PHONOLOGY

**ARTICLE**

# Phonetically incomplete neutralisation can be phonologically complete: evidence from Huai'an Mandarin

Naiyan Du and Karthik Durvasula ⓘ

Department of Linguistics, Languages, and Cultures, Michigan State University, East Lansing, MI 48824, United States of America

**Corresponding author:** Karthik Durvasula; Email: karthikd@msu.edu

**Abstract**

The phenomenon of incomplete neutralisation describes a situation where a putative case of categorical phonological neutralisation is observed to be phonetically non-neutralising. This has been argued to be a problem for phonological theories that employ categorical features. Here, we use two distinct feeding orders of tone sandhi processes from Huai'an Mandarin to show that incomplete phonetic neutralisation is compatible with categorical phonological phenomena. Therefore, incomplete phonetic neutralisation does not automatically inform us of gradient phonological representations. We further show that incomplete phonetic neutralisation can in fact have a large effect size. Such results are not surprising from a classic generative view of phonology where linguistic performance is argued to be a multi-factorial problem, and linguistic knowledge (i.e., competence) is only one of the many factors involved. Furthermore, our results suggest that the observed incompleteness or gradience may have a source outside phonological knowledge.

## Contents

## 1 Introduction

As has long been recognised in discussions of linguistic competence as abstract knowledge, there are multiple potential interacting factors in performance (Chomsky 1964, 1965; Valian 1982; Schütze 1996, *inter alia*). Similarly, there are also multiple interacting sources that affect speech production, for example, lexical knowledge, phonological knowledge and memory and processing constraints (Whalen 1991, 1992; Warner *et al.* 2004; Wright 2004). Consequently, gradience in phonetic manifestations cannot automatically be used as a diagnostic of gradient phonological representations. In this article, we explore the phenomenon of incomplete neutralisation to argue that incomplete phonetic neutralisation does not automatically inform us about phonological representations or phonological knowledge, more generally. We present relevant data from Tone 1 and Tone 4 sandhi processes in Huai'an Mandarin (Huai'an hereafter), both of which crucially participate in feeding orders to trigger other tone sandhi processes, to argue that phonetically incomplete neutralisation can still be phonologically complete.

Since at least the mid-1980s, the effect of incomplete neutralisation has been documented in a variety of languages including Catalan (Dinnsen & Charles-Luce 1984), Dutch (Warner *et al.* 2004; Ernestus & Baayen 2006), Japanese (Braver & Kawahara 2016), Polish (Slowiaczek & Dinnsen 1985; Slowiaczek & Szymanska 1989) and Russian (Dmitrieva 2005; Kharlamov 2012; Matsui 2015). For example, in German, it has been described that the phonological voicing contrast for obstruents is neutralised at the right edge of a prosodic word (Wagner 2002). A rule-based mapping of the relevant phonological process is stated in (1). However, careful phonetic and perceptual experimentation has shown that the neutralisation is incomplete phonetically (Port & O'Dell 1985; Roettger *et al.* 2014, *inter alia*). In other words, *underlying voiceless* stops, *derived voiceless* stops and *underlying voiced* stops all have different phonetic distributions.[1]

---

[1] To be succinct, here we use 'underlying voiceless stop' to mean surface voiceless stops that map from underlying voiceless stops, 'derived voiceless stop' to mean surface voiceless stops that map from underlying voiced stops, and 'underlying voiced stops' to mean surface voiced stops that map from underlying voiced stops. We will use the terms 'underlying' and 'derived' in the same fashion to describe stops and tones for the rest of the article.

(1) $[-\text{sonorant}] \rightarrow [-\text{voice}] / \_)_\omega$
where ω represents a prosodic word

The observed effect of incomplete neutralisation has been argued by many researchers to pose a challenge to traditional formal phonology where categorical phonological representation and modular feed-forward model are assumed (Manaster Ramer 1996; Port & Leary 2005; Goldrick & Blumstein 2006; Roettger *et al.* 2014; Braver 2019; McCollum 2019). It is often assumed that under such a feed-forward framework, the phonological representations are discrete elements that do not contain any gradient phonetic information, and phonetics *only* has access to the output of phonology (Kenstowicz 1994; Pierrehumbert 2002). We call this view the *Standard generative view of phonology*.[2] As a result, underlying representations that undergo phonological neutralisation process should not have any consequence on phonetic manifestations. However, in cases of incomplete neutralisation, there are traces of the underlying representation in the phonetic manifestation.

One trivial but widely adopted solution to the puzzle posed by incomplete neutralisation has been to simply deny that such an effect can be caused by grammatical knowledge (Dinnsen & Charles-Luce 1984; Fourakis & Iverson 1984; Manaster Ramer 1996; Warner *et al.* 2004, *inter alia*). To support this claim, several criticisms have been raised against previous experimental designs as well as the interpretation of the results. One main criticism is whether the observed phenomenon of incomplete neutralisation is due to task effects. Among these effects, the most discussed one is orthography. It has been noticed by many researchers (Fourakis & Iverson 1984; Manaster Ramer 1996, *inter alia*) that in the seminal work of Port & O'Dell (1985), participants were presented stimuli orthographically where minimal pairs were always in contrast. Native speakers of German may have hypercorrected and produced unnatural speech to match the forms of orthography. This suspicion becomes especially disturbing when Warner *et al.* (2004) showed a significant production difference in words that are identical in underlying representations but differ in orthography in Dutch.[3] To circumvent the interference of orthography, two methods have been employed, namely changing the experimental paradigm and looking at languages where the relevant phonological contrast is not reflected orthographically.

In line with changing the experimental paradigm, Fourakis & Iverson (1984) employed a unique strategy aimed at concealing the morphological forms that native speakers of German are supposed to produce, so the influence of orthography is expected to decrease. The participants were instead presented auditorily with the conjugated form where the contrast is maintained in both underlying and surface representations and asked to decompose and produce the bare form where the incomplete neutralisation is expected to happen. Through this paradigm, Fourakis and Iverson camouflaged the task as a morphological exercise to distract the participants to elicit

---

[2]While this conception of phonology is relatively standard in our opinion, it is actually quite different from the *Classic generative view of phonology* where phonology is seen as *knowledge*. We return to this issue in §7, where we point out that the latter notion has no trouble in accounting for facts related to incomplete neutralisation.

[3]In support of our larger point in this article, we would like to point out that their observation in fact shows how powerful performance (i.e., non-phonological) factors can be in accounting for phonetic manifestations.

more natural pronunciations. Interestingly, the effect of incomplete neutralisation was not observed, and Fourakis and Iverson concluded that the previously found incomplete neutralisation was actually a task effect. Using a different strategy, Jassem & Richter (1989) asked participants to answer questions designed to elicit the target words in Polish, and observed no evidence of incomplete neutralisation. However, by implementing the same strategy as Fourakis & Iverson (1984) and increasing the statistical power with more speakers and more test minimal pairs, Roettger *et al.* (2014) found an effect of incomplete neutralisation. However, it is worth noting that, as Roettger *et al.* (2014) themselves pointed out, the strategy employed by Fourakis and Iverson can incur a potential artifact of phonetic accommodation. In the experimental paradigm, the participants hear the conjugated form where neutralisation cannot happen and the voicing contrast is present, and have to produce the form where neutralisation does happen. In such a paradigm, the observed effect of incomplete neutralisation may be due to the participants mirroring vowel duration differences in the stimulus recordings they heard of the conjugated forms. When Roettger *et al.* (2014) controlled for this confound in one of their experiments, they found only a very small, non-significant effect (<3 ms) in the right direction. This suggests that there might indeed be no clear evidence for incomplete neutralisation even in their well-powered study. To sum up, this general strategy to solve the problem of orthography by changing the task performed by the participants leads to very weak evidence (if that) for the presence of incomplete neutralisation.

A second method employed to overcome task effects related to orthography has been to use a language where the crucial contrast is not marked in the orthography. For example, Catalan has been claimed to have a devoicing process but no orthographic marking of an underlying voicing contrast under any phonological conditions, and Dinnsen & Charles-Luce (1984) did not observe any evidence of incomplete neutralisation in Catalan devoicing. However, later Charles-Luce & Dinnsen (1987) reanalysed their data and found incomplete neutralisation in the cue of voicing into closure. Here, it is worth noting that in Catalan, quite a few words actually maintain the underlying voicing contrast in orthography, so the real situation is more complicated and Catalan cannot simply be treated as a language that does not mark underlying voicing contrast orthographically (Badia Margarit 1962; Manaster Ramer 1996).

In another case, Braver & Kawahara (2016) observed a putative case of incomplete neutralisation in Japanese. Most of their stimuli were presented in Chinese characters (kanji), which is an orthographic system that is commonly used in Japanese but only has a very weak connection with pronunciation.[4] Although most Chinese characters were originally created by combining a part that indicates pronunciation and a part that indicates meaning (Yang 1995), the connection between characters and pronunciation is largely obscured by historical sound change and character change (Huang & Liao 2017). In Japanese, most Chinese characters are used to represent both borrowed

---

[4]In Braver and Kawahara's experiment, 10 out of 13 sets of stimuli were presented only in Chinese characters, while the other three sets included some kana (two in katakana and one in hiragana). In current usage, kana refers to two syllabaries where each character represents a mora, which in Japanese phonology may be an onset–vowel combination, a coda, or the second half of a long vowel. It is also worth noting that geminates and some long vowels are marked with diacritics in kana. Braver and Kawahara observed incomplete phonetic neutralisation for each set of stimuli.

words from China (the Sino-Japanese lexicon) and words that are originated in Japan (the Yamato lexicon) (Japan Broadcasting Corporation 1998; Itô & Mester 1999), and the resulting multiple pronunciations (*onyomi* and *kunyomi*) of many Chinese characters can only further weaken the connection between Chinese characters and pronunciations. So it is hard to imagine that Japanese speakers hypercorrected based on Chinese characters, and Braver and Kawahara still appeared to observe incomplete neutralisation in monomoraic prosodic word lengthening process.

To sum up, although the case of Catalan is controversial, the case of Japanese provides good evidence that at least in some languages, the observed incomplete neutralisation is not caused by orthographic knowledge.

Another source of criticism of incomplete neutralisation is that the observed effect size is typically quite small. Small effect sizes have been argued to likely not be functionally significant and therefore not to need a grammatical explanation (Dinnsen & Charles-Luce 1984; Mascaró 1987; Warner *et al.* 2004).[5] For example, among the phonetic cues examined by Port & O'Dell (1985), preceding vowel duration before underlying voiced stops was only about 15 ms longer than that before underlying voiceless stops, voicing into closure of derived voiceless stops was only 5 ms longer than that of underlying voiceless stops and duration of aspiration noise before underlying voiceless stops was only 15 ms longer than that of derived voiceless stops. Similar effect sizes were also found in Polish (Slowiaczek & Dinnsen 1985; Jassem & Richter 1989), Dutch (Warner *et al.* 2004) and two other studies on German (Piroth & Janker 2004; Roettger *et al.* 2014). To summarise the discussion on the criticisms on incomplete neutralisation, the debate on the existence of incomplete neutralisation is still pretty much ongoing, especially with respect to the issue of effect size.

In this article, as mentioned above, we will argue using data from Huai'an that incomplete phonetic neutralisation can stem from phonologically complete neutralisation. By using Huai'an, we avoid the orthographic confound discussed above, as the stimuli can be presented in Chinese characters, an orthographic system that has only a weak connection with pronunciation (this is similar to the Japanese case discussed above). Furthermore, the language allows us to argue that effect sizes are tangential to the issue of phonological neutralisation. Anticipating our results, we show that although there is a rather large phonetic difference with respect to incomplete phonetic neutralisation, there is clear evidence that the relevant processes are phonologically categorically neutralising, as evidenced by the fact that their outputs feed other sandhi processes.

## 2 The issue of phonological neutralisation versus phonetic implementation

As introduced in §1, the definition of incomplete neutralisation is twofold, being a combination of phonological neutralisation and phonetically incomplete neutralisation. An issue with many previous studies of incomplete neutralisation is that

---

[5]To be clear, we are not claiming that the effect size must be large for incomplete neutralisation that is caused by grammatical knowledge. We are only recognising here that it is a reasonable concern that incomplete neutralisation with a small effect size may not even be captured by grammatical knowledge and therefore may not be able to pose any challenges to traditional formal phonology.

researchers do not typically show evidence that the examined processes are truly phonological neutralisation, as opposed to phonetic implementation (Cohn 1993; Dunbar 2013). Under the categorical view of phonological representations, phonological neutralisation entails a change from one phonological category to another phonological category, while phonetic implementation does not result in any categorical changes. To give an example, it is assumed by Port & O'Dell (1985) and other previous studies on German that the devoicing process results in a voiceless obstruent category in the phonological surface form.[6] However, there is no clear evidence, especially evidence from phonological behaviour, that shows a derived voiceless obstruent is actually neutralised with the underlying voiceless obstruent in the phonology. If Dunbar's (2013) suspicion that word-final devoicing in German is actually a phonetic implementational process turns out to be valid, then the so-called 'devoiced obstruent' at the right edge of prosodic word remains phonologically unchanged and still belongs to the same 'voiced' category with voiced obstruents in other positions. As a result, it would not be surprising according to the Standard generative view of phonology that the so-called 'devoiced obstruent' is phonetically different from an underlying voiceless obstruent since they are phonologically different, that is, different in the surface representations.

To the best of our knowledge, the only careful previous study that attempted to establish phonological neutralisation using evidence from phonological behaviour is Braver and Kawahara's (2016) study on the lengthening of prosodic words in Japanese. In Japanese, since a prosodic word has been argued to be at least bimoraic (Itô 1990; Mester 1990; Poser 1990; Mori 2002; Itô & Mester 2003), an *underlying* monomoraic prosodic word has been argued to lengthen to be bimoraic. Braver & Kawahara (2016) showed that this neutralisation is incomplete phonetically, that is, a derived bimoraic prosodic word is still shorter than an underlying bimoraic prosodic word.

