



## Identity Avoidance in Turkish Partial Reduplication: Feature Specificity and Locality

**Kevin Tang\***, Department of English Language and Linguistics, Institute of English and American Studies, Heinrich-Heine-University, Düsseldorf, 40225, Germany; Department of Linguistics, University of Florida, Gainesville, Florida, 32611-5454, United States of America, [kevin.tang@hhu.de](mailto:kevin.tang@hhu.de)

**Faruk Akkuş**, Department of Linguistics, University of Massachusetts Amherst, Amherst MA 01003, United States of America, [fakkus@umass.edu](mailto:fakkus@umass.edu)

\*Corresponding author.

---

This study investigates the Turkish partial reduplication phenomenon, in which the reduplicant is derived by prefixing  $C_1VC_2$  syllable, where  $C_1V$  are identical to the word-initial CV of the base and the  $C_2$  ends in one of the four linking consonants:  $-p$ ,  $-m$ ,  $-s$ ,  $-r$ . This study re-examines the factors conditioning the choice of the linking consonant, by focusing the nature of the (dis)similarity (feature specificity) and the proximity (locality) between the consonants in the base and the linking consonant, using an acceptability rating task with over 200 participants and a diverse set of stimuli in terms of length and word shapes. Results indicate a gradient identity avoidance effect that extends over all consonants in the base. Crucially, the effect of all consonants is not uniform, with the strength of the effect decreasing further into the base. The study also uncovers an elusive interplay between the distance-based decay effect and the syllable position effect, both of which turn out to play a role in these non-categorical patterns with multiple features. Furthermore, results indicate that identity avoidance operates over both individual features as well as whole segments. Overall, the study argues that locality-sensitive feature-specific identity avoidance constraints are part of the grammar.

---



## 1. Introduction

Many languages exhibit a process in which a consonant is altered in order to reduce its resemblance partially or fully to another consonant (see e.g., Stanton 2017; Suzuki 1998 for some discussion of (long-distance) consonant dissimilation). A representative example of dissimilation comes from Latin, in which certain suffixes exhibit dissimilatory effects with respect to the consonants in the stem. For example, when the stem does not contain any lateral consonant [l], the adjectival suffix *-alis* surfaces as is (1a-b). However, when the stem contains a lateral consonant at any position in the stem constituent, the suffix surfaces as [-aris] (1c-d).

- (1)
- |    |                |              |                               |
|----|----------------|--------------|-------------------------------|
| a. | /nav-alis/     | nav-alis     | ‘naval’                       |
| b. | /episcop-alis/ | episcop-alis | ‘episcopal’                   |
| c. | /sol-alis/     | sol-aris     | ‘solar’                       |
| d. | /lun-alis/     | lun-aris     | ‘lunar’ (Suzuki 1998:12, (3)) |

In this regard, reduplication is a particularly fruitful process to examine since although in many cases the material attached to the base in reduplication resembles the base phonologically, exact copying of the base is not always achieved (see e.g., Alderete, Beckman, Benua, Gnanadesikan, McCarthy & Urbanczyk, 1999; Frampton, 2009 for some illustrations). As such, similarity between the reduplicated form and the (relevant part of the) base is reduced. Against this background, we analyze a pattern of reduplication from Turkish with a partial dissimilation, which turns out to be informative as to several open questions related to the effects of the identity avoidance,<sup>1</sup> such as at what granularity it applies (e.g., at the level of segments or phonological features), the extent to which proximity (locality) plays a role, or its interaction with syllable roles, especially in instances with multiple features involved.

Cross-linguistic studies have found that features may differ in their strength in phonological patterns that involve identity avoidance (see an extensive survey of 46 phenomena by Bye (2011)). While many feature classes have been found to participate in identity avoidance, such as place of articulation, laryngeal state, manner (continuancy, liquid, nasality), vowel height, and suprasegmental properties (length and tone), these typological analyses are based mostly on patterns that involve only one or two featural changes. Analysing such patterns means that the researchers are able to clearly identify the relevant features. For example, only the [lateral] feature participates in identity avoidance in the Latin liquid dissimilation. Looking beyond phonological patterns, several studies (Berkley, 2000; Frisch et al., 2004; Graff & Jaeger, 2009) have examined the statistical asymmetries in the lexicon which suggests that there is preference for phonetically dissimilar consonants within a word, such that fewer similar sounds co-occur than would be expected *a priori*. This suggests that identity avoidance plays a role in shaping the phonological

---

<sup>1</sup> In this paper, we use ‘identity avoidance effect’ for the most part, instead of the Obligatory Contour Principle (OCP) as a more neutral term although we still refer to OCP as well occasionally.

lexicon. This type of analyses provides a fruitful ground for better understanding the typology of participating features in identity avoidance. This is because one must examine a larger set of features simultaneously, owing to the logically possible combinations of consonants in a root, as opposed to just one or two features that participate in a phonological pattern. However, there exists another type of phonological phenomena that involve identity avoidance with a larger set of phonological features, such is the case of Turkish partial reduplication, which involves four consonants, potentially allowing multiple features to participate in dissimilation. Looking at these phenomena would further inform our understanding of the typology of identity avoidance.

The partial reduplication in Turkish also helps further our understanding of *syllable position* in assimilation (and dissimilation), a phenomenon examined in Suzuki (1998); Rose and Walker (2004); Bennett (2012); Bennett (2013). Focusing on R-dissimilation and L-assimilation in Sundanese as well as nasal assimilation (i.e., nasal agreement/harmony) in Kikongo, these studies find the following patterns, which are also categorical as they usually involve one or two phonological features: R-dissimilation occurs whenever the affix /r/ and an /r/ in the root have different syllable roles. In the case of both L-assimilation and nasal assimilation, however, matching syllable roles contribute to segments' similarity allowing assimilation to take place between segments that share the same syllable role. For example, in the case of Sundanese L-assimilation, the initial /l/ of the root and the /r/ of the affix /ar/ end up as the onsets of the first two syllables of the stem. In fact, these (mis)match patterns are given as predictions in Bennett (2013, p. 536) using the constraint CC-SROLE, which limits correspondence based on structural position. Accordingly, Bennett argues that while harmony/assimilation is predicted between consonants with matching syllable roles, dissimilation is predicted for consonants with mismatching syllable roles (but not those with matching syllable roles). Various questions arise in light of these patterns: First, is it possible to find dissimilation that would favor matching syllable roles? Secondly, can the effect of syllable role still be observed in non-categorical patterns with multiple features, and relatedly, what are the effects (if any) of syllable role in gradient assimilation/dissimilation?

Moreover, little is known about the interaction of the effect of syllable role with the *distance-decay effect*. Zymet (2014) illustrated that in multiple languages the likelihood of application of assimilation or dissimilation decreases as transparent distance increases. The patterns that were examined involve usually one or two phonological features, such as rounding dissimilation in Malagasy, liquid dissimilation in Latin and in English, and vowel harmony in Hungarian. It remains an open question whether distance-based decay can still be observed in patterns that involve with multiple features, and whether there is an interplay between the distance-based decay and syllable position.

Turkish is apt to address – at least provide significant insights to – these open questions because the partial reduplication in Turkish is a gradient dissimilation phenomenon, which involves four distinct consonants, potentially allowing multiple features to participate in dissimilation.

Furthermore, it shows a strong effect in matching syllable roles, and exhibits an interesting interplay between the role of syllable position and the distance-based decay.

As will be shown below, Turkish has a type of emphatic reduplication, in which the emphatic variants are derived by prefixing a  $C_1VC_2$  syllable, where  $C_1V$  are identical to the word-initial CV of the base, and the  $C_2$  (known as *linking consonant*) ends in one of the four consonants ( $-p$ ,  $-m$ ,  $-s$ ,  $-r$ ), effectively reducing the resemblance with the base. The effect of consonant dissimilation has largely been attributed to the first two consonants of the base (Kelepir 1999, 2000; Wedel 1999). This study examines the nature of the factors behind the choice of the  $C_2$  using an acceptability rating task with over 200 participants and a more diverse set of stimuli in terms of length and word shapes than the ones employed in prior literature. While partially confirming the conclusions of some prior studies, the current study reveals novel findings. The effect beyond the first two consonants have been mentioned but they are not often found to be statistically significant or required when the first two consonants have already been considered (e.g., Yu 1999). This pattern suggests that locality in which dissimilation operates in a categorical manner. This is particularly significant in light of the recent discussions of distance-based decay effect (Zymet 2014, 2018), which, as mentioned above, suggests that the likelihood for the application of a phonological process decreases as the distance increases, and this can happen in a gradient and non-linear manner. This study establishes that the effect of identity avoidance spreads across all consonants in the base, which crucially is not uniform but exhibits a distance-based decay effect. This effect also interacts with the syllable position effects. In particular, the effect would decay with the distance from the linking consonant and would be enhanced if the consonant in the base matches the constituency of the linking consonant (coda vs onset). For instance, in the case of Turkish, the linking consonant is always a coda (in the C-initial forms, which is the focus of this study). When a base form has the word shape  $C_1VC_2V$ , a distance-based decay effect was found, as hypothesised, with  $C_2$  being less influential than  $C_1$ . However, the word shape  $C_1VC_2$  exhibits the opposite pattern with  $C_1$  being less influential than  $C_2$ . This surprising effect can be attributed to the fact that the LC and the  $C_2$  in  $C_1VC_2$  are both codas. Accordingly, the Turkish phenomenon is particularly insightful in this regard too: It showcases the interplay of distance-based decay effect and syllable position effect, which in most studies are treated separately (due to the data involved). Additionally, in line with findings in cross-linguistic typology, features do not participate equally in identity avoidance processes, with some being more influential than others. Our study confirms the importance of only some of the features employed in previous studies, such as [strident], [labial], and [nasal], but not others, such as [coronal], [sonorant], or [continuant].

These new findings contribute to our understanding of the nature of locality, particularly arguing for the view that locality-sensitive feature-specific identity avoidance constraints are part of the grammar. This study overall highlights the value of revisiting a long-debated topic

with a new lens to address any remaining unsolved puzzles, replicate existing findings, and generate new hypotheses that can contribute to the future studies.

The rest of the Introduction section introduces the phenomenon itself, and situates its implications within the larger literature with a focus on locality/syllable role and feature-specificity.

## 1.1. The phenomenon and the implications

Turkish is an agglutinative language, with the majority of derivation achieved through suffixation. A rare instance where prefixation is observed in the language is the so-called *partial reduplication* (or *emphatic reduplication*) (Demircan, 1987; Lewis, 1967, i.a.). It is the only prefix present in the language except for foreign prefixes in borrowed forms. The partial reduplication is found with modifiers, namely adverbs and adjectives. As shown in (2),<sup>2</sup> emphatic variants are derived by prefixing a  $C_1VC_2$  syllable. The initial  $C_1V$  are identical to the word-initial CV of the base. However,  $C_2$  ends in one of the four consonants: *-p*, *-m*, *-s*, *-r* (Lewis, 1967), generally referred to as *linking consonants (LC)*.<sup>3</sup>

(2)	Base	Gloss	Reduplication	Gloss
	<i>kara</i>	‘black’	<i>kap-kara</i>	‘very black’
	<i>beyaz</i>	‘white’	<i>bem-beyaz</i>	‘very white’
	<i>ma:vi</i>	‘blue’	<i>mas-ma:vi</i>	‘fully blue’
	<i>temiz</i>	‘clean’	<i>ter-temiz</i>	‘completely clean’

On the other hand, if the base is V-initial, the LC surfaces as *-p*.

(3)	Base	Gloss	Reduplication	Gloss
	<i>açık</i>	‘open’	<i>ap-açık</i>	‘very open’
	<i>ince</i>	‘thin’	<i>ip-ince</i>	‘very thin’
	<i>ansızın</i>	‘suddenly’	<i>ap-ansızın</i>	‘very suddenly’

Partial reduplication in Turkish has been the subject of many studies, which have aimed to explain the conditions underlying the choice of the LCs (see e.g., Demir, 2018; Demircan, 1987; Dobrovolsky, 1987; Foster, 1969; Hatiboğlu, 1973; Kaufman, 2014; Köylü, 2020; Lewis, 1967; Sofu, 2005; Sofu & Altan, 2008; Taneri, 1990; Wedel, 1999; Yavaş, 1980; Yu, 1998, 1999). This phenomenon is particularly suited for shedding light on the nature of the identity avoidance effect in terms of its locality and feature specificity. It is a type of dissimilation in reduplication with *fixed segmentism* (Alderete et al. 1999). Crucially, prototypical instances of reduplication with fixed segmentism involve only a single fixed unit (e.g., [a] in Javanese habitual reduplication (Yip

<sup>2</sup> In the transcription, the following correspondences hold between the orthographic forms and IPA:  $ı \rightarrow [u]$ ,  $ç \rightarrow [tʃ]$ ,  $c \rightarrow [dʒ]$ ,  $ş \rightarrow [ʃ]$ ,  $ü \rightarrow [y]$

<sup>3</sup> In terms of the morphological analysis, we assume that the reduplication is RED + LINKER + BASE where RED is only CV and only the LINKER can dissimilate.

1997:p. 18), or [m] in Turkish, (4), in which an invariant segment appears whatever the features of the base). However, as seen in (2), the number of fixed segments is four in Turkish emphatic reduplication. This allows for many more features to participate in the identity avoidance effect.

(4)	<b>Base</b>	<b>Gloss</b>	<b>Reduplication</b>	<b>Gloss</b>
	<i>sarı</i>	‘yellow’	<i>sarı marı</i>	‘yellow or similar colors’
	<i>kapı</i>	‘door’	<i>kapı mapı</i>	‘door or the like’ (Turkish)

Another interesting property of the Turkish emphatic reduplication is that, it differs from most crosslinguistic examples in not just allowing multiple possible fixed segments, but also having variability in which LC would appear for a given stem, as exemplified in (5). As such, although it is not always reported in the literature (though see Müller, 2003; Wedel, 1999; Yu, 1999), there is indeed variability with respect to the choice of the LC with a substantial number of bases.

(5)	<b>Base</b>	<b>Gloss</b>	<b>Reduplication</b>	<b>Gloss</b>
	<i>yeşil</i>	‘green’	<i>yem/p-yeşil</i>	‘completely green’
	<i>başka</i>	‘different’	<i>bam/p-başka</i>	‘very different’
	<i>buruşuk</i>	‘creased’	<i>bum/s-buruşuk</i>	‘very creased’
	<i>yırtık</i>	‘torn’	<i>yıs/p-yırtık</i>	‘completely torn’

Moreover, with items that allow multiple LCs, variation exists as to which LC is preferred for a given item across speakers. Even the existence of a varying degree of preference for the choice of a particular linking consonant (and between multiple LCs) differentiates the Turkish phenomenon from the classic reduplication with fixed-segmentism. Crucially, in Turkish partial reduplication, as we will argue in this paper, the choice of the linking consonant is sensitive to the features of the base. In some instances of reduplication with fixed segmentism (e.g., Alderete et al. 1999; McCarthy & Prince, 1993), the identity of that fixed segment is attributed to the emergence of the unmarked, default form. However, as has been noted by Wedel (1999) and Yu (1999), this is not the case in Turkish, and these LCs appear in the output despite clear markedness violations and are not the unmarked segments in the language. The consonants [p, m, s] are not considered the default segments in Turkish, and instead [n, j] are considered to be so (Wedel, 1999). As such, the latter are known as ‘buffer consonants’ in traditional Turkish grammars, appearing in various contexts (e.g., breaking up vowel hiatus as in *keđi-y-i* ‘cat-[j]-Accusative’ ‘the cat’).

While the seminal work by Demircan (1987) and subsequent work all identified the identity avoidance effect as a major factor, their accounts often involve heuristically chosen features that operate over the first two consonants in the base.<sup>4</sup> The current study (re)-examines the factors

<sup>4</sup> While the heuristic choice of the features might be a well-established way of doing traditional phonological analysis (Kenstowicz & Kisseberth, 2014, Ch. 2), which we also appreciate, this study demonstrates that a more statistical approach that does not rely on heuristic choices reveals properties that would otherwise be missed. Thus, our goal is not to critique the use of heuristics per se, but to highlight the point that in some cases we can learn more from a statistical approach.

conditioning the choice of the LC, focusing on the C-initial forms since the V-initial forms are consistently reduplicated with the LC *-p*. We extend the previous studies by examining the nature of the proximity (locality) and the (dis)similarity (feature specificity) between the consonants in the base and the LC. Starting with the former, we examine both topics in turn, particularly how they have been handled in the Turkish literature on emphatic reduplication as well as their implications for the broader literature beyond Turkish.

## 1.2. Locality avoidance in Turkish and beyond

With respect to locality, most studies emphasise the importance of  $C_1$  and  $C_2$ , yet some of them (implicitly) assume a cut-off after  $C_2$ . One of the questions this study addresses is whether the consonants beyond  $C_2$  do not play a role in the choice of the LC. Let us begin by taking a closer look at the classic study on the phenomenon by Demircan (1987), which observes that the selection of the linker is subject to various dissimilation constraints. The primary observations are given in (6), in the format they are succinctly summarized in Yu (1999, p. 5 & p. 18) and with slight modifications.<sup>5</sup>

- (6) Demircan's (1987) observations (adapted from Yu, 1999):
- (i) The linker cannot be identical with any of the consonants in the base.
  - (ii) No gemination: The linker should not be identical to the  $C_1$  of the base.
  - (iii) Avoid full reduplication: The linker cannot be identical with  $C_2$  of the base with  $C_1VC_2$  items.
  - (iv) Featural identity avoidance: Avoid a linker that shares similar features, such as [labial], [strident] & [approximant], with any segment in the base.

Note that the observations (ii) and (iii) by Demircan focus on the contrast between the linking consonant and the  $C_1$  and  $C_2$  of the base, while the observations also make reference to the whole base. However, this latter point has been underappreciated in some subsequent studies. For example, based on the minimal pair in (7), Kelepir (1999, 2000) argues that not only  $C_1$  of the base, but also  $C_2$  matters for the choice of the LC.

- (7) a. *yeni* 'new' *yep-yeni* 'completely new'  
 b. *yeşil* 'green' *yem-yeşil* 'completely green' (Kelepir, 2000 p. 11)

Similar to Demircan's (1987) observation (i), Kelepir's (1999) study also has a constraint/restriction, \*Repeat [strident], which makes reference to the whole base (specifically, that rules

---

<sup>5</sup> The summaries are Yu's recasting of Demircan's observations about dissimilation constraints, and thus has some slight modifications which do not make a difference for the content. For example, Yu (1999) uses the term *linker*, while Demircan (1987) uses the term 'closer' for the LC. The observation (i) is given as "Avoid closers identical with any of the base consonants to rule out..." in Demircan (1987), while it is given as "The linker cannot be identical with the final consonant of the base" in Yu (1999) with some other qualifications.

out the strident linker [s] if there is a strident in the base). Yet, Kelepir's constraint system centers around the comparison of the LC with respect to  $C_1$  and  $C_2$  of the base.

