

Naturalistic speech supports distributional learning across contexts

Kasia Hitczenko^{a,1} and Naomi H. Feldman^b

^aLaboratoire de Sciences Cognitives et Psycholinguistique, Département d'Études Cognitives, ENS, EHESS, CNRS, PSL Research University, Paris, France; ^bDepartment of Linguistics and UMIACS, University of Maryland, College Park, USA

This manuscript was compiled on March 26, 2022

At birth, infants discriminate most of the sounds of the world's languages, but by age one, infants become language-specific listeners. This has generally been taken as evidence that infants have learned which acoustic dimensions are contrastive, or useful for distinguishing among the sounds of their language, and have begun focusing primarily on those dimensions when perceiving speech. However, speech is highly variable, with different sounds overlapping substantially in their acoustics, and after decades of research, we still do not know what aspects of the speech signal allow infants to differentiate contrastive from non-contrastive dimensions. Here, we show that infants could learn which acoustic dimensions of their language are contrastive, despite the high acoustic variability. Our account is based on the cross-linguistic fact that even sounds that overlap in their acoustics differ in the contexts they occur in. We predict that this should leave a signal that infants can pick up on, and show that acoustic distributions indeed vary more by context along contrastive dimensions compared with non-contrastive dimensions. By establishing this difference, we provide a potential answer to how infants learn about sound contrasts, a question whose answer in natural learning environments has remained elusive.

phonetic learning | language acquisition | distributional learning

Languages differ in the speech sound inventories they use to reliably convey meaning. For example, Hindi has a distinction between unaspirated retroflex [ɖ] vs. dental [t] that is used to distinguish meanings (e.g., [ɖal] means 'postpone', while [tal] means 'beat'), but English does not. Adult speakers are generally tuned to the properties of the language(s) they speak. For example, while most adult Hindi speakers can hear the difference between [t] and [ɖ], most English-speaking adults cannot (1). Because speech sound inventories differ across languages, listeners must learn about the speech sounds of their language from the input they hear.

The first signs of this phonetic learning appear within the first year of life. During their first couple of months, infants can discriminate most sounds of the world's languages, showing similar perceptual abilities regardless of their language experience. For example, both newborn English-learning and Hindi-learning infants can hear the distinction between retroflex [ɖa] vs. dental [ta], a speech contrast that exists in Hindi, but not English. However, over the course of the first year of life, this changes. Infants become language-specific listeners, starting to more closely resemble adults in their discrimination abilities. Their ability to discriminate non-native contrasts (e.g. between retroflex [ɖa] and dental [ta] for English-learning infants) declines (2–4), whereas their ability to discriminate contrasts in their own language improves (5).

These perceptual changes have generally been taken as evidence that infants are learning which acoustic dimensions are contrastive in their language: that is, which acoustic

dimensions have multiple categories along them (6). Speech sounds differ in how they are acoustically produced and one or more acoustic dimensions will be used to signal differences between sound contrasts. The idea, then, is that infants become aware of which acoustic dimensions are used to contrast the meaningful sounds in their language, and begin primarily focusing on those dimensions when perceiving speech.

Decades of research into how infants learn about contrastiveness in their first year of life has built a wealth of knowledge in this area; however, we still do not know what aspects of the speech signal allow infants to make these inferences from the acoustically variable speech they hear in their daily lives. One of the most well-studied current proposals for how infants learn which dimensions of their language are contrastive is known as distributional learning (6). It proposes that infants learn the contrastive dimensions of their language by tracking the frequency distribution of sounds along acoustic cue dimensions. If an infant observes a bimodal (two-peaked) distribution along a dimension, then they learn that the dimension is contrastive, whereas if an infant observes a unimodal (one-peaked) distribution, then they learn that the dimension is not contrastive. This account has experimental support: distributions of sounds affect infants' discrimination in the lab (6–9). In addition, when bimodality is present in the input, computational models successfully learn correct contrasts (10, 11). However, a key assumption underlying this proposal is that contrastive dimensions do indeed exhibit bimodality, and while this is the case for some contrasts, recent work looking at naturalistic speech corpora has shown that this is

Significance Statement

Languages differ in the speech sounds they use and humans need to learn which sounds their language uses. This learning starts early. By one, infants have already tuned into their language(s): their ability to hear sound distinctions from their language improves, while they often lose the ability to hear other sound distinctions. Understanding how this early learning proceeds is important as it serves as a foundation for later development; however, it has proven difficult to identify a learning mechanism that works on the true input infants hear. Here, we present an account for how infants learn the speech contrasts of their language and show that the necessary signal is present in naturalistic speech, advancing our understanding of early language learning.

K.H. and N.H.F. designed the research; K.H. implemented the analyses; K.H. and N.H.F. wrote the manuscript

The authors declare they have no conflict of interest.

¹To whom correspondence should be addressed. E-mail: kasia.hitczenko@ens.psl.eu

58 not a universal property of child-directed speech (12, 13).

59 For example, in Japanese, vowel length is contrastive (14),
60 meaning that two different words like /toko/ (*bed*) and /toko:/
61 (*travel*) can be distinguished solely by how long a vowel is.
62 However, analyses of a spontaneous corpus of Japanese child-
63 directed speech reveal that the distribution along the duration
64 dimension is unimodal despite being contrastive (12) (Fig. 1a;
65 note that infants do not have access to the individual color-
66 coded short and long vowel distributions shown in this figure,
67 only the combined overall distribution). A similar finding has
68 been reported for Dutch vowel length (13), as well as many
69 other contrasts (5, 15, 16). That is, although infants are able
70 to use distributional information for learning when available,
71 it is not available for all of the contrasts they learn about, so
72 distributional learning is not sufficient.

73 Many follow-up theories have been proposed to explain
74 how infants learn in cases where bimodality is not present.
75 This has included theories arguing that bimodality might be
76 present when considering only the most prominent sounds (e.g.,
77 stressed vowels) (17), when normalizing for effects of neigh-
78 boring sounds or other factors (18), or when using word-level,
79 visual, or referential information (13, 19–23). While many
80 of them have experimental support and work on controlled
81 lab speech, over the past 40 years, it has proven difficult to
82 identify a learning mechanism that works on the true speech
83 infants hear. (13) takes an important step in that direction
84 by showing that, in Dutch, average vowel durations by word
85 type are often longer in word types with long vowels than
86 word types with short vowels (and, thus, that short and long
87 vowels may be separable). However, we still do not have a
88 measure that consistently separates vowels with a contrast
89 from vowels without a contrast across corpora, languages, and
90 vowel qualities. This problem is so extreme that recent work
91 has suggested that infants might not actually be learning how
92 many phonetic categories there are along a dimension at all,
93 because this signal is not present in their input in a way that
94 they have access to (24).

95 In this paper, we show that the necessary signal to learn
96 which acoustic dimensions are contrastive may be present in
97 naturalistic input and accessible to infants. Our proposal takes
98 advantage of the contextual information of a sound, which
99 infants are sensitive to (20–22, 25–29). In this work, we take
100 the context of a sound to include factors like its neighboring
101 sounds, its prosodic position in a word/utterance (i.e., if it
102 immediately borders a word or utterance boundary), and its
103 word frame; however, we think of context more broadly as any
104 information that listeners track about where a sound occurs or
105 who spoke it. When an acoustic dimension is contrastive, there
106 are multiple categories along it and the relative proportion
107 of those categories may differ across contexts (e.g., if two
108 categories are present, one context may be 50% category 1 and
109 50% category 2, whereas another context may be 90% category
110 1 and only 10% category 2). We show that such differences
111 in category frequency—which are extremely common across
112 languages (30–32)—can help infants distinguish contrastive
113 from non-contrastive dimensions.

114 We test our proposal on two test cases, Japanese and Dutch,
115 which have been most problematic for both distributional
116 learning and additional previous theories, and show that our
117 proposal explains how infants could nonetheless learn the
118 contrast from information available to them within their first

119 year of life. Complemented by previous findings that (1) infants
120 are sensitive to distribution shapes and contextual information,
121 and (2) changes in the relative proportion of sounds across
122 contexts is a cross-linguistically widespread property of sound
123 categories, these results are promising and suggest that infants
124 may be able to learn about contrastiveness from naturalistic
125 speech input, thus pointing towards a possible answer to a
126 long-standing question in the field.

127 Distributional Learning Across Contexts

128 The inspiration for our proposal comes from a finding show-
129 ing that the context a sound occurs in (neighboring sounds,
130 prosodic position, speaker, etc.) is predictive of its identity:
131 just knowing what context a Japanese vowel appears in can
132 predict its length with around 95% accuracy (33). This means
133 that short and long vowels appear in different proportions in
134 different contexts. Most contexts have almost all short vowels
135 (e.g. Context 1 in Fig. 1b), whereas some contexts have almost
136 all long vowels (e.g. Context 2 in Fig. 1b), and some are in
137 between (e.g. Contexts 3–4 in Fig. 1b). Figure 1b reveals
138 that these changes in the relative proportion of short and long
139 vowels can change the overall shape of the frequency distribu-
140 tion in the context. All of the distributions in Figure 1b are
141 unimodal, despite the fact that there are two categories. Thus,
142 they would not be conducive to the distributional learning
143 theory proposed by (6). However, this is only one aspect of
144 a distribution’s shape, and across contexts, the distributions
145 differ in how wide or peaky they are, where they peak, and so
146 forth. This arises because of two facts: (i) when a dimension is
147 contrastive, the overall frequency distribution in each context
148 is the sum of the short vowel distribution and the long vowel
149 distribution, (ii) short and long vowels have different distribu-
150 tions, as can be seen in Figure 1. Taken together, this means
151 that in a language like Japanese, where there is a contrast,
152 we would expect different relative proportions of short vowels
153 and long vowels across different contexts, and since short vow-
154 els and long vowels have different acoustic distributions, we
155 would expect the overall distribution to change across different
156 contexts. On the other hand, in a language like French (where
157 there is no length contrast), shape changes cannot arise from
158 different relative proportions of short and long vowels because
159 there is no short vs. long vowel distinction.