The current article utilises a different strategy of examining rules in feeding orders to establish phonological neutralisation. The fact that the output of a process can trigger another process provides evidence that the first process results in complete phonological neutralisation. Yet, despite the categorical neutralisation in the phonology, we will show that there is incomplete neutralisation in the phonetics for each of the feeding processes in Huai'an tone sandhi processes.[7] We will elaborate the feeding orders in Huai'an in §3 with more background information. And, §§4 and 5 will present the two experiments we have run based on two feeding orders in Huai'an.

---

[6]To be fair to them, this assumption is consistent with, and likely based on, what many phonologists have claimed in the prior literature.

[7]As pointed out by an anonymous reviewer, Ernestus *et al.* (2006) showed that Dutch has optional progressive voicing assimilation across word boundaries for obstruents. A word-final devoicing process feeds into the assimilation process and causes the initial obstruent of the next word to become voiceless, which suggests that the devoiced obstruent in the preceding word belongs to the same phonological category as its underlyingly voiceless counterpart. In relation to this, Ernestus *et al.* (2006) showed in a separate experiment that there is incomplete phonetic neutralisation in word-final devoicing process in Dutch (see also Warner *et al.* 2004). Further study is needed to show that a devoiced obstruent that can trigger another devoicing process is actually incompletely neutralised in the phonetics. If this is indeed the case, Dutch will serve as another clear case of incomplete phonetic neutralisation under the condition of complete phonological neutralisation.

## 3 Background

Huai'an belongs to the Jianghuai Guanhua Group (Lower Yangtze Mandarin) of the Mandarin language family. Native speakers are mainly from Huai'an city, which is located in the northern part of Jiangsu Province (Li 1989). Huai'an has four phonemic tones, labelled as Tone 1, Tone 2, Tone 3 and Tone 4 (Jiao 2004; Wang & Kang 2012). Following the tradition of tone description in Chinese languages, in Table 1, the four tones are given in tone letters using a scale of 1–5 where 1 is the lowest f0 and 5 is the highest f0 and followed by a contour description in words (Chao 1930).[8] The tonal contours of phonemic tones in isolation are given in Figure 1. The speaker (male, age: 53) pronounced four repetitions of four monosyllabic morphemes that share the same segmental content [sɔ] and only contrast in the tone on the vowel. f0 was extracted only from the vowel at 5% steps with a script in Praat (Boersma & Weenink 2021). However, it is worth noting that the tonal contours in isolation for Mandarin tones have been noticed to be quite different when compared with their counterparts in context (Shen 1990; Xu 1994, 1997; Jongman *et al.* 2006, *inter alia*). So, we expected the same kinds of differences in our experiments where tones are pronounced in sentences.

**Table 1.** *Description of phonemic tones in Huai'an.*

| Phonemic tone | Tone letter | Contour description |
|---|---|---|
| Tone 1 | 42 | high falling |
| Tone 2 | 24 | high rising |
| Tone 3 | 312 | low/low rising |
| Tone 4 | 55 | high level |



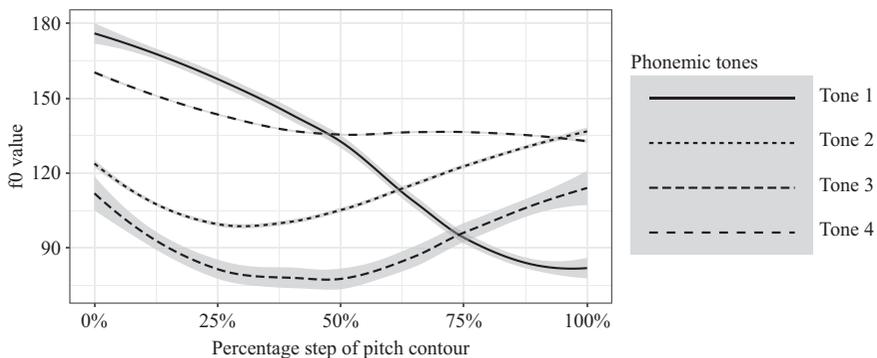**Figure 1.** *Tonal contour of phonemic tones in Huai'an.*

---

[8]Note that Huai'an phonemic tones are different from those in Standard Mandarin. In Huai'an, Tone 1 is a high falling tone and Tone 4 is a high level tone. While in Standard Mandarin, Tone 1 is a high level tone (tone letter: 55) and Tone 4 is a falling tone (tone letter: 51). Tone 2s and Tone 3s in Huai'an and Standard Mandarin are largely similar.

In subsequent examples, tones will be identified with just a T before the tone number, as in T3 for Tone 3; we will, however, continue to use full forms such as Tone 3 in the text.

The three tone sandhi rules relevant for this article are shown in (2). At the post-lexical level, the low-register Tone 3 sandhi is mandatory in some contexts and optional in others (we will elaborate at a later point in this article). In contrast, the high-register Tone 1 and Tone 4 sandhis are always optional. Furthermore, Tone 3 undergoes tone sandhi to become Tone 2, and this Tone 3 sandhi process can only happen when immediately preceding Tone 3 (underlying or derived).[9] As dissimilation processes, the tone sandhis in Huai'an can be straightforwardly explained by the Obligatory Contour Principle (Leben 1973; McCarthy 1986; Yip 2002, *inter alia*). However, some researchers reject the Obligatory Contour Principle as the motivation for tone sandhi processes in Mandarin languages (Duanmu 1994, 2007, *inter alia*). We will not address this debate since it is tangential to the main argument of this article.[10]

(2) *Relevant tone sandhi rules in Huai'an Mandarin* (Wang & Kang 2012)

    a. *Low-register tone sandhi*[11]

        i. Tone 3 sandhi: T3 + T3 → T2 + T3

    b. *High-register tone sandhi (optional processes)*

        i. Tone 1 sandhi: T1 + T1 → T3 + T1

        ii. Tone 4 sandhi: T4 + T4 → T3 + T4

Crucially, the Tone 3 output of the high-register tone sandhi processes feeds the low-register Tone 3 sandhi process as in (3). Since high-register tone sandhis are optional and Tone 3 sandhi is also optional for trisyllabic utterances in (3), multiple surface representations are possible for both examples.

(3) *Feeding Order in Huai'an Mandarin (boldface represents the locus of a potential tonal change due to the relevant tone sandhi process; the data are from the authors)*

    a. *Tone 1 sandhi feeds Tone 3 sandhi*

        u        ku       fən
        Mr. Wu estimate score

        'Mr. Wu estimates scores.'

---

[9]Consistent with our use of the terms 'underlying' and 'derived' throughout this article, 'underlying Tone 3' means surface Tone 3 that maps from underlying Tone 3, and 'derived Tone 3' means surface Tone 3 that maps from underlying Tone 1 or Tone 4.

[10]Depending on different proposed representations of the tones (Chen 2000; Duanmu 2007; Yip 2002, to name a few), the motivation of tone sandhis can be different. As we will discuss in §6.1, our assumption is that each Mandarin tone is represented as a single phonological unit.

[11]Register is a tonal feature first proposed by Yip (1980) and then widely adopted in the literature of Chinese tonal phonology. Here, we simply use the feature to distinguish Tone 3 sandhi from other tone sandhi processes in Huai'an.

| UR | T3 | T1 | T1 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Tone 1 Sandhi | T3 | **T3** | T1 | | | | | (or) | T3 | **T1** | T1 |
| Tone 3 Sandhi | **T2** | T3 | T1 | (or) | **T3** | T3 | T1 | | **T3** | T1 | T1 |
| SR | T2 | T3 | T1 | (or) | T3 | T3 | T1 | (or) | T3 | T1 | T1 |

b.  *Tone 4 sandhi feeds Tone 3 sandhi*

u        to        zəɯ
Mr. Wu chop meat

'Mr. Wu chops meat.'

| UR | T3 | T4 | T4 | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Tone 4 Sandhi | T3 | **T3** | T4 | | | | | (or) | T3 | **T4** | T4 |
| Tone 3 Sandhi | **T2** | T3 | T4 | (or) | **T3** | T3 | T4 | | **T3** | T4 | T4 |
| SR | T2 | T3 | T4 | (or) | T3 | T3 | T4 | (or) | T3 | T4 | T4 |

The feeding relationships between each of the high-register tone sandhis and Tone 3 sandhi suggest that the high-register tone sandhis result in a Tone 3 category that is phonologically the same as an underlying Tone 3. This interpretation of the data remains the same given a parallel approach to phonology such as Optimality Theory (Prince & Smolensky 1993). A usually employed markedness constraint for low-register tone sandhi in Mandarin languages is *33, which is based on the Obligatory Contour Principle. This constraint penalises adjacent Tone 3 syllables (Zhang 1997; Wang & Lin 2011; see also Chen 2000 for an implicit use of this constraint). For this constraint to trigger the structural change in the first syllable (namely, Tone 3 → Tone 2), the second syllable in an underlying /Tone 3 Tone 4 Tone 4/ or /Tone 3 Tone 1 Tone 1/ sequence must surface with Tone 3. Consequently, an Optimality Theory analysis would also maintain the crucial categorical aspects of the feeding order that are focal for the current article.

With regard to this interpretation of phonological identity between derived and underlying Tone 3s, concerns may be raised about application rates, especially when an underlying Tone 3 mandatorily triggers Tone 3 sandhi while a derived Tone 3 can only optionally trigger Tone 3 sandhi when the two types of Tone 3 syllables are the middle syllable of a trisyllabic utterance. The comparison is shown in (3) and (4). And, some researchers may want to ascribe the difference in application rates to a difference between derived Tone 3 and underlying Tone 3 in the phonology, either as different phonological representations or as the same representations indexed to different variable processes. However, intervening factors are not controlled for when compared in this way. For Tone 3 sandhi to mandatorily apply before an underlying Tone 3, the established planning window only needs to include the two Tone 3 syllables. Therefore, in (4), the established planning window only needs to include the first two syllables. In contrast, for Tone 3 sandhi to apply before a derived Tone 3, the established planning window needs to include at least three syllables to ensure both high-register Tone 1/Tone 4 sandhi and low-register Tone 3 sandhi occur. It is well recognised in the previous literature that a larger planning window has more planning

difficulty and is therefore less likely (Ferreira & Swets 2002; Wagner *et al.* 2010; Kilbourn-Ceron & Goldrick 2021, *inter alia*). The reason is the increasing burden on working memory, which can lead to speech errors or delays. Huai'an turns out to not be an exception. Previous experimental study on Tone 3 sandhi in Huai'an does support the existence of the effect of planning difficulty (Du & Lin 2021). Due to such an effect, a planning window that extends three syllables long is less likely to be established in (3), which means Tone 3 sandhi is less likely to apply before a derived Tone 3. Overall, the difference in application rates comes naturally from the planning difficulty effect and does not need to be accounted for in the phonology.[12]

As pointed out by two anonymous reviewers, proponents of gradient phonological representation may argue that although both underlying and derived Tone 3s can trigger Tone 3 sandhi, they may still have different phonological representations. By this analysis, the difference in application rates would be explained by the difference in the phonological representations. First, we would like to point out that any analysis that predicts application rates based on gradient phonological representations or phonetic similarity would have to be precise in accounting for not only cases where the process is triggered but also cases where the process is *not* triggered; namely, it would have to explain why only the derived Tone 3 shows a variation in application rates and not the underlying Tone 3, and not the other way around. Furthermore, it would have to account for the fact that any other tones that are phonetically similar (along the relevant dimensions) do not trigger the process. While an evaluation of such an analysis is not possible without a concrete specification of the proposal, we suspect that, to explain the difference between derived Tone 3 and underlying Tone 3, one will have to make reference to performance factors anyway. Relatedly, we appeal to the need to prioritise relatively simple categorical phonological representations when they are sufficient to account for the observed patterns (per Occam's razor/the law of parsimony); in our case, the difference in application rates can be accounted for by independently needed performance factors, namely planning, and therefore we need not complicate our understanding of the relevant phonological (tonal) representations. For this reason, we see the feeding rule interaction as evidence of complete phono-logical neutralisation of the derived Tone 3 from Tone 1 and Tone 4 sandhi processes. Furthermore, we use the processes to probe the phonetic (acoustic) consequences of the neutralising processes in the case of the derived Tone 3 that in turn trigger Tone 3 sandhi.

(4)    *Application of Tone 3 sandhi before underlying Tone 3 in trisyllabic utterances*

    a.  *Tone 1 sandhi feeds Tone 3 sandhi*

        u       pɔ    tɕi
        Mr. Wu protect car

        'Mr. Wu protects cars.'

---

[12]It has been noted in previous literature that tone sandhi patterns in many Mandarin languages are sensitive to prosodic structures (Zhang 1997; Chen 2000; Duanmu 2007; Wang & Lin 2011, *inter alia*). Therefore, the observed difference in application rates between derived and underlying Tone 3s can also be cashed out in terms of planning difficulty related to prosodic phrases.

| UR | | T3 | T3 | T1 |
|----|----|----|----|----|
| Tone 3 Sandhi | | **T2** | T3 | T1 |
| SR | | T2 | T3 | T1 |

b. *Tone 4 sandhi feeds Tone 3 sandhi*

u tɛ ɕiæ̃

Mr. Wu catch elephant

'Mr. Wu catches elephants.'

| UR | T3 | T3 | T4 |
|----|----|----|----|
| Tone 3 Sandhi | **T2** | T3 | T4 |
| SR | T2 | T3 | T4 |

To further ensure the phonological equivalence of derived Tone 3 and underlying Tone 3, we only analyse the derived Tone 3 tokens that actually trigger Tone 3 sandhi in this article, which allows us to have perfect surface minimal pairs in each of our experiments. By doing so, we also exclude the possibility that any identified incomplete phonetic neutralisation patterns arise as a result of averaging the outcomes of an optional phonological process, since we only look at the cases where we have reason to believe that the process has applied. Despite the categorical phonological behaviour of the derived Tone 3 in Huai'an, in the next two sections, we will show that there is substantial incomplete phonetic neutralisation of derived Tone 3 and underlying Tone 3 for the feeding orders involving Tone 1 sandhi and Tone 4 sandhi.