Wedel (1999, 2000) is another study that argues for the presence of dissimilatory phonological constraints in partial reduplication, yet lacks a constraint of the sort proposed by Demircan and Kelepir, who allowed for the scanning of the whole base to avoid the choice of an LC that shares all or some of the features of the LC. In fact, the author explicitly mentions that there should be a cut-off after  $C_2$ . This is reflected in the generalizations and the OT constraints that Wedel proposes. Consider (8):

- (8) (i) [p] is not selected if  $C_1$  is labial: \*PLOSIVE- $\alpha$ PL  
 (ii) The interpolated consonant [LC] must be non-identical to  $C_1$ : \*GEM  
 (iii) The interpolated consonant [LC] must be non-identical to  $C_2$ : \*REPEAT (Wedel 2000: 550)

The assumption that consonants beyond  $C_2$  have no significant effect has led some researchers that do nonce-word studies to only create nonce-words consisting of two-consonant stems. This is most clearly seen in the study of Köylü (2020), whose all 48 nonce-words consists of bases with a maximum of two consonants.

In fact, examining the features labial, coronal, and strident, Yu (1999) found that the strident feature has an effect with respect to  $C_3$ , and corroborates Kelepir's \*Repeat [strident] constraint. More generally, this supports the observation that the LC is not restricted to identity avoidance effect with respect to  $C_1$  and  $C_2$ . Crucially, Yu did not find a significant effect from segments beyond  $C_3$ , which he attributed to the limited number of adjectives and adverbs with more than three consonants in his corpus data.

To sum up, the role of  $C_1$  and  $C_2$  has been the focus of many previous studies of the Turkish phenomenon.<sup>6</sup> There is nonetheless evidence from other languages that suggests all the consonants in most dissimilation phenomena could play a role in identity avoidance, being also subject to a *distance-based decay effect* (Zymet 2014, 2018). This effect states that the likelihood for the application of a phonological process decreases as transparent distance increases. Arabic is the poster child for this kind of effect in identity avoidance (see e.g., Coetzee & Pater, 2008; Frisch & Zawaydeh, 2001; Frisch et al., 2004; McCarthy, 1986; McCarthy & Prince, 1994). For example, OCP-Place is a single gradient constraint that restricts consonant co-occurrence in Arabic based on (i) similarity (see below for the discussion of similarity) and (ii) proximity (Frisch & Zawaydeh, 2001; Frisch et al., 2004). Particularly, the influence of similarity on consonant co-occurrence is affected by distance, as the constraint is weaker for non-adjacent consonants. Accordingly, the fact that many prior studies on Turkish partial reduplication consider only  $C_1$  and  $C_2$ , and the null effect of segments beyond  $C_3$  in Yu (1999) might also be due to the distance-based decay

---

<sup>6</sup> But see Demircan, 1987; Yu, 1999.



effect. The present study indeed uncovers that the identity avoidance in Turkish is subject to such an effect, where it holds for all the segments, with the effect being strongest from  $C_1$  and being weakened as a function of its distance from the target segment.

### 1.3. (Dis)similarity of features in Turkish and beyond

The (dis)similarity of features plays an even more prominent role in the identity avoidance literature. For example, studies usually aim to probe which specific features identity avoidance constraints are sensitive to whether they are categorical (Bennett, 2013), or gradient, (see e.g., Frisch et al., 2004; Gallagher & Coon, 2009; McCarthy, 1986), or whether features matter to the same degree. If the answer to the latter is no, then what is the extent to which specific features matter as opposed to other features in a given language and cross-linguistically? (see e.g., Bye, 2011; Coetzee & Pater, 2008; Frisch et al., 2004; Gallagher & Coon, 2009; Graff & Jaeger, 2009) It turns out, typologically, not all features participate equally in identity avoidance. Based on an extensive survey of 46 phenomena by Bye (2011), the following phonological dimensions have been found to play a role: The place of articulation, the laryngeal state, the manner of articulation (continuancy, liquid, nasality), vowel height, and suprasegmental properties such as length and tone. However, major class features such as [consonantal], [sonorant], and [approximant] do not play much of a role. Concerning the place feature, [labial] is relatively common, while [coronal] is rare and [dorsal] is unattested. Furthermore, alternations involving laterals and rhotics are relatively common.

When we look at previous studies on Turkish, we observe that they tend to focus on specific (set of) features. For instance, the analysis in Keleşir (1999) is built on avoidance constraints that use specific features. In particular, base consonants contrast with their correspondents in the reduplicant in place and sonorancy (i.e., the features [coronal], [sonorant], [labial], or [continuant]). Demircan (1987) identifies [coronal], [labial], and [nasal] as features of importance, and Yu (1999) adds that [strident] and [approximant] features also play a role in the interaction of segments.

Some remarks are in order regarding these previous studies on Turkish and their connection to the larger literature. Firstly, the choice of the relevant features is usually heuristic. For instance, Wedel (1999, 2000) uses PLOSIVE- $\alpha$ PL to explain why [p] is not selected if  $C_1$  in the base is labial, while Keleşir (2000) uses the exact consonant [b] directly through the constraint, \*-pb-. The heuristic nature mostly stems from the fact that researchers posit features, OT-constraints, or perception-related restrictions that can explain the respective datasets. In fact, a lot of the proposed constraints are correlated and refer to overlapping issues, in that a constraint might be implicated or even entailed by another constraint. For example, the rule (i) in (6) encompasses the rule (ii). Secondly, and relatedly, it is unclear whether these features are the only features that matter, or perhaps there are some other features that may also trigger identity avoidance, but have been missed since most studies rely on the researcher's informed observations rather than a

systematic examination of all possible features. Thirdly, cross-linguistic studies have found that features may differ in their strength in identity avoidance. For example, Gallagher and Coon (2009) find that [+strident] and [+ejective] have greater effects than others in Chol, while non-coronal place features were found to play a greater role in Arabic and Muna (Coetzee & Pater, 2008; Pierrehumbert, 1993). As such, various features have been shown to trigger the identity avoidance effects of varying strengths in different languages. Several proposals have been put forth with the aim of capturing the varying strength of the features. For example, Frisch et al. (2004) proposes a *similarity metric of natural classes*, which are shown to strongly correlate the observed-over-expected ratios (O/E) of consonant pairs co-occurring within Arabic roots, in a way that relies on the feature inventory of a language. However, data from other languages have not been able to replicate the effect of this metric. Coetzee and Pater (2008) failed to establish a correlation using this metric for Muna and Rotuman. Their discussion of place of articulation patterns in Muna and Arabic also demonstrates that the relative strength of place co-occurrence patterns cannot be due to inventory structure alone. Similarly, Graff and Jaeger (2009) found no evidence for such a correlation using this metric by any of the three languages examined (Dutch, Aymara, and Javanese). An alternative approach was proposed by Graff and Jaeger (2009) which did find effects of feature-specific effects in identity avoidance of the three language using a more complex model. This model allows each individual feature from each consonant in the base being weighted freely.

Additionally, previous analyses involved using a specific set of features, but their formulation of these features are often correlated, for example, a constraint is entailed by another. The choice of relevant features, therefore, tends to be heuristic. This is partly due to feature redundancy, a design property of standard versions of feature theory. Multiple sets of relevant features that could play a role in consonant dissimilation. For example, one could use any of following features [nasal], [labial], [sonorant], [voice] to model the dissimilation with /m/. That is not to say redundant feature values do not matter. In fact, Keyser and Stevens (2006), and Stevens and Keyser (1989) have proposed that redundant feature values can enhance the phonetic interpretation of contrastive values (see Clements & Ridouane, 2006 for an overview). On top of feature identity avoidance, previous analyses also make reference to total identity of  $C_2$ , such as  $C_2$  cannot be identical to any of the consonants in the base. Again, these two types of factors are correlated. Previous analyses have also not fully addressed the question of whether the effect of total identity can be reduced to the effect of feature identity. These considerations call for revisiting the phenomenon with a different theoretical angle,<sup>7</sup> as well as new methodological tools that can jointly evaluate identity avoidance of all features and linking consonants, therefore allowing us to tease apart their relative contributions.

---

<sup>7</sup> For example, this study is not concerned with the particular angle some prior studies take (i.e., whether the Turkish partial reduplication fits into a phonological analysis (Emergence of the Unmarked), or morphological analysis (melodic overwriting) of fixed segmentism) (see e.g., Kelepir, 1999, 2000; Yu, 1999).

Accordingly, besides probing the locality effect, this study also aims to better understand the inventory of features that may be involved in triggering the identity avoidance effect in Turkish emphatic reduplication, and if certain features have a greater effect in this phenomenon in relation to other features. In so doing, we also aim to elucidate the question of how the similarity between the consonants in the base and the LC should be specified, for instance, at the level of the total identity of the consonants, or at the level of the individual features? As will be discussed later in Section 2, this study adopts the methodological approach by Graff and Jaeger (2009) given its ability to help us rigorously examine the nature of both locality and feature specificity, and establish whether locality-sensitive feature-specific identity avoidance constraints are part of the grammar.

To address the questions that revolve around locality and feature-specificity, in this study we conducted an acceptability rating judgement experiment (which has rarely been used for the study of the partial reduplication phenomenon) as opposed to researcher's intuitions or experimentally obtained forced-choice task responses. Among other things, our findings provide support for the view that speakers' grammars have active identity avoidance constraints that operate on specific features (e.g., Gallagher & Coon, 2009) and the strength of their effects is a function of the distance between similar consonants (e.g., Pierrehumbert, 1993; Zymet, 2014, 2018). These are in line with previous findings on Arabic (Frisch et al., 2004), Dutch, Aymara, and Javanese (Graff & Jaeger, 2009) which highlighted a gradient feature-similarity-based restriction that is also subject to locality (see also Suzuki, 1998).

In terms of the nature of feature specificity, our results reveal that in the Turkish emphatic reduplication process, the similarity between the consonants in the base and the LCs operates at the segmental level (total identity) as well as the level of individual phonological features (partial identity). Furthermore, we found that not all features participate equally by allowing individual features to be free parameters. Our study also confirms the importance of only some of the features employed in the previous studies, such as [strident], [labial], and [nasal], but not others, such as [coronal], [sonorant], or [continuant]. The important features like [strident] and [labial] and the unimportant features such as [coronal] and [sonorant] are more in line with the cross-linguistic tendencies (see e.g., Bye, 2011; Pierrehumbert, 1993).

In terms of the effect of locality, the study demonstrates that not only the first two consonants, but all the consonants in the base form contribute to the identity avoidance effect (Zymet, 2014), with the strength of the effect decreasing further into the base. Another locality-related factor that studies on assimilation (and dissimilation) have noted is the role of *syllable position* (Rose & Walker, 2004) or *syllable-role specific correspondence* (Bennett, 2012). This approach states that matching syllable roles might contribute to segments' (dis)similarity. Given this aspect has not been investigated in the context of Turkish emphatic reduplication, our study also examines whether syllable position is a significant factor. Strikingly, it turns out the identity avoidance

effect in Turkish is also sensitive to the syllable position. In this regard, it lends support to the view that syllable position might play a role in contributing to segments' (dis)similarity (Bennett, 2012; Rose & Walker, 2004; Suzuki, 1998). In particular, the effect would decay with the distance from the linking consonant and would be enhanced if the consonant in the base matches the constituency of the linking consonant (coda vs onset). For instance, in the case of Turkish, the linking consonant is always a coda.<sup>8</sup> When a base form has the word shape  $C_1VC_2V$ , a distance-based decay effect was found, as hypothesised, with  $C_2$  being less influential than  $C_1$ . However, the word shape  $C_1VC_2$  exhibits the opposite pattern with  $C_1$  being less influential than  $C_2$ . This surprising effect can be attributed to that the LC and the  $C_2$  in  $C_1VC_2$  are both codas (see Section 3.2.2 for the complete result). Accordingly, the Turkish phenomenon is particularly insightful in this regard too: It showcases the interplay of distance-based decay effect and syllable position effect, which in most studies are treated separately (due to the data involved). Moreover, unlike other examples that illustrate the syllable position effect, the presence of this effect in Turkish is not immediately clear both due to the just-mentioned interplay with distance-based decay effect and also because the dissimilation is not categorical (cf. Bennett, 2013), but much more gradient.

This study also replicates some of the previous Turkish-specific findings. For example, the preference hierarchy regarding the choice of the LC ([p] > [s] > [m] > [r]) still broadly holds. Methodologically we demonstrate that the precise nature of the identity avoidance effect can be revealed using hierarchical regression and statistical model comparisons (Graff & Jaeger, 2009; Zymet, 2019).

The rest of the paper is organised as follows: Section 2 outlines the details of our study. We first introduce the logic of this study before we present the materials (Section 2.1), followed by the methodological details of the experiment (Section 2.2). Section 2.3 introduces the variables of interest. Section 2.4 outlines the modelling procedure. In this section, we will introduce the model specification and evaluation as well as the modelling details of Study I on feature specificity (Section 2.4.2) and Study II on locality (Section 2.4.3). Section 3 presents our results for the two studies. Section 4 contextualises our findings of the Turkish case into the broader literature, focusing on feature specificity and locality, particularly the interaction between distance-based decay and the syllable position. This section also discusses issues such as factors beyond identity avoidance and representation of speakers' knowledge. Section 5 summarises and concludes the paper.

## 2. The present study

The present work consists of two studies. Study I addresses the level of similarity on which the identity avoidance effect operates. Study II addresses whether the proximity between the consonants and the linking consonant plays a role in the identity avoidance effect.

---

<sup>8</sup> Note again that the linking consonant is always a coda in the C-initial forms.

To examine these research questions, we conducted a large-scale rating study of 162 real base forms of Turkish sampled from previous studies. **Table 1** summarises most of the previous studies along the lines of various criteria that will be referred to in the current study: (i) Whether the study relies on the researcher’s intuition or an experiment, (ii) what type of experiment was conducted and the number of participants, and (iii) whether the items used are real words or nonce-words.

Sources	Intuition	Experiment	Type of Exp.	# of Participants	# of Items and Types
Hatiboğlu (1973)	Yes	No	–	–	Real (142)
Demircan (1987)	Yes	Yes	FC	100	Real (110), Nonce (20)
Dobrovolsky (1987)	Yes	No	–	–	Real (9)
Taneri (1990)	–	Yes	FC	32	Real (300)
Wedel (1999)	Yes	Yes	FC	3–8	Real (125 + 80)
Yu (1998)	No	Yes	Rating	4	Real (101), Nonce (56)
Yu (1999)	Yes	No	–	–	Real (152)
Kelepir (2000)	Yes	No	–	–	Real (89)
Sofu (2005)	–	Yes	FC	25 adults, 89 children	Nonce (38)
Sofu & Altan (2008)	–	Yes + Corpus	FC	80	Real (132)
Kaufman (2014)	–	Yes	1-FC, 2-Rating	1–16, 2–50	1-Nonce (44), 2-Real (45)
Demir (2018)	–	Yes	OSR	125	Real (10), Nonce (34)
Köylü (2020)	–	Yes	FC	14	Nonce (48)

**Table 1:** A summary of the data examined in 13 previous studies. FC stands for *forced choice*, and OSR stands for *open-set response*.

The decision of using a rating task was motivated by our desire to fill a methodological gap that was found in previous studies. As summarised in **Table 1**, in many studies, the judgements were often based on the researcher’s intuitions only and whenever there is an experimental component, the task is almost exclusively a forced-choice task. The use of acceptability rating

might have an advantage over a forced-choice task in that the forced-choice task might be masking potential variability within participants that we noted in (5).

Note also that previous studies on Turkish either completely ignored the variable nature of the LC, or incorporated into their analyses only the most dominant LC (even when they have empirical basis for the variability). For instance, Wedel (2000) reported that most participants responded with only one LC for each base form, even though the participants as a group chose multiple LCs. In a post experiment interview, the participants reported that the other forms that they did not choose were also possible; however, they simply selected the first one that they had in mind. A rating task would therefore be able to better examine base forms that have a high level of variability, allowing for multiple LCs to apply to the same base form by a given participant.

Our study reports on a large-scale rating experiment with 162 real base forms with at least 40 participants rated each form. It is a methodological improvement over existing studies that made use of a rating task, since they either tested a small number of participants such as Yu (1998) with four participants, or tested a small number of items such as Kaufman (2014) with 10 real base forms.

## 2.1. Materials

The experimental items were 162 real base forms taken from previous studies to enhance the comparability and the replicability of this work. Most of the items were taken from the classic study by Hatiboğlu (1973) because many of the later studies also examined a subset of these items, and the rest were sampled from the other studies.<sup>9</sup> The 162 items were then evenly divided into five lists (three lists have 33 items; one list has 32 items; and one list with 31 items) with each list containing roughly the same distribution of dominant linking consonants as well as variable items.<sup>10</sup>

## 2.2. Methods

Each participant was asked to perform both a rating task and a forced-choice task (not reported here).<sup>11</sup> For each base form, all four of its reduplicated forms (each with a different linking consonant (LC)) were shown on the same screen orthographically. Each participant was randomly

---

<sup>9</sup> In addition to the items from Hatiboğlu (1973), we included a few more items from the list in Stachowski (2014), which compiled the list of mostly overlapping items from Hatiboğlu (1973), Demircan (1987), Müller (2003), and others. Out of the 178 items in Stachowski's (2014) list, we left out 16 items since some of them were nouns (e.g., *buz* 'ice', or *çevre* 'environment'), and others were items not available in Turkey Turkish, but other Turkic varieties such as Azeri Turkish (e.g., *deyirmi* 'circle').

<sup>10</sup> The expected dominant linking consonant of each item was based on previous studies as well as the linguistic intuition of Faruk Akkuş, the co-author who is a native Turkish speaker. The expected distribution of the items consists of 55% P-items, 9% M-items, 17% S-items, 15% R-items and 24% variable items. This information was not used in our statistical analyses and served merely for the purpose of the experiment design.

<sup>11</sup> A simple item-level correlation analysis suggests the judgements from the forced-choice task are similar to those obtained from the rating task ( $R^2: 0.85$ ).

assigned to one of the five lists. The order of the two tasks and the order of the four reduplicated forms were also randomised for each participant. Each reduplicated form was rated on a scale of naturalness: DOĞAL DEĞİL ‘not natural’ [1 to 7] DOĞAL ‘natural’. The experiment was programmed using Experigen (Becker & Levine, 2013), hosted at <http://db.phonologist.org/>.