160 In light of this, we propose that infants might learn that
161 a dimension is contrastive by tracking the acoustic distribu-
162 tion along that dimension across different contexts. They
163 could compare the shapes of the distributions across those
164 contexts, and infer that a dimension is contrastive if the shape
165 varies substantially across contexts, but infer that it is not
166 contrastive if the shape is largely the same across contexts.
167 We operationalize a sound’s context as (i) its (immediately)
168 neighboring sounds, its prosodic position (whether it falls at a
169 word or utterance boundary), and its quality (learned before
170 length), or (ii) its word frame, due to evidence that infants are
171 sensitive to this information in their input (20–22, 25–29, 34).
172 However, we are not tied to these particular factors. Any con-
173 textual factors that infants track, and that change the relative
174 proportion of sound category membership, could work.

175 It is important to note that the learning outcome of this
176 proposal is the same as in (6), but differs from the learning out-
177 comes of some phonetic learning theories that have arisen since
178 then (19). In particular, the learning outcome here is whether

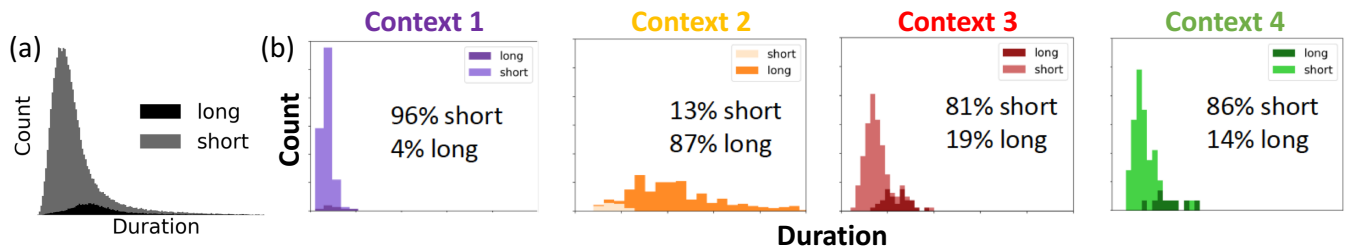


Fig. 1. (a) The frequency distribution of sounds along the duration dimension in Japanese is unimodal, despite vowel length being contrastive. (b) Vowel frequency distributions along duration, for four Japanese contexts (defined by prosodic position, neighboring sounds, and vowel quality). The relative proportion of phonemically short and long vowels changes substantially across contexts, which results in differently shaped distributions. The short vs. long categories are color-coded for the reader's benefit. Infants (and our analyses) do not have access to this color information when learning, only the overall distributions.

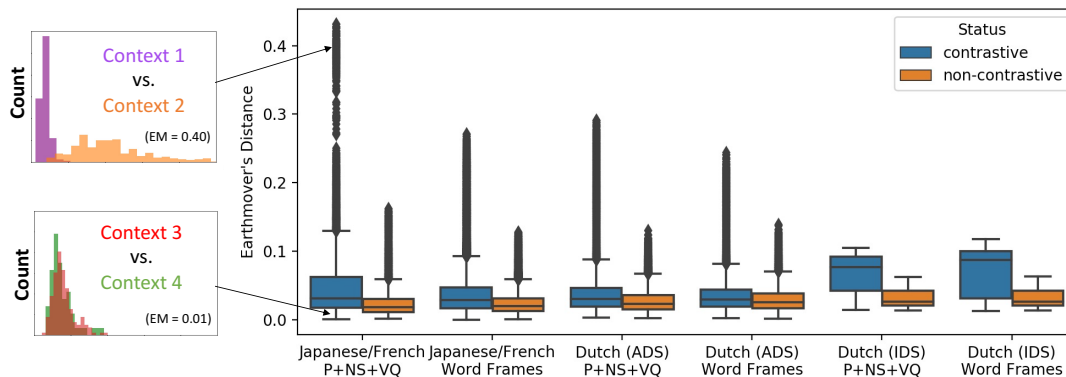


Fig. 2. Distribution of Earthmover's distances by test case. Each datapoint represents the pairwise Earthmover's distance (EM) between distributions from two different contexts (e.g., we show the comparison for Fig. 1b's Context 1 vs. Context 2, which has a high Earthmover's distance, and Fig. 1b's Context 3 vs. Context 4, which has a small Earthmover's distance). Across all test cases, the tail of the contrastive boxplot (left) is longer than that of the non-contrastive boxplot (right), suggesting that there are more extreme distribution shape changes across contexts when the acoustic dimension is contrastive. P+NS+VQ = prosodic position + neighboring sounds + vowel quality.

179 or not an acoustic dimension is contrastive - i.e. whether it is
 180 used to distinguish multiple categories. The learning outcome
 181 of some other theories included this knowledge implicitly, but
 182 often also included knowledge about what those categories
 183 were. Certainly listeners eventually learn about the categories,
 184 and a full learning account would need to eventually explain
 185 how that happens. However, the discrimination behavior infants
 186 exhibit in (2) does not require them to have learned
 187 categories (2, 24, 35), so we follow the original literature and
 188 focus on how infants learn which dimensions are contrastive
 189 in their language.

190 In what follows, we ask whether the necessary signal for
 191 this learning account is present in naturalistic speech; that
 192 is, whether there are larger distribution shape changes across
 193 contrastive dimensions than non-contrastive dimensions. We
 194 focus on three test cases, which each involve some data in
 195 which vowel length is contrastive, and some data in which
 196 vowel length is not contrastive. We look at vowel length for
 197 two reasons. First, it has a largely agreed upon primary cue
 198 (duration) that can be easily extracted from any annotated
 199 corpus. Second, it is possibly the best known case of extreme
 200 overlapping categories that cannot be explained by previous
 201 theories like distributional learning (12). We focus on the test
 202 cases that have been problematic for past phonetic learning
 203 theories, but argue in the General Discussion that this same
 204 approach to phonetic learning is likely to be effective across a
 205 wide range of languages and contrasts.

206 Results

207 Our results confirm that in spontaneous speech there are more
 208 extreme distribution shape changes across contexts when a

dimension is contrastive than when it is not.

209 For all of the corpora we study, we extract the acoustic
 210 distributions across a number of contexts, and compare them
 211 pairwise, using Earthmover's distance (36), a commonly used
 212 metric of distribution shape difference (see Supplementary
 213 Materials for discussion of results using an alternative metric,
 214 KL divergence, instead). We operationalize 'context' in two
 215 different ways, both of which rely on information that infants at
 216 the relevant age are sensitive to: (i) a combination of prosodic
 217 position, neighboring sounds, and vowel quality (P+NS+VQ)
 218 and (ii) word frames (WF).

219 We first compare a spontaneous speech corpus of Japanese
 220 (which has a vowel length contrast) against a spontaneous
 221 speech corpus of French (which does not). We then test two
 222 spontaneous Dutch corpora. Dutch has the property that a
 223 subset of its vowels has a length contrast, whereas a different
 224 subset does not. Comparing the subset that has a contrast
 225 against the subset that does not allows us to control for any
 226 effects that may arise due to differences in how the French and
 227 Japanese corpora were collected and annotated. Two of our
 228 tests examine adult-directed speech (ADS) corpora because
 229 they allow us to test this proposal on large-scale, spontaneous
 230 speech corpora which do not exist for infant-directed speech
 231 (IDS), but we include results from a small corpus of infant-
 232 directed Dutch as well.

233
 234 **A. Japanese vs. French ADS.** We first compared Japanese and
 235 French, defining context as a combination of prosodic position,
 236 neighboring sounds, and vowel quality (Fig. 2). Each data-
 237 point contributing to the boxplot represents the Earthmover's
 238 distance between a pair of contextual acoustic distributions.
 239 For example, the comparison between Context 1 and Context

Context	Percent Long	Count	Frequency Rank
Phrase-initial, word-final /e/	64.7	1357	18
Phrase-initial, phrase-final /a/	56.7	255	95
Phrase-initial, phrase-final /e/	87.9	244	100

Table 1. Information about the Japanese contexts that drive the tail in the case of the P+NS+VQ analysis, including what percentage of vowels in that context are long, how many times that context occurred (Count), as well as its frequency rank out of all contexts that occurred.

2 in Fig. 1b has a high Earthmover’s distance, whereas the distance for Contexts 3 and 4 is much small because they are very similar.

The boxplot corresponding to Japanese (where vowel length is contrastive) has a much larger tail, extending upwards towards large Earthmover’s distances, than the boxplot corresponding to French (where vowel length is not contrastive). This means that, as predicted, there are many more pairs of contexts that have substantially different shapes (like Context 1 vs. Context 2 in Fig. 1b) when there is a contrast than when there is not. The maximum distance, the mean distance, and the distance variance are all larger for Japanese than French (max = 0.43 vs. 0.16; mean = 0.05 vs. 0.02; variance = 0.003 vs. 0.0004). Analyzing the contents of the tail in Japanese reveals that the tail is driven by contexts that have a much higher percentage of long vowels than observed overall and that occur frequently in the input (see Table 1 for frequency counts and ranks of the contexts that drive the signal).

These same patterns hold when we continue looking at French vs. Japanese, but instead use word frames as contexts. As before, there are more contexts with more extreme distribution shape changes in Japanese than French (i.e. along contrastive than non-contrastive dimensions), as seen by the longer tail in the second pair of boxplots in Fig. 2. As before, the maximum distance, the mean distance, and the distance variance are all larger for Japanese than French (max = 0.27 vs. 0.12; mean = 0.04 vs. 0.02; variance = 0.001 vs. 0.0002).