## 4 Experiment 1: Tone 1 sandhi

### 4.1 Participants

We recruited 11 native speakers of Huai'an Mandarin via personal relationships in Huai'an City. The age range was from 37 to 55 years. Among them, eight self-identified as female, and three as male. Due to the language standardisation trend in mainland China (Ramsey 1989), young speakers in Huai'an are generally bilingual and are native speakers of both Huai'an and Standard Mandarin. To minimise the influence of Standard Mandarin, we recruited older speakers who are only fluent in Huai'an. All the participants were born and raised in Huai'an City. None of them had participated in any linguistic studies before or heard about the concept of incomplete neutralisation.

### 4.2 Stimuli

The stimuli were composed of trisyllabic sentences with each syllable forming a separate word, to ensure that the tone sandhi processes observed are post-lexical and completely productive. Also, only right-branching utterances as in (3) are employed, simply because not enough left-branching utterances could be constructed that would have all the other characteristics required by the experimental design. The stimuli were divided into four sets as shown in (5). Furthermore, the third syllable was always

Tone 1. The second syllable was one of the following: (a) an underlying Tone 1 that optionally underwent Tone 1 sandhi to become Tone 3 or (b) an underlying Tone 3 that did not undergo any tone sandhi in this context. The first syllable was underlyingly Tone 3 or Tone 2. As a consequence of the possibilities in the second syllable, there were a few different possibilities for the first syllable, including (a) an underlying Tone 3 that could undergo Tone 3 sandhi to become Tone 2 with reference to the second syllable and (b) an underlying Tone 2 that did not undergo any tone sandhi in this context. The four sets differed only in tonal patterns and not in segmental content. Furthermore, the crucial second syllable was always a sequence of a voiceless unaspirated stop followed by a vowel. Voiceless unaspirated stops were chosen to make sure that there would be a consistent way to identify the acoustic onset of the vowel by referring to the burst of the stop. The full stimulus list is summarised in Appendix A. It is worth noting that one character, 搭, may be pronounced with the only checked tone in Huai'an (Jiao 2004; Wang & Kang 2012), which is an allophone of Tone 4 and appears only on monomoraic syllables ending with glottal stop. We excluded all checked tone productions when extracting f0 information.

(5)  *Four sets of stimuli in Experiment 1 (the syllables crucial for the current comparison are in boldface)*

    a.  underlying T3 following underlying T2:
       /T2 T3 T1/ → [T2 T3 T1]

    b.  underlying T3 following underlying T3:
       /T3 T3 T1/ → [T2 **T3** T1]

    c.  derived T3 following underlying T2:
       /T2 T1 T1/ → [T2 T3 T1] or [T2 T1 T1]

    d.  derived T3 following underlying T3:
       /T3 T1 T1/ → [T2 **T3** T1] or [T3 T1 T1]

Out of the above set of possibilities, the most crucial comparison is between two tones in the second syllable, namely, an underlying Tone 3 as in (5b) and a derived Tone 3 as in the first possibility in (5d). This particular comparison controls for the preceding surface context (derived Tone 2) and the following surface context (underlying Tone 1) and is therefore a perfect minimal pair. Furthermore, the two cases also show evidence that both tones are in fact categorically Tone 3, as they trigger Tone 3 sandhi on the preceding tone. Finally, as mentioned previously, the comparison allows us to exclude the possibility that any identified incomplete phonetic neutralisation pattern arises as a result of averaging the outcomes of an optional phonological process. This is the crucial pair we will focus on in this experiment.

The set of possibilities also allows us to visually compare the derived Tone 3 against an underlying Tone 1 in the same surface context, as in the second possibility in (5c) (although the preceding syllable in this case is an underlying Tone 2 instead of a derived Tone 2).

Each participant produced four repetitions of 24 test and 27 filler sentences at a natural speech rate, which means each participant read a total of 204 sentences. All stimuli were randomised for each participant.
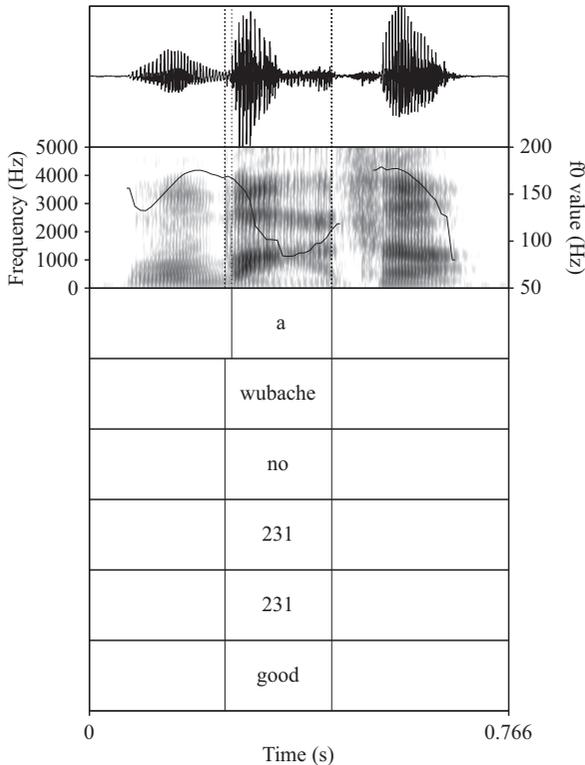
**Figure 2.** *Annotation scheme of Experiment 1 (Tone 1).*

### 4.3 Procedure

The experiment was conducted entirely in Huai'an city. Each participant was recorded by a trained research assistant using Audacity (Audacity Team 2019) and a Popu Line BK USB microphone on a Lenovo laptop in a quiet room that was located in the participant's home or workplace. The participants were told that the purpose of the study was to collect some general information on Huai'an. In post-experiment interviews, none of the participants reported noticing the minimal pairs, or that tones were the real focus of the study. The participants were instructed to read at a normal speech rate using their everyday voice, and the stimuli were presented in Chinese characters. The participants were also encouraged to read through the stimulus list to be familiar with the reading materials before producing them.

### 4.4 Measurement

Using Praat (Boersma & Weenink 2021), the recordings were manually annotated by the first author, who is a native speaker of Huai'an. An example is shown in Figure 2. Only the second syllable was marked, and the annotation file had six tiers in total. The first tier marked the vowel of the second syllable for phonetic analysis.

The first zero-crossing at the beginning of the voicing of the target vowel and after the burst of the unaspirated stop was identified as the vowel onset, and the zero-crossing immediately following the vowel's final glottal pulse was identified as the vowel offset. All other tiers marked the whole second syllable to index phonological information and recording quality. The onset of the second syllable was marked just before the release burst of the initial stop, and the offset of the second syllable corresponded with the offset of the nuclear vowel. The second tier indicated the whole sentence in pinyin, which is the official romanisation system for Chinese characters in China. The third tier was the tone sandhi condition where 'yes' meant the second syllable had undergone tone sandhi and 'no' meant it had not; the fourth and fifth tiers indicated the underlying tones and surface tones, respectively; and the last tier had the quality of the recording. We only used productions that were marked 'good' in the analysis. The reasons that productions were not marked as 'good' included background noise, speech errors, any long delay while producing the utterance, and checked tone pronunciation. f0 was extracted only from the vowel at 5% steps with a script in Praat.

To compare across different speakers and different vowels, *z*-score transformation was performed for each vowel of each speaker based on Hz scale (Laplace 1820; Lobanov 1971).

### 4.5  *Results and statistical modelling*

All data analyses in this article were performed in R (R Core Team 2021) using the **tidyverse** suite of packages (Wickham *et al.* 2019). The statistical modelling was done using the **lme4** package (Bates *et al.* 2021).[13]

The number of tokens for each possible combination of underlying representation and surface representation is summarised in Table 2. The application rate of Tone 3 sandhi before underlying Tone 3 is 97.2%, while the application rate before derived Tone 3 is 74.0%.[14] Seventy-one tokens were not marked as 'good' and excluded, which accounts for 6.7% of all test stimuli.

The *z*-score transformed f0 contours on the crucial second syllable are shown in Figure 3. As a reminder, the crucial comparison is between a derived Tone 3 and an underlying Tone 3 after derived Tone 2s in the same surface context – the context establishes that both the Tone 3s are categorically Tone 3 phonologically, as they trigger Tone 3 sandhi. We also present the tone contour for an underlying Tone 1 in the same surface context for visual comparison with the two crucial Tone 3s.
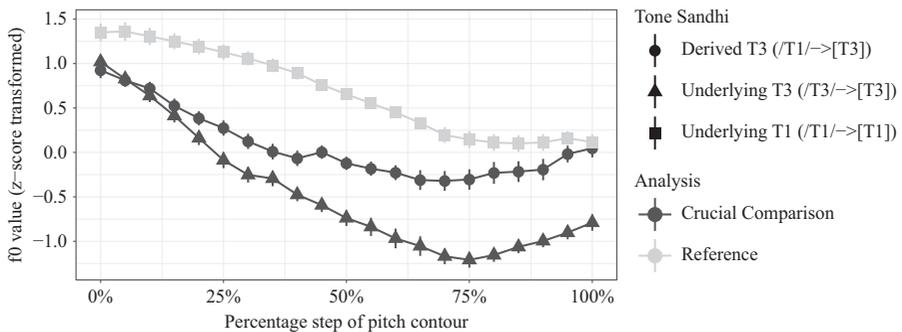
Based on visual inspection of the data, the derived Tone 3 seems to start like an underlying Tone 3 and end like an underlying Tone 1. And, the contour shape of the derived Tone 3 is close to that of an underlying Tone 3. Furthermore, the comparison

---

[13] All the data presented in this article are available in the form of csv files with the extracted measurements at the permanent link https://osf.io/9qcbr. The repository also includes the R script used to analyse and plot the data.

[14] As mentioned before, the difference in application rates does not necessarily inform us of any differences in phonological representations, since the difference is explicable by independently needed mechanisms, namely the difficulty in planning longer utterances.

***Table 2.*** *Number of tokens for each UR and SR combination in Experiment 1.*

| UR | SR | Number of tokens |
|---|---|---|
| T2T3T1 | T2T3T1 | 259 |
| T3T3T1 | T3T3T1 | 7 |
| | T2T3T1 | 242 |
| T2T1T1 | T2T1T1 | 74 |
| | T2T3T1 | 167 |
| T3T1T1 | T3T1T1 | 59 |
| | T3T3T1 | 46 |
| | T2T3T1 | 131 |



**Figure 3.** *Comparison of second-syllable contours in Experiment 1 (Tone 1; error bars indicate standard error).*

between underlying Tone 3 and derived Tone 3 clearly shows that the neutralisation is incomplete.[15]

For the purposes of statistical modelling, we used just the two-group factor (underlying Tone 3 vs. derived Tone 3), and ignored underlying Tone 1, in order to simplify the modelling and address only the crucial question of whether or not the underlying and derived Tone 3s have incompletely neutralised. The results turn out to support the observation that the neutralisation is indeed incomplete phonetically.

In dealing with time-course data, traditional techniques like *t*-tests and ANOVA have to divide continuous time into multiple time bins and therefore have to make

---

[15] To further address the concern that incomplete neutralisation patterns identified in this study may arise as a result of averaging the outcomes of an optional phonological process, the distributions of underlying Tone 1, derived Tone 3 and underlying Tone 3 are shown for each time step in the Supplementary Material. Crucially, the derived Tone 3 distribution is generally unimodal, and distinct from the other two distributions, across the time steps. Thus, there is no evidence of an averaging artifact over optional surface representations for the derived Tone 3 case.

multiple comparisons. This method has been argued by Mirman (2017) to be problematic for increasing the risk of 'false positives'. Since each time bin incurs the nominal 5% false positive rate implied by '$p < 0.05$', overall, the false positive rate with multiple time bins and multiple comparisons will be much higher than a single comparison.

To solve this problem, multiple analysis methods have been developed, including Smooth Spline Analysis of Variance (Wang 1998), generalised additive model (Hastie & Tibshirani 1990) and growth curve analysis (Mirman *et al.* 2008; Mirman 2017). In this article, we follow Chen *et al.* (2017) in modelling f0 contours using growth curve analysis. Growth curve analysis uses multilevel linear regression to avoid multiple comparisons, and has been argued to be a useful modelling technique in different fields (Baldwin & Hoffmann 2002; McArdle & Nesselroade 2003, *inter alia*). To apply growth curve analysis in Huai'an tones, we started with a simple model as in (6) (Mirman *et al.* 2008).

