

# Mandarin third tone sandhi may be incompletely neutralizing in perception as well as production

Stephen Politzer-Ahles<sup>1\*</sup>, Katrina Connell<sup>2</sup>, Lei Pan<sup>1</sup>, Yu-Yin Hsu<sup>1</sup>

<sup>1</sup>Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic University, Hong Kong,

<sup>2</sup>Department of Spanish, Italian, and Portuguese; Pennsylvania State University, United States

*Submitted to Journal:*  
Frontiers in Psychology

*Specialty Section:*  
Language Sciences

*Article type:*  
Original Research Article

*Manuscript ID:*  
447184

*Received on:*  
09 Jan 2019

*Frontiers website link:*  
[www.frontiersin.org](http://www.frontiersin.org)

In review

---

### *Conflict of interest statement*

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest

### *Author contribution statement*

SP and YH conceived the experiment. SP, KC, LP and YH designed the experiment. KC and LP created the stimuli. KC programmed the experiment and LP collected the data. SP and KC analyzed the data. All authors wrote and approved the manuscript.

### *Keywords*

incomplete neutralization, Mandarin tone, third tone sandhi, visual world eye-tracking, Speech Perception

### *Abstract*

Word count: 131

Mandarin third tone sandhi is traditionally assumed to be incompletely neutralizing in production but completely neutralizing in perception, based on metalinguistic judgment tasks in which participants cannot reliably identify the underlying tone of syllables neutralized by tone sandhi. We performed a visual world eye-tracking study to see if implicit sensitivity to the differences between the surface forms influences participants' eye movement patterns, even if they cannot consciously access this for identification tasks. We found a slight trend in this direction, with participants looking more towards orthographic representations that match the underlying form of the neutralized syllable they hear. The results are statistically inconclusive but suggest that this paradigm may be able to provide evidence that Mandarin neutralized tones are indeed incompletely neutralized, and that further research along these lines is warranted.

### *Funding statement*

This research was supported by grant G-UACX from the Department of Chinese and Bilingual Studies to SP.

### *Ethics statements*

(Authors are required to state the ethical considerations of their study in the manuscript, including for cases where the study was exempt from ethical approval procedures)

*Does the study presented in the manuscript involve human or animal subjects:* Yes

*Please provide the complete ethics statement for your manuscript. Note that the statement will be directly added to the manuscript file for peer-review, and should include the following information:*

- Full name of the ethics committee that approved the study
- Consent procedure used for human participants or for animal owners
- Any additional considerations of the study in cases where vulnerable populations were involved, for example minors, persons with disabilities or endangered animal species

*As per the Frontiers authors guidelines, you are required to use the following format for statements involving human subjects: This study was carried out in accordance with the recommendations of [name of guidelines], [name of committee]. The protocol was approved by the [name of committee]. All subjects gave written informed consent in accordance with the Declaration of Helsinki.*

*For statements involving animal subjects, please use:*

*This study was carried out in accordance with the recommendations of 'name of guidelines, name of committee'. The protocol was approved by the 'name of committee'.*

*If the study was exempt from one or more of the above requirements, please provide a statement with the reason for the exemption(s).*

*Ensure that your statement is phrased in a complete way, with clear and concise sentences.*

This study was carried out in accordance with the recommendations of the Human Subjects Ethics Sub-committee at the Hong Kong Polytechnic University with written informed consent from all subjects. All subjects gave written informed consent in accordance with the Declaration of Helsinki. The protocol was approved by the Human Subjects Ethics Sub-committee at the Hong Kong Polytechnic University.

*Data availability statement*

Generated Statement: All datasets generated for this study are included in the manuscript and the supplementary files.

In review

# Mandarin third tone sandhi may be incompletely neutralizing in perception as well as production

1 Stephen Politzer-Ahles<sup>1\*</sup>, Katrina Connell<sup>2</sup>, Yu-Yin Hsu<sup>1</sup>, Lei Pan<sup>1</sup>

2 <sup>1</sup>Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic University, Hong Kong,  
3 Hong Kong

4 <sup>2</sup>Department of Spanish, Italian & Portuguese; Pennsylvania State University, State College,  
5 Pennsylvania, United States

6 **\* Correspondence:**

7 Stephen Politzer-Ahles

8 sjpolit@polyu.edu.hk

9 **Keywords: incomplete neutralization, Mandarin tone sandhi, third tone sandhi, visual world**  
10 **eye-tracking, speech perception**

11 **Abstract**

12 Mandarin third tone sandhi is traditionally assumed to be incompletely neutralizing in production but  
13 completely neutralizing in perception, based on metalinguistic judgment tasks in which participants  
14 cannot reliably identify the underlying tone of syllables neutralized by tone sandhi. We performed a  
15 visual world eye-tracking study to see if implicit sensitivity to the differences between the surface  
16 forms influences participants' eye movement patterns, even if they cannot consciously access this for  
17 identification tasks. We found a slight trend in this direction, with participants looking more towards  
18 orthographic representations that match the underlying form of the neutralized syllable they hear. The  
19 results are statistically inconclusive but suggest that this paradigm may be able to provide evidence  
20 that Mandarin neutralized tones are indeed incompletely neutralized, and that further research along  
21 these lines is warranted.

22 **1 Introduction**

23 A common phenomenon in language is *neutralization*, whereby two or more sounds that normally  
24 would be different, are instead changed to have the same surface pronunciation in a certain context.  
25 For example, standard Dutch has a distinction between phonologically voiced and voiceless stops  
26 like /d/ and /t/, but at the end of the word both are realized as [t] (the words *baat* “benefit” and *baad*  
27 “bathe” are both pronounced [baat]) (Warner, Jongman, Sereno, & Kemps, 2004). In many cases,  
28 however, these contrasts are not completely neutralized. For example, in the Dutch example above,  
29 there are still slight pronunciation differences between *baat* and *baad*, and these differences are  
30 perceptible to listeners (Warner et al., 2004). The same is true in English: for a minimal pair like *beat*  
31 and *bead*, while /t/ and /d/ may not differ in aspiration like they would at the onset of a word, there  
32 are still several subtle acoustic differences between them, as well as large acoustic differences in the  
33 vowels that precede them (see, e.g., Flege, Munro, & Skelton, 1992). Such a pattern is known as  
34 *incomplete neutralization* (for review see Nicenboim, Roettger, & Vasisth, 2018, among others).

35 Incomplete neutralization has important implications for psycholinguistic accounts of speech  
 36 perception and production, which must account for the consequences of incompletely neutralized  
 37 representations in lexical storage and in speech comprehension. It also has implications for  
 38 phonological theory—for example, struggles to theoretically account for phonological opacity  
 39 (situations where a surface form exists in a context that looks like it should have triggered a change  
 40 to a different surface form, but somehow does not) may be obviated if one assumes that an  
 41 incompletely neutralized surface form does not actually trigger the relevant phonological alternation  
 42 and thus is not actually opaque (see, e.g., Zhang, Lai, & Sailor, 2011). Finally, incomplete  
 43 neutralization is pervasive; it may even be the norm, as many putatively neutralized contrasts have  
 44 been found to have some remaining acoustic correlates that distinguish them, and patterns where two  
 45 neutralized forms are actually indistinguishable in both acoustic measurements and perception seem  
 46 to be rare (Kim & Jongman, 1996).<sup>1</sup>

47 Some neutralization patterns are argued to be incomplete in production but complete in perception:  
 48 that is to say, two putatively neutralized sounds have reliable acoustic differences which can be  
 49 detected with computer-assisted measurements and statistical methods, but human listeners cannot  
 50 reliably detect them when tested. A classic example of this case is third tone sandhi in Mandarin.  
 51 Standard Mandarin has lexical tones, such that segmentally identical syllables with different tones  
 52 may correspond to different morphemes (e.g., [kwa<sup>1</sup>]<sup>2</sup>, with a High tone, is the pronunciation of the  
 53 morpheme written 瓜, "melon", among others; whereas [kwa<sup>˥</sup>], with a Falling tone, is the  
 54 pronunciation of the morpheme written 挂, "to hang"). Crucially, two tones, Rising and Low, are  
 55 putatively neutralized in a certain context. When a Low tone (also known as a third tone or "tone  
 56 three", hence the name "third tone sandhi") is preceded by another within the same intonational  
 57 domain, it is instead pronounced as a Rising tone (see Kuo, Xu, and Yip, 2007, and Zhang & Lai,  
 58 2010, among others, for further details regarding other phonological and morphosyntactic constraints  
 59 on this pattern). For instance, the morpheme written 雨 is normally pronounced [y<sup>˩</sup>], with Low tone;  
 60 but when it appears before another Low tone, as in the compound word 雨伞 [y<sup>˨˩</sup> san<sup>˩</sup>] "umbrella", it  
 61 is instead produced with a Rising tone, homophonic with the morpheme written 鱼 [y<sup>˨˩</sup>], "fish". This  
 62 alternation causes the distinction between Low and Rising tones to be neutralized in pre-Low  
 63 positions that license third tone sandhi. Many acoustic studies, though, have found that the distinction  
 64 is not completely neutralized: a "Rising" tone derived via tone sandhi from an underlying Low tone  
 65 tends to be slightly lower and have a slightly later turning point in its tonal contour (Chen, Wiltshire,  
 66 & Li, 2017; Cheng, Chen, & Gubian, 2013; Liu, 2013, Peng, 2000; Yuan & Chen, 2014; Zhang &  
 67 Lai, 2010; Zhang & Peng, 2013; among others). On the other hand, several studies have also shown  
 68 that speakers asked to identify the underlying tone (i.e., if they are presented with an ambiguous  
 69 Rising-Low disyllable and asked to identify it as an underlyingly Low-Low disyllable or an