The experiment was advertised via social media. Participants were invited to take part in the experiment voluntarily and given informed consent. The inclusion criteria of our target population were native Turkish speakers, born in Turkey without language-related disorders. A total of 283 participants completed the experiment. 207 participants who met our inclusion criteria were included in the analyses. All items were rated by at least 40 participants.<sup>12</sup>

We evaluated the results of the rating study by adapting Graff and Jaeger’s (2009) methodological approach. Graff and Jaeger (2009) examined the feature specificity and locality of the identity avoidance effect in the lexical organisation of Aymara, Dutch, and Javanese. Methodologically they made use of a regression approach to allow for individual identity factors as well as nuisance factors to act as free parameters. They compared the different types of identity factors (total identity and featural identity between two segments) and whether these factors are strictly local by using a model comparison approach. Their key findings were that i) both the total identity and the featural identity of two segments influence the formation of lexical roots, ii) the identity avoidance effect operates over individual features which have their individual weights, and iii) the identity avoidance effect operate both locally and non-locally. Our study aims to ask the same research questions concerning the identity avoidance effect. Using this approach we were able to quantitatively examine the effect of feature specificity and locality rather than heuristically as was done in many previous studies. Unlike the previous studies which modelled only the dominant linking consonant of each base form, we modelled trial-level responses from each participant using linear mixed-effects regression.

While the individual features included as predictors are also the same set of features that suffers from the same feature redundancy issue discussed in Section 1.3, the statistical regression approach enables researchers to deal with this issue by disentangling the direct effects of specific variables. The coefficient of each variable in a multiple regression model represents the variable’s influence while statistically controlling for the influence of other variables (Winter, 2019). Thus this study makes use of a novel approach of using statistical reasoning to determine which features are relevant, rather than preselecting a particular set of features using more traditional phonological means.

The 162 items were chosen to ensure that words with more than two consonants are well represented (see **Table 3** for the breakdown of the selected words by the number of consonants). This was motivated by how previous studies were restricted in their ability to properly examine the

---

<sup>12</sup> Their mean age was 27.44 years, ranged between 18 and 63. 146 were women and 61 were men.

potential contribution from consonants further away from the reduplicant due to limited number of stimuli with more than two/three consonants limit (see Section 1.1 for the full discussion).

These items are part of the existing vocabulary of the Turkish lexicon. Therefore, they are likely to have an uneven representation of consonants and their features. Due to this uneven representation in the items, it is possible that, for example,  $C_1$  and its features might not have the same effect across the stems with different numbers of consonants. The nature of the stimuli therefore imposes a limit as to how strong an effect of a given feature can exert; for example, [strident] can be shown to be significant across items but [labial] may not be because its representation happens to be low. Not separating the items by the number of consonants could mask these potential differences (as we will see in the result section, there are indeed notable differences; see Section 4 for a discussion).<sup>13</sup> We therefore chose to model our items separately for stems with different numbers of consonants (see Sections 2.4.2 and 2.4.3 for details). Furthermore, this analytical approach is necessitated by the regression modelling approach by Graff and Jaeger (2009) that we are adopting. For instance, if we were to combine both items with two and three consonants in a single regression model, then the predictors concerning the third consonants would not be able to be specified for the items with only two consonants. Relatedly, given the nature of the dataset and the modelling procedure, it is worth noting what can be inferred from the results. First, we can infer the overall importance of the features in Turkish emphatic reduplication by examining how prevalent they are across models in terms of their role in capturing the data. Second, we can also infer the nature of feature specificity by comparing between models (using model comparisons) with different formulations of identity avoidance predictors. Third, we can also infer the nature of locality by comparing models (again with model comparison) with and without the features associated with a specific consonant. Please see Section 4.1 for a detailed discussion.

## 2.3. Variables

This section describes the fixed effect variables and the random effect variables we included in our analyses. Variables that we excluded can be found in Appendix B.

### 2.3.1 Fixed effect variables

**Total identity:** Each consonant in the base is encoded for whether it is identical to the LC with non-identical being the reference level. For instance, a base with two consonants have two binary variables. This variable is a version of the previously proposed constraint such as

---

<sup>13</sup> The use of existing words raises the issue of lexicalization, which is addressed in Section 4.4 and argued to not be the case. Ultimately, even if this was a matter of lexicalization, our study would provide a more delicate measure of how much lexicalization affects people's sensitivity in their choice due to rating task we employ.



full reduplication or no-repeat which checks the total identity of  $C_2$  and the LC; however, it generalises across all consonants and not only  $C_2$ .

**Partial identity:** Each consonant in the base is encoded for whether the feature value of each of its phonological features is identical to that of the LC with non-identical being the reference level. It is important to note that only positive featural values are compared. Following the phonological system outlined in Erguvanlı Taylan (2015), eight binary consonantal features were used: Sonorant, voice, continuant, strident, anterior, coronal, labial, and nasal (see Appendix A for the feature chart).<sup>14</sup> The features, high, back, and lateral, were excluded because all of the four linking consonants have a negative value for these features. The partial identity was modelled using two approaches. The first approach is to allow each of the eight matched features to contribute differently by using them as individual variables. This approach would create eight binary variables for each of the consonants in the base form. We refer to these variables as *individual feature identity*. The second approach is to sum up the number of matched features, thus assuming that all features have the same weight. This would yield one continuous variable for each of the consonants in the base form. We refer to this variable as *sum featural identity*.

**Transitional phonotactic probability:** Demircan (1987) proposed speakers might avoid selecting a linking consonant that would lose or change the distinctive features of  $C_1$  due to the principle of least effort. For instance, consonantal sequences across syllables such as [p.b] might undergo devoicing, which makes the base less intelligible; therefore, [p] is unlikely to be selected in that context. Similarly, Wedel (2000) observed that plosives are generally dispreferred in Turkish phonotactics if followed by a homorganic consonant since they are articulatorily or perceptually marked. We confirmed this observation by conducting a corpus search using an online Turkish lexicon (TELL) (Inkelas et al., 2000).

Following Wedel (2000), lexical statistics were used to quantify the degree of junctural markedness. The assumption is that speakers are unlikely to produce consonant clusters that are articulatorily difficult or perceptually less distinctive. The token frequencies of all intervocalic heterosyllabic two-consonant clusters beginning with one of the four LCs ([Vp.CV], [Vm.CV], [Vs.CV], and [Vr.CV]) were extracted from a large subtitle-based corpus of Turkish. The written corpus was compiled using over 40,000 subtitle texts of Turkish. The corpus contains approximately 200 million word tokens and over 200,000 word types. The use of a subtitle corpus was motivated by the fact that lexical frequencies derived from subtitle texts have consistently shown to outperform those from other genres in capturing behavioural responses in psycholinguistic tasks across languages (Brysbaert & New, 2009; Keuleers et al. 2010; Tang,

---

<sup>14</sup> Following Erguvanlı Taylan (2015), the so-called ‘soft-g’ ğ was treated as a voiced velar fricative /ɣ/ in the calculation of partial identity. Since only nine out of 162 items contain the ‘soft-g’, we anticipate that a different methodological decision should have a negligible effect on our findings. We encourage readers to experiment with different feature systems using our data and scripts on [osf.io](https://osf.io).

2012; Tang & de Chene, 2014). The expectation is that the higher the transitional phonotactic probability (estimated using token frequency) of a juncture sequence, the higher the acceptability rating (Albright, 2007; Bailey & Hahn, 2001; Goldrick, 2011).

### 2.3.2. Coding illustration of identity variables

**Table 2** illustrates how the identity variables are coded for  $C_1$  [s] of the stimulus *sarı* ‘yellow’. The total identity of  $C_1$  matches with the LC [s] and not with the LCs [p,m,r]; therefore, the total identity variable with LC [s] is coded as 1 and the others are coded as 0s. Concerning the individual featural identity variables,  $C_1$  and the LC [r] both have a positive feature value for the features [continuant], [anterior], [cororal], therefore these features are coded as 1s, while the other features ([sonorant], [voice], [strident], [labial], and [nasal]) are coded as 0s. The sum featural identity variable for the LC [r] is the sum of the number of matched features which is 3 ([continuant], [anterior], and [cororal]).

<i>sarı</i> ‘yellow’			Matching $C_1$									
$C_1$	$C_2$	LC	son	voice	cont	strid	ant	cor	lab	nas	sum	total
s	r	p	0	0	0	0	1	0	0	0	1	0
s	r	m	0	0	0	0	1	0	0	0	1	0
s	r	s	0	0	1	1	1	1	0	0	4	1
s	r	r	0	0	1	0	1	1	0	0	3	0

**Table 2:** Illustration of identity variables of  $C_1$  for the stimulus *sarı* ‘yellow’ with each of the four linking consonants. Abbreviations: son: sonorant, cont: continuant, strid: strident, ant: anterior, cor: coronal, lab: labial, lat: lateral, and nas: nasal. The individual feature identity variables are son, voice, cont, strid, ant, cor, lab, and nas. sum stands for the sum featural identity; total stands for the total identity variable.

### 2.3.3. Random effect variables

As is typical of psycholinguistic research, participant and base form were included as random effects to allow for idiosyncrasies of individual participants and items.

Prior studies have established that only the consonants of the base trigger the identity avoidance effect in Turkish. Therefore, we analysed any base words that have the same number of consonants together. However, to recognise that these base words do in fact have different shapes (different number of syllables, and closed vs open syllables in different positions), we encoded the word shape of the base forms (such as CVC and CVCV, etc. as shown in **Table 3**) as a random effect to capture its potential effect. These models thus allow us to focus on the factors

of interest (i.e., feature mismatch).<sup>15</sup> The linking consonant of the reduplicated forms was also included as a random effect to capture the general preference for specific linking consonants.

## 2.4. Modelling procedure

The 162 items consist of 27 vowel-initial items and 135 consonant-initial items. Vowel-initial items were not analysed since they have an overwhelming preference for the LC [p] (Demircan, 1987; Kelepir, 1999; Sofu, 2005; Sofu & Altan, 2008, i.a.). This preference for [p] is supported by the descriptive statistics of the ratings of the 27 items. The mean and median ratings (on a scale of 1 to 7) are 6.402 and 7 for [p], 1.502 and 1 for [m], 1.533 and 1 for [s], and 1.197 and 1 for [r]. The 135 consonant-initial items were divided up into four groups based on the number of consonants they contain in the base and their word shapes as shown in **Table 3**. The five-consonant items were filtered out because there were only six of them and they might not provide enough statistical power for the analyses.

$C_1C_2$	$C_1C_2C_3$	$C_1C_2C_3C_4$	$C_1C_2C_3C_4C_5$
$C_1VC_2$ (23)	$C_1VC_2VC_3$ (37)	$C_1VC_2C_3VC_4$ (20)	$C_1VC_2VC_3C_4VC_5$ (5)
$C_1VC_2V$ (19)	$C_1VC_2C_3V$ (14)	$C_1VC_2VC_3VC_4$ (8)	$C_1VC_2VC_3VC_4C_5$ (1)
	$C_1VC_2C_3$ (4)	$C_1VC_2VC_3C_4V$ (1)	
	$C_1VC_2VC_3V$ (2)	$C_1VC_2C_3VC_4V$ (1)	

**Table 3:** A summary of the stimuli categorised by the number of consonants in the base and their word shapes. The number in parentheses indicates the number of items for each word shape.

### 2.4.1. Model specification and evaluation

Linear mixed-effects regression models were fit to the rating responses conducted using the *lme4* package in R (Bates et al., 2015; R Core Team, 2013). Following standard practice in regression modelling, the continuous variables were z-score normalised (e.g., Baayen, 2008, Sec. 2.2). Z-score normalization allows us to compare the relative strength of our continuous predictors directly. As per standard practice with token frequency, the transitional phonotactic probability was log-transformed (base 10) before z-score normalization. Our categorical predictors was summed (Wissmann et al., 2007) with non-identical as the base level.

The statistical significance of the individual predictors in all the models was evaluated by bootstrapping. Bootstrapping was carried out using the *bootmer* function in the *lme4* library. One

<sup>15</sup> Note that the inclusion of word shape as a random effect was only for the models for which word shape was not the variable of interest. In Study II, when examining the effect of word shape, we constructed a separate regression model for each word shape, therefore word shape was not included as a random effect.

thousand bootstrap simulations were performed for each model. Bootstrapped p-values and confidence intervals at 95% were computed for each predictor in each model. We follow the conventional alpha-level of 0.05 for significance. Model comparisons were performed using Akaike information criterion (AIC).<sup>16</sup> All models underwent the process of model criticism. For each model, the residuals were extracted and data points that were 2.5 standard deviations above or below the mean residual value were excluded. No more than 1% of the data points were excluded in any of the models.

To evaluate potential collinearity issues, we computed the Variance Inflation Factor (VIF) of the predictor variables in each of the models. The variables in all but four of the models have  $VIF < 5$ . These four models have in total 14 variables that have  $VIF > 5$  but  $< 10$ , and two variables have VIF slightly above 10. These values are mostly below the typical critical values of 5 or 10 (Chatterjee & Hadi, 2015; Tomaschek et al., 2018), which indicates no serious issues of collinearity.<sup>17</sup>

In this study, we make use of zero-order correlations to address the issues of interpretability due to potential collinearity, even if the potential of serious collinearity is low as suggested by the VIF analyses.<sup>18</sup> The correlation amongst the identity predictors could still cause their effects to be counterintuitive and hard to interpret. Collinearity between two predictors can cause the reduction or sign reversal in one of the model estimate. For instance, some of the identity predictors might behave in the opposite direction of identity avoidance with their regression coefficients being positive. One diagnostic of a suppressor effect is whether the model estimate is in the same or opposite direction as the correlation between the dependent and independent variable. Model estimates in the opposite direction of the correlation suggest a suppressor effect. Follow-up inspections were performed by examining the pairwise zero-order association between each of individual predictors and the response variable. Pairwise zero-order association is to estimate the effect of individual independent variables has on the dependent variable (see Appendix D for discussion). This was done by fitting multiple mixed-effects regression models with only one independent variable at a time. A regression model with the same random effect

---

<sup>16</sup> Model comparisons were also performed using Bayesian information criterion (BIC). The penalty term for the number of parameters is larger in BIC than in AIC. BIC is useful when the number of parameters between the two models being compared is particular different. All the results were the same using AIC or BIC, therefore only the results using AIC were reported.

<sup>17</sup> Two of these four models have only two variables with a VIF above 5 but below 6 with the maximum VIF of 5.79 and 5.87. One of the models has five variables with a VIF above 5 but below 10 and one variable with a VIF of 10.2. The remaining model has seven variables with a VIF above 5 but below 10 and one variable with a VIF of 10.5. The variables associated with these higher VIFs are all individual featural identity variables. It is worth noting that recommendations of VIF's cutoff vary depending on the literature (e.g., Hepworth et al. (2007) suggest 4 to be the cutoff). It is not clear what a meaningful boundary is for a low versus a high value. No fixed set of guidelines can guarantee the correct analysis of collinear data (Tomaschek et al., 2018). We acknowledge that the best practice of dealing with collinearity has not been established.

<sup>18</sup> We have made our data and analysis scripts available on [osf.io](https://osf.io) and we encourage readers to evaluate the data and our procedures themselves and to examine with different statistical techniques.

structure as the above models was fitted with only each of the predictors for each of the three groups of base forms. A null model with no fixed effect variables was fitted to compare with each of these models with one fixed effect variable. The drop in AIC values was used as a measure of the importance of each variable ( $AIC_{null} - AIC_{superset}$ ). A drop in AIC of more than 2 indicates statistical significance. When interpreting the direction of the identity effects in the full models, these pairwise associations would assist us in identifying cases of a sign reversal due to collinearity.

The distribution of the variables (both the response variable and the predictors) for each of the three item groups with two, three, and four consonants in the base form can be found in Appendix C. Means and standard deviations of the by-participant standardised ratings of each item (a base form with one of the four linking consonants) can be found in Appendix H. The complete report of the statistical analyses (regression tables, model comparisons, and figures) can be found on the osf.io repository (see Section 5: *Data accessibility statement*).

#### 2.4.2 Study I: Feature specificity

Study I focuses on three groups of items: 42 two-consonant base forms, 57 three-consonant base forms, and 30 four-consonant base forms. The analyses were conducted separately for each of the three item groups. Two full models were initially fitted. The two models differ in the type of partial identity variables. One model has individual feature identity variables (eight binary variables per consonant in the base), while the other has sum featural identity variables which are the number of identical features and are computed by summing up the number of identical features, thus assuming that all features have the same weight (one continuous variable per consonant in the base). To assess the level of feature specificity, a series of model comparisons was performed by removing each type of identity variables in bulk. Three more subset models were therefore fitted: a) A model with total identity variables without partial identity variables, b) a model with individual featural identity variables without total identity variables, c) a model with sum featural identity variables without total identity variables. These models were fitted with the predictor variables outlined in Section 2.3.1 as fixed effects and four random intercepts with the variables outlined in Section 2.3.3.

The regression structures of the two full models are shown below. Note that the identity variables have  $C_i$  in parentheses and the index  $i$  is referring to a specific consonant in the base form. If the base form contains  $N$  consonants, then there would be  $N$  sets of identity variables.<sup>19</sup>

Model with total identity and partial identity using individual features:

$$\text{Rating} \sim \text{Total identity } (C_i) + \text{Sonorant identity } (C_i) + \text{Voice identity } (C_i) + \text{Continuant identity } (C_i) + \text{Strident identity } (C_i) + \text{Coronal identity } (C_i) + \text{Labial identity } (C_i) + \text{Nasal identity } (C_i) + \text{Transitional phonotactic probability} + (1 \mid \text{Participant}) + (1 \mid \text{Base form}) + (1 \mid \text{Word shape}) + (1 \mid \text{Linking consonant})$$


---

<sup>19</sup> See Section 2.2 for the rationale behind the modelling approach.

Model with total identity and partial identity using the sum of the matched features:

$$\text{Rating} \sim \text{Total identity } (C_i) + \text{Sum featural identity } (C_i) + \text{Transitional phonotactic probability} + (1 \mid \text{Participant}) + (1 \mid \text{Base form}) + (1 \mid \text{Word shape}) + (1 \mid \text{Linking consonant})$$

### 2.4.3. Study II: Locality: Distance-based decay and syllable role

Study II consists of two analyses. The first analysis aims to address the importance of the consonants beyond  $C_2$  (namely  $C_3$  and  $C_4$ ). The second analysis focuses on examining how the identity avoidance effect would be affected if the consonant in the base matches the constituency of the linking consonant.

In the first analysis, the initial models were the best models found in Study I using model comparisons. Model comparisons were performed by removing identity variables (total and partial) that are associated with each consonant position in bulk. The drop in AIC values were used as a measure of the importance of the consonant ( $AIC_{\text{subset}} - AIC_{\text{superset}}$ ).

