

Phonotactic probabilities in Mandarin syllables

Anonymous ACL submission

Abstract

In this study, a spoken data sample is computed via frequency-based probabilistic phonotactics in Mandarin syllables by categorizing the spoken dataset into 12 syllable structure types. Phonotactic probabilities are measured by the bigram or biphone frequencies with which phonological segments and phone sequences occur in word types in Mandarin. Spoken data drawn from 2,384,567 lexical items show that correlations between speech production measures and phonological/articulatory complexity are not found. Instead, phonotactic probabilities influence speech production processes in Mandarin speakers independent of phonological complexity.

Keywords: Mandarin spoken corpus, phonotactic probability, bigram/biphone frequency, speech production

1 Introduction

It is generally believed that speakers can process certain sound sequences faster than others. The possible sound sequences in languages are not all equiprobable as some are more frequent than others. Some researchers suggested that certain sound sequences have attributed similar behavioral effects that are easier to articulate (i.e., less phonological complexity), but others attributed the patterning to varying degrees of probabilistic constraints. Such constraints can be referred to as phonotactic probabilities where phonological phones and sound sequences are legally arranged in lexical items (Jusczyk et al., 1994). For example, in English, the initial sequence [str] is allowable whereas the sequence [stn] does not form a legal arrangement. Or, in Mandarin, the initial sequence [kwa] is permissible while the sequence [kja] or

[kwn] is not. Additionally, the single phone unit in the above phone sequences does not distribute evenly. The glide [w] or [j] occurs more frequently than the consonant [k] in Mandarin due to the fact that glides have a wider distribution (i.e., syllable-initially, syllable-medially, and syllable-finally) than the true consonant [k] (syllable-initially exclusively) (Wan, 2022). In Goldrick and Larson's (2008) experiments, English speakers were sensitive to variations in frequency, demonstrating that phonotactic probabilities are encoded by speech production processes. Such novel phonotactic constraints were found to be correlated with the phonotactic probability of specific phonological structures. However, other research has shown a highly correlated association between speech production and phonological/articulatory complexity such as markedness in phones or syllable structure (e.g., Jakobson, 1941/1968; Romani and Calabrese, 1998). Evidence from these studies presented a rather small number of structures that yielded mixed and uncertain findings. The increasing variety of approaches to probability in phonology indicates a growing agreement that phonological analysis needs to incorporate probability and frequency into the theoretical framework (Alderete and Finley, 2023). Therefore, determining whether phonological/articulatory complexity or probabilistic constraints having a strong correlation in natural languages is not straightforward.

A number of studies further found that phonotactic probabilities exhibit a strong correlation with neighborhood density, which refers to the quantity of lexical items that share phonological similarity with a target (e.g., Goldrick and Rapp, 2007; Vitevitch et al., 2004). These effects manifest at separate and independent levels within the spoken production system. In this study, we are going to compute frequency-based probabilistic phonotactics in Mandarin syllables by

82 categorizing the spoken dataset into 12 syllable
 83 structure types via Biphone/Phone or
 84 Bigram/Gram frequencies (i.e., segment-to-
 85 segment co-occurrence probability of sounds
 86 within the lexical items; Vitevitch and Luce, 2004)
 87 and tone is omitted in the calculation. In addition,
 88 the effects of phonotactic probabilities, the
 89 likelihood of occurrence of a phone sequence, will
 90 be measured along with the syllable structure types.

91 2 Methodology

92 Spoken data in the study that has been collected
 93 over decades were drawn to be discussed from
 94 the first author's lab, Phonetics and
 95 Psycholinguistics Laboratory (N= 2,384,567
 96 syllables, 202 hours). The topics of the spoken
 97 content that were recorded in a naturalistic setting
 98 varied from lecture notes, class discussions,
 99 interviews, presentations, conversations of daily
 100 lives, etc., between multiple speakers in Taiwan.
 101 The sound files collected after 2020 were sent to a
 102 Speech-to-Text (STT) system for transcription into
 103 the International Phonetic Alphabet (IPA) via
 104 Chinese characters; the STT Package was
 105 developed from the application pyTranscriber
 106 (<https://github.com/raryelcostasouza>
 107 /[pyTranscriber](https://github.com/raryelcostasouza)) in the Phonetics and
 108 Psycholinguistics lab. The transcription in Chinese
 109 characters of a 60-minute audio file took 80
 110 seconds. The accuracy of the output can vary a lot,
 111 depending on factors such as voice quality, noise
 112 clarity, gender, age, and/or speech speed of
 113 speakers. The accuracy rate varied between 70%
 114 and 90% depending on the combination of these
 115 factors. The output of the STT system was then
 116 manually checked. Afterward, the entire transcript
 117 was automatically segmented by the CKIP parser
 118 (Ma and Chen, 2003) and POS tagged by the CKIP
 119 tagger from the Chinese Knowledge and
 120 Information Processing group (CKIP, 1998). The
 121 parsed and tagged transcription was also checked
 122 manually according to the criteria of word
 123 segmentation and POS tagging of the Academia
 124 Sinica Corpus (CKIP, 1998), which are commonly
 125 applied in corpora such as the Linguistic Data
 126 Consortium (Ma and Huang, 2006) and the Peking
 127 University corpus (Huang et al, 2008).