<sup>1</sup> Alternatively, though, it is possible that the abundance of attested incomplete neutralization patterns could also be due to methodological practices. In short, if one contrast has many possible acoustic correlates that could be tested and could be quantified in different ways, and these are all checked for evidence of incomplete neutralization, this greatly increases the chance of finding spurious differences (see Nicenboim et al., 2018, for discussion; and Roettger, under review, for a general introduction to this statistical issue). Meta-analytic methods, as well as frequent replication, provide evidence that some classic incomplete neutralization patterns like Dutch and German final devoicing are indeed reliable (Nicenboim et al., 2018); it would be valuable to extend these methods, as well as other meta-analytic methods like the p-curve (Simonsohn, Nelson, & Simmons, 2014) to other less studied cases of putatively incomplete neutralization.

<sup>2</sup> Transcriptions in brackets are representations in the International Phonetic Alphabet. We use the following symbols to indicate tones: 1 for High tone, ˨˩ for Rising tone, ˩ for Low tone, and ˥ for Falling tone.

70 underlyingly Rising-Low one) cannot reliably identify the correct tone (Liu, 2013; Peng, 2000; Wang  
 71 & Li, 1967; Zhang & Peng, 2013). Thus, it is generally believed that while Mandarin third tone  
 72 sandhi is not completely neutralizing in production (i.e., a Rising tone resulting from tone sandhi is  
 73 acoustically different from an underlyingly Rising tone), it is completely neutralizing in perception  
 74 (i.e., native listeners cannot reliably hear the contrast).

75 Nonetheless, some results in the extant literature do suggest that the subtle contrast between  
 76 underlying and sandhi-derived Rising tones may influence listeners' perception. These come from a  
 77 priming experiment, a visual world eye-tracking experiment, a discrimination task, and an  
 78 identification task, described below.

79 Zhou and Marslen-Wilson (1997) report two auditory-auditory priming experiments in which  
 80 participants heard primes and targets that were each disyllabic words, and in which the first syllable  
 81 was sometimes an underlyingly or a sandhi-derived Rising tone. In their first experiment, the critical  
 82 targets began with a sandhi-derived Rising tone (e.g., [ts<sup>h</sup>ai<sup>1</sup> tɛ<sup>h</sup>y<sup>1</sup>] 采取 "to carry out", the first  
 83 syllable of which has Low tone in its citation form but has Rising tone in this word because it is  
 84 followed by another Low tone), and reaction time in performing lexical decision to these items was  
 85 sped up (relative to a control condition with unrelated primes) when they were immediately preceded  
 86 by primes with a Low tone in the first syllable, like [ts<sup>h</sup>ai<sup>1</sup> xoŋ<sup>1</sup>] 彩虹 "rainbow". On the other hand,  
 87 in their second experiment (and in a direct replication of their second experiment), the critical targets  
 88 began with an underlyingly Rising tone (e.g., [ts<sup>h</sup>ai<sup>1</sup> pan<sup>1</sup>] 裁判 "referee", the first syllable of which  
 89 is pronounced with Rising tone even in its citation form), and reaction time to these targets was  
 90 slowed down (relative to a control condition with unrelated primes) when they were preceded by  
 91 primes with a Low tone in the first syllable, like [ts<sup>h</sup>ai<sup>1</sup> na<sup>1</sup>] 采纳 "to adopt [a measure]". These  
 92 patterns appear to be opposite: a sandhi target is facilitated by a Low prime whereas a Rising target is  
 93 inhibited by a Low prime. This might be taken as evidence that sandhi-derived and underlying Rising  
 94 tones were treated differently by the parser. However, the experiment was not designed to test this  
 95 issue, and this result could be accounted for without recourse to incomplete neutralization. For  
 96 example, the critical syllables were embedded in meaningful two-syllable words, which are sufficient  
 97 to allow participants to identify the underlying form even if they can't hear the difference between the  
 98 incompletely neutralized tones. E.g., the form [ts<sup>h</sup>ai<sup>1</sup> tɛ<sup>h</sup>y<sup>1</sup>] corresponds to an existing Low-Low  
 99 disyllable word 采取 and does not correspond to any existing Rising-Low disyllable word. Thus, a  
 100 participant does not need to be able to hear the difference between sandhi-derived and underlying  
 101 tones to be able to recognize that the first syllable in this word has undergone tone sandhi; the  
 102 participant only needs to recognize the following syllable and to recognize that this combination only  
 103 corresponds to one existing word, one where the first syllable is underlyingly Low.

104 Speer and Xu (2008) report two priming experiments and a visual world eye-tracking experiment,  
 105 specifically designed to investigate the comprehension of ambiguous neutralized forms. Participants  
 106 heard Mandarin sentences in which the critical Rising-tone word was in a sandhi-licensing context,  
 107 and thus could be interpreted as either underlyingly Rising or underlyingly Low. For example, one  
 108 sentence (今天海边\_\_很多) meant something like "Today by the sea there was/were<sup>3</sup> a lot of [qy<sup>1</sup>]."  
 109 In the Mandarin sentence (which has a different word order than English), the critical syllable [qy<sup>1</sup>] is

<sup>3</sup> Mandarin verbs are not morphologically marked for number, so the verb was the same in either sentence and thus could not be a cue to the identity of the critical ambiguous word.

110 in a context that licenses tone sandhi. Thus, this word could either be interpreted as meaning "fish" (  
 111 鱼, the citation form of which is [ɥʏ˥] with Rising tone), or as meaning "rain" (雨, the citation form  
 112 of which is [ɥʏ˩] with low tone, but which would change to [ɥʏ˥] with Rising tone in this sandhi-  
 113 licensing context). The sentence stimuli were cross-spliced such that sometimes the critical  
 114 ambiguous syllable was one produced as a natural Rising tone, with its concomitant acoustic cues,  
 115 and sometimes the critical ambiguous syllable was one produced with a sandhi-derived Rising tone,  
 116 with its concomitant acoustic cues. In the visual world eye-tracking experiment, participants heard  
 117 these sentences while looking at a screen with both potential characters (the one corresponding to the  
 118 interpretation with Rising tone and the one corresponding to the interpretation with Low tone) shown  
 119 on the screen, along with distractor characters. An eye tracker monitored where the participants' eyes  
 120 were looking as they heard the critical ambiguous syllable. Counterintuitively, when participants  
 121 heard the sentence with a sandhi-derived Rising tone in it, they initially looked at the character  
 122 corresponding to the Rising tone more than the one corresponding to the Low tone; also  
 123 counterintuitively, when they heard the sentence with an underlyingly Rising tone in it, they initially  
 124 looked at the Low-tone character more than the Rising tone character. These results are opposite  
 125 what one would initially predict (it makes the most sense to assume that if participants can recognize  
 126 the difference, they would look at the appropriate character more than the inappropriate character).  
 127 Nonetheless, they are potentially consistent with the notion that participants are subtly sensitive to  
 128 the difference between the incompletely neutralized tones (although they are also consistent with  
 129 other possibilities, of course; e.g., they may be a Type 1 error or they may be due to some other  
 130 unexpected factor).

131 Liu (2013) tested six listeners' explicit metalinguistic identification and discrimination of underlying  
 132 and sandhi-derived Rising tones. In the identification task they heard ambiguous disyllables like  
 133 [tɕʰi˥ ma˩], in which the first syllable was either underlyingly Rising or was changed to Rising via  
 134 tone sandhi, and had to select the appropriate orthographic representation (for this example, either 起  
 135 码 /tɕʰi˥ ma˩/ "at least" or 骑马 /tɕʰi˥ ma˩/ "ride a horse"). In the discrimination they heard disyllable  
 136 pairs, where either both disyllables were the same, or one disyllable was a production with an  
 137 underlying Rising and one a production with a sandhi-derived Rising tone. They then judged whether  
 138 they were the same or different. The participants were not significantly more accurate than chance  
 139 level in identifying the tones, but they were significantly above chance in discriminating them.  
 140 Above-chance discrimination does not indicate that the tones were incompletely neutralized, though,  
 141 because there are many other ways that different acoustic cues can allow participants to perform well  
 142 on within-category discrimination (not just in tone perception and not just in incomplete  
 143 neutralization contexts). Evidence that the tones are incompletely neutralized should take the form of  
 144 either above-chance ability to identify the appropriate tones, or better discrimination across tone  
 145 categories (underlying and sandhi-derived Rising tones) than within either category; above-chance  
 146 discrimination in of itself is not necessarily indicative of sensitivity to a categorical difference  
 147 between these tones, as opposed to sensitivity to any other differences between any pair of tokens.  
 148 (We could make an analogy to dishes from different culinary traditions. For example, if a person eats  
 149 dishes at a Hunan-style restaurant and a Sichuan-style restaurant, that person might be able to notice  
 150 the difference between a dish from one restaurant and a dish from the other restaurant, as well as the  
 151 difference between two dishes from the same restaurant, without necessarily being sensitive to the  
 152 difference between Hunan cuisine and Sichuan cuisine. Likewise, being able to tell the difference  
 153 between two stimuli—which may differ in many other respects other than their tones—is not  
 154 necessarily evidence that participants are sensitive to the general difference between derived and  
 155 underlying Rising tones.)

