0.975
    - consonant high similarity score = 0.95
- consonant medium similarity score = 0.925
- consonant low similarity score = 0.9
- vowel diff "h" penalty = 0.975
  - vowel high similarity score = 0.95
  - vowel medium similarity score = 0.925
- vowel low similarity score = 0.9