B. Dutch ADS and IDS. To test our proposal using a within-language comparison, we compare the subset of Dutch vowels that do contrast in length and the subset of Dutch vowels that do not. We find that the predicted pattern still holds – and it holds for both ways of defining context and both the ADS and IDS corpora (Fig. 2). This confirms that the results are not merely an artifact of using different corpora, as in the French vs. Japanese case, but seem to reflect something inherent to the existence or nonexistence of categories along an acoustic dimension. It is worth noting that Dutch-learning infants would not be able to perform this exact analysis to learn whether there is a length contrast, because they would not yet know enough to separate the vowels into contrastive and non-contrastive subsets. We return to the issue of what a learning account might look like in the Discussion. Meanwhile, we conclude from this analysis that the signal our account predicts exists in Dutch: contrastive dimensions differ from non-contrastive dimensions.

Despite the qualitative similarity in results across all test cases, the scale of the difference in tail length varies. For example, in the French vs. Japanese P+NS+VQ case, the maximum Earthmover’s distance in the contrastive Japanese case is 0.43, whereas for the other ADS cases, the maximum is less than 0.3. In the Dutch IDS corpus, which only has 284

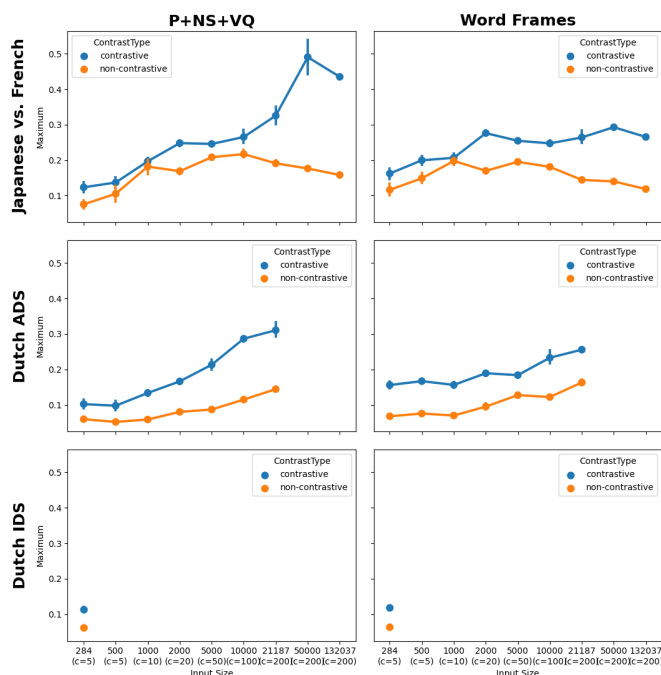


Fig. 3. Bootstrapped analyses reveal that observed differences between contrastive (top/blue line) vs. non-contrastive (bottom/orange line) dimensions are meaningful, but that input size does matter. “C” refers to the number of contexts included in the analysis. The maximum input size for which data is shown depends on the corpus size: 284 for Dutch IDS, 21187 for Dutch ADS, and 132037 for Japanese vs. French.

vowel tokens, the maximum is only around 0.1. One possibility is that these differences arise because of the large differences in corpus size. The Japanese vs. French corpora considered 132,037 tokens and the Dutch ADS corpus considered 21,187 tokens, but the Dutch IDS corpus only considered 284 tokens.

C. Corpus size analyses. To test how corpus size impacts results, we used bootstrap samples to run each analysis 50 times for 10 different corpus sizes ranging from the size of smallest corpus (284) to the size of the largest corpus (132,037). This also allowed us to test how much the size of the tail varied, and whether differences observed between contrastive vs. non-contrastive cases were meaningful. Fig. 3 shows these results when calculating the maximum Earthmover’s distance across all 50 runs; analogous plots for mean are provided in the Supplementary Materials. First, this analysis reveals that the differences observed are meaningful: across many runs, at large enough corpus sizes, the contrastive line is higher than the non-contrastive line. That being said, in the Japanese vs. French case, the difference does not emerge until around 2000 vowel tokens have been observed, so input size does matter. Second, this analysis reveals that differences in scale may be partially, but are not entirely, due to corpus size. When subsetting to the size of the Dutch ADS corpus, the Japanese vs. French word frame maximum matches the remaining ADS results. However, the results are less clear for Dutch IDS: subsetting the Dutch ADS corpus to the size of the Dutch IDS corpus yields results more in line with each other for the P+NS+VQ analysis, but less so for the word frames analysis.

From a learning perspective, this means that an ideal learner would need to observe around 2000 vowel tokens and track the acoustic distribution within the 20 most frequent contexts in order to observe the difference (though we discuss potential ways to reduce the memory demands of the proposal next as

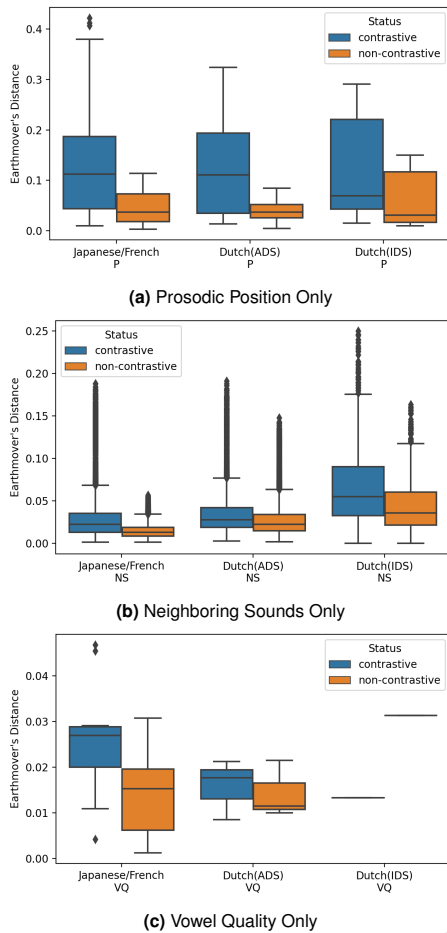


Fig. 4. Results are similar when we relax our assumptions about infants' knowledge. Here, instead of studying a combination of factors, we study (a) prosodic position (P), (b) neighboring sounds (NS), and (c) vowel quality (VQ) individually (left: contrastive dimension; right: non-contrastive dimension). While less clear for vowel quality, we see the same pattern of results in the prosodic position and neighboring sounds only cases.

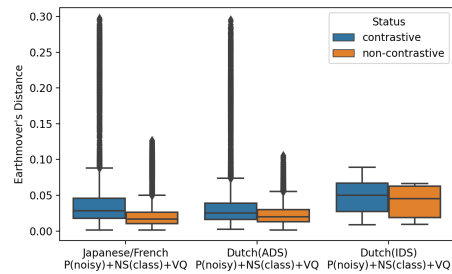


Fig. 5. Results are similar when we relax our assumptions about infants' knowledge. Here, we define neighboring sounds by their broad class (stop, fricative, etc.) and add noise to the prosodic position with a noise rate of 20% (left: contrastive dimension; right: non-contrastive dimension).

distributions across combinations of contexts, these results suggest that this need not be the case in order for our account to be successful. It suffices to track distributions across individual contexts (e.g., word frames or prosodic position). From a learning perspective, this means that infants would only need to be tracking the acoustic distribution across ~ 4 -10 contexts (rather than ~ 20 as observed before), and these contexts include extremely prominent contexts (e.g., utterance-final, utterance-initial tokens).

The next two assumptions we revisit are whether infants can perfectly encode the identity of neighboring sounds, and whether they have a solid enough grasp on word segmentation to have access to the prosodic position information we use. On the one hand, prosodic boundaries are one of the first signals that infants are sensitive to (26, 27); however, there is concurrent evidence that infants make missegmentation errors (37-39). To address these two assumptions, we test what happens when we re-run the P+NS+VQ simulations, with neighboring sounds defined by their broad class (i.e., stop, fricative, vowel, approximant, etc.) rather than their particular identity (e.g., /k/, /g/, /b/) and with noise added to the prosodic information (we simulate a 20% error rate here, such that 20% of the time, the infant misrepresents the prosodic position of the vowel, but the results generalize across error rates). Fig. 5 shows that these differences do not qualitatively change the results, suggesting that even with a more rudimentary grasp on contextual factors, infants could still use this method to learn the sound contrasts.

Overall, we show that the necessary linguistic knowledge and capabilities can be considerably reduced and yet the correct finding still emerges. This suggests that this finding is a robust one that immature learners could learn from even in noisy learning environments.

E. Analyses with long vowels removed. Finally, to test whether these results arose because of the contrast, we removed all vowels labeled as long from the corpora and reran the same analyses. We predicted that removing the long vowels would cause the tail for the contrastive dimension to disappear, such that the results for the contrastive dimension with long vowels removed (i.e. with the contrast artificially removed) would resemble those of the non-contrastive dimension. As can be seen in Fig. 6, in Japanese and in Dutch IDS (the two hand-annotated corpora we use), the tail disappears or is reduced once the long vowels are removed, suggesting that it is at least partially the presence of the long vowels that causes these large changes in distribution shape. However, this is not case in Dutch ADS: the contrastive dimension still

well as in the General Discussion).

D. Relaxing our assumptions about infants' knowledge. Our analyses so far have been conducted assuming that (a) infants can track acoustic distributions across combinations of three contexts (prosodic position, neighboring sounds, and vowel quality), (b) infants can perfectly identify neighboring consonants, and (c) infants can perfectly segment words from speech. Although some of these assumptions have yet to be tested (e.g., we don't know whether infants can track distributions along multiple contextual dimensions), we know that others are likely overestimating infants' prior knowledge when learning about contrastive dimensions. Here, we show that the same qualitative results still emerge even when we weaken these three assumptions.