156 Finally, in a recent perception study, Lin and Hsu (2018) presented natural tokens with underlying or  
157 sandhi-derived Rising tones to four listeners who were instructed to indicate whether the tone they  
158 heard was Rising, Low, or something in between. When hearing a sandhi-derived Rising tone,  
159 participants identified it as Rising 53% of the time and as Low 29% of the time. On the other hand,  
160 when hearing an underlyingly Rising tone, they identified it as Rising 79% of the time, and as Low  
161 8% of the time. While the reliability of these observations is limited by the small scale of the study,  
162 this is another piece of evidence that listeners may perceive derived and underlying Rising tones  
163 slightly differently, even if they are not completely accurate in identifying them.

164 Overall, there is some evidence that listeners may be sensitive to the difference between incompletely  
165 neutralized Mandarin tones, but the evidence is weak and spotty, with some of it coming from studies  
166 not designed to test this issue and some coming from counterintuitive results that are difficult to  
167 explain. Something that is noteworthy is that, to our knowledge, all the studies suggesting that  
168 listeners are not sensitive to this difference are based on explicit metalinguistic judgment tasks. On  
169 the other hand, some of the potential evidence that listeners are sensitive to the difference comes  
170 from online measures like eye movements and reaction times, which participants do not have direct  
171 control over and which may reflect processes that participants are not consciously aware of. Because  
172 of these conflicting sets of results, we hypothesized that listeners may be able to hear the difference  
173 between sandhi-derived and underlying Rising tones at the unconscious, automatic level, but not able  
174 to consciously access that for a metalinguistic judgment.

175 We test this with a visual world eye-tracking experiment, using a design that is essentially a  
176 simplification of that used by Speer and Xu (2008). Participants heard ambiguous disyllabic words,  
177 without a sentence context, like [du<sup>1</sup> pən<sup>4</sup>], which might correspond to the word 读本 "reading  
178 book", where the citation form of the first syllable has Rising tone, or to the word 赌本 "bookie",  
179 where the citation form of the first syllable has Low tone. Like in the design of Speer and Xu (2008),  
180 we predict that, if participants are somewhat sensitive to the difference between the incompletely  
181 neutralized tones, they should look more at the word with an underlyingly Low first syllable when  
182 they hear a token that was spoken as a production of the word with underlying Low tone, compared  
183 to when they hear a token that was spoken as a production of the word with underlying Rising tone.  
184 By using single words rather syllables embedded in sentences, we aimed to minimize potential  
185 complications related to sentence processing and the plausibility or semantic fit between critical  
186 words and the rest of the sentence—although, of course, there are still other differences between the  
187 critical targets, like their frequencies, stroke counts, and other lexical properties. Crucially, however,  
188 in both conditions we are comparing (auditory stimulation with the underlyingly Rising-tone word,  
189 and stimulation with the sandhi-derived Rising-tone word), the Low and Rising target words on the  
190 screen are the same, so any lexical differences between them would not cause a difference between  
191 conditions. If the Low target is looked at more when hearing the sandhi-derived Rising form than the  
192 underlyingly Rising form, as far as we know this can only be due to listeners' sensitivity to the  
193 phonetic difference between these forms.

## 194 **2 Methods**

195 All experiment materials, data, de-identified participant demographic information, and analysis  
196 scripts are available at <https://osf.io/ursh9/>. Experiment methods were pre-registered at  
197 <https://osf.io/35ang/register/5771ca429ad5a1020de2872e>.

### 198 **2.1 Participants**



199 60 native speakers of Mandarin (mean age 23.55, age range 18-34, 48 women and 12 men)  
 200 participated in the study. Participants were required to filled in the Language Background  
 201 Questionnaire before the experiment. According to the questionnaire, all participants came from  
 202 homes where Mandarin was the primary language in use between birth and 5 years old and they  
 203 listed Mandarin as the most dominant language used in their daily life. They had normal or corrected-  
 204 to-normal vision and hearing. They provided informed consent to participate and were compensated  
 205 in cash. All experiment procedures were approved by the Human Subjects Ethics Sub-committee at  
 206 the Hong Kong Polytechnic University (project reference # HSEARS20171012002). Two additional  
 207 volunteers participated in the experiment but their data were not included in the analysis due to  
 208 technical issues relating to the eye tracker.

## 209 2.2 Materials

210 The experiment consisted of two auditory conditions: a sandhi-derived Rising tone and an  
 211 underlyingly Rising tone. The critical stimuli consisted of 14 pairs of disyllabic words that were  
 212 identical in their segmental structure and differed in the underlying tone of the initial syllable. In the  
 213 sandhi-derived Rising tone condition, these words had an initial Low tone and a final Low tone,  
 214 which causes the initial low tone to be phonetically realized as a rising tone (i.e., 土改 /t<sup>h</sup>u<sup>l</sup> kai<sup>l</sup>/ [t<sup>h</sup>u<sup>l</sup>  
 215 kai<sup>l</sup>]). In the Underlying Rise condition, the initial tone was a Rising tone and the final tone was a  
 216 Low tone (i.e., 涂改 /t<sup>h</sup>u<sup>1</sup> kai<sup>l</sup>/ [t<sup>h</sup>u<sup>1</sup> kai<sup>l</sup>]). The words across conditions were maximally similar in  
 217 terms of frequency, homophone density, neighborhood density, and neighborhood frequency for the  
 218 word as a unit, as well as for each syllable individually, summarized in Table 1.

219 A female native speaker of Mandarin Chinese recorded the words in isolation in a sound attenuated  
 220 room with a Telefunken M-80 dynamic microphone and a Focusrite Scarlett 2i2 sound interface.  
 221 Three repetitions were recorded at slow, normal, and fast speeds. The words were cut from the main  
 222 audio recording and no further manipulation was done. Tokens were selected from the normal-speed  
 223 repetition only.

224 In the eye-tracking experiment, a critical display presented 4 words in characters. Two of the words  
 225 formed the target pair, as described above. The other two words served as distracters. The words in  
 226 the target pairs matched in their segments and differed only in the initial tone (i.e., 涂改 /t<sup>h</sup>u<sup>1</sup> kai<sup>l</sup>/  
 227 versus 土改 /t<sup>h</sup>u<sup>l</sup> kai<sup>l</sup>/). The words in the distracter pairs also matched in their segments and differed  
 228 only in the initial tone (i.e., 安保 /an<sup>l</sup> pau<sup>l</sup>/ - 暗堡 /an<sup>l</sup> pau<sup>l</sup>/), but the target and distracter pairs  
 229 differed from each other in all of their segments.

230 Given the salience of the sandhi contrast, there was a concern that participants would realize that the  
 231 focus of the experiment was on sandhi, which may then impact their eye-movements during the  
 232 tasks. To avoid this issue, several steps were taken to distract away from the prominence of the  
 233 critical sandhi trials. First, the structure of a single display was carefully controlled. All words in the  
 234 experiment had a low tone as the second syllable. In every display, each of the four tones was present  
 235 as an initial syllable once. As such, a single display would have four words present: one High+Low,  
 236 one Rising+Low, one Low+Low, and one Falling+Low. In this way, it stands to reason that there  
 237 would occasionally be sandhi items that may be difficult to tell apart, but since every item has a low  
 238 tone second syllable and each tone appears as the initial tone once in a display, it was hoped that the  
 239 sandhi items would be considered an coincidence, as opposed to the focus of the experiment. Second,  
 240 forty-two filler trials were included. The most crucial of these fillers was a set of 14 fillers that were  
 241 identical to the critical trials in all respects, except that which pairs served as targets and distractors

242 was reversed, with the sandhi pair now serving as the distractor pair and one of the other words (with  
243 High or Falling tone on the first syllable) serving as the target. This filler set is crucial, in that it  
244 establishes that a sandhi pair can be present in the display but not targeted. This prevents participants  
245 from being able to know *a priori* if the sandhi pair will be targeted without hearing any acoustic input  
246 and prevents baseline effects. The listener must wait for segmental information to determine which of  
247 the two pairs is being targeted, and then for the tone to determine which item in the pair is being  
248 targeted.

249 The remaining fillers were identical to the critical trials and crucial filler set except that what tone  
250 patterns that constituted the pairs of words was switched. While the Rising+Low and Low+Low  
251 items matched in their segments in the critical and first filler sets, the Rising+Low and Falling+Low  
252 items matched in their segments and the Low+Low and High+Low items matched in their segments.