(6)     $Y_{ij} = (\gamma_{00} + \zeta_{0i}) + (\gamma_{10} + \zeta_{1i}) * \text{Time}_{ij} + \varepsilon_{ij}$

Here, $i$ is the $i$th f0 ($z$-score transformed) contour and $j$ is the $j$th time point; $Y_{ij}$ is the f0 ($z$-score transformed) value for the $i$th contour at the $j$th time point; $\gamma_{00}$ is the population average value for the intercept, $\zeta_{0i}$ is individual variation on the intercept, $\gamma_{10}$ is the population average value for the fixed effect of time, $\zeta_{1i}$ is individual variation on the fixed effect of time and $\varepsilon_{ij}$ is the error term.[16] To optimise the model for the data, we employed higher-order polynomial functions, and allowed individuals to vary on each term only when those terms reached significance according to chi-square likelihood ratio tests (Chen *et al.* 2017; Chen & Li 2021, *inter alia*). In Mandarin languages, a tone-bearing unit, which is assumed to be the syllable, the rhyme or the nucleus, has been widely argued to be associated with at most three tonal targets (Bao 1990, 1992; Duanmu 1994, *inter alia*). Therefore, the most complex tones can only have one change of direction, which will produce U-shaped contours, such as high-low-high and low-high-low. To conform to this observation, we only considered up to second-order functions to ensure that the final model is no more complex than a U-shaped contour. Also, orthogonal polynomials were used to make sure that the linear and quadratic terms were not correlated (Mirman 2017). After optimising the model by including all significant terms, we first treated underlying Tone 3 and derived Tone 3 as the same and modelled them as one single contour to get Model 1. Then we built models that treat them as different, namely, models that include a tone sandhi condition (underlying Tone 3 vs. derived Tone 3) to do model comparison. Based on Model 1, tone sandhi condition is first allowed to affect only intercept to get Model 2. Then tone sandhi condition is allowed to affect both intercept and linear term to get Model 3. Finally, tone sandhi condition is allowed to affect all fixed effects, which include intercept, linear term and quadratic term, and the outcome is Model 4.

---

[16]The individual variation terms, $\zeta_{0i}$ and $\zeta_{1i}$, are akin to the random intercept by participant and random slope of time by participant.
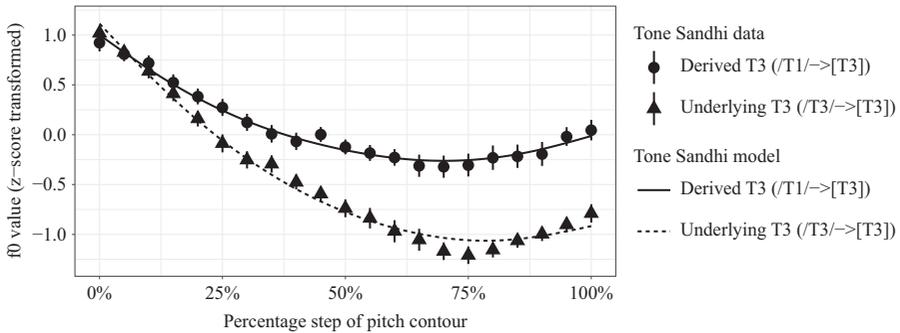
**Figure 4.** *Observed data and growth curve model fits for derived and underlying Tone 3 (error bars indicate standard error).*

**Table 3.** *Parameter estimates of the full model with the assumption of tone sandhi affecting every fixed effect (baseline: derived Tone 3).*

|  | Estimate | Std. error | $t$ | $p$ |
|---|---|---|---|---|
| Intercept | 0.07 | 0.06 | 1.13 | 0.28 |
| Linear | −16.01 | 2.14 | −7.47 | <0.01 |
| Quadratic | 9.57 | 1.59 | 6.04 | <0.01 |
| Tone Sandhi: Intercept | −0.50 | 0.03 | −19.68 | <0.01 |
| Tone Sandhi: Linear | −13.64 | 1.24 | −10.99 | <0.01 |
| Tone Sandhi: Quadratic | 4.80 | 1.23 | 3.90 | <0.01 |

A chi-square likelihood ratio test was used to determine whether two minimally different models differ significantly.

The result shows that the difference between underlying Tone 3 and derived Tone 3 is in fact supported by model comparisons. The addition of a tone sandhi condition improves the model on the intercept as shown by comparing Model 1 and Model 2 ($\chi^2(1) = 331.81$, $p < 0.01$), on the linear term as shown by comparing Model 2 and Model 3 ($\chi^2(1) = 118.34$, $p < 0.01$) and on the quadratic term as shown by comparing Model 3 and Model 4 ($\chi^2(1) = 14.99$, $p < 0.01$). Figure 4 shows how the full model (Model 4) fits the observed data. The parameter estimates for the full model are summarised in Table 3.

Moreover, the effect size of incomplete neutralisation is large in Tone 1 sandhi. The mean difference in f0 between underlying Tone 3 and derived Tone 3 across all steps is 18 Hz, which is more than two times the just-noticeable difference (JND) of f0 value (7 Hz) for Mandarin speakers (Jongman *et al.* 2017). Furthermore, across the last 10 steps (steps 11–20), the f0 difference is over 22 Hz, which is more than three times the JND. The f0 difference (f0 of derived Tone 3 minus f0 of underlying Tone 3) of each step is summarised in Table 4. Recall that the underlying premise of those who criticise the small effect size of incomplete neutralisation is that only if the differences were robust and large in size, the existence of such an effect should be accepted as

**Table 4.** *f0 difference of each step in Experiment 1 (Tone 1).*

| Step | f0 difference (Hz) | Step | f0 difference (Hz) |
|------|-------------------|------|-------------------|
| 0 | −5 | 11 | 22 |
| 1 | −2 | 12 | 25 |
| 2 | 2 | 13 | 26 |
| 3 | 4 | 14 | 29 |
| 4 | 7 | 15 | 30 |
| 5 | 12 | 16 | 32 |
| 6 | 12 | 17 | 30 |
| 7 | 11 | 18 | 28 |
| 8 | 15 | 19 | 29 |
| 9 | 22 | 20 | 25 |
| 10 | 21 | | |

functionally relevant.[17] According to that standard, Huai'an Tone 1 sandhi is clearly a case of phonetically incomplete neutralisation.

## 5 Experiment 2: Tone 4 sandhi

To show that the pattern is not unique to Tone 1 sandhi, and to extend the scope of the current study, we ran a second experiment on Tone 4 sandhi process in Huai'an.

### 5.1 Participants

We recruited 20 native speakers of Huai'an Mandarin, again via personal relationships in Huai'an City. The age range was from 33 to 57 years old. Again, to minimise the influence of Standard Mandarin, we avoided younger speakers in this study. Among them, 16 self-identified as female, and 4 as male. Five of these speakers had also participated in Experiment 1. The interval between the two experiments was about 7 months; the five participants from Experiment 1 failed to guess, and were not told, the purpose of Experiment 2. As in Experiment 1, all the participants were born and raised in Huai'an City. Other speakers had not participated in any linguistic studies before or heard about the concept of incomplete neutralisation.

### 5.2 Stimuli

The stimuli were organised in the same way as in Experiment 1. The four sets of trisyllabic sentences are shown in (7), and the full stimulus list is summarised in Appendix B.

---

[17]We acknowledge that it is not entirely clear to us what is intended by the use of phrases such as 'functional relevance', since many aspects of linguistic behaviour might be important to the speaker/listener while not stemming from the grammar *per se*; however, we retain the phrase here to reflect the terminology in the subfield.

(7) *Four sets of stimuli in Experiment 2 (the syllables crucial for the current comparison are in boldface)*

   a. underlying T3 following underlying T2:
   /T2 T3 T4/ → [T2 T3 T4]

   b. underlying T3 following underlying T3:
   /T3 T3 T4/ → [T2 **T3** T4]

   c. derived T3 following underlying T2:
   /T2 T4 T4/ → [T2 T3 T4] or [T2 T4 T4]

   d. derived T3 following underlying T3:
   /T3 T4 T4/ → [T2 **T3** T4] or [T3 T4 T4]

As with Experiment 1, the crucial comparison is between two tones in the second syllable, namely the underlying Tone 3 in (7b) and the derived Tone 3 as in the first possibility in (7d). This comparison allows us to control for the surface context, while also establishing that the two tones are indeed categorical Tone 3s, since they trigger Tone 3 sandhi on the preceding tone. Furthermore, as mentioned previously, the comparison allows us to exclude the possibility that any identified incomplete phonetic neutralisation pattern arises as a result of averaging the outcomes of an optional phonological process.

The set of possibilities also allows us to look at an underlying Tone 4 in roughly the same surface context, as in the second possibility in (7c), for visual comparison.

Each participant produced four repetitions of 20 test sentences at a natural speech rate with 20 fillers, which means that each participant read a total of 160 sentences.

## 5.3 Procedure

The procedure was identical to that of Experiment 1.

## 5.4 Measurement

The recordings were manually annotated by the first author but with a somewhat different scheme. For this experiment, both the first and second syllables were marked. The first syllable was marked to confirm that derived Tone 3 can in fact trigger Tone 3 sandhi on this syllable. An example is shown in Figure 5. The annotation file had five tiers in total. The criteria for marking vowels and syllables remained the same. The first tier marked the vowel of the syllable. All other tiers marked the whole second syllable to index phonological information and recording quality. The second tier indicated the position of the syllable inside the sentence, where a first syllable was marked '1' and a second syllable was marked '2'. The third tier contained the pinyin romanisation of the whole sentence followed by the underlying tone of the syllable. The fourth tier marked whether the syllable underwent tone sandhi. And, the last tier indicated the quality of the recording. Similar to the previous experiment, we only used productions from recordings that were marked 'good'. The f0 extraction, normalisation and visualisation processes are identical to those in the previous experiment.
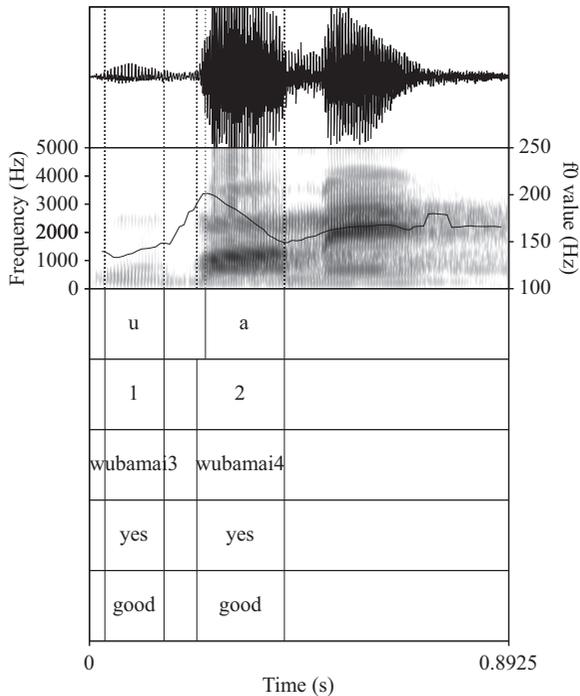
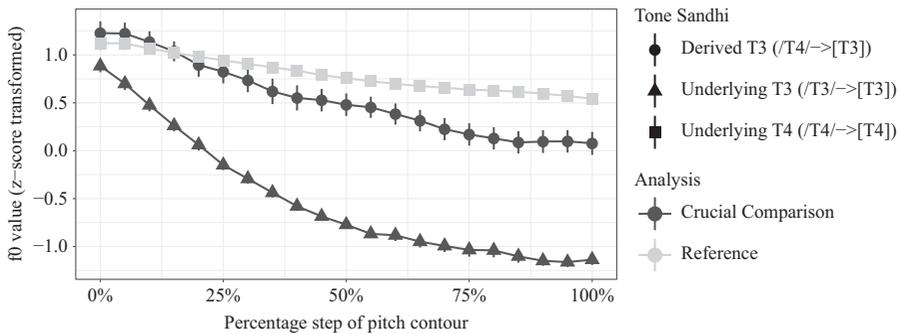***Figure 5.*** *Annotation scheme of Experiment 2 (Tone 4).*

### 5.5  Results and statistical modelling

The number of tokens for each possible combination of underlying representation and surface representation is summarised in Table 5. The application rate of Tone 3 sandhi before underlying Tone 3 is 94.8%, while the application rate before derived Tone 3 is 24.2%.[18] Seventy-nine tokens were not marked as 'good' and excluded, which accounts for 5.0% of all test stimuli.

---

[18] Again, we point out that the difference in application rates does not necessarily inform us of any differences in phonological representations. The difference can be accounted for by the difficulty in planning longer utterances. We also recognise that the application rate of Tone 3 derived from Tone 4 here is different from Tone 3 derived from Tone 1 in Experiment 1 (74.0%), and as one anonymous reviewer pointed out, such a difference could be used to argue for a phonological difference between the Tone 3s from different underlying sources. However, such a difference may simply be due to different groups of speakers in two independent experiments, or even more simply, the observed difference in effect sizes may simply be random variation, as would be expected between any two experiments measuring the same phenomenon. Future study is needed to compare Huai'an Tone 1 sandhi and Tone 4 sandhi on the same group of speakers in one single experiment. If the application rates are indeed replicable, an intriguing possibility that we note for the readers is that the phonetic difference between derived Tone 3 from Tone 1 and derived Tone 3 from Tone 4 may itself serve as a performance factor (related to the differential difficulty in implementing different tones, which in turn effects planning) that can account for the difference of application rate outside phonology. At the moment, the explanations based on planning difficulty are simply speculative, since it is not clear how phonological planning can affect tone sandhi application rate, and future study is needed to quantify the size of variation caused by planning difficulty.

**Table 5.** *Number of tokens for each UR and SR combination in Experiment 2.*

| UR | SR | Number of tokens |
|---|---|---|
| T2T3T4 | T2T3T4 | 386 |
| T3T3T4 | T3T3T4 | 20 |
|  | T2T3T4 | 368 |
| T2T4T4 | T2T4T4 | 156 |
|  | T2T3T4 | 212 |
| T3T4T4 | T3T4T4 | 98 |
|  | T3T3T4 | 213 |
|  | T2T3T4 | 68 |



**Figure 6.** *Comparison of second-syllable contours in Experiment 2 (Tone 4; error bars indicate standard error).*

The *z*-score transformed f0 contours on the crucial second syllable are shown in Figure 6. Again, the crucial comparison is between a derived Tone 3 and an underlying Tone 3 after a derived Tone 2 in the same surface context. We also present the tone contour for an underlying Tone 4 in the same surface context for visual comparison with the two crucial Tone 3s.