The first assumption we revisit is whether infants can track distributions across combinations of contexts (i.e., prosodic position, neighboring sound, and vowel quality). We test what happens when we study each of these three factors individually. Fig. 4 shows that, especially for prosodic position, but also for the other factors, the same patterns, for the most part, emerge. While it is still quite conceivable that infants track

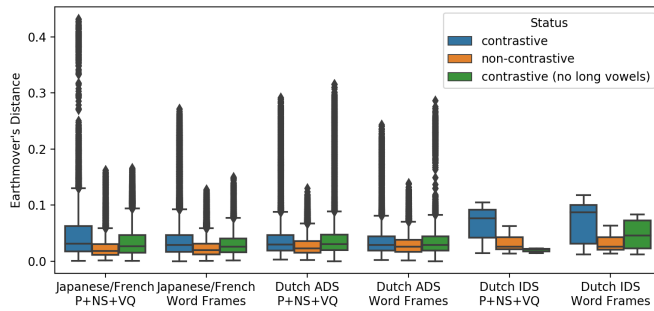


Fig. 6. Results from artificially removing long vowels. Within each case, the left boxplot corresponds to the contrastive dimension, the middle boxplot corresponds to the non-contrastive dimension, and the right boxplot corresponds to the contrastive dimension with long vowels removed. We observe that the tail length is reduced for the Japanese/French and Dutch IDS analyses; however, against predictions, not for the Dutch ADS analysis.

they learn that an acoustic dimension is contrastive if the distribution shape along that dimension varies substantially across different contexts. For this account to work, it needs to be the case that the distribution shape varies more across contexts when a dimension is contrastive than when it is not. We tested this prediction in three test cases, with two different definitions of context. Across the board, our results show that the distribution shape along an acoustic dimension changes more across different contexts when that dimension is contrastive than when it is non-contrastive. This is a signal that differentiates contrastive and non-contrastive dimensions, and it is the kind of signal that listeners are likely to be sensitive to. As such, this is one of the first phonetic learning accounts that has been shown to work on spontaneous data and suggests that infants could be learning which acoustic dimensions are contrastive after all. In the remainder of the paper, we discuss the promise and open questions of this proposal, including its generalizability, evidence on whether infants have the necessary sensitivities, and how infants could use this signal to learn.

Generalizability. We considered the test cases of Japanese and Dutch vowel length because they are famous problem cases for many of the phonetic learning theories that already exist (12, 41). However, they are unique contrasts in a number of ways. They have *low functional load* in that they are not frequently used to distinguish different meanings and they have particularly overlapping acoustic distributions. In addition, the Japanese contrast is primarily signaled by one acoustic cue (i.e. duration), while most contrasts are signaled by more (42, 43) and, as around 90% of Japanese vowels are short, it is less balanced than many other contrasts. Future work will need to test the generalizability of this proposal. Nonetheless, we think it is likely that this signal will generalize to other contrasts because the signal we illustrate in Fig. 1b does not stem from these idiosyncrasies. Rather it arises because of a handful of contexts that have particularly different frequency distributions and we think these exist because of properties of language that phonologists have argued are universal (30).

All languages are thought to have phonotactic or co-occurrence constraints. For example, for the English contrast [n]-[ŋ] (e.g. *sin* vs. *sing*), [n] can occur at the beginning or end of a syllable, while [ŋ] can only occur at the end of syllables. This means that the relative proportion of [n]'s and [ŋ]'s will change across those contexts, and could show patterns like in Fig. 1. Similarly, sometimes sounds will be pronounced differently based on their context. For example, if we consider the contrast between [n]-[m], the [n] sound never occurs before the sounds [p,b], but is instead pronounced as [m] in those contexts (as in *impossible*). These types of phonological alternations also create differences in which sounds occur in which contexts and are cross-linguistically widespread (30–32). Finally, there are systematic regularities based on the words in the language. For example, for the [p]-[b] contrast in English, [b] is more likely in the word frames *_aby* and *_ar* (*baby* is a word, but *paby* is not and *bar* is more frequent than *par*), whereas [p] is more likely in the word frames *_lay* and *_in* (*play* is a word, but *blay* is not and *pin* is more frequent than *bin*). Taken together, phonotactics, phonological alternations, and word regularities all create systematic regularities in the contexts that sounds occur in, such that different sounds occur in different contexts and different contexts are made up of

has a longer tail even when long vowels are removed. While this is not predicted by our account, there are a number of reasons why we may observe this result here. One possibility is that the annotations are imperfect. The remaining datasets studied here include hand-corrected segmental annotations, while the Dutch ADS data included force-aligned annotations which were not validated for their duration. Another possibility is that it has to do with the fact that these data were phonetically annotated, rather than phonemically annotated. That is, some phonemically long vowels were marked as being short. In fact, (40) reports that as many as 20% of word tokens that had long vowels underwent a shortening process. Especially combined with the fact that these data were automatically annotated, this could mean that we are unable to actually remove all long vowels and that some long vowels remain which are driving the differences observed. Finally, another possibility is that the presence of long vowels in a context changes the distribution of short vowels. For example, if a context is 50% short vowels and 50% long vowels, then the short vowels may be pronounced with shorter durations than in a context with 90% short vowels and only 10% long vowels (i.e. to better differentiate the vowel types). If this is the case, then even if we remove the long vowels, we should expect to see differences in the short vowel distributions across contexts. However, this finding is nonetheless different from what our account predicts and more work should be done with hand-annotated corpora to see whether this finding remains and, if so, what is driving it.

Nonetheless, across the board, we do observe that the contrastive dimensions have a longer tail than the non-contrastive case, suggesting that contrastive dimensions exhibit more extreme distribution shape changes across contexts (even though all distributions are unimodal) than non-contrastive dimensions. These results show that contrastive dimensions look different than non-contrastive dimensions, and that infants would learn the correct generalization about their language by using this signal. As such, this is one of the first pieces of signal that has been shown to successfully differentiate contrastive and non-contrastive dimensions using spontaneous speech.

Discussion

In this paper, we proposed a new account for how infants could learn which acoustic dimensions of their language are contrastive. The idea is that infants track the frequency distributions of sounds in different contexts, and that

497 different ratios of sounds. This is currently thought to be true
498 of all contrasts - even those that are equally balanced, have
499 higher functional load, or are multidimensional - and will be
500 the case in adult- as well as infant-directed speech. Indeed, (33)
501 showed the context is similarly predictive of which category a
502 sound belongs to in both adult-directed and infant-directed
503 speech.

504 Nonetheless, while our results are promising, they ultimately
505 come from one acoustic dimension, one contrast type, and only
506 three languages. Furthermore, there are additional complexities
507 that we have not considered here, like how extralinguistic
508 factors (e.g., speaking style/register, emotional content of
509 speech, speech rate) could affect the signal and how this
510 proposal interfaces with other learning strategies that infants
511 have been shown to use (e.g., using visual and referential
512 information).

513 To test generalizability, future work should replicate these
514 findings on other data sets as they become available (we have
515 made our code publicly available to facilitate this effort). Particular
516 focus should be placed on replicating these results on a large
517 ecological corpus of infant-directed speech, replicating these
518 results on other contrasts whose acquisition is difficult to explain
519 (e.g., Filipino nasals (44)), studying Dutch to understand why
520 removing long vowels does not always change the result, as well
521 as replicating these results in languages where the identity of a
522 particular sound cannot be predicted as well from its context, due
523 to having fewer phonotactic and other systematic restrictions (45).
524 It will also be important to test this proposal on contrasts that
525 are signaled by multiple acoustic dimensions (as the contrasts we
526 study here are unique in primarily being signaled by one: duration).
527 In order for this account to work in those cases, we would need to
528 observe the signal we report here along at least one of the acoustic
529 dimensions that signal the contrast (e.g., along VOT or F0 for stop
530 contrasts). Alternatively, rather than operating over individual
531 acoustic dimensions, infants could search for this signal along
532 composite acoustic dimensions that are discovered from the input
533 by combining dimensions that are highly correlated in the input
534 (e.g., as discussed in (35, 46)). Finally, we note that Earthmover's
535 distance can be calculated over multidimensional distributions, so
536 it should be possible to scale this approach up to multidimensional
537 contexts if so desired; however, for that to work, the infant would
538 have to first identify the relevant combinations of acoustic dimensions
539 they should focus on. Overall, it will be important to replicate
540 these findings across highly variable corpora that adequately
541 represent the full range of speech types that infants could
542 encounter.

543 Finally, it will also be important to test that these results
544 do not *overgeneralize* and wrongly label non-contrastive dimensions
545 as contrastive. Here, it will be particularly important to test
546 behavior on allophonic variation, where a particular sound is
547 realized differently depending on the context it occurs in. While
548 this pattern is similar to that of different phonemes occurring
549 in different contexts, our analyses provide preliminary evidence
550 that allophonic variation is not labeled as contrastive. In particular,
551 in French, vowel length varies allophonically (vowels are
552 lengthened depending on the following consonant) (47), yet our
553 analyses reliably treated French (allophonic variation in vowel
554 length) differently from Japanese (phonemic variation in vowel
555 length). This suggests that this method

556 may correctly differentiate contrastive and allophonic variation,
557 though it will be important to study this further and understand
558 how it does so (if it does). We offer two speculative reasons
559 why allophonic variation may not be detected to the same degree
560 as phonemic variation, though they will need to be tested. First,
561 this proposal relies on contrasting sounds having sufficiently
562 different acoustic distributions (so that changing the relative
563 proportion of the sounds changes the shape of the distribution).
564 It is possible that allophonic variation changes how a sound is
565 produced to a lesser degree than phonemic variation, though this
566 may be difficult to assess. A second possibility is that allophony
567 often, though not always, affects all of the sounds produced in
568 a particular context (e.g., the duration of all French vowels is
569 affected by neighboring consonants). This may lead to a shift
570 in distribution between contexts, without a change in shape,
571 which may lead to smaller distribution shape changes as measured
572 by Earthmover's distance. Certainly, more research studying how
573 this proposal handles allophony will be crucial.