253 All pair types served as the targeted pair an equal number of times throughout the experiment. Each  
254 of the four tones served as the initial syllable for a target word an equal number of times through the  
255 experiment. Each of the four tones were present on the screen as the initial syllable an equal number  
256 of times through the experiment.

257 Additionally, the characters of the second syllables of the targeted pairs were controlled so that for  
258 half of the trials the second characters matched (i.e., 夹板 - 甲板) and for half they mismatched even  
259 though their pronunciations were the same (i.e., 宣纸- 选址). All critical trials had matching second  
260 characters.

261 Critical items were not repeated across audio conditions. The audio conditions for each item were  
262 separated out onto 2 lists in a Latin square design, and participants were randomly assigned to a list.  
263 In List 1, participants heard a token that was produced as an underlyingly Rising tone for half of the  
264 critical trials and a token that was produced as a sandhi-derived Rising tone for the other half of the  
265 critical trials. This pattern was reversed for List 2. The displays were identical across the lists and  
266 conditions, only the audio was manipulated.

### 267 **2.3 Procedure**

268 The experiment was compiled using Experiment Builder software (SR Research). Participants' eye-  
269 movements were recorded with a desktop EyeLink 1000 Eye Tracker recording at 1000 Hz (1 gaze  
270 position sample recorded every millisecond). The experiment began with a calibration of the  
271 participants' pupil and corneal reflection. This calibration was followed by the practice session of 4  
272 trials. After any questions were answered, the experiment began. A trial began with four words  
273 appearing on the screen in Times New Roman size 80 font in white on a black background in a non-  
274 displayed 2x2 grid. The words remained on the screen for 3,000ms (preview time). This time allowed  
275 participants to pre-activate the pronunciations of each of the words and to familiarize themselves  
276 with their locations. No auditory stimulus was heard during this presentation. After the 3,000ms  
277 preview, the images disappeared, and a fixation cross appeared in the middle of the screen for 500ms  
278 to return the participant's gaze to a neutral starting point. As the fixation cross disappeared, the words  
279 reappeared on the screen in the same locations as during the preview, and an auditory stimulus was  
280 heard through headphones. This auditory stimulus was the target word for that trial, heard in isolation.  
281 Participants were instructed to click on the word spoken as quickly as possible. Once the participant  
282 clicked, a blank screen appeared for 700ms, after which the next trial began. Both eye-movements  
283 (recorded from the target-word onset in the auditory stimulus) and selection accuracy were recorded.

284 **2.4 Data analysis**

285 The data were exported using SR Research DataViewer software in a samples report with an interest  
 286 period from the word onset until the participant clicked on an item. The samples report provides  
 287 information on the individual samples of gaze position and does not pre-determine groupings of  
 288 samples into fixations. The resulting file was pre-processed using a Python script written by SR  
 289 Research that primarily bins time into 8ms bins, calculates the proportion and count data for each of  
 290 the four interest areas for each bin, excludes blinks and saccades, and removes several unnecessary  
 291 columns. This script also excludes fixations that fell outside of the interest areas. The interest areas  
 292 were set to a 300x200 pixel rectangle around the approximately 150x75 pixel words (The width of  
 293 this region varied slightly depending on the characters of the word). This allows for slight human or  
 294 tracker error in the accuracy of the tracking. The interest areas did not overlap and were separated by  
 295 at least 200 pixels of blank space. The size and arrangement of the interest areas was pre-  
 296 programmed into the experiment before data collection began.

297 After this preprocessing, visualization and statistical analysis were carried out using R (R Core Team,  
 298 2016). Trials in which the participant did not click on either the Low- or the Rising-tone visual target  
 299 word were excluded from further analysis. The hypothesis that there were more looks to the Low  
 300 target when hearing the sandhi-derived than the underlying Rising stimulus was evaluated with a  
 301 cluster-based permutation test (Maris & Oostenveld, 2007). This test, originally developed for  
 302 analyzing event-related potential data, controls for multiple comparisons and removes some  
 303 researcher degrees of freedom (albeit creating others) by eliminating the need for the researcher to  
 304 choose a particular time window for analysis. Instead, in this test, an uncorrected comparison is done  
 305 at every sample, clusters of adjacent samples that pass an *a priori* threshold are formed (to address  
 306 the fact that consecutive samples of data are not statistically independent), and a Monte Carlo  
 307 permutation test is performed to evaluate the statistical significance of the difference. For our test, we  
 308 used a paired-sample t-test at each time bin to test the hypothesis that the looks to the Low target  
 309 were greater when hearing the sandhi-derived Rising stimulus than when hearing the underlying  
 310 Rising stimulus.<sup>4</sup> Any series of one or more adjacent time bins with uncorrected  $p < .1$  in this test was  
 311 marked as a temporal cluster (we used a relatively liberal clustering threshold based on our *a priori*  
 312 interest in effects that may be weak but long-lasting, which are best detected with a loose threshold,  
 313 per Maris & Oostenveld, 2007). The t-statistics in each cluster were summed, and the greatest t-  
 314 statistic sum across all clusters was saved as the test statistic. Next, 1000 Monte Carlo permutation  
 315 statistics were drawn. For each permutation, the condition labels were randomly permuted in each  
 316 participant (since there were only two conditions, this essentially means that for each participant we  
 317 flipped a coin to randomly determine whether to keep the original pairing between data and condition  
 318 label or to switch them), then a paired t-test was conducted on each time bin, the t-statistics were  
 319 summed within each cluster (where a "cluster" is the same bin or series of bins identified as a cluster  
 320 in the original analysis of the observed data), and the largest of the t-statistic sums was taken as the  
 321 permutation statistic. The distribution of these 1000 permutation statistics formed the permutation  
 322 distribution, and the final corrected p-value for the cluster-based permutation test was the proportion

8

---

<sup>4</sup> In our pre-registered plan we said that we would logit-transform the data first, but we ended up not doing this because the individual participant data included many zeros, the logit of which is undefined. While our data do not meet the distributional assumptions required for getting a p-value from a t-test, it is important to note that neither these t-tests nor the p-values derived from them were used to draw statistical inferences or conclusions in this study. Rather, they are only used to put time bins into clusters for the purpose of creating a test statistic. The final p-value, and the inferences and conclusions drawn with its help, are based on a non-parametric permutation test, as described by Maris & Oostenveld (2007). Thus, this non-parametric test is still valid, regardless of how the initial sample-wise tests for clustering are performed.

323 of permutation statistics that were greater than or equal to the original test statistic. This test does not  
324 license inferences about exactly when a difference between conditions starts and ends; rather, it  
325 licenses an inference about whether there is a difference between conditions overall, anywhere in the  
326 curve, after controlling for multiple comparisons.

### 327 **3 Results**

#### 328 **3.1 Behavioural**

329 Overall, participants clicked on the character corresponding to the Low-tone target on 56% of trials  
330 when they heard the stimulus with the sandhi-derived Rising tone, and on 52.6% of trials when they  
331 heard the stimulus with the underlyingly Rising tone. This is in the direction one might expect if  
332 listeners are sensitive to the difference, but this difference was not significant in a logistic mixed-  
333 effects model with maximal random effects for participants and items ( $b=0.28$ ,  $z=1.15$ ,  $p=.249$ ). As  
334 shown in Figure 1, while some items tended to be preferentially identified with the Low or Rising  
335 word target (likely due to differences between the two targets' frequencies or other lexical  
336 characteristics), the choice of which word to click was not strongly or reliably influenced by the  
337 auditory stimulus.

338 ---Figure 1 about here---

#### 339 **3.2 Eye-tracking**

340 The left side of Figure 2 shows the proportion of looks to the Low target over time, as a function of  
341 which auditory stimulus was heard. The right side shows the proportion of looks to the Rising target;  
342 while we did not perform statistical analyses on these, since we only pre-registered analysis for the  
343 looks to the Low target, we show the data here for completeness. For each condition, the proportion  
344 of fixations peaks around 50% rather than reaching 100%; this is not surprising, since the stimuli are  
345 ambiguous. Similar to what Speer and Xu (2008) observed, the early portion of the time windows  
346 shows a counterintuitive pattern, with participants looking more at the putatively inappropriate  
347 target—i.e., looking at the Low tone more when hearing an underlying Rising tone word than when  
348 hearing a sandhi-derived Rising tone, and looking at the Rising tone word more when hearing a  
349 sandhi-derived Rising tone than when hearing an underlying Rising tone. Crucially, in the latter  
350 portion of the window, participants appear to look more at the putatively appropriate target: i.e., for  
351 targets corresponding to a word whose first syllable is Low tone, participants look more at these  
352 targets when hearing a Rising tone that was derived from a Low tone via tone sandhi, compared to  
353 when hearing an underlyingly Rising tone. This trend was not significant, however, in our pre-  
354 registered statistical analysis ( $p=.351$ ).