In the second analysis, the best model structure in the first analysis was fitted over base forms with each of the word shapes separately without the random variable (word shape). The five most frequent word shapes were selected since they have a relatively higher number of base forms (at least 14) to enhance the generalisability of our findings. The structures are  $C_1VC_2$ ,  $C_1VC_2VC_3$ ,  $C_1VC_2C_3VC_4$ ,  $C_1VC_2V$ , and  $C_1VC_2C_3V$ . Identity variables (total and partial) that are associated with each consonant position were dropped in bulk and the drop in AIC values were computed. The variable importance values between each of the consonants within each word shape were compared to enable an examination of how syllabification of the consonants plays a role, specifically whether the consonant in the base matches the constituency of the linking consonant.

## 3. Results

### 3.1. Study I: Feature specificity

In the following sections, we present the results from the model comparisons of different levels of feature specificity in Section 3.1.1. Given the model comparisons, the selected best models were then evaluated further individually for each of the three item groups. To enhance the interpretability of each of the predictors in the best models, all pairwise associations between the response and each of the predictors were computed and can be found in Appendix D. Finally, we report the detailed model evaluations of each of the predictors in Sections 3.1.2, 3.1.3, and 3.1.3 for the two-consonant, three-consonant, and four-consonant groups, respectively. The regression tables of all of the models can be found on the osf.io repository (see Section 5: *Data accessibility statement*).

### 3.1.1. Model comparison

To evaluate the level of specificity of our identity variables, five models were fitted and evaluated for their AIC levels for each of the three item groups (base forms with either two ( $C_1C_2$ ), three ( $C_1C_2C_3$ ) or four ( $C_1C_2C_3C_4$ ) consonants). These five models include the two full models (total identity and sum featural identity; total identity and individual featural identity) and three subset models without either the total identity variables or the partial identity variables (total identity, sum featural identity, individual featural identity). The AIC values of all models are summarised in Table 15 (See Section F in the Appendix).

The model structure with both the total identity variables and the individual featural identity variables consistently yielded lower AICs (the best model fit) across the three item groups (two-consonant: 28853.63, three-consonant: 39131.32, and four-consonant: 20298.78). This finding suggests that both total identity and partial identity play a role in identity avoidance in Turkish partial reduplication. This supports many of the previous analyses which take into account of both total and partial identity avoidance, for instance, Demircan (1987)'s observations that the linking consonant (LC) cannot be identical with any of the consonants in the base and the features of the LC should not be identical to the  $C_2$  of the base. This finding is also in line with how both total and partial identity influence consonant co-occurrence patterns within lexical roots (Gallagher & Coon, 2009; Graff & Jaeger, 2009).

Given that both total identity and partial identity are important, their relative importance is also examined. Their relative importance can be evaluated by comparing the drop in AIC values when either of these variable types were dropped from a full model. This was computed separately from the two full model structures (total identity and sum featural identity; total identity and individual featural identity). The drop in AIC is summarised in Table 14 (See Section F in the Appendix). The individual featural identity has a bigger drop in AIC than total identity in the two-consonant group (1016.18 vs 87.07), the three consonant-group (1791.52 vs 291.94), and the four-consonant group (1094.96 vs 249.08). Similarly, the sum featural identity has a bigger drop in AIC than total identity in the two-consonant group (461.79 vs 150.88), the three consonant-group (1028.13 vs 376.39), and the four-consonant group (473.5 vs 339.42). This finding suggests that partial featural identity (sum or individual) has a stronger effect on the linking consonant than total identity. This, as far as we know, has not been formally established in previous studies of Turkish partial reduplication.

### 3.1.2. Model evaluation: Two-consonant group

The fixed and random effects estimates of the two-consonant model are summarised in Table 4 and Table 11 (see Appendix E), respectively. First of all, the transitional phonotactic probability was not statistically significant ( $p = 0.514$ ). We turn now to the identity variables, starting with those in  $C_1$ . All but the continuant identity variable ( $p = 0.750$ ) and the anterior identity

variable ( $p = 0.082$ ) were statistically significant. An identity avoidance effect was found in five of the seven significant identity variables with a negative coefficient – total identity ( $\hat{\beta} = -1.6748$ ), sonorant identity ( $\hat{\beta} = -1.4779$ ), strident identity ( $\hat{\beta} = -1.7381$ ), coronal identity ( $\hat{\beta} = -0.6159$ ), and labial identity ( $\hat{\beta} = -1.2808$ ). The remaining two significant variables, voice identity ( $\hat{\beta} = 1.3000$ ) and nasal identity ( $\hat{\beta} = 2.2870$ ), have a positive coefficient, suggesting the opposite effect of identity avoidance. To clarify these two variables, we turn to the pairwise

		$\hat{\beta}$	SE	$t$	CI <sub>Lower 95%</sub>	CI <sub>Upper 95%</sub>	$p_{\text{Bootstrapped}}$
	(Intercept)	-0.7366	0.5939	-1.2402	-1.9217	0.4272	0.252
C <sub>1</sub>	<b>Total identity</b>	-1.6748	0.2209	-7.5826	-2.1176	-1.2124	<.001***
	<b>Sonorant identity</b>	-1.4779	0.2098	-7.0449	-1.8832	-1.0523	<.001***
	<b>Voice identity</b>	1.3000	0.1485	8.7545	1.0002	1.5936	<.001***
	Continuant identity	0.0378	0.1297	0.2915	-0.2043	0.2868	0.750
	<b>Strident identity</b>	-1.7381	0.2278	-7.6310	-2.2085	-1.2861	<.001***
	Anterior identity	0.5291	0.2961	1.7870	-0.0577	1.1192	0.082
	<b>Coronal identity</b>	-0.6159	0.1130	-5.4507	-0.8271	-0.3867	<.001***
	<b>Labial identity</b>	-1.2808	0.1327	-9.6542	-1.5404	-1.0173	<.001***
	<b>Nasal identity</b>	2.2870	0.3653	6.2600	1.5439	3.0079	<.001***
C <sub>2</sub>	<b>Total identity</b>	0.7255	0.1217	5.9596	0.4904	0.9635	<.001***
	<b>Sonorant identity</b>	-0.7207	0.1258	-5.7281	-0.9809	-0.4719	<.001***
	<b>Voice identity</b>	-1.7778	0.1211	-14.6801	-2.0168	-1.5447	<.001***
	<b>Continuant identity</b>	-0.8210	0.0996	-8.2377	-1.0144	-0.6254	<.001***
	<b>Strident identity</b>	-1.2272	0.1399	-8.7738	-1.4799	-0.9418	<.001***
	Anterior identity	0.3377	0.2854	1.1833	-0.2387	0.9052	0.262
	Coronal identity	-0.1035	0.1054	-0.9810	-0.3146	0.1180	0.336
	<b>Labial identity</b>	-3.6485	0.2078	-17.5598	-4.0599	-3.2392	<.001***
	Nasal identity	0.1762	0.2763	0.6375	-0.3666	0.7075	0.504
	Transitional phonotactic probability	-0.0408	0.0594	-0.6859	-0.1596	0.0773	0.514

**Table 4:** Fixed effects summary for Study I (two-consonant base forms).  $\hat{\beta}$ : coefficient; SE: standard error;  $t$ : t-value; CI<sub>Lower 95%</sub> and CI<sub>Upper 95%</sub>: 95% confidence intervals of the coefficient from bootstrapping;  $p_{\text{Bootstrapped}}$ : p-value from bootstrapping simulations. Significant variables ( $p \leq 0.05$ ) are in bold.

Number of observations: 6,883; number of participants: 207; number of base forms: 42; number of word shapes: 2; number of linking consonants: 4.

Level of significance:  $\cdot$  ( $p \leq 0.1$ ), \* ( $p \leq 0.05$ ), \*\* ( $p \leq 0.01$ ), \*\*\* ( $p \leq 0.001$ ).



association analysis in **Table 10**. The voice identity variable has a genuine identity *preference* effect ( $\hat{\beta} = 1.3000$ ) since the zero-order association is also positive (**Table 10**,  $\hat{\beta} = 0.4045$ ). The nasal identity variable shows a suppressor effect since the coefficient is positive ( $\hat{\beta} = 2.2870$ ) even though the zero-order association is negative (**Table 10**,  $\hat{\beta} = -2.2434$ ).

We now focus on the identity variables in  $C_2$ . All but three variables were statistically significant. The insignificant variables are anterior identity ( $p = 0.262$ ), coronal identity ( $p = 0.336$ ), and nasal identity ( $p = 0.504$ ). An identity avoidance effect was found in five of the six significant identity variables with a negative coefficient – sonorant identity ( $\hat{\beta} = -0.7207$ ), voice identity ( $\hat{\beta} = -1.7778$ ), continuant identity ( $\hat{\beta} = -0.8210$ ), strident identity ( $\hat{\beta} = -1.2272$ ), and labial identity ( $\hat{\beta} = -3.6485$ ). The total identity variable shows a suppressor effect since the coefficient is positive ( $\hat{\beta} = 0.7255$ ) even though the zero-order association is negative (**Table 10**,  $\hat{\beta} = -0.8473$ ).

### 3.1.3. Model evaluation: Three-consonant group

The fixed and random effects estimates of the three-consonant model are summarised in **Table 5** and **Table 12** (see Appendix E), respectively. First of all, the transitional phonotactic probability was statistically significant in the positive direction ( $\hat{\beta} = 0.1821$ ,  $p = 0.514$ ), suggesting a preference for an articulatorily or perceptually unmarked heterosyllabic cluster. We turn now to the identity variables, starting with those in  $C_1$ . An identity avoidance effect was found in all of the five significant identity variables with a negative coefficient – total identity ( $\hat{\beta} = -1.8850$ ), voice identity ( $\hat{\beta} = -0.4231$ ), continuant identity ( $\hat{\beta} = -0.5254$ ), strident identity ( $\hat{\beta} = -1.6904$ ), and labial identity ( $\hat{\beta} = -2.2846$ ). The remaining three variables were not significant and they are sonorant identity ( $p = 0.312$ ), anterior identity ( $p = 0.350$ ), and coronal identity ( $p = 0.092$ ). Nasal identity was excluded in this model because the base forms have no nasals in  $C_1$ .

We now focus on the identity variables in  $C_2$ . All but one variable were statistically significant. The insignificant variable is sonorant identity ( $p = 0.118$ ). An identity avoidance effect was found in six of the eight significant identity variables with a negative coefficient – voice identity ( $\hat{\beta} = -2.0183$ ), continuant identity ( $\hat{\beta} = -0.6266$ ), strident identity ( $\hat{\beta} = -1.3080$ ), coronal identity ( $\hat{\beta} = -1.1598$ ), labial identity ( $\hat{\beta} = -2.2197$ ), and nasal identity ( $\hat{\beta} = -1.0419$ ). The total identity variable shows a suppressor effect since the coefficient is positive ( $\hat{\beta} = 1.2713$ ) even though the zero-order association is negative (**Table 10**,  $\hat{\beta} = -0.9739$ ). The positive coefficient ( $\hat{\beta} = 0.4863$ ) of the anterior identity variable is unlikely to be genuine because its level of significance is weak with a *p-value* of 0.024 and while the zero-order association is also positive ( $\hat{\beta} = 0.0160$ ), it was insignificant with a small effect size (a featural match increases the rating by only 0.016 on a scale from 1 to 7).

We now turn to the identity variables in  $C_3$ . All but the total identity variable ( $p = 0.1$ ) and the anterior identity variable ( $p = 0.1$ ) were statistically significant. An identity avoidance effect

		$\hat{\beta}$	SE	<i>t</i>	CI <sub>Lower 95%</sub>	CI <sub>Upper 95%</sub>	<i>P</i> <sub>Bootstrapped</sub>
	(Intercept)	-2.6151	0.4537	-5.7641	-3.5669	-1.7126	<.001***
C <sub>1</sub>	<b>Total identity</b>	-1.8850	0.1361	-13.8534	-2.1466	-1.6117	<.001***
	Sonorant identity	0.1596	0.1492	1.0697	-0.1345	0.4529	0.312
	<b>Voice identity</b>	-0.4231	0.1118	-3.7850	-0.6471	-0.2043	<.001***
	<b>Continuant identity</b>	-0.5254	0.0933	-5.6289	-0.7022	-0.3440	<.001***
	<b>Strident identity</b>	-1.6904	0.1259	-13.4288	-1.9466	-1.4399	<.001***
	Anterior identity	-0.2035	0.2223	-0.9155	-0.6491	0.2353	0.350
	Coronal identity	0.1700	0.0971	1.7504	-0.0237	0.3658	0.092
	<b>Labial identity</b>	-2.2846	0.1264	-18.0688	-2.5308	-2.0290	<.001***
	Nasal identity	-	-	-	-	-	-
	C <sub>2</sub>	<b>Total identity</b>	1.2713	0.1158	10.9754	1.0434	1.5075
Sonorant identity		0.1525	0.0980	1.5552	-0.03312	0.3422	0.118
<b>Voice identity</b>		-2.0183	0.1149	-17.5622	-2.2374	-1.8008	<.001***
<b>Continuant identity</b>		-0.6266	0.0951	-6.5911	-0.8161	-0.4352	<.001***
<b>Strident identity</b>		-1.3080	0.1154	-11.3340	-1.5395	-1.0837	<.001***
<b>Anterior identity</b>		0.4863	0.2242	2.1691	0.0412	0.9340	0.024*
<b>Coronal identity</b>		-1.1598	0.0976	-11.8878	-1.3522	-0.9646	<.001***
<b>Labial identity</b>		-2.2197	0.1142	-19.4466	-2.4463	-1.9995	<.001***
<b>Nasal identity</b>		-1.0419	0.1487	-7.0058	-1.3387	-0.7497	<.001***
C <sub>3</sub>		Total identity	-0.2254	0.1453	-1.5512	-0.5064	0.0480
	<b>Sonorant identity</b>	-1.3457	0.1424	-9.4516	-1.6230	-1.0667	<.001***
	<b>Voice identity</b>	0.9853	0.1337	7.3689	0.7202	1.2473	<.001***
	<b>Continuant identity</b>	0.2059	0.0807	2.5491	0.0438	0.3747	0.016*
	<b>Strident identity</b>	-1.2153	0.1232	-9.8688	-1.4585	-0.9716	<.001***
	Anterior identity	0.3722	0.2261	1.6464	-0.07387	0.8059	0.1
	<b>Coronal identity</b>	-0.8877	0.0899	-9.8803	-1.0608	-0.7172	<.001***
	<b>Labial identity</b>	0.8738	0.1722	5.0740	0.5124	1.2165	<.001***
	<b>Nasal identity</b>	-1.6972	0.1410	-12.0394	-1.9702	-1.4165	<.001***
		<b>Transitional pho- notactic probability</b>	0.1821	0.0420	4.3398	0.1017	0.2628

**Table 5:** Fixed effects summary for Study I (three-consonant base forms).  $\hat{\beta}$ : coefficient; SE: standard error; *t*: t-value; CI<sub>Lower 95%</sub> and CI<sub>Upper 95%</sub>: 95% confidence intervals of the coefficient from bootstrapping; *P*<sub>Bootstrapped</sub>: p-value from bootstrapping simulations. Significant variables ( $p \leq 0.05$ ) are in bold.

Number of observations: 9,374; number of participants: 207; number of base forms: 57; number of word shapes: 4; number of linking consonants: 4.

Level of significance: · ( $p \leq 0.1$ ), \* ( $p \leq 0.05$ ), \*\* ( $p \leq 0.01$ ), \*\*\* ( $p \leq 0.001$ ).

was found in four of the seven significant identity variables with a negative coefficient – sonorant identity ( $\hat{\beta} = -1.3457$ ), strident identity ( $\hat{\beta} = -1.2153$ ), coronal identity ( $\hat{\beta} = -0.8877$ ), and nasal identity ( $\hat{\beta} = -1.6972$ ). The voice identity variable and the labial identity variable both show a suppressor effect since the coefficients are positive (voice:  $\hat{\beta} = 0.9853$ ; labial:  $\hat{\beta} = 0.8738$ ) even though the zero-order associations are negative (Table 10, voice:  $\hat{\beta} = -0.5918$  and labial:  $\hat{\beta} = -1.5442$ ). The continuant identity variable potentially has a genuine identity preference effect ( $\hat{\beta} = 0.2059$ ) since the zero-order association is also positive but only near-significant (Table 10,  $\hat{\beta} = 0.1514$ ,  $\Delta\text{AIC} = 1.63$ ).

### 3.1.4. Model evaluation: Four-consonant group

The fixed and random effects estimates of the four-consonant model are summarised in Table 6 and Table 13 (see Appendix E), respectively. First of all, the transitional phonotactic probability was statistically significant in the expected positive direction ( $\hat{\beta} = 0.4553$ ). We turn now to the identity variables, starting with those in  $C_1$ . Only three identity variables were significant and they are all in the negative direction, suggesting an identity avoidance effect – total identity ( $\hat{\beta} = -3.1510$ ), strident identity ( $\hat{\beta} = -0.8103$ ), and labial identity ( $\hat{\beta} = -3.0314$ ). Nasal identity was excluded in this model because the base forms have no nasals in  $C_1$ .

We now focus on the identity variables in  $C_2$ . All but one variable were statistically significant. The insignificant variable is nasal identity ( $p = 0.73$ ). An identity avoidance effect was found in four of the eight significant identity variables with a negative coefficient – voice identity ( $\hat{\beta} = -2.7016$ ), continuant identity ( $\hat{\beta} = -1.5544$ ), coronal identity ( $\hat{\beta} = -0.8464$ ), and labial identity ( $\hat{\beta} = -2.1668$ ). The total identity variable and the sonorant identity variable both show a suppressor effect since the coefficients are positive (total:  $\hat{\beta} = 1.1706$ ; sonorant:  $\hat{\beta} = 0.4956$ ) even though the zero-order associations are negative (Table 10, total:  $\hat{\beta} = -0.2840$  and sonorant:  $\hat{\beta} = -0.5162$ ). The strident identity variable also shows a suppressor effect since the coefficient is negative ( $\hat{\beta} = -0.8427$ ) even though the zero-order association is positive (Table 10,  $\hat{\beta} = 0.4654$ ). The anterior identity variable has a genuine identity preference effect ( $\hat{\beta} = 1.3207$ ) since the zero-order association is also positive (Table 10,  $\hat{\beta} = 0.4449$ ).

We now turn to the identity variables in  $C_3$ . All but the anterior identity ( $p = 0.618$ ) were statistically significant. An identity avoidance effect was found in four of the eight significant identity variables with a negative coefficient – total identity ( $\hat{\beta} = -0.7889$ ), voice identity ( $\hat{\beta} = -0.6932$ ), labial identity ( $\hat{\beta} = -0.7807$ ), and nasal identity ( $\hat{\beta} = -1.1454$ ). The strident identity variable shows a suppressor effect since the coefficient is negative ( $\hat{\beta} = -2.3943$ ) even though the zero-order association is positive (Table 10,  $\hat{\beta} = 0.9814$ ). Two of the significant identity variables (continuant and coronal identity variables) show a genuine identity preference effect since the coefficients are positive (continuant:  $\hat{\beta} = 0.8278$ ; coronal:  $\hat{\beta} = 0.4544$ ) and the zero-order associations are also positive (Table 10, continuant:  $\hat{\beta} = 1.2100$  and coronal:  $\hat{\beta} = 0.6786$ ).