Could infants do this? This proposal places higher computational
574 and memory demands on infant learners than many past theories
575 have (e.g., distributional learning). Infants would need to be
576 able to track distribution shapes across many contexts and then
577 compare their shapes pairwise.

(1) *Could infants track distributions across different contexts?*
578 Though this has not been tested, we know that listeners track
579 the shape of frequency distributions overall, and use this for
580 phonetic learning and processing. For example, infants make
581 different phonetic inferences depending on whether the distribution
582 they hear is bimodally- or unimodally- shaped (6). Another
583 study showed that adult listeners are sensitive to the variance of
584 the sound distributions they are exposed to, another property of
585 a distribution's shape (48). In that experiment, adults categorized
586 sounds differently depending on whether they heard a distribution
587 with high or low variance, though this has not been tested in
588 infants.

589 In addition, we know that listeners are sensitive to the context
590 of a sound and use it for phonetic learning and processing. Both
591 toddlers 12 months and older and adults have been argued to
592 track acoustic distributions across speakers (which can be thought
593 of as a context), can adapt to speakers who have different accents
594 (i.e. different distributions of sounds) (49-55), and mirror the
595 speech of their interlocutors. In addition, infants are sensitive
596 to phonotactics (56), as well as phonological alternations - the
597 fact that sounds tend to be pronounced differently in different
598 contexts (57, 58). Additionally, multiple studies have shown
599 that infants use the word frame of a sound in phonetic learning
600 (20, 22). That is, infants seem to assign acoustically similar
601 sounds to different categories if they occur in different word
602 frames, suggesting that infants can track the context that a
603 sound occurs in and use it for phonetic learning. Adult speech
604 perception is affected by contextual factors, like neighboring
605 sounds. For example, (59) showed that Japanese perception of
606 whether the final vowel in CoC'V was phonemically short or
607 long depended on the identity of both C and C'. Furthermore,
608 (60) showed that adults continually track how informative
609 particular acoustic cues are and will selectively reweight these
610 cues in some contexts but not others (e.g., when categorizing
611 /b-p/ in beer-pier, but not for /d-t/ in deer-tier and vice
612 versa). Indeed, listeners may even track information across
613 contexts defined by multiple factors.

619 For example, in English, to learn whether a stop consonant
620 will be aspirated or not, the speaker must track whether the
621 sound is voiceless or voiced (akin to Vowel Quality), whether
622 it occurs at the start of a stressed syllable or not (akin to
623 Prosodic Position), and what the neighboring sounds are and
624 then notice the change in pronunciation/distribution across
625 those different contexts.

626 Taken together, while these results do not provide direct
627 evidence that infants track distribution shapes across contexts,
628 they suggest that listeners can track complex statistical reg-
629 ularities across complex contexts and use them in real-time
630 phonetic learning and processing.

631 (2) *Could infants compare distribution shapes?* Finally, the
632 last skill necessary for the proposal is for infants to be able to
633 compare distributions. While this has again not been directly
634 tested, one possible clue is that listeners seem to reweight
635 acoustic cues depending on how variable/informative they
636 are, with cues that have narrower distributions being more
637 informative than cues with wider distributions (48, 61, 62).
638 Another possible clue is that toddlers and adults are able
639 to identify when they need to adapt their representations to
640 speech they hear. Being able to identify an accent implies
641 that listeners can identify when the speech they are hearing
642 differs from the speech they usually hear, a computation that
643 is likely to involve tracking at least some properties of the
644 distributions (63).

645 Overall, given infants' demonstrated sensitivity to distribu-
646 tion shape and to changes across context, there is good reason
647 to believe that infants could be sensitive to the type of distribu-
648 tional information that our account assumes, but future work
649 should test whether infants/listeners can track distributions
650 across different contexts (defined by one or more contextual
651 factors) and compare distributions' shapes. In addition to test-
652 ing whether infants can, in theory, perform the computations
653 this account requires, future work should also test whether
654 infants actually use them to learn about contrastiveness in
655 the way we propose here. One approach would be to test this
656 experimentally, by exposing infants to acoustic distributions
657 that differ or remain the same across contexts and seeing if
658 this affects their learning/behavior. Another approach would
659 be to use cross-linguistic corpora to identify contrasts that
660 should be easier/harder to learn according to our proposal
661 and compare that against age of acquisition and speech per-
662 ception/production data. For example, controlling for degree
663 of acoustic overlap, this theory would predict that contrasts
664 that have stronger phonotactic restrictions or that are more
665 predictable (i.e., it is easier to predict which member of the
666 contrast occurred based on the context it occurred in) should
667 be easier to learn through this method. These approaches
668 will allow us to overcome the next big hurdle for this account,
669 which is determining whether infants use this signal to learn.

670 **Reducing the computational complexity of the proposal.** It is
671 also possible that the memory and computational restrictions
672 of the proposal could be reduced. On the one hand, we saw
673 that considering individual contexts (e.g., just the most fre-
674 quent word frames or just prosodic position) was still effective,
675 as was introducing some parsing errors and considering broader
676 segment classes rather than individual neighboring sounds. On
677 the other hand, online approximations or metrics that do not
678 require the whole distribution to be tracked in order to get a
679 measure of distribution shape distance could also reduce the

680 computational and memory complexity of the proposal. For
681 example, rather than exactly representing the distribution,
682 this proposal could operate over a compressed representation
683 of the distribution that keeps track of how many points fall
684 within larger bins/bands (similar to reducing the number of
685 bins in a histogram). It is likely that we would still observe
686 the critical pattern even with this less detailed representation
687 of distribution shape, and, as this only requires keeping track
688 of one number per bin (its count), it could reduce the size
689 of the representation of a distribution to just 5-10 numbers.
690 Even considering all 200 of the most frequent contexts, this
691 could involve storing as few as 1000 numbers. In addition,
692 there may be a way to zero in on the contexts that yield
693 the necessary signal without doing all of the pairwise com-
694 parisons represented in the boxplots. It is possible that the
695 "key" contexts that drive the signal are overrepresented in the
696 outliers of the overall distribution (i.e., particularly short or
697 particularly long vowels could be more likely to occur in a
698 context that drives the tail). If this were true, infants could
699 arrive at the same signal we observe here, by focusing in
700 on the contexts of outlier sounds, rather than tracking the
701 distribution across all contexts. Finally, another possibility
702 is that infants could compare contextual distributions using
703 higher-order measures of distribution shape (e.g., variance)
704 rather than tracking the entire distribution. The fact that the
705 distribution shape changes across contexts could also mean
706 that the variance of the distributions changes across contexts.
707 If so, infants could pick up on this difference without encoding
708 the entire distribution across contexts.

Moving from signal to learning account. The data we report
709 plots Japanese and French side-by-side, but most Japanese
710 infants do not get French input to compare against. Assuming
711 that this pattern generalizes to other contrasts and that infants
712 have the necessary sensitivities to detect this signal, how could
713 infants actually use it to learn?
714

715 One possibility is that infants use a built-in threshold to
716 determine whether a dimension is contrastive: if the metric
717 (this could be something like the average, range, variance, or
718 maximum Earthmover's distance) exceeds the threshold, they
719 learn the dimension is contrastive; otherwise, they learn that
720 it is not. Another possibility is that infants compare against
721 other acoustic dimensions of their own language (instead of
722 against other languages, as we did). If these metrics turn out to
723 be larger for all contrastive dimensions than all non-contrastive
724 dimensions, infants could easily separate contrastive vs. non-
725 contrastive dimensions.

726 One complication for this possibility, however, is that the
727 metrics we report are sensitive to the scale of a dimension,
728 making it difficult to compare across dimensions with different
729 scales (e.g., formants vs. duration). To overcome this problem,
730 we tried z-scoring the dimensions, but found that the key
731 effect partially disappeared: the effect was retained for French
732 vs. Japanese and when considering only the most frequent
733 contexts, but when more contexts were considered, the pattern
734 reversed in Dutch ADS. This happened because z-scoring is
735 sensitive to variance, and as there was more variability along
736 the contrastive than non-contrastive dimensions, z-scoring
737 led to artificially lowered Earthmovers' distances along the
738 contrastive (more variable) dimension. Nonetheless, a method
739 that standardizes the scales with less sensitivity to the overall
740 variance could allow for comparison across dimensions.

741 Finally, given the overall distribution along an acoustic 802
742 dimension, infants could have a probabilistic model of how 803
743 different they should expect distributions to be across contexts, 804
744 depending on whether the overall distribution is made up of one 805
745 vs. two categories. With this, they could compare how likely a 806
746 one-category vs. two-category solution is to have generated the 807
747 observed Earthmover’s distances (controlling for phonological
748 processes, it would be unlikely for one category to produce
749 extremely different distribution shapes as in Fig. 1b).