355 ---Figure 2 about here---

356 Figure 3 shows the pattern across participants and items. While the effect was statistically significant  
357 in the direction we predicted, it is also clear that it varies substantially across participants and also  
358 varies somewhat across items. The pattern is much clearer for items, given that each item had  
359 observations from as many as 30 participants per condition, whereas each participant had  
360 observations from only 7 or fewer items per condition. We performed an additional exploratory  
361 cluster-based permutation test, relaxing the cluster threshold to  $p<.3$  (as looser thresholds are more  
362 sensitive for detecting weak but long-lasting effects [Maris & Oostenveld, 2007], like the one  
363 observed here) and calculating the fixation proportions by item rather than by participant, since the

364 effects for items were more stable. This yielded a  $p=.098$  effect, although  $p$ -values are not  
 365 interpretable for this test since it is not confirmatory.

366 ---Figure 3 about here---

367 Overall, the results are inconclusive, consistent with both the presence and the absence of a  
 368 difference between conditions. On the one hand, the pre-registered analysis did not yield a  
 369 statistically significant effect, so we are not able to reject the possibility that the tones are completely  
 370 neutralized in perception. On the other hand, visual inspection of the data suggest that it may not be  
 371 reasonable to conclude that they are completely neutralized either (keeping in mind that failure to  
 372 reject the null hypothesis is not the same as acceptance of the null hypothesis), and exploratory  
 373 analysis suggests that our pre-registered analysis plan may not have been as sensitive as it could  
 374 have; i.e., if we replicate the study with more items, and perform confirmatory analysis on a more  
 375 stable proportion measure and use a more appropriate clustering threshold, then it may be possible to  
 376 detect a significant effect. None of these facts justify concluding that there is a reliable difference in  
 377 *this* experiment, as these are all post-hoc decisions about researcher degrees of freedom (e.g.,  
 378 Roettger, ms.; Simmons et al., 2011) and could just reflect capitalizing on noise patterns. At the very  
 379 least, however, they suggest that this paradigm *may* be capable of providing evidence for incomplete  
 380 neutralization in the perception of Mandarin Low and Rising tones, and they demonstrate that this  
 381 issue is worth further investigation and replication.

## 382 4 Discussion

383 We used visual world eye-tracking to examine whether participants are somewhat sensitive to the  
 384 underlying form of underlying vs. sandhi-derived Rising tones in Mandarin. Using eye-tracking, we  
 385 examined whether participants were more likely to look at an orthographic form representing an  
 386 underlying Low tone (e.g., 赌本 "bookie", /du˩ pən˩/) when they heard a token of [du˩ pən˩] that was  
 387 produced as a reading of this underlying form (with the tone of its first syllable changed to Rising  
 388 tone because of tone sandhi) and less likely to look at this orthographic form when they heard a token  
 389 of [du˩ pən˩] that was produced as a reading a form with Low tone underlyingly (such as 读本  
 390 "reading book", /du˩ pən˩/). The results of our eye-tracking experiment were inconclusive: there was  
 391 a numerical trend in the direction we predicted, which was not statistically significant but also not  
 392 consistent with there being no difference. Statistically speaking, with a result like this we do not have  
 393 evidence to conclude either that there is a difference or that there is no difference; rather, we can only  
 394 conclude that the result is indeterminate between these possibilities and that further research with  
 395 higher power (especially with more items, and perhaps also with more participants) is needed to  
 396 resolve the issue. At least, however, the results suggest that the present paradigm may be a way to  
 397 observe partially intact discrimination in the perception of Mandarin tone sandhi, challenging the  
 398 previously accepted view that this alternation is completely neutralizing in the perception domain.

399 If perception of Mandarin Low and Rising tones is indeed not completely neutralized by tone sandhi  
 400 (a speculation which needs further research to confirm), it is worth noting that this clearly is a weak  
 401 trend. Even if participants are found to reliably identify sandhi-derived Rising tones as underlyingly  
 402 Low slightly more often than they identify underlyingly Rising tones as Low, this trend is far from  
 403 100%; if anything, we expect that further research will show participants to be either at chance or just  
 404 slightly above chance. This is very different, then, than usual phonological contrasts that speakers  
 405 can very reliably distinguish. This would put Mandarin third tone sandhi on a similar level as other  
 406 classical incomplete neutralization cases, such as Dutch final devoicing, where participants also show  
 407 a slight (better than chance but still small and very imperfect) sensitivity to incompletely neutralized

408 sounds (e.g., Warner et al., 2004). Thus, this perceptual difference may be of little practical use to  
409 interlocutors during real-time language comprehension.

410 To further examine why the experiment may have failed to observe a significant sensitivity to  
411 incompletely neutralized Mandarin tone sandhi, we checked the F0 tracks of the stimuli to confirm  
412 whether the speaker actually did incompletely neutralize their tones as is typically reported. Figure 4  
413 shows the pitch tracks for each item. While some items, such as 1, 7, and 12, do indeed show the  
414 typical pattern, with underlying Rising tone (solid red line) being higher overall than sandhi-derived  
415 Rising tones (dashed blue line), some show little difference (e.g., items 5 and 11) or a difference in  
416 the opposite direction (e.g., items 3 and 8). Thus, the incomplete neutralization of Rising and Low  
417 tones in production cannot be said to be reliable in this speaker.

418 ---Figure 4 about here---

419 Did listeners more accurately distinguish Rising from Low tones in those items where the production  
420 distinguished them better? This does not appear to be the case, as shown by comparing Figures 4 and  
421 5. The items which showed the largest or most reliable difference in eye-tracking are not necessarily  
422 the ones which showed the largest difference in production.

423 ---Figure 5 about here---

424 Overall, then, we can see that the incomplete neutralization of Low and Rising tones in production is  
425 real (i.e., statistically significant and widely replicated) but is not entirely reliable across or within  
426 speakers (i.e., it does not happen all the time for every item produced by every speaker), and that  
427 perception of these neutralized tones is also not very reliable: either it is completely neutralized, as  
428 previous studies have argued, or even if it is not completely neutralized, it is still far less than perfect.  
429 In fact, the unreliability of the contrast in production may play a role in why the distinction is not  
430 reliably used in perception. We can think of three possible mechanisms for why perception of  
431 incompletely neutralized tones may be poor. First, it may be that the acoustic difference between  
432 underlying and sandhi-derived Rising tones is imperceptible or difficult to perceive, even when it is  
433 present. If that is the case, we expect that perception of these tones will be inaccurate on all items, in  
434 experiments such as the present one. Secondly, it is possible that the acoustic difference between the  
435 surface forms is indeed perceptible when the surface forms are not completely neutralized, but that  
436 this incomplete neutralization does not always occur (as seen in our own stimuli, shown in Figure A).  
437 In that case, we would expect that participants are fairly reliable at identifying the tones, but only on  
438 the items that show good incomplete neutralization in production; this does not seem very consistent  
439 with our findings. Finally, a last possibility is that the difference between incompletely neutralized  
440 tones is perceptible when it is present, but that participants know it is not a reliable cue they do not  
441 regularly use it to drive their perception. If that is the case, we would expect perception to be  
442 inaccurate on all items, the same as in the first possibility outlined above. A potential piece of  
443 evidence for this last possibility would be if participants show above-chance accuracy at perceiving  
444 these tones in an on-line, implicit measure like eye movements, but not in an explicit metalinguistic  
445 measure like end-state judgments; that would be consistent with the notion that participants can  
446 perceive the difference somewhat but end up not using it to inform their eventual categorization of  
447 the sound. This is indeed a pattern suggested by the present results (in our case both the eye  
448 movements and the clicks do show non-significant trends in the direction consistent with incomplete  
449 neutralization, but the trend for eye movements seems stronger than that for clicks; although, given  
450 that these are very different measures with different properties, it is difficult to directly compare the

451 sizes of these measures, so future research will be valuable to confirm whether perception of tones  
452 really is less neutralized in eye movements than in end-state responses).

453 A surprising aspect of our results is that, as shown in Figure 2, there was a period in the early time  
454 window at the beginning of the syllable where participants showed a trend in the opposite of the  
455 expected direction: they looked at the Low target slightly more when hearing underlying Rising tones  
456 compared to when hearing sandhi-derived (underlyingly Low) tones. We had not predicted this, and  
457 it was not significant in our analysis (in fact it could not have been, as our analysis used one-sided  
458 tests). This pattern was also observed, however, by Speer and Xu (2008). Thus far we have no  
459 explanation for it. Given that it has shown up in two experiments by now, however, it is highly  
460 subjective, and in future research it will be important to examine whether this pattern is replicated,  
461 now that there is substantial *a priori* reason to expect it. If future research finds this pattern to be  
462 reliable, an explanation will be needed for why this counterintuitive eye movement pattern emerges  
463 in tone perception.

464 One limitation of the study is that, while we used a large number of listeners and items, the  
465 experiment only used one speaker. Thus, we cannot be sure that the results will generalize to  
466 perception of tokens from other speakers; this is an open question for future research. We do note  
467 that using one speaker is currently standard for studies in this area (Peng [2000], Liu [2013], and  
468 Zhang & Peng [2013] each used one speaker in their perceptual tasks, and Wang and Li [1967] used  
469 two). Nonetheless, acoustic differences between underlying and sandhi-derived Rising tones may not  
470 be constant. For instance, the acoustic difference between sandhi-derived and underlyingly Rising  
471 tones in real words reported in several studies (Peng, 2000; Liu, 2013; Zhang & Lai, 2010) was not  
472 significant in the study by Zhang and Peng (2013); and even when a difference is significant, it is not  
473 necessarily present for all speakers (see, e.g., Liu, 2013, and Zhang and Peng, 2013). Therefore,  
474 examining the extent to which incompletely neutralized perception generalizes across speakers (if it  
475 occurs at all) is a valuable question for future study.