		$\hat{\beta}$	SE	<i>t</i>	CI <sub>Lower 95%</sub>	CI <sub>Upper 95%</sub>	<i>p</i> <sub>Bootstrapped</sub>
	(Intercept)	-5.4181	0.9340	-5.8006	-7.2647	-3.5117	<.001***
C <sub>1</sub>	<b>Total identity</b>	-3.1510	0.2061	-15.2876	-3.5550	-2.7555	<.001***
	Sonorant identity	-0.4286	0.2290	-1.8715	-0.8854	0.0268	0.07
	Voice identity	0.4343	0.2426	1.7906	-0.0347	0.9120	0.074
	Continuant identity	0.0156	0.1672	0.0930	-0.3133	0.3449	0.942
	<b>Strident identity</b>	-0.8103	0.2130	-3.8041	-1.2120	-0.3986	<.001***
	Anterior identity	0.3652	0.4687	0.7791	-0.5533	1.2584	0.420
	Coronal identity	-0.0721	0.1758	-0.4098	-0.4376	0.2848	0.714
	<b>Labial identity</b>	-3.0314	0.2253	-13.4560	-3.4557	-2.5791	<.001***
	Nasal identity	-	-	-	-	-	-
C <sub>2</sub>	<b>Total identity</b>	1.1706	0.1662	7.0455	0.8619	1.4836	<.001***
	<b>Sonorant identity</b>	0.4956	0.1743	2.8434	0.1621	0.8409	0.006**
	<b>Voice identity</b>	-2.7016	0.2063	-13.0954	-3.1069	-2.2986	<.001***
	<b>Continuant identity</b>	-1.5544	0.1463	-10.6265	-1.8314	-1.2641	<.001***
	<b>Strident identity</b>	-0.8427	0.1906	-4.4205	-1.2156	-0.4716	<.001***
	<b>Anterior identity</b>	1.3207	0.5778	2.2855	0.1535	2.4554	0.016*
	<b>Coronal identity</b>	-0.8464	0.2171	-3.8995	-1.2629	-0.4309	<.001***
	<b>Labial identity</b>	-2.1668	0.3002	-7.2173	-2.7372	-1.5935	<.001***
	Nasal identity	-0.0815	0.2503	-0.3257	-0.5714	0.4146	0.73
C <sub>3</sub>	<b>Total identity</b>	-0.7889	0.1881	-4.1947	-1.1574	-0.4198	<.001***
	<b>Sonorant identity</b>	1.5489	0.2462	6.2909	1.0774	2.0444	<.001***
	<b>Voice identity</b>	-0.6932	0.2390	-2.9008	-1.1736	-0.2361	0.002**
	<b>Continuant identity</b>	0.8278	0.2077	3.9866	0.4211	1.2334	<.001***
	<b>Strident identity</b>	-2.3943	0.2254	-10.6203	-2.8504	-1.9424	<.001***
	Anterior identity	0.2458	0.4789	0.5131	-0.7456	1.2054	0.618
	<b>Coronal identity</b>	0.4544	0.1658	2.7407	0.1197	0.7935	0.012*
	<b>Labial identity</b>	-0.7807	0.2630	-2.9688	-1.3117	-0.2481	0.006**
	<b>Nasal identity</b>	-1.1454	0.2120	-5.4039	-1.5793	-0.7288	<.001***
C <sub>4</sub>	<b>Total identity</b>	-0.9657	0.2389	-4.0422	-1.4505	-0.4703	<.001***
	<b>Sonorant identity</b>	1.9742	0.2551	7.7381	1.4611	2.4885	<.001***
	<b>Voice identity</b>	-0.6780	0.2347	-2.8890	-1.1401	-0.2026	0.004**
	Continuant identity	0.1308	0.2288	0.5718	-0.3241	0.6022	0.616
	<b>Strident identity</b>	-3.4709	0.2483	-13.9762	-3.9699	-2.9808	<.001***
	Anterior identity	-0.1093	0.4723	-0.2314	-1.0578	0.8355	0.806
	Coronal identity	0.2520	0.1974	1.2768	-0.1212	0.6165	0.184
	Labial identity	-0.3079	0.2828	-1.0886	-0.8769	0.2499	0.296
	<b>Nasal identity</b>	-1.7796	0.2056	-8.6571	-2.1915	-1.3756	<.001***
	<b>Transitional pho- notactic probability</b>	0.4553	0.0803	5.6728	0.3006	0.6109	<.001***

**Table 6:** Fixed effects summary for Study I (four-consonant base forms).  $\hat{\beta}$ : coefficient; SE: standard error; *t*: t-value; CI<sub>Lower 95%</sub> and CI<sub>Upper 95%</sub>: 95% confidence intervals of the coefficient from bootstrapping; *p*<sub>Bootstrapped</sub>: p-value from bootstrapping simulations. Significant variables ( $p \leq 0.05$ ) are in bold. Number of observations: 4,900; number of participants: 207; number of base forms: 30; number of word shapes: 4; number of linking consonants: 4. Level of significance: · ( $p \leq 0.1$ ), \* ( $p \leq 0.05$ ), \*\* ( $p \leq 0.01$ ), \*\*\* ( $p \leq 0.001$ ).

The positive coefficient ( $\hat{\beta} = 1.548$ ) of the sonorant identity variable is unlikely to be genuine because while the zero-order association is also positive ( $\hat{\beta} = 0.0436$ ), it was insignificant with a small effect size (a featural match increases the rating by only 0.0436 on a scale from 1 to 7).

Finally we turn to the identity variables in  $C_4$ . Four of the variables were not statistically significant – continuant identity ( $p = 0.616$ ), anterior identity ( $p = 0.806$ ), coronal identity ( $p = 0.184$ ), and labial identity ( $p = 0.296$ ). Five of the variables were statistically significant. An identity avoidance effect was found in four of the five significant identity variables with a negative coefficient – total identity ( $\hat{\beta} = -0.9657$ ), voice identity ( $\hat{\beta} = -0.6780$ ), strident identity ( $\hat{\beta} = -3.4709$ ), and nasal identity ( $\hat{\beta} = -1.7796$ ). The sonorant identity variable has a genuine identity preference effect ( $\hat{\beta} = 1.9742$ ) since the zero-order association is also positive (Table 10,  $\hat{\beta} = 0.2887$ ).

### 3.2. Study II: Locality: Distance-based decay and syllable role

In the following sections, we examine which consonants play a role in the identity avoidance effect and whether a distance-based decay effect can be found through a series of model comparisons in Section 3.2.1. We then repeat the same analyses with the five most frequent word shapes in our dataset in Section 3.2.2 to examine the role of syllable position in the sense of Bennett (2012) and Rose and Walker (2004). The regression tables of all of the models can be found on the *osf.io* repository (see Section 5: *Data accessibility statement*).

#### 3.2.1. Identity avoidance beyond $C_2$

To examine whether the identity avoidance effect operates beyond  $C_2$ , model comparisons were performed by comparing a full model with identity variables from all consonants with models without any identity variables of a specific consonant. The drop in AIC was used as a measure of variable importance.

All  $\Delta$ AIC values are above 2 (a typical significance threshold for AIC values), therefore the identity avoidance effect operates over every single consonant, including  $C_3$  and  $C_4$ . This is a surprising finding since most of the previous studies did not find an effect from beyond  $C_2$  (e.g., Keleşir, 2000; Wedel, 1999, 2000).

To examine whether there is a distance-based decay effect such that the importance of each consonant decreases as the distance from the linking consonant increases, we compare the relative level of importance across consonant positions. In the two-consonant group,  $C_2$  ( $\Delta$ AIC = 872.30) is more important than  $C_1$  ( $\Delta$ AIC = 803.71). In the three-consonant group,  $C_1$  ( $\Delta$ AIC = 1679.46) is the most important consonant, followed by  $C_2$  ( $\Delta$ AIC = 1201.62), and  $C_3$  ( $\Delta$ AIC = 549.03). In the three-consonant group,  $C_1$  ( $\Delta$ AIC = 1679.46) is the most important consonant, followed by  $C_2$  ( $\Delta$ AIC = 1201.62) which in turn is more important than  $C_3$  ( $\Delta$ AIC = 549.03). In the four-consonant

group, the consonants from the most important to the least are  $C_1$  ( $\Delta AIC = 628.49$ ),  $C_2$  ( $\Delta AIC = 379.64$ ),  $C_4$  ( $\Delta AIC = 370.35$ ), and  $C_3$  ( $\Delta AIC = 179.20$ ). The  $\Delta AIC$ s of  $C_3$  and  $C_4$  are lower than those of  $C_2$  and  $C_1$ . However, the exact rankings do not clearly suggest a distance-based decay effect. While a distance-based decay effect can be seen with the three-consonant group ( $C_1 > C_2 > C_3$ ), the order of importance diverges with the other two groups – 1)  $C_4$  of the four-consonant group is more important than  $C_3$ , and 2)  $C_2$  of the two-consonant group is more important than  $C_1$ . We speculate that the divergence is attributable to the syllable position of the consonants. In the next section we examine the effect of syllable position by separating the items by word shape.<sup>20</sup>

### 3.2.2. Syllable position

The method of evaluating variable importance of each consonant remains the same as section 3.2.1. The difference is that in Section 3.2.1, the data was divided up by the number of consonants the items contain, while in the current section, the data was divided up by the five word shapes to enable an examination of the syllable position. These five word shapes are  $C_1VC_2V$ ,  $C_1VC_2$ ,  $C_1VC_2C_3V$ ,  $C_1VC_2VC_3$ , and  $C_1VC_2C_3VC_4$ . The word shapes with the same number of consonants were compared.  $C_1VC_2$  and  $C_1VC_2V$  were compared because  $C_2$  is a coda consonant in  $C_1VC_2$  but an onset consonant in  $C_1VC_2V$ . Similarly,  $C_1VC_2C_3V$  and  $C_1VC_2VC_3$  were compared because  $C_2$  and  $C_3$  are syllabified differently across the two shapes. Finally,  $C_1VC_2C_3VC_4$  was examined because it could reveal the combined effect of syllable position and the distance decay effect, since it contains two sets of onsets ( $C_1$  and  $C_3$ ) and codas ( $C_2$  and  $C_4$ ) which differ only in their distance from the linking consonant.

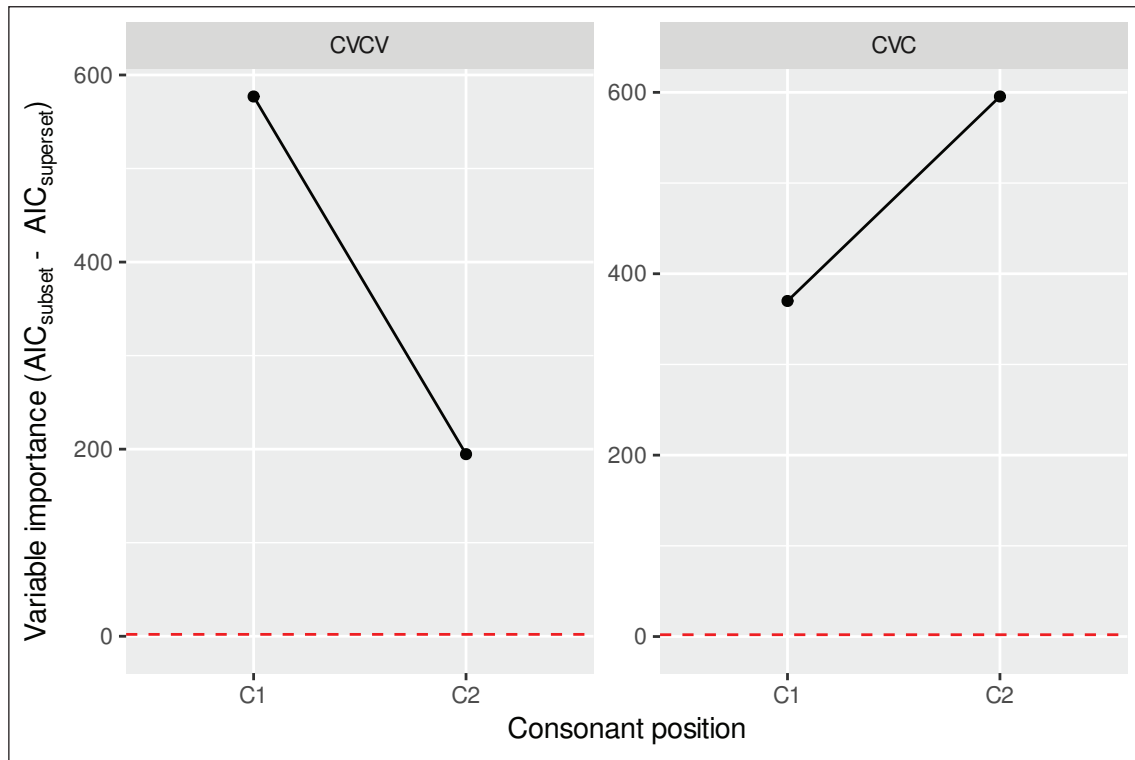
AIC is being used only to compare models with the same set of data. For example, a fully-specified model was fitted over the subset of the data with the  $C_1VC_2V$  items. To evaluate the importance of  $C_1$ , we fitted a model on the same subset of the data but without any identity variables of  $C_1$  (i.e., the model regression structure has the identity variables of  $C_2$  and not of  $C_1$ ). AIC was computed for both of these models (the initial model with both identity variables of  $C_1$  and  $C_2$ , and the new model with only identity variables of  $C_2$ ). The two AIC values were used to compare these two models.

The word shapes with the same number of consonants were compared and visualised. **Figure 1** visualises the variable importance of each consonant in  $C_1VC_2V$  and  $C_1VC_2$  base forms.  $C_1VC_2V$  shows a distance-based decay effect with  $C_2$  being less important than  $C_1$ , while the  $C_1VC_2$  shows the opposite pattern with  $C_1$  being less important than  $C_2$ . These differences indicate that the divergence observed in the two-consonant group (as described in Section 3.2.1) was driven

---

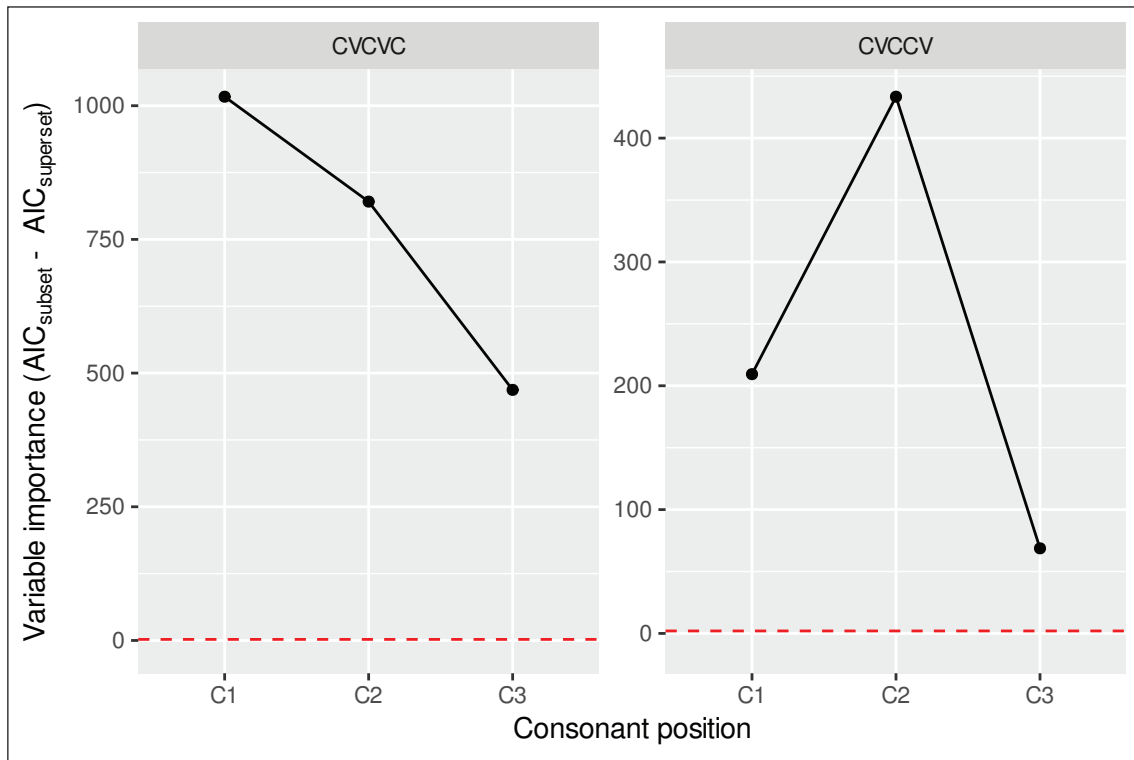
<sup>20</sup> We acknowledge that one could test the syllable position effect more directly by comparing two models, one with identity avoidance predictors and make references to the syllable constituency of their corresponding consonant, and one without making such references. We will leave such alternative analyses for the future.

by  $C_1VC_2$  which has slightly more base forms than  $C_1VC_2V$  (23  $C_1VC_2$  base forms compared to 19  $C_1VC_2V$  base forms). The primary difference between  $C_1VC_2$  and  $C_1VC_2V$  is that  $C_2$  is a coda consonant in  $C_1VC_2$  but an onset consonant in  $C_1VC_2V$ . The extra vowel/syllable in  $C_1VC_2V$  is also a possible cause of the divergence. To further evaluate these observations, we turn to the two word shapes in the three-consonant group.



**Figure 1:** Variable importance of each consonant in two-consonant base forms by word shape:  $C_1VC_2V$  and  $C_1VC_2$ .

**Figure 2** visualises the variable importance of each consonant in  $C_1VC_2VC_3$  and  $C_1VC_2C_3V$  base forms.  $C_1VC_2VC_3$  and  $C_1VC_2C_3V$  both have three consonants and two syllables but they differ in terms of the constituency of their  $C_2$  and  $C_3$ . The  $C_2$  is an onset in  $C_1VC_2VC_3$ , but a coda in  $C_1VC_2C_3V$ . The  $C_3$  is a coda in  $C_1VC_2VC_3$ , but an onset in  $C_1VC_2C_3V$ .  $C_1VC_2VC_3$  shows a distance-based decay effect with a decrease in importance from  $C_1$  to  $C_3$ . However,  $C_1VC_2C_3V$  shows a different pattern. A general distance-based decay can still be seen with  $C_1$  being more important than  $C_3$ , but  $C_2$  diverges from the pattern being more important than both  $C_1$  and  $C_3$ .  $C_1VC_2C_3V$  matches the distance-based decay pattern observed in the three-consonant group (as described in Section 3.2.1). This is again not surprising since  $C_1VC_2VC_3$  is the dominant word shape in the three-consonant group with 37 base forms, while  $C_1VC_2C_3V$  has 14 base forms.