Based on visual inspection of the data, the pattern seems to be different from the case of Tone 1 sandhi. The derived Tone 3 seems to start like an underlying Tone 4,[19] instead of like an underlying Tone 3 as in Experiment 1. Furthermore, the derived Tone 3 gradually deviates from underlying Tone 4 through the whole contour; note that this is in contrast to Experiment 1, where the derived Tone 3 ended up at a value almost identical to the underlying Tone 1. However, the contour shape of the derived

---

[19] As observed by one anonymous reviewer, derived Tone 3 from Tone 4 actually starts higher than underlying Tone 4, although the difference in raw pitch between derived Tone 3 and underlying Tone 4 is very small (less than 5 Hz for the first four steps). As suggested by the reviewer, this pattern may be due to some effort to maintain contrast between derived Tone 3 and underlying Tone 4, since the remainder of derived Tone 3 is relatively high and close to underlying Tone 4.
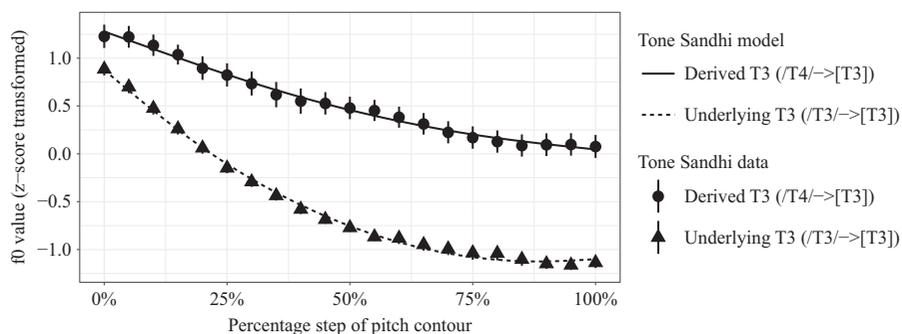
**Figure 7.** *Observed data and growth curve model fits for derived and underlying Tone 3 (error bars indicate standard error).*

Tone 3 is again close to that of an underlying Tone 3, as in Experiment 1. Despite the difference, incomplete phonetic neutralisation is again clearly observed in the comparison between underlying Tone 3 and derived Tone 3.[20]

The modelling method remains the same as in Experiment 1, and four models are generated. The observation of incomplete phonetic neutralisation is again supported by model comparisons. The addition of a tone sandhi condition improves the model on the intercept as shown by comparing Model 1 and Model 2 ($\chi^2(1) = 1,429.23$, $p < 0.01$), the linear term as shown by comparing Model 2 and Model 3 ($\chi^2(1) = 66.22$, $p < 0.01$) and the quadratic term as shown by comparing Model 3 and Model 4 ($\chi^2(1) = 32.67$, $p < 0.01$). Figure 7 shows how the full model with the assumption of tone sandhi affecting every fixed effect fits the observed data. And, the parameter estimates for full model are summarised in Table 6.

Again, the effect size of incomplete neutralisation is also large in Tone 4 sandhi. The mean difference in f0 between underlying Tone 3 and derived Tone 3 across all steps is 17 Hz, which is more than two times the JND of f0 value (7 Hz) for Mandarin speakers (Jongman *et al.* 2017). Also, across the last 11 steps (steps 9–20), the f0 difference is over 21 Hz, which is more than three times the just noticeable difference. The f0 difference (f0 of derived Tone 3 minus f0 of underlying Tone 3) of each step is summarised in Table 7. Therefore, the case of Huai'an Tone 4 sandhi can also be

---

[20]We recognise that the underlying Tone 3 in Experiment 2 appears to be slightly phonetically different from its counterpart in Experiment 1, especially at the tonal offset position. Tone 3, which is usually used to represent low tone in Mandarin languages, generally involves creakiness. The small difference between underlying Tone 3s in Experiments 1 and 2 may be caused by inconsistency related to the Praat f0 estimation algorithms for creaky sounds. Whether it is due to this issue or to simple random variation is something we leave for further inquiry.

Again, to further address the concern that incomplete neutralisation patterns identified in this study may arise as a result of averaging the outcomes of an optional phonological process, the distributions of underlying Tone 4, underlying Tone 3, and derived Tone 3 are shown for each time step in the Supplementary Material. Again, crucially, the derived Tone 3 distribution is generally unimodal, and distinct from the other two distributions, across the time steps. Thus, there is no evidence of an averaging artifact over optional surface representations for the derived Tone 3 case.

**Table 6.** *Parameter estimates of the full model with the assumption of tone sandhi affecting every fixed effect (baseline: derived Tone 3).*

|  | Estimate | Std. error | $t$ | $p$ |
|---|---|---|---|---|
| Intercept | 0.50 | 0.06 | 8.26 | <0.01 |
| Linear | −22.09 | 2.23 | −9.89 | <0.01 |
| Quadratic | 4.38 | 1.33 | 3.29 | <0.01 |
| Tone Sandhi: Intercept | −1.01 | 0.02 | −43.58 | <0.01 |
| Tone Sandhi: Linear | −10.43 | 1.26 | −8.30 | <0.01 |
| Tone Sandhi: Quadratic | 7.17 | 1.25 | 5.75 | <0.01 |

**Table 7.** *f0 Difference of each step in Experiment 2 (Tone 4).*

| Step | f0 difference (Hz) | Step | f0 difference (Hz) |
|---|---|---|---|
| 0 | −4 | 11 | 24 |
| 1 | −1 | 12 | 23 |
| 2 | 1 | 13 | 24 |
| 3 | 3 | 14 | 23 |
| 4 | 5 | 15 | 23 |
| 5 | 11 | 16 | 23 |
| 6 | 14 | 17 | 24 |
| 7 | 16 | 18 | 26 |
| 8 | 19 | 19 | 27 |
| 9 | 22 | 20 | 27 |
| 10 | 22 | | |

safely identified as phonetically incomplete neutralisation, and not susceptible to the criticism of a small effect size.

The coding in Experiment 2 also allowed us to answer another question that we did not answer for Experiment 1. In Experiment 1, we impressionistically coded whether or not the first syllable was in fact subject to Tone 3 sandhi. One could have argued that this impressionistic coding could have been inaccurate, and was based on a perceptual bias of the annotator (first author). To address this concern, it would have been optimal if we could have shown through phonological behaviour that the derived Tone 2 is indeed phonologically identical to underlying Tone 2. Although historically Tone 2 sandhi (Tone 2 + Tone 2 → Tone 3 + Tone 2) existed in Huai'an (Wang & Kang 2012), this tone sandhi rule was not observed in our fieldwork in early 2020, probably because of influence from the standard language, as is generally observed in other languages (Labov 1963; Milroy 2001, *inter alia*). And, no researchers before have tested if derived Tone 2 can trigger another tone sandhi process. Therefore, we cannot verify if the derived Tone 2 can trigger Tone 2 sandhi like an underlying Tone 2. Furthermore, we are not aware of any other phonological processes in the language
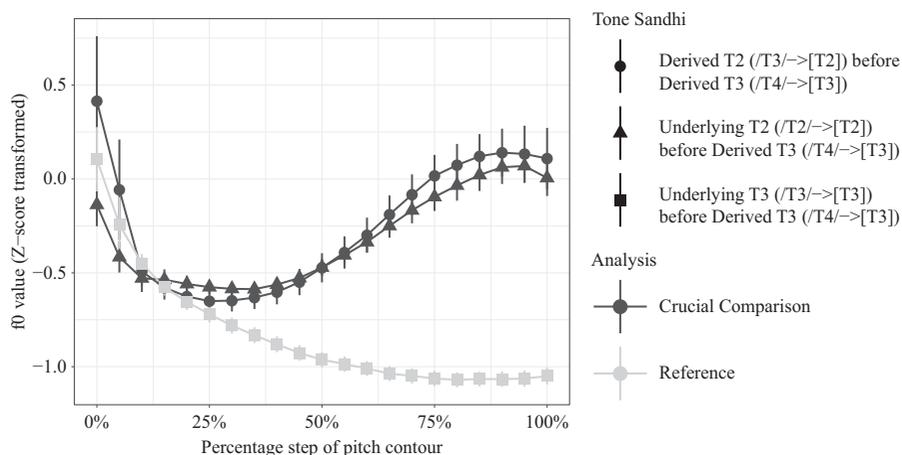
**Figure 8.** *Comparison of first-syllable contours in Experiment 2 (Tone 4; error bars indicate standard error).*

that are triggered by Tone 2. As a result, it is not possible to establish Tone 2 category by phonological behaviour in Huai'an and we turn to provide phonetic evidence for the Tone 2 identity of the derived rising tone.

To make some inroads into the question of the phonological nature of the (putatively) derived Tone 2 in initial position, in Experiment 2, we also annotated the first syllable, and are therefore able to observe the f0 contours for derived Tone 2 (from underlying Tone 3) and compare it to an underlying Tone 2 to see if the impressionistic coding was appropriate. The tone contours of the *z*-score transformed f0 for the relevant first syllables are shown in Figure 8. For comparison, we also present the tone contour for an underlying Tone 3 on the first syllable that comes from a derived Tone 3 failing to trigger Tone 3 sandhi on the preceding syllable. By doing so, a three-way visual comparison is possible at the position of the first syllable under the same phonological environment, that is, before derived Tone 3.

Based on the visual inspection of the data, the derived Tone 2 that undergoes Tone 3 sandhi with reference to the following derived Tone 3 is phonetically highly similar to an underlying Tone 2 with regard to the f0 contour. Both derived Tone 2 and underlying Tone 2 f0 contours are phonetically very different from underlying Tone 3. Furthermore, as with the other tone sandhi processes discussed in this article, there is incomplete phonetic neutralisation of the derived Tone 2 (from an underlying Tone 3) and the underlying Tone 2 in the first syllable. With the modelling method introduced in §4.5, the addition of a tone sandhi condition improves the model on the quadratic term as shown by comparing Model 3 and Model 4 ($\chi^2(1) = 4.96, p = 0.03$), but not on the intercept as shown by comparing Model 1 and Model 2 ($\chi^2(1) = 2.16, p = 0.14$) or the linear term as shown by comparing Model 2 and Model 3 ($\chi^2(1) = 1.10$, $p = 0.29$). Figure 9 shows how the full model (Model 4) with the assumption of tone sandhi affecting every fixed effect fits the observed data. And, the parameter estimates for the full model are summarised in Table 8.
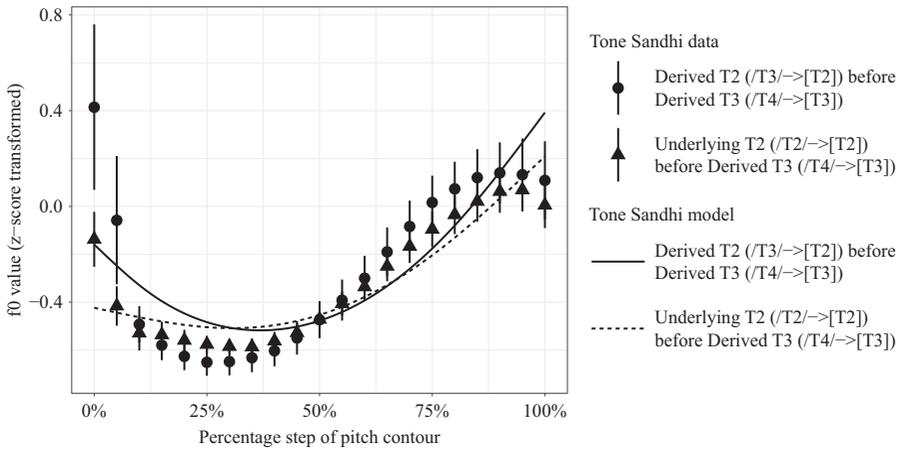
**Figure 9.** *Observed data and growth curve model fits for derived and underlying Tone 2 before derived Tone 3 (error bars indicate standard error).*

**Table 8.** *Parameter estimates of the full model with the assumption of tone sandhi affecting every fixed effect (baseline: derived Tone 2).*

|                         | Estimate | Std. error | $t$    | $p$    |
|-------------------------|---------:|-----------:|-------:|-------:|
| Intercept               | $-0.26$  | 0.06       | $-4.07$ | <0.01 |
| Linear                  | 10.00    | 2.69       | 3.72   | <0.01  |
| Quadratic               | 9.63     | 1.89       | 5.11   | <0.01  |
| Tone Sandhi: Intercept  | $-0.04$  | 0.03       | $-1.51$ | 0.13  |
| Tone Sandhi: Linear     | $-1.32$  | 1.30       | $-1.02$ | 0.31  |
| Tone Sandhi: Quadratic  | $-2.88$  | 1.29       | $-2.23$ | 0.02  |

However, consistent with our larger claim, this should not be interpreted as incomplete phonological neutralisation. The mean difference in f0 between underlying Tone 2 and derived Tone 2 across all steps is only 1 Hz, which is *much* lower than the JND of f0 value (7 Hz) for Mandarin speakers (Jongman *et al.* 2017). The f0 difference (f0 of underlying Tone 2 minus f0 of derived Tone 2) of each step is summarised in Table 9. This indicates that native speakers of Huai'an may not be able to distinguish underlying *versus* derived Tone 2s and therefore are likely to analyse them as belonging to the same phonological category. It is worth noting that an assumption has been made here that a phonetic difference that is *much* smaller than or around JND means phonologically complete neutralisation, and a phonetic difference that is *much* bigger than the JND is compatible with both phonologically complete neutralisation (as in Huai'an Tone 1 and Tone 4 sandhis) and phonologically incomplete neutralisation. We acknowledge that some previous studies on incomplete neutralisation have shown that phonetic differences that are smaller than the relevant JND are still perceptually distinguishable (Port & O'Dell 1985; Warner *et al.* 2004,

**Table 9.** *f0 Difference of each step for first syllable in Experiment 2 (Tone 2).*

| Step | f0 difference (Hz) | Step | f0 difference (Hz) |
|------|--------------------|------|--------------------|
| 0    | −8                 | 11   | 2                  |
| 1    | −8                 | 12   | 2                  |
| 2    | 0                  | 13   | 1                  |
| 3    | 2                  | 14   | 1                  |
| 4    | 2                  | 15   | 1                  |
| 5    | 3                  | 16   | 1                  |
| 6    | 3                  | 17   | 1                  |
| 7    | 3                  | 18   | 2                  |
| 8    | 3                  | 19   | 3                  |
| 9    | 2                  | 20   | 3                  |
| 10   | 2                  |      |                    |

*inter alia*). However, the substantial phonetic difference between derived Tone 2 and underlying Tone 3 and the phonetic similarity between derived Tone 2 and underlying Tone 2 are difficult to account for by any mechanism known to us other than Tone 3 sandhi – it cannot simply be random variation or a coarticulatory change. Therefore, the impressionistic coding was in our opinion appropriate.