## 476 **5 Conflict of Interest**

477 *The authors declare that the research was conducted in the absence of any commercial or financial*  
478 *relationships that could be construed as a potential conflict of interest.*

## 479 **6 Author Contributions**

480 SP and YH conceived the experiment. SP, KC, LP and YH designed the experiment. KC and LP  
481 created the stimuli. KC programmed the experiment and LP collected the data. SP and KC analyzed  
482 the data. All authors wrote and approved the manuscript.

## 483 **7 Funding**

484 This research was supported by grant G-UACX from the Department of Chinese and Bilingual  
485 Studies to SP.

## 486 **8 Ethics statement**

487 This study was carried out in accordance with the recommendations of the Human Subjects Ethics  
488 Sub-committee at the Hong Kong Polytechnic University with written informed consent from all  
489 subjects. All subjects gave written informed consent in accordance with the Declaration of Helsinki.

490 The protocol was approved by the Human Subjects Ethics Sub-committee at the Hong Kong  
491 Polytechnic University.

## 492 **9 Acknowledgments**

493 The authors would like to thank Jueyao Lin for assistance with stimulus creation.

## 494 **10 References**

495 Cai, Q., & Brysbaert, M. (2010). SUBTLEX-CH: Chinese word and character frequencies based on  
496 film subtitles. *PLoS ONE*, *5*, e10729.

497 Chen, S., Wiltshire, C., & Li, B. (2017). Statistical modeling of Mandarin tone sandhi: neutralization  
498 of underlying pitch targets. *Proceedings of the World Academy of Science, Engineering, and*  
499 *Technology*, *11*.

500 Cheng, C., Chen, J., & Gubian, M. (2013). Are Mandarin sandhi tone 3 and tone 2 the same or  
501 different? The results of functional data analysis. *The 27th Pacific Asia Conference on Language,*  
502 *Information, and Computation*.

503 Cousineau, D. (2005). Confidence intervals in within-subject designs: a simpler solution to Loftus  
504 and Masson's method. *Tutorials in Quantitative Methods for Psychology*, *1*, 42-45.

505 Flege, J., Munro, M., & Skelton, L. (1992). Production of the word-final English /t-/d/ contrast by  
506 native speakers of English, Mandarin, and Spanish. *Journal of the Acoustical Society of America*, *92*,  
507 128-143.

508 Kim, H., & Jongman, A. (1996). Acoustic and perceptual evidence for complete neutralization of  
509 manner of articulation in Korean. *Journal of Phonetics*, *24*, 295-312.

510 Kuo, Y., Xu, Y., & Yip, M. (2007). The phonetics and phonology of apparent cases of iterative tone  
511 change in Standard Chinese. In *Phonology and Phonetics: Tones and Tunes, volume 2*, 212-237.  
512 Berlin: Mouton de Gruyter.

513 Lin, Y., & Hsu, Y. (2018). Whether and how do Mandarin sandhied tone 3 and underlying tone 2  
514 differ? *The 32nd Pacific Asia Conference on Language, Information and Computation*.

515 Liu, X. (2013). 上声变调的声学与感知实验研究 [Acoustic and perceptual research on third tone  
516 sandhi]. 文教资料 [Culture and Education Data], *29*, 139-142. (In Chinese)

517 Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-data.  
518 *Journal of Neuroscience Methods*, *164*, 177-190.

519 Morey, R. (2008). Confidence intervals from normalized data: a correction to Cousineau (2005).  
520 *Tutorials in Quantitative Methods for Psychology*, *4*, 61-64.

521 Neergard, K., Xu, H., & Huang, C. (2016). Database of Mandarin neighborhood statistics.  
522 *Proceedings of the Conference on Language Resources and Education*.



- 523 Nicenboim, B., Roettger, T., & Vasisht, S. (2018). Using meta-analysis for evidence synthesis: The  
524 case of incomplete neutralization in German. *Journal of Phonetics*, 70, 39-55.
- 525 Peng, S. (2000). Lexical versus phonological representations of Mandarin sandhi tones. In *Language*  
526 *acquisition and the lexicon: Papers in laboratory phonology V*, Michael Broe and Janet  
527 Pierrehumbert [Eds.]. pp. 152-167. Cambridge, UK: Cambridge University Press.
- 528 R Core Team (2016). R: A language and environment for statistical computing. R Foundation for  
529 Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.
- 530 Roettger, T. (ms.). Researcher degrees of freedom in phonetic research. <https://psyarxiv.com/fp4jr>
- 531 Simonsohn, U., Nelson, L., & Simmons, J. (2014). P-curve: a key to the file drawer. *Journal of*  
532 *Experimental Psychology: General*, 143, 534-547.
- 533 Speer, S., & Xu, L. (2008). Processing lexical tone in third-tone sandhi. Talk presented at *Laboratory*  
534 *Phonology* 11.
- 535 Wang, W., & Li, K. (1967). Tone 3 in Pekinese. *Journal of Speech and Hearing Research*, 10, 629-  
536 636.
- 537 Warner, N., Jongman, A., Sereno, J., & Kemps, R. (2004). Incomplete neutralization and other sub-  
538 phonemic durational differences in production and perception: evidence from Dutch. *Journal of*  
539 *Phonetics*, 32, 251-276.
- 540 Yuan, J., & Chen, Y. (2014), 3rd tone sandhi in standard Chinese: a corpus approach. *Journal of*  
541 *Chinese Linguistics*, 42, 218-237.
- 542 Zhang, J., & Lai, Y. (2010). Testing the role of phonetic knowledge in Mandarin tone sandhi.  
543 *Phonology*, 27, 153-201.
- 544 Zhang, J., Lai, Y., & Sailor, C. (2011). Modeling Taiwanese speakers' knowledge of tone sandhi in  
545 reduplication. *Lingua*, 121, 181-206.
- 546 Zhang, C., & Peng, G. (2013). Productivity of Mandarin Third Tone sandhi: A wug test. In Peng G.  
547 and Shi F. (Eds.) *Eastward Flows the Great River: Festschrift in Honor of Prof. William S-Y. Wang*  
548 *on his 80th Birthday*, pp. 256-282. Hong Kong: City University of Hong Kong Press.
- 549 Zhou, X., & Marslen-Wilson, W. (1997). The abstractness of phonological representation in the  
550 Chinese mental lexicon. In *Cognitive Processing of Chinese and Related Asian Languages*, Hsuan-  
551 Chih Chen [Ed.], 32-27. Hong Kong, Hong Kong: The Chinese University Press.

## 552 1 Data availability statement

553 The datasets generated for this study can be found in the Open Science Foundation repository  
554 [<https://osf.io/ursh9/>].

555

556

557 **Tables**

558 Table 1: Summary statistics about lexical properties of the materials. Frequency is listed in natural  
 559 log of words/characters per million; for words that did not appear as the corpus, we coded them as  
 560 having 1 occurrence (0.03 per million) to avoid undefined log values. For each cell, the first row  
 561 shows the measure for the whole word, and the second row shows the measure for the first and  
 562 second syllable, in parentheses. Frequency, measures are from the SUBTLEX-CH corpus (Cai &  
 563 Brysbaert, 2010) and homophone/neighbourhood density measures are from the Database of  
 564 Mandarin Neighborhood Statistics (Neergard, Xu, & Huang, 2016).

	Frequency	Homophone Density	Neighbourhood Density	Neighbourhood Frequency	Stroke Count
Underlyingly	-0.95	0.86	0.97	-0.56	16.43
Rising	(4.49, 6.14)	(8.43, 3.57)	(13.93, 13.43)	(2.94, 3.32)	(8.43, 8.07)
Sandhi- derived	-0.17	1.00	0.93	-0.16	15.50
Rising	(5.09, 6.14)	(5.50, 3.57)	(15.50, 13.43)	(3.96, 3.93)	(7.50, 8.07)

565

566 **Figure captions**

567 Figure 1. Proportion of clicks on the visual character target corresponding to the word with Low tone  
 568 in its first syllable, depending on what audio stimulus was heard. The x-axis shows the proportion of  
 569 clicks to this target when participants heard an auditory stimulus in which the first syllable's Rising  
 570 tone was derived through tone sandhi, and the y-axis shows the proportion when the auditory  
 571 stimulus' first syllable's Rising tone was underlyingly Rising. A diagonal line from the lower left to  
 572 the upper right corner indicates where the two proportions are equal. The rest of the plot shows a  
 573 cloud of points, either purple squares or gray circles; each purple square indicates these proportions  
 574 for one item, and each gray circle indicates these proportions for one participants. Because some  
 575 participants' points are directly on top of one another, the gray circles are opaque; darker gray circles  
 576 indicate more participants' data appearing at this same point, and lighter gray circles indicate fewer  
 577 participants. If the Low target is reliably clicked more after hearing a sandhi-derived Rising stimulus  
 578 than after hearing an underlyingly Rising stimulus, then most of the observations should be below the  
 579 diagonal. If most observations are on the diagonal or symmetrically distributed around the diagonal,  
 580 this indicates that there is little effect of the auditory stimulus.