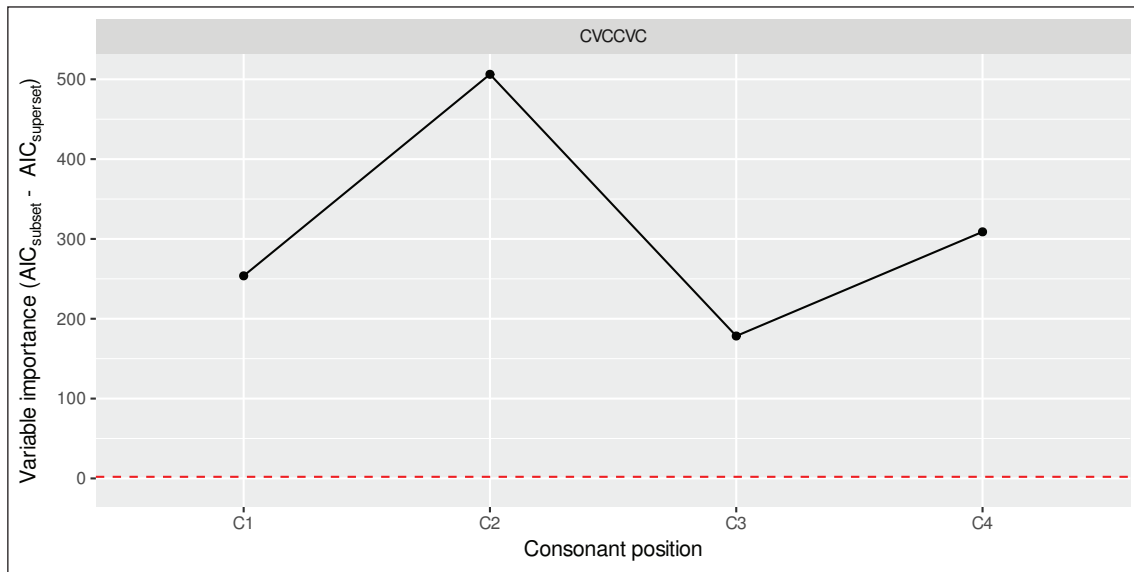


**Figure 2:** Variable importance of each consonant in three-consonant base forms by word shape:  $C_1VC_2VC_3$  and  $C_1VC_2C_3V$ .

One explanation of the patterns of these four word shapes is to consider two effects operating in tandem. The first effect is the distance-based decay effect which predicts a decrease in importance as distance increases. The second effect is that the syllable position effect which predicts coda consonants to be more important than onset consonants. The pattern of  $C_1VC_2C_3V$  can be explained if we consider both effects together.  $C_1$  is more important than  $C_3$  because of the distance-based decay effect, while the divergence of  $C_2$  being more important than  $C_1$  is because  $C_2$  is a coda consonant. The  $C_2$  of  $C_1VC_2VC_3$  conforms to the distance-based decay effect and did not diverge because it is not a coda. Similarly, the divergence of  $C_2$  from the distance-based decay effect in  $C_1VC_2$  but not in  $C_1VC_2V$  can be explained since the  $C_2$  in  $C_1VC_2$  is a coda but an onset in  $C_2VC_2V$ . To further evaluate the two hypothesised effects, we turn to the four-consonant group.

**Figure 3** visualises the variable importance of each consonant in  $C_1VC_2C_3VC_4$  base forms.  $C_1$  and  $C_3$  are both onsets, while  $C_2$  and  $C_4$  are both codas. A distance-based decay effect can be observed with the two onsets and the two codas;  $C_1$  is more important than  $C_3$ , and  $C_2$  is more important than  $C_4$ . The syllable position effect can also be observed with  $C_2$  and  $C_4$  (codas) being more important than  $C_1$  and  $C_3$  (onsets).





**Figure 3:** Variable importance of each consonant in four-consonant base forms by word shape:  $C_1VC_2C_3VC_4$ .

## 4. Discussion

This section presents our results for the two studies. We discuss and situate our findings of the Turkish case into the broader literature, focusing on feature specificity and locality. The section also connects our findings to higher-level topics such as factors beyond identity avoidance and representation of speakers' knowledge.<sup>21</sup>

### 4.1. Feature specificity

The model comparisons in Study I revealed two key findings. The first finding concerns the importance of total identity and partial identity. While previous studies have evaluated both types of identity in their formal analyses, it is entirely conceivable that total identity might not be needed after partial identity has been taken into account. Partial identity is a function of individual featural identity. If two segments are totally identical, both total identity and partial identity would be able to capture the same degree of similarity. Given that the consonants in the base forms might not be one of the four linking consonants, partial identity is expected to be much more important factor than total identity. In other words, partial identity can in theory capture what total identity can but not vice-versa. Our result suggests that partial identity is more important than total identity across all three item groups but total identity still contributes above and beyond partial identity.

<sup>21</sup> See Appendix G for the discussion of the findings regarding the Turkish specific preference hierarchy relating to the linking consonant.

The second finding concerns the nature of partial identity. Recall that cross-linguistic studies have found that features that participate in identity avoidance processes are not weighed equally, in that they may differ in the extent to which they influence the phenomenon in question. Gallagher and Coon (2009) found that features such as [+strident] and [+ejective] trigger greater OCP effects than others in Chol. Moreover, non-coronal place features have been shown to trigger OCP effects of varying strengths (see e.g., Frisch et al., 2004; Pierrehumbert 1993). In light of this background, we addressed the question of whether or not the individual features contribute equally in Turkish. It was found that the models with individual featural identity variables consistently outperformed the models with sum featural identity variables, therefore partial identity operates on an individual featural level. This finding echoes many of the formal analyses from previous studies which formulated a number of identity-avoidance constraints using specific features. In particular, Demircan (1987) identified coronal, labial, and nasal as features of importance and Kelepir (2000)'s analyses involve strident, labial, continuant, and sonorant.

To determine whether these six features are of particular importance in Turkish partial reduplication in general, we inspect the model summary of each of the three item groups in **Tables 4, 5, and 6.**<sup>22</sup> Strident appears to be the most prevalent feature since it was significant in all of the consonant positions across all three models. Labial was similarly prevalent since it was only insignificant once ( $C_4$  of the four-consonant model). Nasal was insignificant twice ( $C_2$  of the two-consonant model and the four-consonant model). Continuant was insignificant three times ( $C_1$  of the two-consonant model, and  $C_1$  and  $C_4$  of the four-consonant model) and it was in the opposite direction of identity avoidance twice ( $C_3$  of the three-consonant and the four-consonant model). Sonorant was insignificant three times ( $C_1$  of the three-consonant model, and the four-consonant model, and  $C_2$  of the three-consonant model) and it was in the opposite direction in  $C_3$  of the four-consonant model. Coronal was insignificant four times ( $C_1$  of the three-consonant model and the four-consonant model,  $C_2$  of the two-consonant model and  $C_4$  of the four-consonant model) and it was in the opposite direction in  $C_3$  of the four-consonant model. In terms of the size of the coefficients, strident and labial were in generally higher than the other features. The features, continuant, sonorant, and coronal, were not as consistent in terms of their statistical significance and the direction of their effects.

Overall, these models suggest that strident and labial were the most prevalent features. This finding is in line both with two of the previous studies on Turkish (i.e., Kelepir, 1999; Yu, 1999), which emphasised the importance of [strident], as well as cross-linguistic studies that found a similar effect (e.g., Bye, 2011; Gallagher & Coon, 2009). The feature [labial] has also been found to be influential as opposed to some other features, such as [nasal], which again are in line with

---

<sup>22</sup> Note that we are not comparing the effect sizes of the features (the coefficients of the features) across models.

the cross-linguistic typology, in that non-coronal place features were argued to play a larger role in the OCP literature (e.g., Bye, 2011; Pierrehumbert, 1993). Another finding of our study concerns the status of [coronal]: While both Demircan, 1987 and Kelepir, 1999 argue for the importance of coronal feature in their systems, our results found no such effect. In this regard, our finding is more consistent with the cross-linguistic generalization that coronal is not an influential feature (e.g., Bye 2011; Coetzee & Pater, 2008; Pierrehumbert, 1993).

While our findings are in line with cross-linguistic typology, it is important to remind ourselves that these cross-linguistic tendencies might be triggered by language-specific phenomena. In the case of Turkish, the alternating segments being [p, m, s, r] largely determine which features could *potentially* participate in identity avoidance. As such, although [strident] is influential in both Chol and Turkish, this importance is due to different factors.

We also believe that it is unsurprising that strident, labial, and nasal are particularly prevalent. Strident and nasal decrease the preference for [s] and [m], respectively, while labial decreases the preference for [p] and [m]. The three features together influence the preference for three of the four linking consonants except for [r]. The preference for [r] can be captured by using a markedness constraint since, as we have discussed earlier, it is the least preferred linking consonant even when other fixed and random effects were taken into account. In contrast to strident, for example, continuant was not a prevalent feature in our study, contrary to Kelepir's analyses. Similarly, the nasal feature is more important than the sonorant feature, also used in the previous literature. These findings, we believe, further confirm the heuristic side of the inventory of features employed by many previous studies. Additionally, we speculate that this state of affairs could be explained by resorting subset-superset relation between the features in question. In particular, the feature sonorant encompasses a larger set of consonants as opposed to nasals; likewise, continuant picks out a larger set of consonant than strident. The results indicate that the feature that is more specific applies first or is more influential, in a way that makes the superset feature redundant.

Our results also reveal that the identity effects of a given feature for a given consonant can differ across different stem types. For example, the continuant identity effect of  $C_1$  is not significant in bi-consonantal and quadri-consonantal stems, but it is significant in tri-consonantal stems.

These apparent inconsistencies call for some *post-hoc* speculative remarks. While it is plausible to assume that the same feature in the same position will exert the same strength across the board for stems of different lengths, we believe this does not need to be the case. The strength of the effect of a feature may very well be not determined once and for all, rather it may be contextually determined, which might include the number of other consonants and the syllable roles in the stem. For example, while a feature F in  $C_1$  might be strong in a two-consonant stem, the same feature F might have a weaker role in a three-consonant stem if a competing feature

is found in  $C_3$  position. In statistical terms, additional consonants could introduce additional identity avoidance effects which could influence the significance of other consonants (e.g.,  $C_1$ ) by taking up variance in the regression model. As an analogy, in phonological theories like OT, constraints are not strictly-ranked, but rather are weighted, where multiple violations of lower-priority constraints are able to overcome the violation of a higher-priority constraint (the “gang effects”) (Pater, 2016). As such, we should not necessarily interpret the mismatches/inconsistencies about the effect of a particular feature as a negative or as a shortcoming of the model, but rather as indicators of the contextual factors that give rise to the apparent mismatch.

Relatedly, what can our statistical models, which were fitted separately over different word shapes, tell us about the Turkish speakers’ grammar? The models were trained on Turkish speakers’ responses to real words which as discussed in Section 2.2, are confined to certain feature combinations. Therefore, our models would allow us to predict how a Turkish speaker would choose  $C_2$  when emphatically reduplicating a novel existing or nonce stem of given length with a C with a given feature specification in a particular position, as long as the stems have comparable feature combinations as the real words that we examined.

If our set of real words are generally representative of the reduplicable items in the Turkish lexicon, then two theoretical stances are conceivable in terms of what aspects of speakers’ grammar our models are capturing. On the one hand, one could take a position with a usage-based lexicon-based approach (e.g., Baayen et al., 2019; Bybee, 2003, 2006, 2010; Chuang & Baayen, 2021). Under this view, Turkish speakers’ grammar is confined to the attested feature combinations in their mental lexicon. Therefore, there are limits on what a Turkish speaker can learn and thus predict from the stems they know, and therefore on what we, as the modelers, can learn and thus predict from our models of what a Turkish speaker can learn from their lexicons. On the other hand, one could take a different position which allows for Turkish speakers’ grammar to not be solely based on their mental lexicon, but rather also under the influence of learning strategies, such as analogical learning (e.g., Arndt-Lappe, 2014; Nosofsky, 1986; Skousen et al., 2002) and discriminative learning (e.g., Baayen et al., 2011). Further research is needed to tease these two stances apart by examining nonce words that contains unattested feature combinations.<sup>23,24</sup>

---

<sup>23</sup> We thank John Kingston and Gaja Jarosz for extensive discussions on this topic.

<sup>24</sup> Another possibility is that our sample of real words are not representative enough of the Turkish lexicon and thus reflecting only an experimental accident with our item selection. In other words, our apparent inconsistent effects could be the result of how our model overfitted the sampled data. Future research could conduct a lexical analysis of how representative our items are, and if they were unrepresentative then one could extend our experiments to a larger set of real words.

Furthermore, yet another possibility is that we have not taken into account factors beyond the identity avoidance effect as will be discussed in Section 4.3, as well as factors that we excluded in this study (as described in Appendix B). Future research could incorporate these additional factors to examine whether and how they would influence our complete findings, including both the apparent inconsistent effects and the consistent effects.

## 4.2. Locality: Distance-based decay and syllable role

The analyses of locality revealed a number of findings, including the domain over which this phenomenon operates, whether its effect is weighed equally across the domain or other factors such as syllable position effects are at play.

The first finding concerns the rightward boundary up to which the identity avoidance effect applies. Recall that most previous studies emphasise the importance of  $C_1$  and  $C_2$  (e.g., Keleşir, 1999; Wedel, 1999), with some of them (explicitly) assuming a cut-off after  $C_2$  (Wedel, 1999). At the same time, Yu (1999) found that the strident feature has an effect with respect to  $C_3$ , which suggests that the LC is not restricted to identity avoidance effect only with respect to  $C_1$  and  $C_2$ . Our results reveal that Yu (1999)'s conclusion was on the right track, but goes further to include all the consonants in the base, and that the effects of the other consonants in the base were not particularly strong or insignificant compared to  $C_1$  and  $C_2$ .

The second finding shows that despite all the consonants contributing to the identity avoidance effect, the effect was not uniform but linear, with the strength of the effect decreasing further into the base. This suggests all the consonants play a role in identity avoidance but they are subject to a distance-based decay effect (Zymet, 2014, 2018), which states that the likelihood for the application of a phonological process decreases as transparent distance increases.

The third finding is that the distance-based decay effect is not completely linear, and is shown to be sensitive to the word shape. Specifically, some consonants diverged from the linear order, with  $C_3$  of the four-consonant base forms is more important  $C_4$ , and  $C_2$  of the two-consonant base forms is more important than  $C_1$ . We found that these divergences can be fully explained, once the constituency of a consonant is taken into account. We postulated that on top of the distance-based decay effect, the syllable position effect is also present with a coda consonant exhibiting a stronger identity avoidance effect than an onset consonant. When considering both effects in tandem, we were able to capture the patterns of importance across a number of frequent word shapes in our dataset. To further illustrate the syllable position effect, we can consider a minimal pair from our dataset – *sık* 'tight' and *sıktı* 'frequent'. In both base forms, the  $C_2$  /k/ should disprefer the linking consonant [p]. Given the syllable position effect,  $C_2$  /k/ in *sık* should disprefer [p] more than that in *sıktı*. This is indeed confirmed by our rating study. *sıp* + *sık* has a mean rating of 3.7 which is less acceptable than *sıp* + *sıktı* which has a mean rating of 5.3. The coda effect was observed in a recent nonce word study. Köylü (2020) asked 14 native speakers of Turkish to reduplicate 48 nonce-words of four word shapes (VCV, CVC, CVCV, and VCCV). It was found that the linking consonant was never identical to  $C_1$  or  $C_2$  with the CVC nonce-words. However, the linking consonant was never identical to  $C_1$  but sometimes identical to  $C_2$  with the CVCV nonce-words. In other words,  $C_2$  did not always affect the linking consonant in terms of total identity in CVCV but not in CVC. This can be explained by how  $C_2$  is a coda in CVC which has a stronger

effect on the linking consonant than when it is an onset in CVCV. Why should a coda consonant outweigh an onset consonant? We speculate that it is due to the fact that the linking consonant itself is also a coda consonant for consonant-initial base forms. Since the identity avoidance effect is a function of the similarity between a consonant in the base and the linking consonant, the similarity of the two consonants would be stronger if they have the same type of constituency. This finding lends support to the view that syllable position might play a role in contributing to segments' (dis)similarity (Bennett, 2012; Rose & Walker, 2004).<sup>25</sup> In particular, the effect would decay with the distance from the linking consonant and would be enhanced if the consonant in the base matches the constituency of the linking consonant (coda vs onset). By uncovering the presence of this effect in the Turkish emphatic reduplication, our study also adds to the typology of syllable position effects. As mentioned earlier in the paper, Bennett (2013) notes that while harmony/assimilation is predicted between consonants with matching syllable roles, dissimilation is predicted for consonants with mismatching syllable roles (but not those with matching syllable roles). Turkish constitutes an example in which dissimilation favors matching syllable roles, with the effect of syllable role observed in non-categorical patterns with multiple features.

All in all, Turkish serves as a fruitful testing ground in showing that the identity avoidance effect holds for all the segments, with the effect being strongest from  $C_1$  and being weakened as a function of its distance from the target segment. Moreover, our study reveals that in addition to the distance, the syllable position also plays a role in the application of the partial reduplication. Given our findings, one could formulate the grammar by deriving the weights of the distance-decay function (e.g., the decay parameter of the decay function which describes the shape of the decay), the weights of syllable position (e.g., the weight of coda), and their possible interactions (McMullin & Burness, 2021; Zymet, 2014). This, however, will likely require a more well-balanced dataset with nonce words in terms of the representation of particular features, word shapes and the number of consonants in the base. This is left for future studies.

### 4.3. Beyond identity avoidance

The current study examined the Turkish partial reduplication phenomenon as a phonological operation with a focus on the identity avoidance effect. Our models were able to explain a sizable portion of the variance of our naturalness judgements.<sup>26</sup> The identity avoidance factors captured around 25% to 35% of the variance (two-consonant group: 26.96%, three-consonant group:

---

<sup>25</sup> It is possible to recast the empirical findings from our study in more formal terms such as Structural Surface Correspondence constraints (e.g., Structural SCorr-CC, which would limit correspondence to only those consonants which have matching syllable roles) as proposed by Bennett (2012) and Rose and Walker (2004).

<sup>26</sup> The proportion of variance captured by fixed effects in the models was computed with the function `r.squaredGLMM()`, part of the `MuMIn` library in R (Bartoń 2022). This function returns both marginal  $R^2$  and conditional  $R^2$ . Marginal  $R^2$  represents the variance explained by fixed factors. Conditional  $R^2$  represents the variance explained by fixed and random factors.