To summarise the results of Experiment 2, we showed, using the feeding interaction between Tone 4 sandhi and Tone 3 sandhi, that the Tone 4 sandhi results in a phonological completely derived Tone 3. Despite this phonologically complete neutralisation, we observed a (rather large) incomplete neutralisation between the derived Tone 3 and underlying Tone 3 in the same surface tonal context. The experiment therefore replicates the results of Experiment 1.

## 6 Discussion

This article offers two clear cases of incomplete neutralisation based on data from Huai'an high-register tone sandhi processes. We observed robust phonetic differences (with large effect sizes) between a derived Tone 3 and an underlying Tone 3 in two independent experiments. This indicates that the observed effect is not likely to be a 'false positive' or functionally unimportant. Moreover, the Huai'an cases avoid any potential interference of orthography by presenting stimuli in Chinese characters. Therefore, some previous criticisms related to experimental design and the interpretation of data do not apply to the current Huai'an evidence.

A crucial aspect of the article is that we first established that the relevant tone sandhi processes are in fact phonological processes. To establish this fact, we look at the phonological behaviour of the derived tones, which to us is the best way of establishing phonological representations. More specifically, we looked at cases of tone sandhi that had feeding interactions, namely high-register tone sandhis including Tone 1 sandhi (Experiment 1) and Tone 4 sandhi (Experiment 2) feeding Tone 3 sandhi

in Huai'an Mandarin. This establishes the fact that the Tone 1 and Tone 4 sandhi processes are indeed cases of phonological neutralisation. Despite this, we observed incomplete phonetic neutralisation between underlying Tone 3 and derived Tone 3s stemming from the two tone sandhi processes. Consequently, our results establish the fact that phonologically complete neutralisation can still be phonetically incomplete.

## 6.1 *The phonological representation of Mandarin tone*

It is worth noting that the interpretation of Huai'an tone sandhi cases as incomplete neutralisation relies on the general consensus that a Mandarin tone is a single phonological unit even though it is realised phonetically as a tonal contour (Yip 1989; Bao 1990, 1992, *inter alia*). Under this view, it is not phonologically possible for part of a Mandarin contour tone to neutralise while another part of the tone remains unchanged. Perhaps the most convincing evidence for this single-phonological-unit representation in Mandarin languages comes from contour tone spreading. The most discussed case is undoubtedly Danyang (Chan 1991; Chen 1991; Yip 1989; data from Lü 1980). The pattern of interest is given in (8):

(8)   a.  2-syllable: hl. lh
      b.  3-syllable: hl. hl. lh
      c.  4-syllable: hl. hl. hl. lh
          where h stands for high tone, l stands for low tone, and. indicates a syllable boundary

 According to Yip's (1989) analysis, in these cases, a falling tone is associated with the first syllable, and a rising tone is associated with the last syllable. Then the falling tone spreads rightwards over the domain as one single unit. If the falling tone is not a unit in phonology, one would not expect the whole contour to spread, but only the low tone at its right edge. A similar phenomenon of tone spreading is also found in Changzhi (Hou 1983). It is worth noting that Duanmu (1994) challenges the above evidence by pointing out that contour tone spreading examples are only found in two languages and restricted to certain morphosyntactic structures. However, since Changzhi City and Danyang City are geographically far away from each other (roughly 734 km apart), tone spreading may be discovered in more languages and potentially more morphosyntactic structures. To summarise, despite dispute, the tone spreading pattern itself offers strong support for phonological contour tone. It is also worth pointing out that despite disagreement with the single-unit analysis of contour tones, Duanmu (1994) claims that tone sandhi results in a categorical change, which is argued in this article to support the interpretation of incomplete neutralisation in Huai'an.

 With the above phonological viewpoint of tonal representations as backdrop, in Huai'an, the fact that both derived Tone 3 and underlying Tone 3 can trigger Tone 3 sandhi suggests that a derived Tone 3 is phonologically identical to underlying tone 3. In fact, to the best of our knowledge, we are not aware of any Mandarin languages where only underlying Tone 3 triggers Tone 3 sandhi, and not derived Tone 3 – this correlation would be accounted for by phonological neutralisation. However, a

phonetic difference on any part of the contour between a derived contour tone and its underlying counterpart indicates phonetic incomplete neutralisation of the whole contour tone unit. In the case of Tone 1 and Tone 4 sandhis in Huai'an, there is a clear phonetic difference at the tonal offset position as shown in Experiment 1 and 2.

Based on the above, we would like to explicitly acknowledge that our claims in the article about incomplete phonetic neutralisation in the face of complete phonological neutralisation are contingent on the phonological representations we have assumed. As we see it, it cannot be any other way. The argument for incomplete neutralisation in any language depends on a certain set of assumed phonological representations. For example, in German, the interpretation of incomplete neutralisation depends on the devoicing rule actually resulting in a [−voice] feature (or equivalent). If the devoicing process results in some other phonological representation with similar phonetics, then the whole issue of incomplete neutralisation vanishes, and there is no need to entertain any more gradience in the phonological system to explain the observed phonetic patterns. In fact, a version of such a featural account is implied by Hale *et al.* (2007), who argue that language-specific phonetics can in fact be accounted for by different phonological feature combinations. Similarly, in Huai'an, it is possible to explain what is observed in the phonetics by changing or adding new phonological representations, but then of course independent evidence of the same representations in the language or in other related languages generally needs to be provided; otherwise it becomes an ad hoc, and therefore unjustified, claim. More generally, any set of representations or computations cannot simply be post hoc accounts of the data/patterns but need to be independently justified claims.

### 6.2  Desiderata for any explanation for incomplete neutralisation

With the two clear cases of incomplete neutralisation, the next step is naturally the explanation for incomplete neutralisation. Due to the limitation of the current study, the exact source of incomplete neutralisation cannot be pinpointed. However, we would like to lay out the desiderata that we think any explanation of incomplete neutralisation must achieve and illustrate the problems with previous explanations alongside.

(9)  Desiderata for a theory of incomplete neutralisation
   a.  The simplest explanation of why incomplete neutralisation exists as a phenomenon
   b.  An explanation for the actual distribution of effect sizes among different phonological processes and within a single phonological process
   c.  An explanation of why 'over-neutralisation' is never observed
   d.  An explanation of how a feeding interaction is possible where the derived representation still incompletely neutralises with the element that triggers the process
   e.  Related to (9d), an explanation of why incompletely neutralised segments can trigger the process, but other phonetically similar segments do not

First, to ensure the priority of a relatively simple theoretical model, explanations that can solve the problem while retaining a relatively simple phonological model

should be considered first (Occam's razor/the law of parsimony). Consequently, if independently needed performance mechanisms have the potential to account for the observation of incomplete phonetic neutralisation, they should be prioritised. Consistent with this principle, in the current study, the difference in Tone 3 sandhi application rates is assigned to independently needed performance factors of phonological planning, and therefore there is no need to complicate our understanding of the relevant phonological (tonal) representations. For the explanation of incomplete neutralisation, beyond previously identified factors such as orthography and task effects, the best performance factors in our opinion that need to be explored further are again phonological planning (Wagner 2012; Tanner *et al.* 2017; Kilbourn-Ceron & Goldrick 2021) and cascaded activation of morphemes during production (Goldrick & Blumstein 2006). If they are able to account for the patterns, we would be able to maintain a much simpler and, consequently, more predictive phonological theory.

The second challenge (9b) facing theories of incomplete neutralisation is the systematic disparity in effect sizes. Any proposed theory should explain among the observed cases why effect sizes of incomplete neutralisation are rather small in devoicing processes (as in German, Dutch, Russian, etc.), but can be quite large as in Huai'an tone sandhis or Japanese vowel lengthening. Moreover, the proposed explanation should also account for the newly found disparity in effect sizes within a single phonological process as in two Huai'an tone sandhis. In Experiment 1, the effect size is very small at the tonal onset position, as shown in Table 4, and becomes quite large as the contour progresses. A similar pattern is also found in Experiment 2, as shown in Table 7. A model that can simply account for a variety of effect sizes misses the systematic nature among different neutralisation processes and within a single time-varying neutralisation process.

The third challenge (9c) is that the proposed explanation should not only predict cases of incomplete neutralisation where the derived category is phonetically close to an underlying category (and in fact, between the phonetic manifestation of two underlying categories – its own UR and the phonological representation it is putatively changing to), but also avoid predicting cases of 'over-neutralisation' where the degree of application is beyond the phonetic distribution of the underlying category it is neutralising to. To return to the case of German devoicing, under the scenario of incomplete neutralisation, the phonetic cues of derived voiceless stops fall between underlying voiceless stops and underlying voiced stops. In the scenario of over-neutralisation, the phonetic cues of underlying voiceless stops would fall between derived voiceless stops and underlying voiced stops. However, only incomplete neutralisation has been observed in examined languages including Huai'an. This observation would be particularly problematic for purely exemplar representations (Brown & McNeill 1966; Bybee 1994; Goldinger 1996, 1997; Port & Leary 2005; Roettger *et al.* 2014, *inter alia*). Many previous theories account for the absence of over-neutralisation by proposing some mechanism whereby phonetically incomplete neutralisation is simply intermediate between two representations as it results from a blend of all phonetic cues of two distinct representations (Gafos & Benus 2006; van

Oostendorp 2008; Smolensky *et al.* 2014; Braver 2019).[21] Either such theories are not specific enough, or other independently needed mechanisms must be incorporated to capture the systematic disparity in effect sizes in (9b).

The fourth challenge (9d) that any theory of incomplete neutralisation faces is to explain how a feeding interaction is possible when the derived representation still incompletely neutralises with the element that triggers the process. In the case of Huai'an, the Tone 3 output of the high-register tone sandhi processes can feed the low-register Tone 3 sandhi process as in (3) despite incompletely neutralising with underlying Tone 3 in the phonetics. Any categorical theory of phonological representations naturally accounts for this as process/rule interactions. Of course, it is possible for a theory of gradient phonological representations to do so too; however, to assess the effectiveness of such a theory, one needs to grapple with the specifics of the representations and computations proposed. To return to the Tone 3 sandhi application rate difference, if one were to propose that the differential application rates are a consequence of gradient phonological representations, where phonetic proximity triggers application of a process, then one has to address two things. First, why do we see the gradience in application rates with the derived category but not with the underlying category, though both vary in terms of phonetic manifestations? Second, we need to ensure that other phonetically similar sounds do not trigger the process too (9e). For example, in German, though both voiced obstruents and sonorants are phonetically voiced, only obstruents devoice at the end of a prosodic word. One may grant that the distinction between obstruents and sonorants is a difference in phonological representations; however, by making use of such distinction, a view of category is implicitly implemented.

We raise these challenges here to move the goalpost in a constructive direction on the debate about incomplete neutralisation. Given the above desiderata, we believe that previous explanations are not perfectly satisfying, and therefore the phenomenon of incomplete neutralisation remains an open problem.

## 7 Conclusion

The primary goal of this article is to offer two clear cases of incomplete neutralisation using data from Huai'an. Our results suggest that incomplete phonetic neutralisation can in fact have a large effect size, and more importantly that the phenomenon does not automatically reflect (gradient) phonological representations. Furthermore, echoing the general advice of Roettger *et al.* (2014), we would like to encourage more work on the topic and on our particular claim, since the acceptance of any phenomenon should not be based on a single study or a single language, and only by accumulating converging evidence from different methodologies can we be more certain of it.

Finally, the phenomenon of incomplete neutralisation highlights a discrepancy between the Standard generative view of phonology (Kenstowicz 1994; Pierrehumbert 2002), wherein the output of phonological computation (the surface phonological

---

[21]Incomplete neutralisation is typically argued to be a blend of the surface representation and the underlying representation, or of the surface representation and a base representation, or of two co-activated surface representations.

representation) uniquely feeds into a phonetics module, and the Classic generative view of phonology, where phonology is seen as *knowledge* (Chomsky 1965; Chomsky & Halle 1965, 1968, *inter alia*). Note that both views represent feed-forward models, where phonological computation feeds into phonetic manifestations, but phonetic manifestations cannot feed into phonological computation. However, per the latter view, linguistic performance is a multi-factorial problem, and linguistic knowledge (i.e. competence) is only one of the many factors involved (Chomsky 1964, 1965; Valian 1982; Schütze 1996; Warner *et al.* 2004, *inter alia*).[22]

Our results from Huai'an tone neutralisations are problematic for the Standard generative view of phonology – if phonetic manifestations depend solely on the output of phonology and nothing else, then it is of course the case that such a view cannot account for cases where phonological neutralisation can still result in distinctness in the phonetics. However, our results are not in conflict with the Classic generative view of phonology. Phonology, per this latter view, is conceived of as grammatical knowledge that is used by a speaker to map a string of lexical items in a specific syntactic structure to articulations, and the use of this knowledge is affected by multiple other performance factors. Consequently, gradience in performance, and more specifically differences in speech production between two identical surface phonological representations, are not surprising. That is, there is no tension between incomplete phonetic neutralisation and categorical phonological neutralisation for the Classic generative view of phonology; instead, the actual mystery as per this view has always been with any observed cases of *complete* phonetic neutralisation stemming from a process of phonological neutralisation.