581 Figure 2. Proportion of looks to each word target over time. The figure consists of two panels, one  
 582 representing proportion of looks to the word associated with an underlying Low tone (left) and one  
 583 representing proportion of looks to the word associated with an underlying Rising tone (right). In  
 584 each panel, the x-axis represents time since the onset of the auditory stimulus, and the y-axis  
 585 represents proportion of looks to the given target. A solid red line indicates the proportion of looks to  
 586 that target when the participant heard a token pronounced with an underlyingly Rising tone, and a  
 587 dashed blue line indicates the proportion of looks when the participant heard a token pronounced  
 588 with a Rising tone derived from tone sandhi. Each line is surrounded by an opaque ribbon  
 589 representing a difference-adjusted by-participant Cousineau-Morey interval (Cousineau, 2005;  
 590 Morey, 2008), such that where one condition's interval does not include the other condition's mean  
 591 and vice versa, the conditions might be significantly different; note that these intervals are only an aid  
 592 for visualization/exploration and will not necessarily correspond to the statistical results reported in  
 593 the prose, as they are two-sided and do not take multiple comparisons into account. In the left panel  
 594 (looks to Low tone), from about 250 to 500 ms there appear to be more looks to the Low target when

595 hearing underlying Rising tones, and from about 750 ms to the end of the window (1500 ms) there  
 596 are more looks when hearing sandhi-derived Rising tones. In the right panel (looks to Rising target)  
 597 the reverse of this pattern is seen: initially more looks to the Rising target when hearing a sandhi-  
 598 derived Rising tone, and then more when hearing an Underlying rising tone.

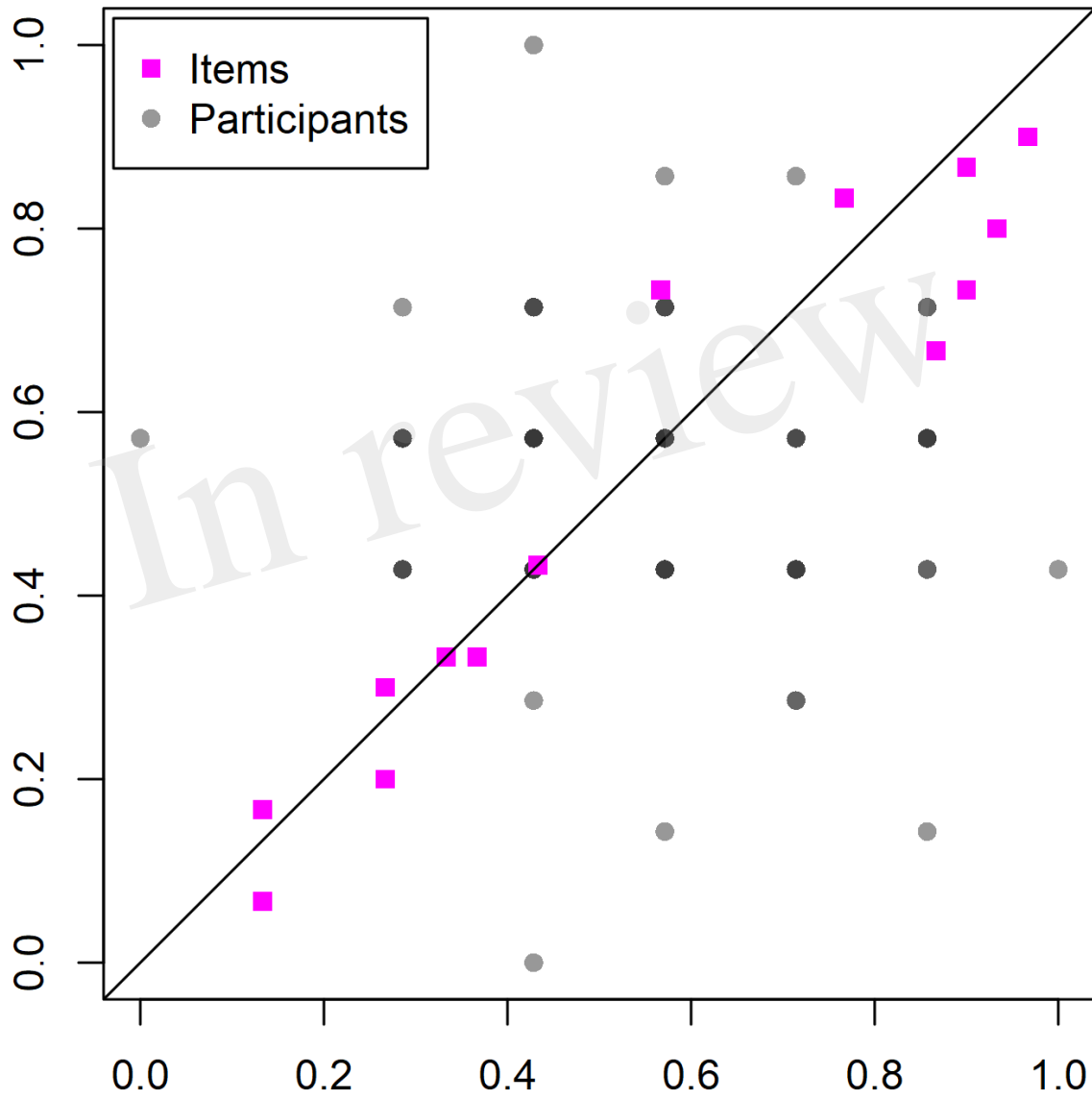
599 Figure 3. Incomplete neutralization effect (difference between looks to Low tone when hearing a  
 600 sandhi-derived Rising tone vs. when hearing an underlying Rising tone) by item (left panel) and by  
 601 participant (right panel). Each panel shows a red line indicating the mean difference over time, and  
 602 multiple gray lines indicating the difference for a given participant or item over time. When a line is  
 603 above zero, this indicates that there were more looks to the Low target when hearing sandhi-derived  
 604 Rising tone than when hearing underlying Rising tone, at this time. The red line is surrounded by a  
 605 shaded red ribbon indicating a two-tailed 95% confidence interval based on the t-statistic; this is just  
 606 a visualization aid and does not necessarily correspond to the statistical results in the prose, as it  
 607 represents two-sided tests and does not take multiple comparisons into account. In each graph, there  
 608 is a region around 1000ms where the mean difference is pretty reliably above zero (consistent with  
 609 this same effect shown in Figure Y and through the statistical analysis). At this time period, 11 out of  
 610 14 items are above zero as well and only three are below zero; whereas for participants, a substantial  
 611 number of participants are below zero.

612 Figure 4. Pitch tracks for the first syllable of each item, as a function of whether it was pronounced  
 613 with an underlying Rising tone (solid red lines) or sandhi-derived Rising tone (dashed blue lines).  
 614 The figure includes fifteen subplots arranged in a 3x5 grid, where the first fourteen subplots each  
 615 represent the first syllable of one item, and the last shows the legend and axis labels. In each subplot,  
 616 the horizontal axis represents time in a syllable, and the vertical axis represents fundamental  
 617 frequency. Each subplot shows two upward-curving lines, representing the fundamental frequencies  
 618 of underlying and sandhi-derived Rising tones over the course of the syllable. Seven of the fourteen  
 619 subplots show the red line (representing underlying Rising tone) mostly higher than the dashed blue  
 620 line (representing sandhi-derived Rising tone), consistent with typically observed incomplete  
 621 neutralization in third tone sandhi. Three plots show the red line representing underlying Rising tone  
 622 mostly below the blue line representing sandhi-derived Rising tone, two show the two lines pretty  
 623 much on top of one another, and two show the two lines crossing each other more or less evenly.

624 Figure 5. Proportion of looks to the Low target over time, for each item. The plot consists of 15  
 625 subplots in a 3x5 grid, as in the previous figure (Figure 4). As in Figure 2, in each subplot shows the  
 626 x-axis represents time since the onset of the auditory stimulus, and the y-axis represents proportion of  
 627 looks to the given target. Solid red lines indicate the proportion of looks to that target when the  
 628 participant heard a token pronounced with an underlyingly Rising tone, and dashed blue lines  
 629 indicate the proportion of looks when the participant heard a token pronounced with a Rising tone  
 630 derived from tone sandhi.