750 One issue that should be considered in future work is 809
751 whether learners consider aggregate distributions across all 810
752 of the vowel qualities when comparing distribution shapes 811
753 across contexts. In our analyses of Japanese and French, 812
754 we have assumed that they do, and in fact, vowel quality 813
755 is one of the contextual factors we analyze. However, this 814
756 creates a potential problem for Dutch, because only a subset 815
757 of Dutch vowels contrast in length. If Dutch infants were to 816
758 initially rely on the aggregated distribution shape comparison 817
759 approach we put forth for Japanese, this might lead them to 818
760 conclude that vowel duration is a contrastive dimension that 819
761 they should tune into (as the context pairs that showed high 820
762 Earthmover’s distance would still be in this analysis), but 821
763 they would not realize that only some of the vowels contrast. 822
764 It is possible that they could later learn which specific vowel 823
765 qualities contrast in length. It is also possible that Dutch 824
766 infants use a different strategy entirely for discovering the 825
767 vowel length contrast: those Dutch vowels that contrast 826
768 in length also contrast in vowel quality (e.g., [a] vs. [a:] 827
769 contrast in Dutch, but [a] vs. [a:] contrast in Japanese). 828
770 Having already separated [a] from [a:] using their vowel 829
771 qualities, Dutch infants could simply notice that these vowels 830
772 differ systematically in their durations, without doing any 831
773 distribution shape comparisons of the type we propose here. 832
774 Ultimately, we remain agnostic as to how exactly infants learn. 833
775 Given the complexity of the task infants are faced with as 834
776 well as past experimental findings showing infants use many 835
777 types of information in phonetic learning, it likely involves 836
778 a combination of strategies (e.g., using word-level, visual, 837
779 referential, and other distributional information in addition to 838
780 the types of analyses we report here). Having established that 839
781 a signal exists in naturalistic data, we hope future research 840
782 will investigate how this signal might best be used for learning 841
783 and how it integrates with other promising accounts. 842

785 Conclusion

786 Infants need to learn which acoustic dimensions of their 843
787 language are contrastive in order to learn the sound system of 844
788 their language. However, we still do not know what aspects 845
789 of naturalistic input provide the necessary signal for them 846
790 to do so. In this paper, we propose a potential account for 847
791 how infants learn this and show that there is a signal about 848
792 whether a dimension is contrastive in noisy, spontaneously 849
793 produced input. This account is particularly promising for 850
794 two reasons. First, the signal that we pick up on is a direct 851
795 consequence of multiple categories exhibiting properties that 852
796 hold true across most languages, so we think it is likely that 853
797 this result will generalize to other contrasts. In addition, 854
798 the signal is something that even infants may be sensitive 855
799 to. Past work has shown that infants track the shapes of 856
800 overall frequency distributions, and know about how sounds 857
801 are likely to sound in different contexts (6, 57, 58). Adults have

802 been shown to track distributions across situations (i.e. across 803
804 different talkers) (49). In conclusion, we show that even when 805
806 two sounds overlap acoustically, the fact that they occur in 807
808 different contexts leaves signal to their contrastiveness. These
809 results provide initial support for a phonetic learning account
810 that works on highly acoustically variable spontaneous speech.

808 Materials and Methods

809 **Methods.** For each test case, one of the datasets (contrastive or 809
810 non-contrastive) was larger than the other. To correct for this, we 811
812 only considered the first N tokens of the larger dataset, where N 813
814 was the size of the smaller dataset. We extracted the duration, 814
815 the primary acoustic cue to length, of each vowel token in seconds, 815
816 rounding to the same degree of precision. In addition, we extracted
817 all contextual information that was available across all of the corpora
818 we study and that infants of the relevant age are sensitive to: 816

- 817 • **Vowel quality:** For Japanese, this was: /a/, /e/, /i/, /o/, or 817
818 /u/. For French, this was: /a/, /e/, /i/, /o/, /u/, /y/, /ø/, 818
819 /ã/, /õ/, or /ô/. For Dutch, this was: /a-a/, /ɔ-o/, /œ-ø/, 819
820 /ɛ-e/, /i/, /u/, /y/, or /i/. The first four listed pairs are 820
821 differentiated by quality and length, but we do not incorporate 821
822 these vowel quality differences into this paper. Vowel quality 822
823 is thought to be learned before vowel length (25). 823
- 824 • **Prosodic position:** We represented prosodic position (a 824
825 vowel’s position relative to prosodic boundaries) with four 825
826 indicator values: (1) whether the vowel was word-initial or not, 826
827 (2) whether the vowel was word-final or not, (3) whether the 827
828 vowel was phrase-initial or not, and (4) whether the vowel was 828
829 phrase-final or not. Infants have been shown to be sensitive to 829
830 prosodic boundaries quite early (26, 27). 830
- 831 • **Neighboring sounds:** We extracted the identity of the 831
832 immediately previous sound and the immediately following sound, 832
833 as labelled by the phonetic transcription, ignoring length infor- 833
834 mation. Again, vowel length contrast is thought to be learned 834
835 later than other types of contrasts (25). 835
- 836 • **Word frame:** We extracted the word frame that the vowel 836
837 occurred in, excluding all length information. For example, 837
838 one word frame could have been [b_i_ru], which would include 838
839 both [biru] and [biru]. We chose to include word frames, as 839
840 infants know and can segment words early (28, 29, 34, 64), and 840
841 use word frames in phonetic learning (20–22). 841

842 We looked at two main ways of defining context, though we 842
843 do not have any commitments about which contexts infants would 843
844 compute over. In the first way of defining context, we used a 844
845 combination of vowel quality, prosodic position, and neighboring 845
846 sounds (e.g. /o/ vowels that follow a /t/ and precede a /k/ that 846
847 are word- and phrase-internal) - though we also consider each of 847
848 these three contextual factors individually. This combined set of 848
849 factors corresponds to the subset of factors considered in (33) that 849
850 were available for the corpora we study and that infants are most 850
851 sensitive to. In the second way of defining context, we used word 851
852 frames, as has been done in (19) and (13) among others. 852

853 Because most contexts occur very infrequently, we looked at 853
854 a subset of all possible contexts. We subsetted the contexts in 854
855 two qualitatively different ways: either by taking the top X most 855
856 frequent contexts, or by taking all contexts that had at least N 856
857 tokens, varying X and N. Results were qualitatively similar in all 857
858 cases, so we present results from including the 200 most frequent 858
859 contexts for the French vs. Japanese and Dutch ADS analyses, and 859
860 only the 5 most frequent contexts for the Dutch IDS analysis due to 860
861 its much smaller size. Once we had the contexts, we extracted the 861
862 vowel duration frequency distributions in each context (examples 862
863 shown in Fig. 1). 863

864 We compared the shape of each pair of contextual frequency 864
865 distributions, using a metric known as Earthmover’s distance or 865
866 Wasserstein distance (65, 66), which is commonly used to measure 866
867 the difference in shape between two distributions (see Supplemen- 867
868 tary Materials for methods and results using KL divergence instead). 868
869 Earthmover’s distance is often talked about in terms of two piles of 869
870 dirt, which represent the two distributions being compared. In this 870

871 context, Earthmover's distance can be thought of as the minimum
872 cost of turning one earth pile into the other, where cost corresponds
873 to a combination of the amount of earth being moved as well as the
874 distance it has to be moved. In other words, the distance is the min-
875 imum average distance a piece of dirt will have to be moved in order
876 to turn one pile into the other. A higher distance means there was a
877 greater shape mismatch. We plot the distribution of Earthmover's
878 distances, and report its mean, variance, and maximum.

879 Having found that the distribution of Earthmover's distances for
880 contrastive dimensions had a longer tail than for non-contrastive
881 dimensions, we qualitatively analyzed the contents of this long tail
882 to determine which individual contexts led to the pattern observed.
883 We identified the contexts that showed up most frequently in the
884 tail and analyzed how frequent they were (both in terms of absolute
885 count and their frequency ranking relative to all contexts), as well
886 as what the relative frequency of short and long vowels was in each
887 of these key contexts.

888 To assess the reliability of these differences between contrastive
889 and non-contrastive dimensions were reliable, we used bootstrap
890 statistics. We sampled (with replacement) particular vowel tokens
891 with their contexts, creating a new contrastive and a new non-
892 contrastive dataset. We then recalculated Earthmover's distances
893 across the contexts using these bootstrap samples, repeating the
894 process 50 times. We plot the maximum Earthmover's distance with
895 the standard deviation, which allows us to observe the reliability of
896 these differences. To study the effect of the input corpus size, we
897 varied the number of vowel tokens sampled from 284 vowel tokens
898 (the size of the Dutch IDS corpus) to 132037 (the size of the French
899 vs. Japanese corpora).

900 In our final simulation, we relaxed our assumptions about infants'
901 prior knowledge. First, while previous analyses used neighboring
902 segment identity directly (e.g., /k/, /t/, /s/), this simulation only
903 used the segment's broad class (e.g., stop, fricative, vowel, etc.),
904 which infants are more sensitive to. Second, to simulate imper-
905 fect segmentation, we added noise to the prosodic position factor.
906 Prosodic position is represented with four indicator values (depend-
907 ing on whether the vowel in question is word/utterance-initial and
908 word/utterance-final). To add noise, we changed 20% of these val-
909 ues (making sure that the resulting prosodic position was real -
910 e.g., sounds considered to be utterance-final were necessarily also
911 considered to be word-final). We then used the same procedure
912 from above with these updated factors.

913 **Corpora.** The French vs. Japanese analysis compared the Corpus
914 of Spontaneous Japanese against the Nijmegen Corpus of Casual
915 French. The Dutch analyses looked at the Ernestus Corpus of
916 Spontaneous Dutch (ADS) and the Levelt/Fikkert corpus (IDS).

917 **Corpus of Spontaneous Japanese (CSJ).** The CSJ is a large corpus
918 of spontaneously produced adult-directed speech (67). Around 90%
919 of the speech consists of spontaneously produced monologues about
920 academic field, their favorite memory, and so forth. The remaining
921 10% consists of spontaneous dialogues either in free conversation
922 with the experimenter or engaged in a task. Our analysis focuses
923 on the core portion of the corpus, which was force-aligned and
924 hand-corrected with the segmental information required for our
925 analyses (see (67) for more details). The core portion consists of
926 811,731 total vowel tokens of which 89.1% are phonemically short
927 and 10.9% are phonemically long, but only the first 132,307 tokens
928 were used to match the size of the French corpus.