128 It is noted that spoken data samples collected in
 129 this study are calculated on the frequency of
 130 occurrence information in various topics which
 131 were recorded in a naturalistic setting. Word counts
 132 are up to date and are not from movie subtitles.

133 Mandarin is analyzed as having a range of possible
 134 phonetic (i.e., surface) syllables: V, CV, GV, VG,
 135 VN, CVG, CVN, CGV, GVG, GVN, CGVG, and
 136 CGVN. The maximal syllable is CGVX, with C a
 137 [+consonantal] segment, G a glide, V the nucleus
 138 vowel, and X either a nasal or a glide (i.e., Wan
 139 1999).

140 Based on the measurement in Vitevitch and
 141 Luce's study (2004) where Bigram/Gram or
 142 Biphone/Phone frequencies are computed by
 143 dividing the sum of the log frequencies of all of the
 144 words with element A at position N (or N & N+1
 145 in the case of Bigrams/Biphones) by the sum of the
 146 log frequencies of all words with any unit in
 147 position N (or N & N+1). Samples of types and
 148 token frequencies of a syllable structure, CGVN, in
 149 Mandarin are shown in Table 1.

IPA	Freq.	IPA	Freq.	IPA	Freq.
ɛjaŋ	18031	swan	1727	ɕwən	337
ɛjən	14461	tɕwan	1498	lwan	324
tɛjaŋ	11453	tɕ ^h wan	1496	tɛ ^h ɥən	307
mjən	10427	tɛ ^h jaŋ	1337	tɛɥən	287
pjən	9907	kwaŋ	1226	k ^h wan	281
tɛ ^h jən	8286	tɕwaŋ	1171	ɕwan	228
tɛjən	8224	k ^h waŋ	1045	twən	228
tjən	8019	ɥən	972	swən	216
njən	7820	xwaŋ	724	tswən	214
t ^h jən	6324	tɕwən	688	njaŋ	118
ljaŋ	5760	xwən	682	tɛ ^h joŋ	115
kwan	5689	ɛjoŋ	641	kwən	55
tɛ ^h ɥən	3705	tɕ ^h waŋ	619	nwan	54
xwan	2990	tɕ ^h wən	576	t ^h wən	49
ljən	2766	tɛɥən	466	ɕwən	39
ɥən	2452	tɕ ^h wən	453	ɕwan	27
twan	2096	t ^h wan	378	tswan	20
lwən	1874	ɕwaŋ	361	tɕ ^h wan	3
p ^h jən	1806	k ^h wən	348	tɛjoŋ	1

151 Table 1: Samples of CGVN in IPA and token frequencies

152
 153 Table 1 lists the spoken dataset in CGVN with all
 154 possible sound sequences of token frequencies in
 155 Mandarin (Note that tone is omitted in the study).

156 The following (1-2) shows the formula for the
 157 bigram/biphone phonotactic probability in
 158 Mandarin with an example of the syllable type
 159 CGVN [ɛjaŋ] (value at 0.0821).

$$\begin{aligned}
162 \text{ PhonProb}_{tokens} = & \\
163 & \left[\frac{\text{Sum of log frequencies of words with [c] in initial biphone position}}{\text{Sum of log frequencies of words with any biphone in initial biphone position}} + \right. \\
164 & \frac{\text{Sum of log frequencies of words with [ja] in initial biphone position}}{\text{Sum of log frequencies of words with any biphone in initial biphone position}} + \\
165 & \left. \frac{\text{Sum of log frequencies of words with [an] in initial biphone position}}{\text{Sum of log frequencies of words with any biphone in initial biphone position}} \right] / \\
166 & [\text{Number of biphone positions}] \quad (1)
\end{aligned}$$

$$\begin{aligned}
167 \text{ PhonProb}_{tokens} = & \\
168 & \left[\frac{\log(18031)+\log(14461)+\log(641)}{\log(18301)+\log(14461)+\log(11453)+\log(10427)+\dots+\log(1)} + \right. \\
169 & \frac{\log(18031)+\log(5760)+\log(118)}{\log(18301)+\log(14461)+\log(11453)+\log(10427)+\dots+\log(1)} + \\
170 & \left. \frac{\log(18301)+\log(14461)+\log(11453)+\log(10427)+\dots+\log(1)}{\log(18031)+\log(5760)+\log(1226)+\log(724)+\log(361)+\log(118)} \right] / 3 = 0.0821 \\
171 & \quad (2)
\end{aligned}$$