# Proportion of clicks on Low target

after underlyingly Rising stimulus



after sandhi-derived Rising stimulus

Figure 2.TIF

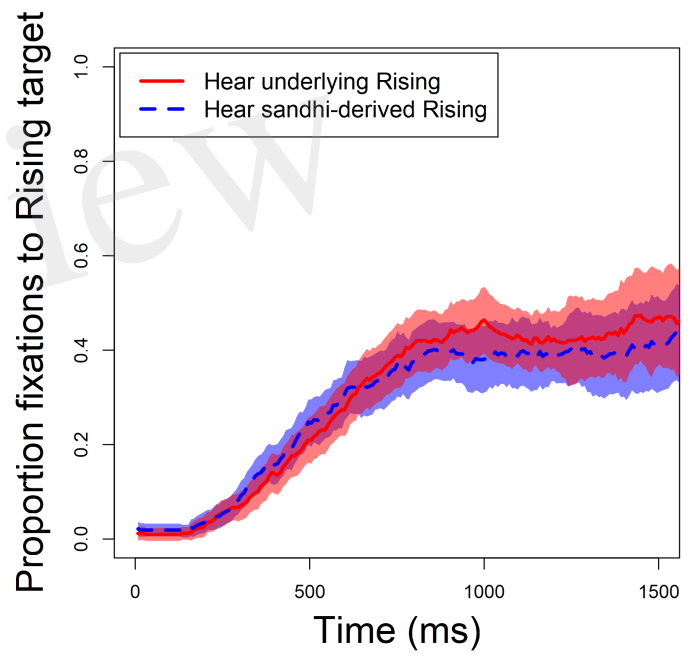
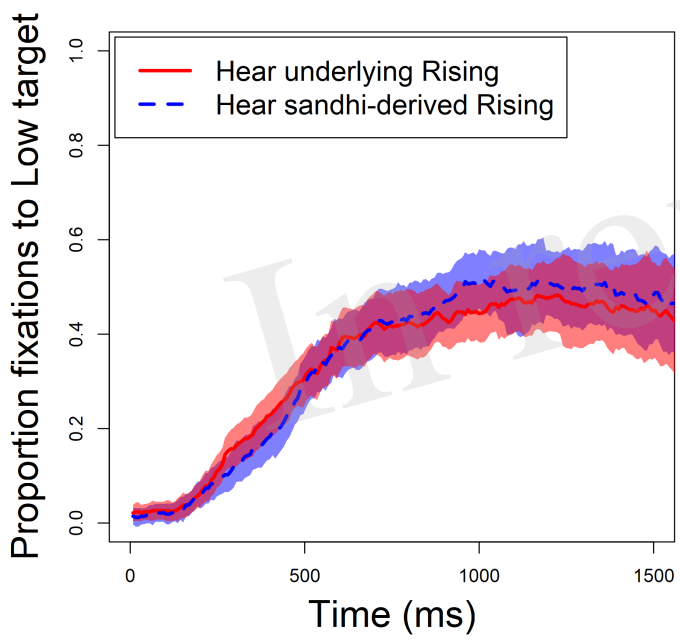


Figure 3.TIF

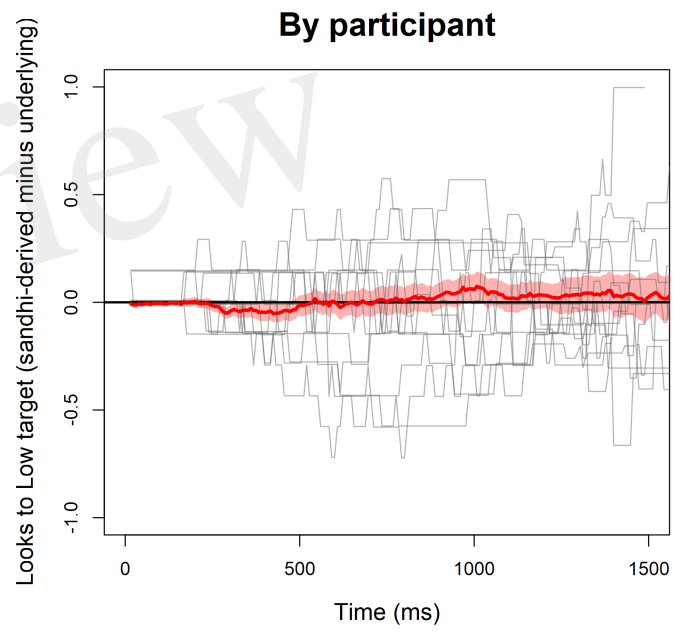
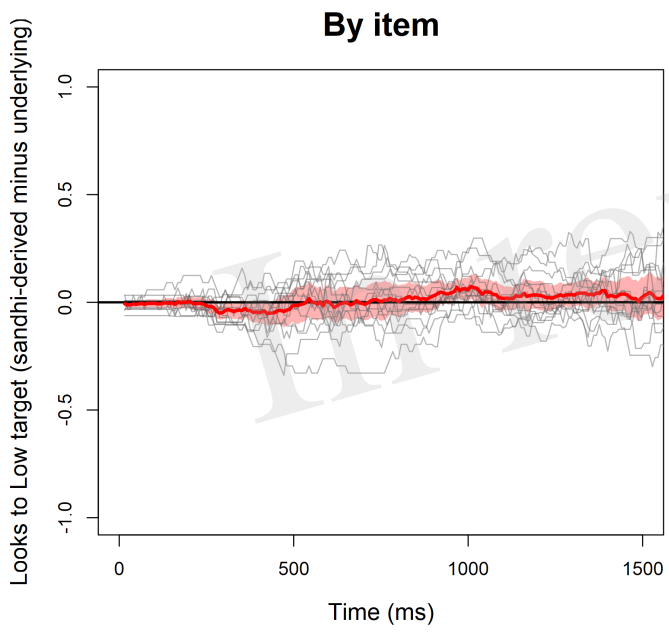


Figure 4.TIF

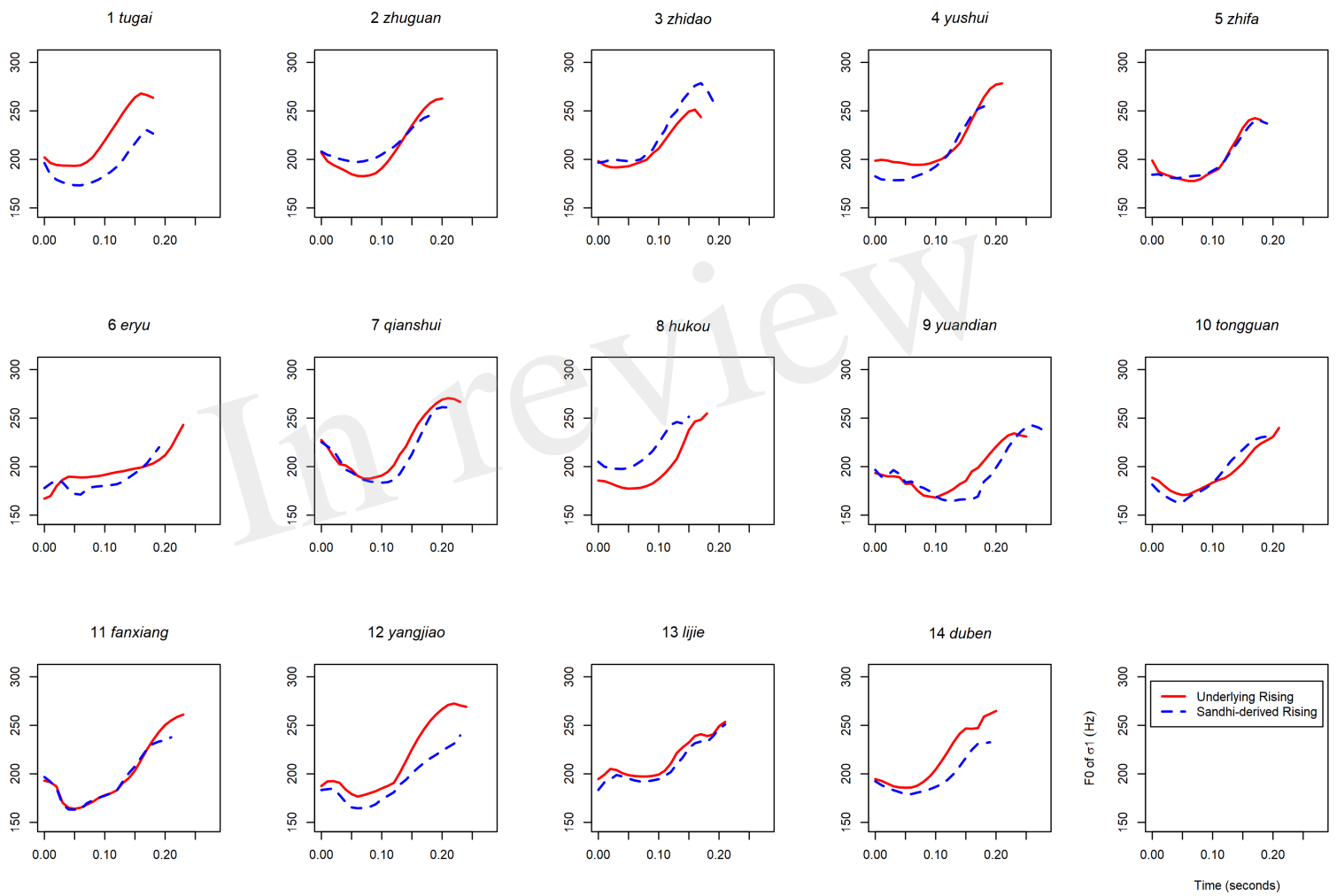


Figure 5.TIF