---

[22] We are not aware of any explicit argumentation that has ever been put forward in support of the Standard generative view over the Classic generative view. So, we are at a loss as to precisely when and, more importantly, *why* this change in viewpoints occurred. Here, we simply note the discrepancy.

## A  Stimuli for Experiment 1 on Tone 1

| Sentence | IPA | Pinyin | Word-by-word gloss | | | Translation of the whole sentence | UR | SR |
|---|---|---|---|---|---|---|---|---|
| 吴把车 | u pa tɕi | wu ba che | 'Mr. Wu' | 'play' | 'car' | 'Mr. Wu plays with cars.' | T2T3T1 | T2T3T1 |
| 吴鼓分 | u ku fən | wu gu fen | 'Mr. Wu' | 'encourage' | 'points' | 'Mr. Wu tries to increase points.' | T2T3T1 | T2T3T1 |
| 吴打车 | u ta tɕi | wu da che | 'Mr. Wu' | 'call' | 'car' | 'Mr. Wu calls for a taxi.' | T2T3T1 | T2T3T1 |
| 吴把虾 | u pa xa | wu ba xia | 'Mr. Wu' | 'play' | 'shrimp' | 'Mr. Wu plays with shrimp.' | T2T3T1 | T2T3T1 |
| 吴摆虾 | u pɛ xa | wu bai xia | 'Mr. Wu' | 'place' | 'shrimp' | 'Mr. Wu places shrimp (in a plate).' | T2T3T1 | T2T3T1 |
| 吴保车 | u pɔ tɕi | wu bao che | 'Mr. Wu' | 'protect' | 'car' | 'Mr. Wu protects cars.' | T2T3T1 | T2T3T1 |
| 吴扒车 | u pa tɕi | wu ba che | 'Mr. Wu' | 'grasp' | 'car' | 'Mr. Wu catches cars.' | T2T1T1 | T2T1T1/T2T3T1 |
| 吴估分 | u ku fən | wu gu fen | 'Mr. Wu' | 'estimate' | 'scores' | 'Mr. Wu estimates scores.' | T2T1T1 | T2T1T1/T2T3T1 |
| 吴搭车 | u ta tɕi | wu da che | 'Mr. Wu' | 'take' | 'cars' | 'Mr. Wu gets a ride.' | T2T1T1 | T2T1T1/T2T3T1 |
| 吴扒虾 | u pa xa | wu ba xia | 'Mr. Wu' | 'smash' | 'shrimp' | 'Mr. Wu smashes shrimp (to eat).' | T2T1T1 | T2T1T1/T2T3T1 |
| 吴掰虾 | u pɛ xa | wu bai xia | 'Mr. Wu' | 'break off' | 'shrimp' | 'Mr. Wu breaks off shrimp (to eat).' | T2T1T1 | T2T1T1/T2T3T1 |
| 吴包车 | u pɔ tɕi | wu bao che | 'Mr. Wu' | 'rent' | 'car' | 'Mr. Wu rents cars.' | T2T1T1 | T2T1T1/T2T3T1 |
| 武把车 | u pa tɕi | wu ba che | 'Mr. Wu' | 'play' | 'car' | 'Mr. Wu plays with cars.' | T3T3T1 | T2T3T1 |
| 武鼓分 | u ku fən | wu gu fen | 'Mr. Wu' | 'encourage' | 'points' | 'Mr. Wu tries to increase points.' | T3T3T1 | T2T3T1 |
| 武打车 | u ta tɕi | wu da che | 'Mr. Wu' | 'call' | 'car' | 'Mr. Wu calls for a taxi.' | T3T3T1 | T2T3T1 |
| 武把虾 | u pa xa | wu ba xia | 'Mr. Wu' | 'play' | 'shrimp' | 'Mr. Wu plays with shrimp.' | T3T3T1 | T2T3T1 |
| 武摆虾 | u pɛ xa | wu bai xia | 'Mr. Wu' | 'place' | 'shrimp' | 'Mr. Wu places shrimp (in a plate).' | T3T3T1 | T2T3T1 |
| 武保车 | u pɔ tɕi | wu bao che | 'Mr. Wu' | 'protect' | 'car' | 'Mr. Wu protects cars.' | T3T3T1 | T2T3T1 |
| 武扒车 | u pa tɕi | wu ba che | 'Mr. Wu' | 'grasp' | 'car' | 'Mr. Wu catches cars.' | T3T1T1 | T3T1T1/T2T3T1 |
| 武估分 | u ku fən | wu gu fen | 'Mr. Wu' | 'estimate' | 'scores' | 'Mr. Wu estimates scores.' | T3T1T1 | T3T1T1/T2T3T1 |
| 武搭车 | u ta tɕi | wu da che | 'Mr. Wu' | 'take' | 'cars' | 'Mr. Wu gets a ride.' | T3T1T1 | T3T1T1/T2T3T1 |
| 武扒虾 | u pa xa | wu ba xia | 'Mr. Wu' | 'smash' | 'shrimp' | 'Mr. Wu smashes shrimp (to eat).' | T3T1T1 | T3T1T1/T2T3T1 |
| 武掰虾 | u pɛ xa | wu bai xia | 'Mr. Wu' | 'break off' | 'shrimp' | 'Mr. Wu breaks off shrimp (to eat).' | T3T1T1 | T3T1T1/T2T3T1 |
| 武包车 | u pɔ tɕi | wu bao che | 'Mr. Wu' | 'rent' | 'car' | 'Mr. Wu rents cars.' | T3T1T1 | T3T1T1/T2T3T1 |

## B Stimuli for Experiment 2 on Tone 4

| Sentence | IPA | Pinyin | Word-by-word gloss | Translation of the whole sentence | UR | SR |
|---|---|---|---|---|---|---|
| 吴保税 | u pɔ suei | wu bao shui | 'Mr. Wu' 'protect' 'tax' | 'Mr. Wu is under bond.' | T2T3T4 | T2T3T4 |
| 吴躲肉 | u to ʐəu | wu duo rou | 'Mr. Wu' 'avoid' 'meat' | 'Mr. Wu avoids eating meat.' | T2T3T4 | T2T3T4 |
| 吴把脉 | u pa mɛ | wu ba mai | 'Mr. Wu' 'touch' 'blood vessel' | 'Mr. Wu diagnoses by touching blood vessels.' | T2T3T4 | T2T3T4 |
| 吴逮象 | u tɛ ɕiã | wu dai xiang | 'Mr. Wu' 'catch' 'elephant' | 'Mr. Wu catches elephants.' | T2T3T4 | T2T3T4 |
| 吴补炮 | u pu pʰɔ | wu bu pao | 'Mr. Wu' 'replenish' 'cannons' | 'Mr. Wu replenishes the stock of cannons.' | T2T3T4 | T2T3T4 |
| 吴报税 | u pɔ suei | wu bao shui | 'Mr. Wu' 'declare' 'tax' | 'Mr. Wu does taxes' | T2T4T4 | T2T4T4/T2T3T4 |
| 吴剁肉 | u to ʐəu | wu duo rou | 'Mr. Wu' 'chop' 'meat' | 'Mr. Wu chops meat.' | T2T4T4 | T2T4T4/T2T3T4 |
| 吴罢卖 | u pa mɛ | wu ba mai | 'Mr. Wu' 'stops' 'sell' | 'Mr. Wu stops selling (to protest).' | T2T4T4 | T2T4T4/T2T3T4 |
| 吴带象 | u tɛ ɕiã | wu dai xiang | 'Mr. Wu' 'take along' 'elephant' | 'Mr. Wu takes along elephants.' | T2T4T4 | T2T4T4/T2T3T4 |
| 吴布炮 | u pu pʰɔ | wu bu pao | 'Mr. Wu' 'deploy' 'cannons' | 'Mr. Wu deploys cannons.' | T2T4T4 | T2T4T4/T2T3T4 |
| 武保税 | u pɔ suei | wu bao shui | 'Mr. Wu' 'protect' 'tax' | 'Mr. Wu is under bond.' | T3T3T4 | T2T3T4 |
| 武躲肉 | u to ʐəu | wu duo rou | 'Mr. Wu' 'avoid' 'meat' | 'Mr. Wu avoids eating meat.' | T3T3T4 | T2T3T4 |
| 武把脉 | u pa mɛ | wu ba mai | 'Mr. Wu' 'touch' 'blood vessel' | 'Mr. Wu diagnoses by touching blood vessels.' | T3T3T4 | T2T3T4 |
| 武逮象 | u tɛ ɕiã | wu dai xiang | 'Mr. Wu' 'catch' 'elephant' | 'Mr. Wu catches elephants.' | T3T3T4 | T2T3T4 |
| 武补炮 | u pu pʰɔ | wu bu pao | 'Mr. Wu' 'replenish' 'cannons' | 'Mr. Wu replenishes the stock of cannons.' | T3T3T4 | T2T3T4 |
| 武报税 | u pɔ suei | wu bao shui | 'Mr. Wu' 'declare' 'tax' | 'Mr. Wu does taxes' | T3T4T4 | T3T4T4/T2T3T4 |
| 武剁肉 | u to ʐəu | wu duo rou | 'Mr. Wu' 'chop' 'meat' | 'Mr. Wu chops meat.' | T3T4T4 | T3T4T4/T2T3T4 |
| 武罢卖 | u pa mɛ | wu ba mai | 'Mr. Wu' 'stops' 'sell' | 'Mr. Wu stops selling (to protest).' | T3T4T4 | T3T4T4/T2T3T4 |
| 武带象 | u tɛ ɕiã | wu dai xiang | 'Mr. Wu' 'take along' 'elephant' | 'Mr. Wu takes along elephants.' | T3T4T4 | T3T4T4/T2T3T4 |
| 武布炮 | u pu pʰɔ | wu bu pao | 'Mr. Wu' 'deploy' 'cannons' | 'Mr. Wu deploys cannons.' | T3T4T4 | T3T4T4/T2T3T4 |

# References

Audacity Team (2019). *Audacity.* Version 2.3.2, retrieved 3 October 2019 from https://www.audacityteam.org/.

Badia Margarit, Antonio M. (1962). *Gramática catalana*, volume 1. Madrid: Editorial Gredos.

Baldwin, Scott A. & John P. Hoffmann (2002). The dynamics of self-esteem: a growth-curve analysis. *Journal of Youth and Adolescence* **31**. 101–113.

Bao, Zhiming (1990). *On the nature of tone*. PhD dissertation, Massachusetts Institute of Technology.

Bao, Zhiming (1992). Toward a typology of tone sandhi. *BLS* **18**. 1–12.

Bates, Douglas, Martin Mächler, Ben Bolker & Steven Walker (2021). Lme4: linear mixed-effects models using 'Eigen' and S4. https://CRAN.R-project.org/package=lme4.

Boersma, Paul & David Weenink (2021). *Praat: doing phonetics by computer.* Version 6.1.41, retrieved 25 March 2021 from http://www.praat.org/.

Braver, Aaron (2019). Modelling incomplete neutralisation with weighted phonetic constraints. *Phonology* **36**. 1–36.

Braver, Aaron & Shigeto Kawahara (2016). Incomplete neutralization in Japanese monomoraic lengthening. In Adam Albright & Michelle A. Fullwood (eds.) *Supplemental proceedings of the 2014 Annual Meeting on Phonology*. Washington, DC: Linguistic Society of America. 12 pp.

Brown, Roger & David McNeill (1966). The 'tip of the tongue' phenomenon. *Journal of Verbal Learning and Verbal Behavior* **5**. 325–337.

Bybee, Joan L. (1994). A view of phonology from a cognitive and functional perspective. *Cognitive Linguistics* **5**. 285–305.

Chan, Marjorie K. (1991). Contour-tone spreading and tone sandhi in Danyang Chinese. *Phonology* **8**. 237–259.

Chao, Yuen-Ren (1930). A system of tone-letters. *Le maître phonétique* **30**. 24–27.

Charles-Luce, Jan & Daniel A. Dinnsen (1987). A reanalysis of Catalan devoicing. *JPh* **15**. 187–190.

Chen, Matthew Y. (1991). An overview of tone sandhi phenomena across Chinese dialects. *Journal of Chinese Linguistics Monograph Series* **3**. 111–156.

Chen, Matthew Y. (2000). *Tone sandhi: patterns across Chinese dialects*. Cambridge: Cambridge University Press.

Chen, Si & Bin Li (2021). Statistical modeling of application completeness of two tone sandhi rules. *Journal of Chinese Linguistics* **49**. 106–141.

Chen, Si, Caicai Zhang, Adam G. McCollum & Ratree Wayland (2017). Statistical modelling of phonetic and phonologised perturbation effects in tonal and non-tonal languages. *Speech Communication* **88**. 17–38.

Chomsky, Noam (1964). The development of grammar in child language: formal discussion. *Monographs of the Society for Research in Child Development* **29**. 35–39.

Chomsky, Noam (1965). *Aspects of the theory of syntax*. Cambridge, MA: MIT Press.

Chomsky, Noam & Morris Halle (1965). Some controversial questions in phonological theory. *JL* **1**. 97–138.

Chomsky, Noam & Morris Halle (1968). *The sound pattern of English*. New York: Harper and Row.

Cohn, Abigail C. (1993). Nasalisation in English: phonology or phonetics. *Phonology* 10. 43–81.