33.67%, and four-consonant group: 28.36%) and, together with the random effects which captured the idiosyncracies of lexical items, the models were able to explain around 50% to 70% of the total variance (two-consonant group: 55.73%, three-consonant group: 55.75%, and four-consonant group: 70.98%). Our study therefore provided ample evidence that phonological factors play a major role in the process behind Turkish partial reduplication. However, there is still variance to be explained by factors beyond phonology, such as the lexicon, morphology, and semantics.

As also noted by an anonymous reviewer, our analysis captures a large portion of the patterns, yet certain minimal contrasts between reduplicated forms such as *köp-kötü* ‘very bad’ vs. *kas-katı* ‘very hard’ still are not fully accounted for. As shown in the Appendix, different linking consonants are selected for the same *k-t* sequence of consonants in the base, and in the same syllabic positions (as argued in Yu 1999, we take it that vowels do not play a role in the reduplication process). Previous analyses have also failed in providing a satisfactory explanation for such minimal contrasts. For example, Demircan (1987) brings up the intuition that speakers might avoid a reduplicant that resembles an existing root and confirms it in a corpus study by Kılıç and Bozşahin (2013), which demonstrated that root-level lexical statistics inversely correlate with the preference of a linking consonant. Speakers’ preference therefore might depend on the knowledge of distributional statistics at the morpholexical level. Speakers might not prefer to use a reduplicated form for the emphatic meaning, but instead prefer to use the word *çok* ‘very’ with the base form. However, in the case of *kas-katı*, it is still unclear why *s* is selected over *p* since both *kas* ‘muscle’ and *kap* ‘container’ are meaningful words in Turkish. In fact, the only linking consonant that does not result in a form resembling another Turkish root is *m* as *kam* is not a Turkish word, whereas *kar* ‘snow’ also is. Therefore, we acknowledge the presence of factors beyond phonological considerations that play a role in the choice of an LC.<sup>27</sup>

In a series of corpus and experimental studies by Kaufman (2014), it was observed that the preference to reduplicate a base form depends on its semantic class and the semantic class of existing base forms that are frequently reduplicated. A low rating of a reduplicated form might not be due to phonological factors but rather due to the participant’s dispreference to reduplicate the base form. Incorporating morpholexical statistics and semantics can therefore provide a more complete picture of the Turkish partial reduplication phenomenon and we leave this for future research.

In light of the phonological generalizations uncovered by our experimental study, in the next section we discuss what speakers’ knowledge of this pattern looks like (i.e., what is the potential representation speakers have in mind when using this phenomenon?)

---

<sup>27</sup> The same anonymous reviewer also notes another issue for the pair *dip-diri* vs. *dup-duru*. Both have the same consonant sequence, the same word shape, and the same linking consonant. However, the acceptability ratings are found to be different for *dip-diri* (mean: 1.43) vs. *dup-duru* (mean: 1.35). This calls for a careful examination of other factors, such as frequency of forms in a corpus of Turkish in future studies.

#### 4.4. Representation of speakers' knowledge

With respect to the representation of speakers' knowledge of the Turkish partial reduplication phenomenon, one point of investigation that researchers have focused on is whether the choice of a particular linking consonant in partial reduplication is simply a matter of *lexicalization*, a term which has been interpreted in more than one way in the literature. One interpretation of this approach is whether the choice of the LC is random/arbitrary, or follows a set of generalizations or rules. If the conclusion is that the choice of the LC is arbitrary, then the phenomenon is considered to be lexicalised. Although earlier studies assumed the choice of the LC to be arbitrary, thus lexicalised (e.g., Foster, 1969; Lewis, 1967; Yavaş, 1980), a number of studies have argued that the choice is not lexicalised, and is indeed conditioned by various rules. These studies include Hatiboğlu (1973), Demircan (1987), Dobrovolsky (1987), Taneri (1990), Wedel (1999), Yu (1999), Kelepir (2000), Sofu and Altan (2008), Kaufman (2014). While these works vary considerably in their implementations of the observations, they converge on the view that the choice of the LC is not arbitrary or lexicalised, and that it is subject to several dissimilation constraints motivated by the OCP, similar to the analyses given for dissimilation processes in other, unrelated languages.

Some other studies have used nonce-words as a diagnostic as to whether Turkish partial reduplication is lexicalised or productive. Under the (often implicit) assumption that real words that participate in this phenomenon obey various rules or generalizations, studies in this line of research probe whether nonce-words are subject to the same generalizations. It turns out the results from these studies are far from clear, and have found varying, and sometimes conflicting results. For example, while Sofu (2005) concludes that speakers seem to extend at least some of the rules to nonce-words, two more recent studies, Demir (2018) and Köylü (2020), reach opposing conclusions as to the status of lexicalization and productivity of the reduplication patterns. While Demir (2018) interprets her results in favour of a lexicalization approach, Köylü (2020) argues that speakers do extend the same strategies they use for real words to nonce-words.

On the side of studies that investigate whether the Turkish emphatic reduplication is arbitrary or rule-governed, Wedel (1999) for example, argues that 'native speakers do abstract some productive phonological generalization from the emphatic forms that exist', as such speakers 'have access to a uniform, constraint-based schema' in using this phenomenon. This is at least the case for the LCs [p, s, m], while the LC [r] may indeed have lost its productivity since it is not used with novel forms and is the least utilised LC as confirmed by other studies including Yu (1999) and this current study.<sup>28</sup> Along the same lines, Yu (1998) also notes that 'modern speakers of Turkish have some grammatical knowledge of emphatic reduplication' on the basis

---

<sup>28</sup> An anonymous reviewer suggested that for Wedel (1999), all reduplicated forms are lexicalised since they are not productive. Given the above statements from Wedel (1999), we believe (see also Köylü, 2020 for this view) that Wedel does not take this phenomenon to be lexicalised, differing from the interpretation of that study by the reviewer.



of the results that speakers do not blindly choose a particular LC (e.g., [p], across the board) and concludes:

Proposals that claim the emphatic construction is unproductive and that all reduplicative closers must be lexically listed with the base form must be taken with great precaution, if not rejected altogether. The experiments reported here clearly suggest that native speakers of Turkish still retain some grammatical knowledge of the selectional restriction of the closer in the emphatic reduplication construction. (Yu, 1998, p. 39)

In this regard, we also add that if partial reduplication was simply a matter of lexicalization, it would be surprising that for a major number of items, speakers exhibit variation, in the sense that more than one linking consonant is permitted. Crucially, in these variable cases, the permissible linking consonants are not identical across items, and even for items that permit multiple identical LCs, the relative acceptability differs from one to another and across speakers. All these suggest that speakers are not blindly memorizing a specific LC for each base form, but instead, as Wedel (1999) and Yu (1999) already suggested, observe certain phonological generalizations, which our study aims to make precise.

Let us now turn to the other interpretation of lexicalization, which focuses on whether nonce-words parallel the behavior of real words. For example, a study by Demir (2018) examines the choice of LC with adults via a comparison of 10 real- and 34 nonce-words using an open-set response task. In one experiment, nonce-words are chosen to have a real word counterpart with identical consonants and word shapes, but different vowels. Demir finds that real words follow the expected observations noted by previous studies, as such using the four common LCs, [p,m,s,r]. On the other hand, in the case of nonce-words, speakers resort to strategies that are not available for real words. For example, they might copy the CVC from the base, as opposed to the CV plus LC strategy, or they might omit a linker altogether. Moreover, nonce-words which were chosen to analogise the real words do not show the same linking consonant as their supposed real word counterparts. Demir (2018) interprets these divergent results between real- and nonce-words to mean that there is not a pattern or set of generalizations for emphatic reduplication of real words, which one would expect to be extended to nonce-words. Therefore, reduplication of real words must be lexicalised.

However, this conclusion might be a bit too hasty. As we just noted, there are other nonce-word studies carried out on Turkish partial reduplication that arrive at different conclusions or interpretations. For example, Sofu (2005) examines the choice of LC with both adults and children using 38 nonce-words. She finds that adults and children conform to the expected patterns in their use of the classic LCs, [p,m,s,r], while children use [p] more than adults. Moreover, children use linking consonants that are different from [p,m,s,r] much more often than adults (e.g., [t,n,f]).

This study shows that both adults and children do extend the patterns observed in real words to nonce-words as well.<sup>29</sup> Similarly, a more recent study by Köylü (2020) also concludes that Turkish native speakers extend the reduplication strategies they employ in real words to nonce-words. These conflicting results call for a careful investigation regarding the causes behind them. Here we speculate on a few potential issues.

There might be various methodological or linguistic factors that lead to this divergence in nonce-word studies. For example, all the previous studies simply present the nonce-word test items by themselves out of a context without assigning them any meaning. This is not a trivial choice in light of the fact that reduplication is applicable only to gradable modifiers. It is likely that these out-of-context nonce-words were not interpreted as such by the participants, and therefore participants resorted to strategies that differ from real words, whose meaning and property of being gradable they are aware of. The restriction regarding the category and property this reduplication requires could also be a factor as to why nonce-word studies might not be replicating or reaching lower scores.

This last point also relates to another concept, productivity, which usually comes up in the discussion of lexicalization. On the point of productivity, it is worth highlighting that it should be approached with caution. This is because emphatic reduplication, as just noted above, applies specifically to a subset of adjectives (i.e., gradable adjectives) and adverbs (which usually are built on adjectives in Turkish) and not to absolute adjectives or modifiers in general. As such, we can make sense of why it is found in a relatively small number of items in the language, due to its nature. Therefore, any statement about productivity should take this important aspect into consideration, and its potential role especially on studies investigating nonce-words.

Moreover, recall that Demir (2018) created certain nonce-words in anticipation of analogy with real words. For example, the nonce-word *mava* was created by Demir on the assumption that the participants would analogise it to the existing adjective *mavi* ‘blue’ while they are attempting to reduplicate *mava*. This is not an innocuous assumption, however. First, analogy does not rely on just a single word, but a number of words (e.g., the Generalised Neighbourhood Model by Bailey and Hahn (2001) considers all the words in the lexicon weighted by lexical frequency and form similarity to the target (nonce-)words). As such this assumption overlooks the complicated aspect of how analogy works. Secondly, some words might have a large number of noun neighbors, which might make the reduplication harder. For example, in the context of the nonce-word *boyuz* from Demir’s study, our native intuition (and those of our consultants)

---

<sup>29</sup> Regarding the use of [p] more frequently and the presence of non-standard linking consonants, it could be showing that children are still in the process of mastering the abstract generalizations. As such, they sometimes revert to the default LC, p, or have a larger set of potential LCs that they have not narrowed down yet to the ones adults use.

analyses *boyuz* to *boyoz*, which is an existing word that refers to a food item. As such, although the experimenter might have a certain real word in mind while designing the nonce-words, participants might have completely different real word analogies. Similarly, with *mava*, native speakers we have consulted brought up the nouns *hava* ‘air’, or *tava* ‘pan’ as the first words that come to their minds, both of which are nouns, rather than the adjective *mavi* that Demir had in mind as a control item. This is significant since the control items that Demir has in mind may not be the items participants are supposedly analogizing to, as such might explain part of the results. A further related note is that analogy may not be solely based on phonological properties, but might be due to semantic resemblance participants establish with real words they might think of.

As just discussed in Section 4.3, there are also other factors beyond identity avoidance, such as lexical (Kılıç & Bozşahin, 2013) and semantic factors (Kaufman, 2014), that play a role but were not consistently considered in these nonce-word studies.

In light of these considerations, the results of our study are in support of the view that partial reduplication is subject to various active phonological rules, particularly speakers have access to locality and feature-based conditions (or generalizations) that they are applying to items that are potentially intensifiable. In this regard, it corroborates the findings/intuitions raised by studies such as Yu (1998), Wedel (1999), and accords with Demircan (1987) and Kelepir (1999) who also conclude that speakers obey various phonological rules when they do partial reduplication. In particular, in this study we have uncovered that the phonological rules that are exhibited by the real words are much more graded than previously thought. The identity avoidance effect is both locality-sensitive (distance-based decay effect, and syllable position effect) and feature-sensitive (individual features with different weights).

With that said, it is important to keep in mind that as Frisch et al. (2004) argues, lexical information, including lexical idiosyncrasies, and rules are not mutually exclusive, in that there does not need to be a categorical choice between the two interpretations. It is possible to have a phenomenon which respects various constraints (e.g., locality- or feature-specificity, or phonotactic ones). As such, it is not necessarily the case that presence of phonological rules conflicts with or rules out the presence of lexical information. The Turkish emphatic reduplication may very well be an example of this sort. Wedel’s (1999) conclusion might also be in line with this interpretation, in that while certain linking consonants [p, s, m] still actively participate in the identity avoidance effect, [r] might have fallen on the memorization side of it.

Having better understood the nature of the phonological grammar that speakers might have, this naturally leads to the second question about lexicalization, (i.e., whether and to what extent these rules are being extended to nonce-words). Concerning nonce-words, we can bring in the insights raised by Becker et al. (2011) who found, on the basis of another phenomenon

in Turkish, that nonce-words are subject to only some of the rules exhibited by real words (see also Harris, Neasom & Tang, 2016; Hayes & White, 2013). As such, there is no *a-priori* reason to expect that nonce-words fully conform to the same generalizations as real words, or reflect all of the generalizations/rules found for real words. Moreover, as discussed above, any study that aims to investigate nonce-words in the Turkish emphatic reduplication must also address a number of methodological and linguistic factors such as framing the nonce-words in a context that signifies its meaning and property of being gradable, as well as controlling for lexical and semantic factors. This can be done by using a combination of careful stimuli design and statistical modelling (Redington & Chater, 1996; Tang & Baer-Henney, 2023). Once these factors are controlled for, we expect (following Becker et al., 2011) that at least some of the real word generalisations would apply to nonce-words (some of which have already been argued to be the case, e.g., Köylü, 2020; Sofu, 2005).

## 5. Conclusions

This paper has re-examined a well-known reduplication phenomenon in Turkish. Modifiers such as adjectives and adverbs can undergo a partial reduplication process to express an emphatic meaning by prefixing a  $C_1VC_2$  syllable. Unlike most instances of reduplication with fixed segmentism which have a single fixed segment, the Turkish emphatic reduplication contains four fixed segments, as such the linking consonant  $C_2$  can be one of the four consonants: [p], [m], [s], and [r]. The study investigates the factors conditioning the choice of the LC, by focusing the nature of the (dis)similarity (feature specificity) and the proximity (locality) between the consonants in the base and the LC. Turkish emphatic reduplication turns out to be well-suited for shedding light on the nature of identity avoidance as it allows multiple possible fixed segments, and the same item itself might be used with multiple LCs.

Using an acceptability rating task conducted with over 200 participants for 162 base forms, the study has uncovered a number of significant findings with implications for the broader research on identity avoidance. Our analyses revealed that the identity avoidance effect is much more graded than it has been previously proposed in both its specificity and locality. Unlike most previous studies which emphasise the importance of  $C_1$  and  $C_2$  of the base as the domain in which dissimilation operates over, our study has found that the effect extends over all consonants in the base. Crucially, despite all the consonants contributing to the identity avoidance effect, the effect was not uniform, with the strength of the effect decreasing further into the base. Therefore, we demonstrate that the phenomenon is subject to a distance-based decay effect (Zymet, 2014, 2018). Moreover, the study uncovers an intricate interplay between the distance-based decay effect and the syllable position effect (Bennett, 2012; Rose & Walker, 2004), which is not as transparent as it is in most other languages. This novel finding for the Turkish emphatic reduplication is made possible due to the methodology and statistical tools adopted in this study.

In terms of the nature of specificity, our results reveal that in the Turkish emphatic reduplication process, the similarity between the consonants in the base and the LCs operates at the segmental level (total identity) as well as the level of individual phonological features (partial identity). Moreover, in line with the cross-linguistic picture, we have found that features that participate in identity avoidance processes may differ in the extent to which they influence the phenomenon in question. Our study confirms the importance of only some of the features employed in the previous studies, such as [strident], [labial], and [nasal], but not others, such as [coronal], [sonorant], or [continuant]. The important features like [strident] and [labial] and the unimportant features such as [coronal] and [sonorant] are more in line with the cross-linguistic tendencies (see e.g., Bye, 2011; Pierrehumbert, 1993). Methodologically we demonstrate that the precise nature of the identity avoidance effect can be revealed using hierarchical regression and statistical model comparisons (Graff & Jaeger, 2009; Zymet, 2019).

---

## Abbreviations

The following abbreviations were used.

**LC:** Linking consonant

**AIC:** Akaike information criterion

**BIC:** Bayesian information criterion

## Data accessibility statement

Given the lack of rating judgements for the Turkish partial reduplication phenomenon, we made our data available in an Open Science Framework repository (<https://www.doi.org/10.17605/OSF.IO/P2JDK>) and in the Appendix H. Furthermore, given the complexity of the analyses used in this study, we made the analysis scripts we used to produce the results available so that readers can evaluate the data and our procedures themselves.

## Additional files

The appendix consists of nine parts: A, B, C, D, E, F, G, and H. In part A, we provide the phonological feature values of Turkish consonants. In part B, we provide a set of excluded variables and the reasons of their exclusion. In part C, we provide the distribution of the variables (both the response variable and the predictors) for each of the three item groups with two, three and four consonants in the base form. In part D, we provide the pairwise association results between the response variable and each of the predictors in Section 2.3.1 for each of the three item groups. In part E, we provide the random effects summaries for all the reported models. In part F, we provide the model comparison for feature specificity. In part G, we provide a discussion of the findings regarding the Turkish specific preference hierarchy relating to the linking consonant. In part H, we provide the by-item acceptability ratings as well as an inter-rater reliability analysis <https://doi.org/10.16995/labphon.6459.s1>.

## Acknowledgements

We dedicate this paper to the residents of all the towns that were impacted by the 2023 earthquakes in southeastern Turkey. We thank our Turkish-speaking participants for their generous contribution to our research. We thank John Kingston, Gaja Jarosz, Ryan Bennett, Joe Pater, Andrew Nevins, Ryan Budnick, Akiva Bacovcin, John Harris, Fabian Tomaschek, Yılmaz Köylü and the audience at the 26th Manchester Phonology Meeting and the 15th Old World Conference in Phonology for their feedback on earlier versions of the study. Special thanks to the anonymous *Laboratory Phonology* reviewers. We thank Anna Sophia Stein and Anh Kim Nguyen for their help with typesetting the paper. All errors remain ours.

## Competing interests

The authors have no competing interests to declare.

## Authors' contributions

We follow the CRediT taxonomy.<sup>30</sup>

Kevin Tang: Conceptualization, Methodology, Data curation, Formal analysis, Writing-Original draft preparation, Visualization, Investigation, Writing-Reviewing and Editing.