929 **Nijmegen Corpus of Casual French (NCCFr).** The NCCFr is a corpus
930 of spontaneously produced adult-directed speech (68). Unlike the
931 CSJ, however, the NCCFr consists exclusively of conversational
932 speech between close friends. Topics included upcoming exams,
933 travel plans, an ongoing strike, and so forth. The corpus consists of
934 speech by 46 French speakers, and includes 132,307 vowel tokens.
935 The corpus was orthographically transcribed by two professional
936 transcribers. The corpus was transcribed at the segmental level by
937 Martine Adda Decker (p.c. with M. Ernestus, January 14, 2019).

938 **Ernestus Corpus of Spontaneous Dutch (ECSD).** The ECSD consists
939 of adult-directed, conversational speech, with speakers talking with
940 a friend, at first freely, and then engaged in a task-oriented dis-
941 cussion (40). The corpus has speech by 20 different speakers, and
942 includes 60,955 tokens with a length contrast and 21,187 tokens

943 without. Professional transcribers created a orthographic transcrip-
944 tion of the interactions, which was manually aligned to the speech.
945 The corpus was also phonetically transcribed using a forced align-
946 ment model (details can be found in (69)). Validations revealed a
947 14% discrepancy between manual annotations and forced-aligned
948 annotations, which is in the range of human disagreement. However,
949 these analyses did not directly validate durational information, so it
950 is unclear how accurate annotations of the start and end points of
951 the phones are. This could introduce some noise into our analyses,
952 as it could affect how accurate the vowel durations are, and how
953 accurately we can determine which word a vowel belonged to.

954 **Fikkert/Levelt/Swingley IDS corpus.** We also tested our account on
955 a corpus of Dutch IDS collected by Fikkert and Levelt (70, 71).
956 The annotated portion of this corpus is small: it contains a total of
957 300 utterances, with a total of only 1296 vowel tokens, but each of
958 the contrastive and non-contrastive datasets had to be subsetted
959 to 284 to make equally sized subsets. The corpus consists of natu-
960 ralistic longitudinal speech interactions with one child (Catootje)
961 aged 1;10. The corpus was transcribed at the word level. Time-
962 aligned phonetic annotations were created by Dan Swingley (DS)
963 (13). Given the transcriptions, the speech toolkit HTK (72) was
964 used to estimate the boundaries of the phones using the HVITE
965 forced-alignment tool. The output of the forced-alignment tool was
966 manually corrected by DS, a speaker of Dutch. KH time-aligned the
967 word-level transcription to the time-aligned phonetic transcriptions
968 based on the location of the phones in Praat (73).

969 Code for all analyses is available on Github:
970 <http://github.com/khitzenko/contextual-dl>. Data are available by
971 request from the researchers who control their distribution.

972 **ACKNOWLEDGMENTS.** This work was supported by NSF grants
973 IIS-1421695, BCS-1734245, and NRT #1449815. We thank Thomas
974 Schatz, Mirjam Ernestus, Dan Swingley, and Francisco Torreira
975 for help acquiring corpora, as well as Micha Elsner, Bill Idsardi,
976 Jeff Lidz, Reiko Mazuka, Rochelle Newman, Thomas Schatz, Dan
977 Swingley and the reviewers for helpful feedback and discussion.

- 978 1. Werker JF, Gilbert JH, Humphrey K, Tees RC (1981) Developmental aspects of cross-
979 language speech perception. *Child development* pp. 349–355.
- 980 2. Werker JF, Tees RC (1984) Cross-language speech perception: Evidence for perceptual re-
981 organization during the first year of life. *Infant Behavior and Development* 7(1):49–63.
- 982 3. Kuhl PK, Williams KA, Lacerda F, Stevens KN, Lindblom B (1992) Linguistic experience alters
983 phonetic perception in infants by 6 months of age. *Science* 255:606–608.
- 984 4. Tsuji S, Cristia A (2014) Perceptual attunement in vowels: A meta-analysis. *Developmental*
985 *psychobiology* 56(2):179–191.
- 986 5. Narayan C, Werker JF, Beddor PS (2010) The interaction between acoustic salience and
987 language experience in developmental speech perception: Evidence from nasal place dis-
988 crimination. *Developmental Science* 13(3):407–420.
- 989 6. Maye J, Werker JF, Gerken L (2002) Infant sensitivity to distributional information can affect
990 phonetic discrimination. *Cognition* 82(3):B101–B111.
- 991 7. Maye J, Weiss DJ, Aslin RN (2008) Statistical phonetic learning in infants: Facilitation and
992 feature generalization. *Developmental science* 11(1):122–134.
- 993 8. Yoshida KA, Pons F, Maye J, Werker JF (2010) Distributional phonetic learning at 10 months
994 of age. *Infancy* 15(4):420–433.
- 995 9. Cristia A (2018) Can infants learn phonology in the lab? a meta-analytic answer. *Cognition*
996 170:312–327.
- 997 10. Vallabha GK, McClelland JL, Pons F, Werker JF, Amano S (2007) Unsupervised learning
998 of vowel categories from infant-directed speech. *Proceedings of the National Academy of*
999 *Sciences* 104(33):13273–13278.
- 1000 11. McMurray B, Aslin RN, Toscano JC (2009) Statistical learning of phonetic categories: insights
1001 from a computational approach. *Developmental science* 12(3):369–378.
- 1002 12. Bion RA, Miyazawa K, Kikuchi H, Mazuka R (2013) Learning phonemic vowel length from
1003 naturalistic recordings of Japanese infant-directed speech. *PLOS ONE* 8(2):e51594.
- 1004 13. Swingley D (2019) Learning phonology from surface distributions, considering Dutch and
1005 English vowel duration. *Language Learning and Development* 15(3):199–216.
- 1006 14. Vance TJ (1987) *An introduction to Japanese phonology*. (SUNY Press).
- 1007 15. Narayan C (2008) The acoustic-perceptual salience of nasal place contrasts. *Journal of*
1008 *Phonetics* 36(1):191–217.
- 1009 16. Swingley D, Alarcon C (2018) Lexical learning may contribute to phonetic learning in infants:
1010 A corpus analysis of maternal spanish. *Cognitive science*.
- 1011 17. Adriaans F, Swingley D (2012) Distributional learning of vowel categories is supported by
1012 prosody in infant-directed speech in *Proceedings of the Annual Meeting of the Cognitive*
1013 *Science Society*. Vol. 34.
- 1014 18. Dillon B, Dunbar E, Idsardi W (2013) A single-stage approach to learning phonological cate-
1015 gories: Insights from inuktitut. *Cognitive Science* 37(2):344–377.
- 1016 19. Feldman NH, Griffiths TL, Goldwater S, Morgan JL (2013) A role for the developing lexicon in
1017 phonetic category acquisition. *Psychological Review* 120(4):751.