172 In this model, the phonotactic probability is
173 calculated in a given syllable using the token
174 frequencies and a dataset of word types involving
175 different syllable types. Initially, the given syllable
176 is segmented into a series of bigrams, representing
177 pairs of adjacent units. Subsequently, for each
178 position within the syllable, the model computes
179 two sums involving one for the logarithm of the
180 same bigram occurring at the position and another
181 one for the logarithm of the frequency of all
182 bigrams at the position. The phonotactic
183 probability of the syllable is determined by
184 summing the ratio of these two sums for each
185 bigram in the syllable and dividing by the total
186 number of bigrams in the syllable. This ratio
187 indicates the relative frequency of each bigram in
188 its position.

189
190
191
192

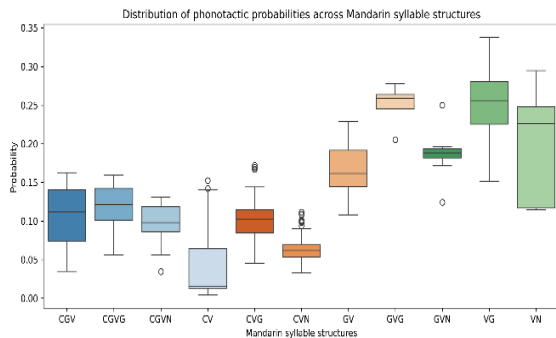


Figure 1: Distribution of phonotactic probabilities across Mandarin syllable structures.

193 In Figure 1, the picture visualizes the distribution
194 of phonotactic probabilities across different
195 syllable structure types, each depicted by a box plot.
196 Each box shows the interquartile range (IQR) of
197 probabilities, with the median marked by a line
198 across the box. The 'whiskers' extend to the furthest
199 points within 1.5 times the IQR, highlighting the

200 range of typical data, while circles denote outliers,
201 representing syllable structures with probabilities
202 outside the typical range. The VN structure has the
203 widest range of phonotactic probabilities,
204 suggesting a high variability in how often it occurs
205 in Mandarin, and it has a higher median probability,
206 indicating the phonotactic constraints in the vowel-
207 nasal sequences in Mandarin are less restricted.
208 The VG and GVG structures have the highest
209 median probability, suggesting that the vowel and
210 the glide sequences can allow a wider range of
211 varieties in Mandarin. On the other hand, CV has
212 the lowest median probabilities, indicating that
213 there is a restricted constraint on the consonant-
214 vowel sequences in Mandarin. The presence of
215 outliers in several structures implies that there are
216 some specific sequences within those structures
217 that are significantly less or more restricted than
218 others.

219 3 Results and limitations

220 In this section, we examine the phonotactic
221 probability distribution calculated in a given
222 Mandarin syllable using the token frequencies and
223 a dataset of word types involving different syllable
224 structures. Type and token frequencies in the
225 current spoken data confirm the studies found in
226 English where the possible sound sequences are
227 not all equiprobable as some are more frequent
228 than others. More importantly, certain sound
229 sequences are related to probabilistic constraints
230 and do not fall in the articulatory complexity since
231 the CV-type structure is supposed to be the easiest
232 pattern at a more flexible range, whereas its
233 phonotactic probability is the lowest.

234 It is interesting to note that among the entire
235 syllable structure types, CV used to be categorized
236 as the easiest type for children to articulate cross-
237 linguistically in acquisition studies. The study
238 suggests that phonotactic constraints in Mandarin
239 disassociate articulatory complexity and
240 phonotactic probabilities influence speech
241 production regardless of the markedness
242 complexity. Our spoken samples via data
243 computation confirm an emerging agreement
244 within the field that phonological theories need to
245 consider phonotactic probabilities. The limitation
246 for the current study is that Levenshtein edit
247 distance will need to be measured in order to
248 further calculate the neighborhood density. The
249 future step will need to investigate the
250 neighborhood density in Mandarin, where the

251 sound-similar words are stored in the mental
252 lexicon.