Faruk Akkuş: Conceptualization, Methodology, Data curation, Formal analysis, Writing-Original draft preparation, Visualization, Investigation, Writing-Reviewing and Editing.

---

## References

- Albright, A. (2007). *Gradient phonological acceptability as a grammatical effect*. Retrieved from <http://web.mit.edu/albright/www/papers/Albright-GrammaticalGradience.pdf>.
- Alderete, J., Beckman, J., Benua, L., Gnanadesikan, A., McCarthy, J., & Urbanczyk, S. (1999). Reduplication with fixed segmentism. *Linguistic Inquiry*, 30(3), 327–364. DOI: <https://doi.org/10.1162/002438999554101>
- Arndt-Lappe, S. (2014). Analogy in suffix rivalry: The case of English-ity and-ness. *English Language & Linguistics*, 18(3), 497–548. DOI: <https://doi.org/10.1017/S136067431400015X>
- Baayen, R. H. [R Harald], Chuang, Y.-Y., Shafaei-Bajestan, E., Blevins, J. P., et al. (2019). The discriminative lexicon: A unified computational model for the lexicon and lexical processing in comprehension and production grounded not in (de) composition but in linear discriminative learning. *Complexity*, 2019. DOI: <https://doi.org/10.1155/2019/4895891>
- Baayen, R. H. [R Harald], Milin, P., Đurđević, D. F., Hendrix, P., & Marelli, M. (2011). An amorphous model for morphological processing in visual comprehension based on naive discriminative learning. *Psychological Review*, 118(3), 438–482. DOI: <https://doi.org/10.1037/a0023851>
- Baayen, R. H. [R Harald]. (2008). *Analyzing linguistic data: A practical introduction to statistics using R*. Cambridge, UK: Cambridge University Press. DOI: <https://doi.org/10.1017/CBO9780511801686>
- Bailey, T. M., & Hahn, U. (2001). Determinants of wordlikeness: Phonotactics or lexical neighborhoods? *Journal of Memory and Language*, 44(4), 568–591. DOI: <https://doi.org/10.1006/jmla.2000.2756>
- Bartoń, K. (2022). *MuMIn: Multi-Model Inference*. R package version 1.47.1. Retrieved from <https://CRAN.R-project.org/package=MuMIn>.
- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, 67(1), 1–48. DOI: <https://doi.org/10.18637/jss.v067.i01>

---

<sup>30</sup> <https://www.ucl.ac.uk/library/research-support/open-access/credit-taxonomy>.

- Becker, M., Ketrez, N., & Nevins, A. (2011). The surfeit of the stimulus: Analytic biases filter lexical statistics in Turkish laryngeal alternations. *Language*, 84–125. DOI: <https://doi.org/10.1353/lan.2011.0016>
- Becker, M., & Levine, J. (2013). Experigen: An online experiment platform. Available at <http://becker.phonologist.org/experigen>.
- Bennett, W. G. (2012). Dissimilation by correspondence in Sundanese. In N. Arnett & R. Bennett (Eds.), *Proceedings of the 30th West Coast Conference on Formal Linguistics* (pp. 76–86). Somerville, MA: Cascadilla Proceedings Project.
- Bennett, W. G. (2013). *Dissimilation, Consonant Harmony, and Surface Correspondence* (Doctoral dissertation, Rutgers University).
- Berkley, D. M. (2000). *Gradient obligatory contour principle effects* (Doctoral dissertation, Northwestern University).
- Brysbaert, M., & New, B. (2009). Moving beyond Kučera and Francis: a critical evaluation of current word frequency norms and the introduction of a new and improved word frequency measure for American English. *Behavior Research Methods*, 41(4), 977–990. DOI: <https://doi.org/10.3758/BRM.41.4.977>
- Bybee, J. (2003). *Phonology and language use*. Cambridge University Press.
- Bybee, J. (2006). From usage to grammar: The mind's response to repetition. *Language*, 711–733. DOI: <https://doi.org/10.1353/lan.2006.0186>
- Bybee, J. (2010). *Language, usage and cognition*. Cambridge University Press. DOI: <https://doi.org/10.1017/CBO9780511750526>
- Bye, P. (2011). Dissimilation. In *The Blackwell Companion to Phonology* (Chap. 60, pp. 1–26). DOI: <https://doi.org/10.1002/9781444335262.wbctp0060>
- Chatterjee, S., & Hadi, A. S. (2015). *Regression analysis by example*. John Wiley & Sons.
- Chuang, Y.-Y., & Baayen, R. H. [R Harald]. (2021). Discriminative learning and the lexicon: NDL and LDL. In *Oxford Research Encyclopedia of Linguistics*. DOI: <https://doi.org/10.1093/acrefore/9780199384655.013.375>
- Clements, G. N., & Ridouane, R. (2006). Distinctive feature enhancement: A review. In A. Botinis (Ed.), *Proceedings of the ISCA tutorial and research workshop on experimental linguistics* (pp. 97–100). Athens, Greece: International Speech Communication Association.
- Coetzee, A. W., & Pater, J. (2008). Weighted constraints and gradient restrictions on place co-occurrence in Muna and Arabic. *Natural Language & Linguistic Theory*, 26(2), 289–337. DOI: <https://doi.org/10.1007/s11049-008-9039-z>
- Demir, N. (2018). Turkish reduplicative adjectives and adverbs. In P. Farrell (Ed.), *Proceedings of Linguistic Society of America* (Vol. 3, pp. 1–14). DOI: <https://doi.org/10.3765/plsa.v3i1.4300>
- Demircan, Ö. (1987). Emphatic reduplication in Turkish. In H. E. Boeschoten & L. T. Verhoeven (Eds.), *Studies on modern Turkish: Proceedings of the third conference on Turkish linguistics* (pp. 24–41). Tilburg, Netherlands: Tilburg University Press.



- Dobrovolsky, M. (1987). Why CVC in Turkish reduplication. In P. Lilius & M. Saari (Eds.), *The Nordic languages and modern linguistics* (Vol. 6, pp. 131–146). Helsinki: Helsinki University Press.
- Erguvanlı Taylan, E. (2015). *The phonology and morphology of Turkish*. İstanbul: Boğaziçi University Press.
- Foster, J. F. (1969). *On some phonological rules of Turkish* (Doctoral dissertation, University of Illinois at Urbana-Champaign).
- Frampton, J. (2009). *Distributed reduplication*. Cambridge, Massachusetts & London: MIT Press. DOI: <https://doi.org/10.7551/mitpress/9780262013260.001.0001>
- Frisch, S. A. [Stefan A.], & Zawaydeh, B. A. (2001). The psychological reality of OCP-Place in Arabic. *Language*, 77(1), 91–106. DOI: <https://doi.org/10.1353/lan.2001.0014>
- Frisch, S. A. [Stefan A.], Pierrehumbert, J. B., & Broe, M. B. (2004). Similarity avoidance and the OCP. *Natural Language & Linguistic Theory*, 22(1), 179–228. DOI: <https://doi.org/10.1023/B:NALA.0000005557.78535.3c>
- Gallagher, G., & Coon, J. (2009). Distinguishing total and partial identity: Evidence from chol. *Natural Language & Linguistic Theory*, 27(3), 545–582. DOI: <https://doi.org/10.1007/s11049-009-9075-3>
- Goldrick, M. (2011). Using psychological realism to advance phonological theory. In *The Handbook of Phonological Theory* (Chap. 19, pp. 631–660). eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781444343069.ch19>. DOI: <https://doi.org/10.1002/9781444343069.ch19>
- Graff, P., & Jaeger, T. (2009). Locality and feature specificity in ocp effects: Evidence from Aymara, Dutch, and Javanese. In M. R. Bochnak, P. Klecha, A. Lemieux, N. Nicola, J. Urban, & C. Weaver (Eds.), *Proceedings from the Annual Meeting of the Chicago Linguistic Society* (Vol. 45, pp. 127–141). Chicago, IL: Chicago Linguistic Society.
- Harris, J., Neasom, N., & Tang, K. (2016). *Phonotactics with [awt] rules: The learnability of a simple, unnatural pattern in English*. 24th Manchester Phonology Meeting, University of Manchester, UK.
- Hatiboğlu, V. (1973). *Pekiştirme ve kuralları*. Türk Dil Kurumu Tanıtım Yayınları.
- Hayes, B., & White, J. (2013). Phonological naturalness and phonotactic learning. *Linguistic Inquiry*, 44(1), 45–75. DOI: [https://doi.org/10.1162/LING\\_a\\_00119](https://doi.org/10.1162/LING_a_00119)
- Hepworth, G., Gordon, I. R., & McCullough, M. J. (2007). Accounting for dependence in similarity data from DNA fingerprinting. *Statistical Applications in Genetics and Molecular Biology*, 6(1). DOI: <https://doi.org/10.2202/1544-6115.1212>
- Inkelas, S., Küntay, A., Sprouse, R., & Orgun, O. (2000). Turkish Electronic Living Lexicon (TELL). *Turkic Languages*, 4, 253–275.
- Kaufman, B. D. (2014). *Learning an unproductive process: Turkish emphatic reduplication* (Master's thesis, University of California Santa Cruz).
- Kelepir, M. (1999). *Emphatic non-identical reduplication in Turkish*. Talk given at the LSA Annual Meeting.

- Keleşir, M. (2000). To be or not to be faithful. In A. Göksel & C. Kerslake (Eds.), *Studies on Turkish and Turkic languages: Proceedings of the ninth international conference on Turkish linguistics* (pp. 11–18). Wiesbaden, Germany: Harrassowitz Verlag.
- Kenstowicz, M., & Kisseberth, C. (2014). *Generative phonology: Description and theory*. Academic Press.
- Keuleers, E., Brysbaert, M., & New, B. (2010). SUBTLEX-NL: A new measure for Dutch word frequency based on film subtitles. *Behavior Research Methods*, 42(3), 643–650. DOI: <https://doi.org/10.3758/BRM.42.3.643>
- Keyser, S. J., & Stevens, K. N. (2006). Enhancement and overlap in the speech chain. *Language*, 82(1), 33–63. DOI: <https://doi.org/10.1353/lan.2006.0051>
- Kılıç, Ö., & Bozşahin, C. (2013). Selection of linker type in emphatic reduplication: Speaker's intuition meets corpus statistics. In M. Knauff, N. Sebanz, M. Pauen, & I. Wachsmuth (Eds.), *Proceedings of the Annual Meeting of the Cognitive Science Society* (Vol. 35, pp. 2722–2727). Berlin, Germany: Cognitive Science Society.
- Köylü, Y. (2020). Abstract knowledge of emphatic reduplication in Turkish. In *Talk given at the 5th workshop on Turkic and languages in contact with Turkic (Tu+ 5)*, University of Delaware. DOI: <https://doi.org/10.3765/ptu.v5i1.4780>
- Lewis, G. L. (1967). *Turkish grammar*. Oxford University Press.
- McCarthy, J. J. (1986). OCP effects: Gemination and antigemination. *Linguistic Inquiry*, 17(2), 207–263.
- McCarthy, J. J., & Prince, A. (1993). Generalized alignment. In G. Booij & J. Van Marle (Eds.), *Yearbook of morphology 1993* (pp. 79–153). DOI: [https://doi.org/10.1007/978-94-017-3712-8\\_4](https://doi.org/10.1007/978-94-017-3712-8_4)
- McCarthy, J. J., & Prince, A. (1994). The emergence of the unmarked: Optimality in prosodic morphology. In M. González (Ed.), *Proceedings of the North East Linguistics Society* (Vol. 24, pp. 333–379). Graduate Linguistics Students Association, University of Massachusetts Amherst.
- McMullin, K., & Burness, P. (2021). Tier-based modeling of gradience and distance-based decay in phonological processes. In H. Björklund & F. Drewes (Eds.), *Proceedings of the 17th meeting on the mathematics of language* (pp. 50–63). Montpellier, France (online): Association for Computational Linguistics.
- Müller, H.-G. (2003). *Morphophonologische untersuchungen an reduplikationen im Türkischen* (Doctoral dissertation, Universität Tübingen).
- Nosofsky, R. M. (1986). Attention, similarity, and the identification-categorization relationship. *Journal of Experimental Psychology: General*, 115(1), 39–57. DOI: <https://doi.org/10.1037/0096-3445.115.1.39>
- Pater, J. (2016). Universal grammar with weighted constraints. In J. McCarthy & J. Pater (Eds.), *Harmonic grammar and harmonic serialism* (pp. 1–46). London: Equinox.
- Pierrehumbert, J. (1993). Dissimilarity in the Arabic verbal roots. In A. J. Schafer (Ed.), *Proceedings of the North East Linguistics Society* (Vol. 23, pp. 367–381). Ottawa, Canada: Graduate Linguistics Students Association, University of Massachusetts Amherst.

- R Core Team. (2013). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. Vienna, Austria. Retrieved from <http://www.R-project.org/>.
- Redington, M., & Chater, N. (1996). Transfer in artificial grammar learning: A reevaluation. *Journal of Experimental Psychology: General*, 125(2), 123–138. DOI: <https://doi.org/10.1037/0096-3445.125.2.123>
- Rose, S., & Walker, R. (2004). A typology of consonant agreement as correspondence. *Language*, 475–531. DOI: <https://doi.org/10.1353/lan.2004.0144>
- Skousen, R., Lonsdale, D., & Parkinson, D. B. (2002). *Analogical modeling: An exemplar-based approach to language*. John Benjamins Publishing. DOI: <https://doi.org/10.1075/hcp.10>
- Sofu, H. (2005). Acquisition of reduplication in Turkish. In B. Hurch & V. Mattes (Eds.), *Studies on reduplication* (pp. 493–509). Berlin & New York: Mouton de Gruyter. DOI: <https://doi.org/10.1515/9783110911466.493>
- Sofu, H., & Altan, A. (2008). Partial reduplication: Revisited. In S. Ay, Ö. Aydin, İ. Ergenç, S. Gökmen, S. İşsever, & D. Peçenek (Eds.), *Essays on Turkish Linguistics: Proceedings of the 14th International Conference on Turkish Linguistics* (pp. 63–73). Wiesbaden, Germany: Harrassowitz Verlag.
- Stachowski, K. (2014). *Standard turkic C-type reduplications*. Jagiellonian University Press.
- Stanton, J. (2017). Segmental blocking in dissimilation: An argument for co-occurrence constraints. In K. Jesney, C. O'Hara, C. Smith, & R. Walker (Eds.), *Proceedings of the 2016 Annual Meetings on Phonology* (Vol. 4), Washington, DC: Linguistic Society of America. DOI: <https://doi.org/10.3765/amp.v4i0.3972>
- Stevens, K. N., & Keyser, S. J. (1989). Primary features and their enhancement in consonants. *Language*, 81–106. DOI: <https://doi.org/10.2307/414843>
- Suzuki, K. (1998). *A typological investigation of dissimilation* (Doctoral dissertation, The University of Arizona).
- Taneri, M. (1990). A Type of Reduplication in Turkish. In I. Lee & S. Schiefelbein (Eds.), *Kansas Working Papers in Linguistics* (Vol. 15, pp. 93–126). Linguistics Graduate Student Association, University of Kansas. DOI: <https://doi.org/10.17161/KWPL.1808.515>
- Tang, K. (2012). A 61 million word corpus of Brazilian Portuguese film subtitles as a resource for linguistic research. *UCL Working Papers in Linguistics*, 24, 208–214.
- Tang, K., & Baer-Henney, D. (2023). Modelling L1 and the artificial language during artificial language learning. *Laboratory Phonology: Journal of the Association for Laboratory Phonology*, 14(1), 1–54. DOI: <https://doi.org/10.16995/labphon.6460>
- Tang, K., & de Chene, B. (2014). *A new corpus of colloquial Korean and its applications. The 14th Conference on Laboratory Phonology*, Tachikawa, Tokyo, Japan.
- Tomaschek, F., Hendrix, P., & Baayen, R. H. (2018). Strategies for addressing collinearity in multivariate linguistic data. *Journal of Phonetics*, 71, 249–267. DOI: <https://doi.org/10.1016/j.wocn.2018.09.004>
- Wedel, A. (1999). Turkish emphatic reduplication. *Phonology at Santa Cruz (PASC)*, 6.

- Wedel, A. (2000). Perceptual distinctiveness in Turkish emphatic reduplication. In R. Billerey-Mosier & B. Lillehaugen (Eds.), *WCCFL 19: Proceedings of the 19th West Coast Conference on Formal Linguistics* (Vol. 19, pp. 546–559). University of California, Santa Cruz. Somerville, MA: Cascadilla Press.
- Winter, B. (2019). *Statistics for linguists: An introduction using R*. Routledge. DOI: <https://doi.org/10.4324/9781315165547>
- Wissmann, M., Toutenburg, H., & Shalabh. (2007). *Role of categorical variables in multicollinearity in the linear regression* (tech. rep. No. 008). Department of Statistics, University of Munich. Munich, Germany. Retrieved from [https://epub.ub.uni-muenchen.de/2081/1/report008\\_statistics.pdf](https://epub.ub.uni-muenchen.de/2081/1/report008_statistics.pdf).
- Yavaş, M. (1980). *Borrowing and its implications for Turkish phonology* (Doctoral dissertation, University of Kansas).
- Yip, M. (1997). Repetition and its avoidance: The case of Javanese. In K. Suzuki & D. Elzinga (Eds.), *Proceedings of the 1995 Southwestern workshop on Optimality Theory (SWOT)* (Vol. 5, pp. 238–262). Tucson: University of Arizona.
- Yu, A. (1998). *Prespecification and dissimilation in Optimality Theory: The Case of Turkish Emphatic Reduplication* (Bachelor's Thesis, University of California, Berkeley).
- Yu, A. (1999). *Dissimilation and allomorphy: The case of Turkish emphatic reduplication*. University of California, Berkeley Ms.
- Zymet, J. (2014). Distance-based decay in long-distance phonological processes. In U. Steindl, T. Borer, H. Fang, A. Pardo, B. H. Peter Guekguezian, C. O'Hara, & I. Ouyang (Eds.), *Proceedings of the 32nd west coast conference on formal linguistics* (pp. 72–81). Somerville, MA: Cascadilla Proceedings Project.
- Zymet, J. (2018). *Lexical propensities in phonology: Corpus and experimental evidence, grammar, and learning* (Doctoral dissertation, University of California, Los Angeles).
- Zymet, J. (2019). Learning a frequency-matching grammar together with lexical idiosyncrasy: Maxent versus hierarchical regression. In K. Hout, A. Mai, A. McCollum, S. Rose, & M. Zaslansky (Eds.), *Proceedings of the annual meetings on phonology* (Vol. 7), Washington, DC: Linguistic Society of America. DOI: <https://doi.org/10.3765/amp.v7i0.4495>