Dinnsen, Daniel A. & Jan Charles-Luce (1984). Phonological neutralization, phonetic implementation and individual differences. *JPh* 12. 49–60.

Dmitrieva, Olga (2005). *Incomplete neutralization in Russian final devoicing: acoustic evidence from native speakers and second language learners*. PhD dissertation, University of Kansas.

Du, Naiyan & Yen-Hwei Lin (2021). Post-lexical tone 3 sandhi domain-building in Huai'an Mandarin: multiple domain types and free application. *University of Pennsylvania Working Papers in Linguistics* 27. 6.

Duanmu, San (1994). Against contour tone units. *LI* 25. 555–608.

Duanmu, San (2007). *The phonology of Standard Chinese*. Oxford: Oxford University Press.

Dunbar, Ewan (2013). *Statistical knowledge and learning in phonology*. PhD dissertation, University of Maryland, College Park.

Ernestus, Mirjam & R. Harald Baayen (2006). The functionality of incomplete neutralization in Dutch: the case of past tense formation. In Louis M. Goldstein, Douglas H. Whalen & Catherine T. Best (eds.) *Laboratory Phonology 8*. Berlin, NY: Mouton de Gruyter. 27–49.

Ernestus, Mirjam, Mybeth Lahey, Femke Verhees & R. Harald Baayen (2006). Lexical frequency and voice assimilation. *JASA* 120. 1040–1051.

Ferreira, Fernanda & Benjamin Swets (2002). How incremental is language production? Evidence from the production of utterances requiring the computation of arithmetic sums. *Journal of Memory and Language* 46. 57–84.

Fourakis, Marios & Gregory K. Iverson (1984). On the 'incomplete neutralization' of German final obstruents. *Phonetica* 41. 140–149.

Gafos, Adamantios I. & Stefan Benus (2006). Dynamics of phonological cognition. *Cognitive Science* 30. 905–943.

Goldinger, Stephen D. (1996). Words and voices: episodic traces in spoken word identification and recognition in memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition* 22. 1166–1183.

Goldinger, Stephen D. (1997). Words and voices: perception and production in an episodic lexicon. In Keith Johnson & John W. Mullennix (eds.) *Talker variability in speech processing*. San Diego, CA: Academic Press. 33–65.

Goldrick, Matthew & Sheila E. Blumstein (2006). Cascading activation from phonological planning to articulatory processes: evidence from tongue twisters. *Language and Cognitive Processes* 21. 649–683.

Hale, Mark, Madelyn Kissock & Charles Reiss (2007). Microvariation, variation, and the features of universal grammar. *Lingua* 117. 645–665.

Hastie, Trevor & Robert Tibshirani (1990). *Generalized additive models*. London: Chapman and Hall.

Hou, Jingyi (1983). Changzhi fangyan jilüe [Notes on the Changzhi dialect]. *Fangyan [Dialect]* 1983. 260–274.

Huang, Borong & Xudong Liao (2017). *Xiandai hanyu [Contemporary Chinese]*. Beijing: Higher Education Press.

Itô, Junko (1990). Prosodic minimality in Japanese. *CLS* 26. 213–239.

Itô, Junko & Armin Mester (1999). The phonological lexicon. In Natsuko Tsujimura (ed.) *The handbook of Japanese linguistics*. Maiden, MA: Blackwell. 62–100.

Itô, Junko & Armin Mester (2003). Weak layering and word binarity. In Takeru Honma, Masao Okazaki, Toshiyuki Tabata & Shin'ichi Tanaka (eds.) *A new century of phonology and phonological theory: a Festschrift for Professor Shosuke Haraguchi on the occasion of his sixtieth birthday*. Tokyo: Kaitakusha. 26–65.

Japan Broadcasting Corporation (1998). *The Japanese Language Pronunciation and Accent Dictionary*. Tokyo: NHK.

Jassem, Wiktor & Lutosława Richter (1989). Neutralization of voicing in Polish obstruents. *JPh* 17. 317–325.

Jiao, Lidong (2004). *Huai'an fangyan de shengdiao fenxi [An analysis of tones in Huai'an dialect]*. Master's thesis, Tianjin Normal University.

Jongman, Allard, Zhen Qin, Jie Zhang & Joan A. Sereno (2017). Just noticeable differences for pitch direction, height, and slope for Mandarin and English listeners. *JASA* 142. EL163–EL169.

Jongman, Allard, Yue Wang, Corinne B. Moore & Joan A. Sereno (2006). Perception and production of Mandarin Chinese tones. In Ping Li, Lihai Tan, Elizabeth Bates & Ovid J. L. Tzeng (eds.) *Handbook of East Asian psycholinguistics*. Cambridge: Cambridge University Press. 209–217.

Kenstowicz, Michael J. (1994). *Phonology in generative grammar*. Cambridge, MA: Blackwell.

Kharlamov, Viktor (2012). *Incomplete neutralization and task effects in experimentally-elicited speech: evidence from the production and perception of word-final devoicing in Russian*. PhD dissertation, Université d'Ottawa/University of Ottawa.

Kilbourn-Ceron, Oriana & Matthew Goldrick (2021). Variable pronunciations reveal dynamic intra-speaker variation in speech planning. *Psychonomic Bulletin & Review* **28**. 1365–1380.

Labov, William (1963). The social motivation of a sound change. *Word* **19**. 273–309.

Laplace, Pierre S. (1820). *Théorie analytique des probabilités*. Paris: Courcier.

Leben, William R. (1973). *Suprasegmental phonology*. PhD dissertation, Massachusetts Institute of Technology.

Li, Rong (1989). Hanyu fangyan de fenqu [The geographic division of Chinese dialects]. *Fangyan [Dialect]* **4**. 19.

Lobanov, Boris M. (1971). Classification of Russian vowels spoken by different speakers. *JASA* **49**. 606–608.

Lü, Shuxiang (1980). Danyang fangyan de shengdiao xitong [The tonal system of the Danyang dialect]. *Fangyan [Dialect]* **1980**. 85–122.

Manaster Ramer, Alexis (1996). A letter from an incompletely neutral phonologist. *JPh* **24**. 477–489.

Mascaró, Joan (1987). Underlying voicing recoverability of finally devoiced obstruents in Catalan. *JPh* **15**. 183–186.

Matsui, Mayuki (2015). *Roshiago ni okeru yuuseisei no tairitsu to tairitsu no jakka: onkyo to chikaku [Voicing contrast and contrast reduction in Russian: acoustics and perception]*. PhD dissertation, Hiroshima University.

McArdle, John J. & John R. Nesselroade (2003). Growth curve analysis in contemporary psychological research. In Wayne F. Velicer & John A. Schinka (eds.) *Handbook of psychology: research methods in psychology*, volume 2. New York: Wiley. 447–480.

McCarthy, John J. (1986). OCP effects: gemination and antigemination. *LI* **17**. 207–263.

McCollum, Adam (2019). Gradient morphophonology: evidence from Uyghur vowel harmony. In Hyunah Baek, Chikako Takahashi & Alex Hong-Lun Yeung (eds.) *Proceedings of the 2019 Annual Meetings on Phonology*, volume 7. Washington, DC: Linguistic Society of America. 12 pp.

Mester, Armin (1990). Patterns of truncation. *LI* **23**. 478–485.

Milroy, James (2001). Language ideologies and the consequences of standardization. *Journal of Sociolinguistics* **5**. 530–555.

Mirman, Daniel (2017). *Growth curve analysis and visualization using R*. Boca Raton, FL: Chapman and Hall/CRC.

Mirman, Daniel, James A. Dixon & James S. Magnuson (2008). Statistical and computational models of the visual world paradigm: growth curves and individual differences. *Journal of Memory and Language* **59**. 475–494.

Mori, Yoko (2002). Lengthening of Japanese monomoraic nouns. *JPh* **30**. 689–708.

van Oostendorp, Marc (2008). Incomplete devoicing in formal phonology. *Lingua* **118**. 1362–1374.

Pierrehumbert, Janet B. (2002). Word-specific phonetics. In Carlos Gussenhoven & Natasha Warner (eds.) *Laboratory Phonology 7*. Berlin, NY: Mouton de Gruyter. 101–139.

Piroth, Hans G. & Peter M. Janker (2004). Speaker-dependent differences in voicing and devoicing of German obstruents. *JPh* **32**. 81–109.

Port, Robert F. & Adam P. Leary (2005). Against formal phonology. *Lg* **81**. 927–964.

Port, Robert F. & Michael L. O'Dell (1985). Neutralization of syllable-final voicing in German. *JPh* **13**. 455–471.

Poser, William J. (1990). Evidence for foot structure in Japanese. *Lg* **66**. 78–105.

Prince, Alan & Paul Smolensky (1993). Optimality Theory: constraint interaction in generative grammar. Technical Report 2, Rutgers University Center for Cognitive Science.

R Core Team (2021). *R: a language and environment for statistical computing*. Version 1.4.1106, retrieved 25 March 2021 from https://www.rstudio.com.

Ramsey, S. Robert (1989). *The languages of China*. Princeton, NJ: Princeton University Press.

Roettger, Timo B., Bodo Winter, Sven Grawunder, James Kirby & Martine Grice (2014). Assessing incomplete neutralization of final devoicing in German. *JPh* **43**. 11–25.

Schütze, Carson T. (1996). *The empirical base of linguistics: grammaticality judgments and linguistic methodology*. Chicago, IL: University of Chicago Press.

Shen, Xiaonan S. (1990). Tonal coarticulation in Mandarin. *JPh* **18**. 281–295.

Slowiaczek, Louisa M. & Daniel A. Dinnsen (1985). On the neutralizing status of Polish word-final devoicing. *JPh* **13**. 325–341.

Slowiaczek, Louisa M. & Helena J. Szymanska (1989). Perception of word-final devoicing in Polish. *JPh* **17**. 205–212.

Smolensky, Paul, Matthew Goldrick & Donald Mathis (2014). Optimization and quantization in gradient symbol systems: a framework for integrating the continuous and the discrete in cognition. *Cognitive Science* **38**. 1002–1138.

Tanner, James, Morgan Sonderegger & Michael Wagner (2017). Production planning and coronal stop deletion in spontaneous speech. *Laboratory Phonology* **8**. 39 pp.

Valian, Virginia (1982). Psycholinguistic experiment and linguistic intuition. In Thomas W. Simon & Robert J. Scholes (eds.) *Language, mind, and brain*. Hillsdale, NJ: Lawrence Erlbaum. 179–188.

Wagner, Michael (2002). The role of prosody in laryngeal neutralization. *MIT Working Papers in Linguistics* **42**. 373–392.

Wagner, Michael (2012). Locality in phonology and production planning. *McGill Working Papers in Linguistics* **22**. 1–18.

Wagner, Valentin, Jörg D. Jescheniak & Herbert Schriefers (2010). On the flexibility of grammatical advance planning during sentence production: effects of cognitive load on multiple lexical access. *Journal of Experimental Psychology: Learning, Memory, and Cognition* **36**. 423–440.

Wang, Chiung-Yao & Yen-Hwei Lin (2011). Variation in Tone 3 Sandhi: the case of prepositions and pronouns. In Zhuo Jing-Schmidt (ed.) *Proceedings of the 23rd North American Conference on Chinese Linguistics*. Eugene, OR: University of Oregon. 138–155.

Wang, Yifeng & Jian Kang (2012). Huai'an nanpian fangyan liangzizu lianxu biaodiao fenxi [An analysis of disyllabic tone sandhi in southern Huai'an dialect]. *Journal of Shenyang Institute of Engineering (Social Sciences)* **8**. 357–359.

Wang, Yuedong (1998). Mixed effects smoothing spline analysis of variance. *Journal of the Royal Statistical Society Series B: Statistical Methodology* **60**. 159–174.

Warner, Natasha, Allard Jongman, Joan Sereno & Rachèl Kemps (2004). Incomplete neutralization and other sub-phonemic durational differences in production and perception: evidence from Dutch. *JPh* **32**. 251–276.

Whalen, Douglas H. (1991). Infrequent words are longer in duration than frequent words. *JASA* **90**. 2311–2311.

Whalen, Douglas H. (1992). Further results on the duration of infrequent and frequent words. *JASA* **91**. 2339–2340.

Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D'Agostino McGowan, Romain François, Garrett Grolemund, Alex Hayes, Lionel Henry, Jim Hester, Max Kuhn, Thomas Lin Pedersen, Evan Miller, Stephan Milton Bache, Kirill Müller, Jeroen Ooms, David Robinson, Dana Paige Seidel, Vitalie Spinu, Kohske Takahashi, Davis Vaughan, Claus Wilke, Kara Woo & Hiroaki Yutani (2019). Welcome to the Tidyverse. *Journal of Open Source Software* **4**. 1686.

Wright, Richard (2004). Factors of lexical competition in vowel articulation. In John Local, Richard Ogden & Rosalind Temple (eds.) *Papers in laboratory phonology VI*. Cambridge: Cambridge University Press. 75–87.

Xu, Yi (1994). Production and perception of coarticulated tones. *JASA* **95**. 2240–2253.

Xu, Yi (1997). Contextual tonal variations in Mandarin. *JPh* **25**. 61–83.

Yang, Jilin (1995). *Zhongguo zhongxiaoxue baikequanshu [encyclopedia for elementary school and middle school students]*. Harbin: Harbin Publishing House.

Yip, Moira J. (1980). *The tonal phonology of Chinese*. PhD dissertation, Massachusetts Institute of Technology.

Yip, Moira J. (1989). Contour tones. *Phonology* **6**. 149–174.

Yip, Moira J. (2002). *Tone*. Cambridge: Cambridge University Press.

Zhang, Ning (1997). The avoidance of the third tone sandhi in Mandarin Chinese. *Journal of East Asian Linguistics* **6**. 293–338.