253 References

254 Chu-Ren Huang, Lung-Hao Lee, Wei-guang Qu, Jia-
255 Fei Hong, and Shiwen Yu. 2008. [Quality Assurance](#)
256 [of Automatic Annotation of Very Large Corpora: a](#)
257 [Study based on heterogeneous Tagging System](#). In
258 *Proceedings of the Sixth International Conference*
259 *on Language Resources and Evaluation (LREC'08)*,
260 pages 2725-2729. Marrakech, Morocco. European
261 Language Resources Association (ELRA).

262 CKIP. 1998. *Academia Sinica Balanced Corpus*
263 (Version 3) [CD-ROM]. Taipei: Chinese
264 Knowledge and Information Processing Group,
265 Academia Sinica.

266 Cristina Romani and Andrea Calabrese. 1998. [Syllabic](#)
267 [constraints in the phonological errors of an aphasic](#)
268 [patient](#). *Brain and language*, 64(1):83-121.
269 <https://doi.org/10.1006/brln.1998.1958>

270 I-Ping Wan. 1999. *Mandarin phonology: Evidence*
271 *from speech errors* (Order No. 9943387) [State
272 University of New York at Buffalo]. Available from
273 ProQuest Dissertations & Theses A&I; ProQuest
274 Dissertations & Theses Global. (304551099)

275 I-Ping Wan. 2022, April 15. [Error analysis in](#)
276 [Mandarin corpus phonology](#), [Invited talk in the
277 Institute of Linguistics at Academia Sinica, Taipei,
278 Taiwan]. Academia Sinica Institute of Linguistics
279 Phonetics Laboratory.
280 [https://drive.google.com/file/d/11giA-](https://drive.google.com/file/d/11giA-12_AIAAH5kkyRPOwK0jra4AZmev/view)
281 [12_AIAAH5kkyRPOwK0jra4AZmev/view](https://drive.google.com/file/d/11giA-12_AIAAH5kkyRPOwK0jra4AZmev/view)

282 John Alderete and Sara Finley. 2023. [Probabilistic](#)
283 [phonology: A review of theoretical perspectives,](#)
284 [applications, and problems](#). *Language and*
285 *Linguistics*, 24(4):565-610.
286 <https://doi.org/10.1075/lali.00141.ald>

287 Matthew Goldrick and Brenda Rapp. 2007. [Lexical](#)
288 [and post-lexical phonological representations in](#)
289 [spoken production](#). *Cognition*, 102(2):219-260.

290 Matthew Goldrick and Meredith Larson. 2008.
291 [Phonotactic probability influences speech](#)
292 [production](#). *Cognition*, 107(3):1155-1164.
293 <https://doi.org/10.1016/j.cognition.2007.11.009>

294 Michael S. Vitevitch and Paul A. Luce. 2004. [A web-](#)
295 [based interface to calculate phonotactic probability](#)
296 [for words and nonwords in English](#). *Behavior*
297 *Research Methods, Instruments, & Computers*,
298 36(3):481-487.
299 <https://doi.org/10.3758/BF03195594>

300 Michael S. Vitevitch, Jonna Armbrüster, and Shinying
301 Chu. 2004. [Sublexical and lexical representations in](#)
302 [speech production: effects of phonotactic](#)
303 [probability and onset density](#). *Journal of*

Experimental Psychology: Learning, Memory, and
304 *Cognition*, 30(2):514-529.
305 <https://doi.org/10.1037/0278-7393.30.2.514>
306

307 Peter W. Jusczyk, Paul A. Luce, and Jan Charles-Luce.
308 1994. [Infants' sensitivity to phonotactic patterns in](#)
309 [the native language](#). *Journal of memory and*
310 *Language*, 33(5):630-645.
311 <https://doi.org/10.1006/jmla.1994.1030>

312 Roman Jakobson. 1941/1968. *Child language aphasia*
313 *and phonological universals*. (A. R. Keiler, Trans.)
314 The Hague: Mouton.

315 Wei-Yun Ma and Chu-Ren Huang. 2006. [Uniform and](#)
316 [Effective Tagging of a Heterogeneous Giga-word](#)
317 [Corpus](#). In *Proceedings of the 5th International*
318 *Conference on Language Resources and Evaluation*
319 *(LREC-5)*, pages 2182-2185, Genoa, Italy.
320 <https://aclanthology.org/L06-1163/>

321 Wei-Yun Ma and Keh-Jiann Chen. 2003. [Introduction](#)
322 [to ckip chinese word segmentation system for the](#)
323 [first international Chinese word segmentation](#)
324 [bakeoff](#). In *Proceedings of SIGHAN*, pages 168-171.
325 <https://aclanthology.org/W03-1726/>

326 A Supplementary Material

327 Supplementary material on the spoken data and
328 coding will be released to the public once the paper
329 is accepted for publication.