- 1018 20. Feldman NH, Myers EB, White KS, Griffiths TL, Morgan JL (2013) Word-level information
1019 influences phonetic learning in adults and infants. *Cognition* 127(3):427–438.
- 1020 21. Swingle D (2009) Contributions of infant word learning to language development. *Philosophical
1021 Transactions of the Royal Society of London B: Biological Sciences* 364(1536):3617–
1022 3632.
- 1023 22. Thiessen ED (2007) The effect of distributional information on children's use of phonemic
1024 contrasts. *Journal of Memory and Language* 56(1):16–34.
- 1025 23. Yeung HH, Chen LM, Werker JF (2014) Referential labeling can facilitate phonetic learning in
1026 infancy. *Child Development* 85(3):1036–1049.
- 1027 24. Schatz T, Feldman NH, Goldwater S, Cao XN, Dupoux E (2021) Early phonetic learning with-
1028 out phonetic categories: Insights from large-scale simulations on realistic input. *Proceedings
1029 of the National Academy of Sciences* 118(7).
- 1030 25. Sato Y, Sogabe Y, Mazuka R (2010) Discrimination of phonemic vowel length by Japanese
1031 infants. *Developmental Psychology* 46(1):106.
- 1032 26. Christophe A, Dupoux E, Bertoni J, Mehler J (1994) Do infants perceive word boundaries?
1033 An empirical study of the bootstrapping of lexical acquisition. *The Journal of the Acoustical
1034 Society of America* 95(3):1570–1580.
- 1035 27. Christophe A, Mehler J, Sebastián-Gallés N (2001) Perception of prosodic boundary corre-
1036 lates by newborn infants. *Infancy* 2(3):385–394.
- 1037 28. Jusczyk PW, Aslin RN (1995) Infants' detection of the sound patterns of words in fluent
1038 speech. *Cognitive Psychology* 29(1):1–23.
- 1039 29. Jusczyk PW, Houston DM, Newsome M (1999) The beginnings of word segmentation in
1040 English-learning infants. *Cognitive Psychology* 39(3-4):159–207.
- 1041 30. Greenberg JH, Ferguson CA, Moravcsik EA (1978) *Universals of human language: phonol-
1042 ogy*. (Stanford University Press) Vol. 2.
- 1043 31. Moreton E (2002) Structural constraints in the perception of english stop-sonorant clusters.
1044 *Cognition* 84(1):55–71.
- 1045 32. Gómez DM, et al. (2014) Language universals at birth. *Proceedings of the National Academy
1046 of Sciences* 111(16):5837–5841.
- 1047 33. Hitczenko K, Mazuka R, Elsner N, Feldman NH (2020) When context is and isn't helpful: A
1048 corpus study of naturalistic speech. *Psychonomic Bulletin & Review* pp. 1–37.
- 1049 34. Bergmann C, Cristia A (2016) Development of infants' segmentation of words from native
1050 speech: A meta-analytic approach. *Developmental science* 19(6):901–917.
- 1051 35. Feldman NH, Goldwater S, Dupoux E, Schatz T (2021) Do infants really learn phonetic cate-
1052 gories? *Open Mind* pp. 1–19.
- 1053 36. Rubner Y, Tomasi C, Guibas LJ (1998) A metric for distributions with applications to image
1054 databases in *Sixth International Conference on Computer Vision*. (IEEE), pp. 59–66.
- 1055 37. Babineau M, Shi R (2011) Processing of french liaisons in toddlers. *BUCLD 35 Proceedings.
1056 Cascadilla Press, Somerville, MA* pp. 25–37.
- 1057 38. Mattys SL, Jusczyk PW (2001) Phonotactic cues for segmentation of fluent speech by infants.
1058 *Cognition* 78(2):91–121.
- 1059 39. Seidl A, Johnson EK (2006) Infant word segmentation revisited: Edge alignment facilitates
1060 target extraction. *Developmental science* 9(6):565–573.
- 1061 40. Ernestus MTC (2000) *Voice assimilation and segment reduction in casual Dutch: A corpus-
1062 based study of the phonology-phonetics interface*. (Holland Institute of Generative Linguistics,
1063 Utrecht).
- 1064 41. Antetomaso S, et al. (2017) Modeling phonetic category learning from natural acoustic data
1065 in *BUCLD 41: Proceedings of the 41st Annual Boston University Conference on Language
1066 Development*.
- 1067 42. Lisker L (1986) "Voicing" in English: A catalogue of acoustic features signaling /b/ versus /p/
1068 in trochees. *Language and Speech* 29(1):3–11.
- 1069 43. Narayan C (2013) Developmental perspectives on phonological typology and sound change.
1070 *Origins of sound change: Approaches to phonologization* pp. 128–146.
- 1071 44. Narayan C, Peters A, Woldenga-Racine V (2017) Fragile phonetic contrasts in longitudinal
1072 infant-directed speech: Implications for infant speech perception in *BUCLD 42: Proceedings
1073 of the 41st Annual Boston University Conference on Language Development*.
- 1074 45. Pimentel T, Roark B, Cotterell R (2020) Phonotactic complexity and its trade-offs. *Transac-
1075 tions of the Association for Computational Linguistics* 8:1–18.
- 1076 46. Roark CL, Plaut DC, Holt LL (2022) A neural network model of the effect of prior experience
1077 with regularities on subsequent category learning. *Cognition* 222:104997.
- 1078 47. Tranel B (1987) *The Sounds of French: An Introduction*. (Cambridge University Press).
- 1079 48. Clayards M, Tanenhaus MK, Aslin RN, Jacobs RA (2008) Perception of speech reflects opti-
1080 mal use of probabilistic speech cues. *Cognition* 108(3):804–809.
- 1081 49. Maye J, Aslin RN, Tanenhaus MK (2008) The weckud wetch of the wast: Lexical adaptation
1082 to a novel accent. *Cognitive Science* 32(3):543–562.
- 1083 50. Cristia A, et al. (2012) Linguistic processing of accented speech across the lifespan. *Frontiers
1084 in Psychology* 3:479.
- 1085 51. van Heugten M, Johnson EK (2012) Infants exposed to fluent natural speech succeed at
1086 cross-gender word recognition. *Journal of Speech, Language, and Hearing Research*.
- 1087 52. White KS, Aslin RN (2011) Adaptation to novel accents by toddlers. *Developmental Science*
1088 14(2):372–384.
- 1089 53. Weatherhead D, White KS (2016) He says potato, she says potahto: Young infants track
1090 talker-specific accents. *Language Learning and Development* 12(1):92–103.
- 1091 54. Kuhl PK (1979) Speech perception in early infancy: Perceptual constancy for spectrally dis-
1092 similar vowel categories. *The Journal of the Acoustical Society of America* 66(6):1668–1679.
- 1093 55. Kuhl PK (1983) Perception of auditory equivalence classes for speech in early infancy. *Infant
1094 behavior and development* 6(2-3):263–285.
- 1095 56. Jusczyk PW, Luce PA, Charles-Luce J (1994) Infants' sensitivity to phonotactic patterns in
1096 the native language. *Journal of Memory and Language* 33(5):630.
- 1097 57. Jusczyk PW, Hohne EA, Bauman A (1999) Infants' sensitivity to allophonic cues for word
1098 segmentation. *Perception & Psychophysics* 61(8):1465–1476.
- 1099 58. White KS, Peperkamp S, Kirk C, Morgan JL (2008) Rapid acquisition of phonological alterna-
1100 tions by infants. *Cognition* 107(1):238–265.
- 1101 59. Moreton E, Amano S (1999) Phonotactics in the perception of Japanese vowel length: Evi-
1102 dence for long-distance dependencies in *EUROSPEECH*.
- 1103 60. Idemaru K, Holt LL (2011) Word recognition reflects dimension-based statistical learning. *Journal of Experimental Psychology: Human Perception and Performance* 37(6):1939.
- 1104 61. Holt LL, Loto AJ (2006) Cue weighting in auditory categorization: Implications for first and
1105 second language acquisition. *The Journal of the Acoustical Society of America* 119(5):3059–
1106 3071.
- 1107 62. Toscano JC, McMurray B (2010) Cue integration with categories: Weighting acoustic cues
1108 in speech using unsupervised learning and distributional statistics. *Cognitive science*
1109 34(3):434–464.
- 1110 63. Liu L, Jaeger TF (2018) Inferring causes during speech perception. *Cognition* 174:55–70.
- 1111 64. Bergelson E, Swingle D (2012) At 6–9 months, human infants know the meanings of many
1112 common nouns. *Proceedings of the National Academy of Sciences* 109(9):3253–3258.
- 1113 65. Rubner Y, Tomasi C, Guibas LJ (1998) A metric for distributions with applications to image
1114 databases in *Sixth International Conference on Computer Vision (IEEE Cat. No. 98CH36271)*.
1115 (IEEE), pp. 59–66.
- 1116 66. Villani C (2008) *Optimal transport: Old and new*. (Springer Science & Business Media) Vol.
1117 338.
- 1118 67. Maekawa K (2003) Corpus of Spontaneous Japanese: Its Design and Evaluation in *ISCA &
1119 IEEE Workshop on Spontaneous Speech Processing and Recognition*.
- 1120 68. Torreira F, Adda-Decker M, Ernestus M (2010) The Nijmegen Corpus of Casual French.
1121 *Speech Communication* 52(3):201–212.
- 1122 69. Schuppler B, Ernestus M, Scharenborg O, Boves L (2011) Acoustic reduction in conversa-
1123 tional Dutch: A quantitative analysis based on automatically generated segmental transcrip-
1124 tions. *Journal of Phonetics* 39(1):96–109.
- 1125 70. Fikkert P (1994) *On the acquisition of prosodic structure*. (Ph.D. Thesis).
- 1126 71. Levelt C (1994) *On the acquisition of place*. (Ph.D. Thesis).
- 1127 72. Young S, et al. (2002) *The HTK Book*. Cambridge University Engineering Department
1128 3(175):12.
- 1129 73. Boersma P (2001) Praat: A system for doing phonetics by computer. *Glott International*
1130 5(9/10):341–345.

1132 **Supplementary Materials**

1133 Figure 3 in the main text reports our bootstrapping analysis results, treating the maximum Earthmover’s distance as the metric of interest
 1134 (i.e., showing how maximum Earthmover’s distance varies across 50 runs). In Figure S1, we provide analogous plots treating the mean
 1135 Earthmover’s distance as the primary metric instead. Results are similar: across input sizes the mean contrastive Earthmover’s distance is
 1136 greater than the mean non-contrastive Earthmover’s distance, suggesting that observed differences between contrastive vs. non-contrastive
 1137 dimensions are meaningful, but that input size does not matter.

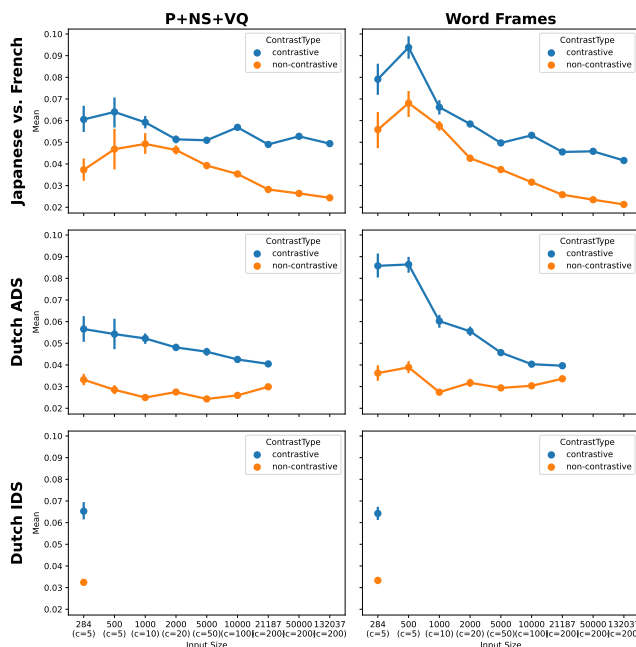


Fig. S1. Results from bootstrapped input size analyses. Across input sizes, the mean contrastive Earthmover’s distance is greater than the mean non-contrastive Earthmover’s distance. “C” refers to the number of contexts included in the analysis. The maximum input size for which data is shown depends on the corpus size: 284 for Dutch IDS, 21187 for Dutch ADS, and 132037 for Japanese vs. French.

1138 The analyses in the main text use Earthmover’s distance to compare distribution shapes across contexts. Figure S2 shows results using an
 1139 alternative metric, namely KL divergence (or Kullback-Leibler divergence). KL divergence is a measure of how different two probability
 1140 distributions are. As we do not have access to the closed-form probability distributions within each context (only samples), we estimated
 1141 KL divergence as follows. We divided the tokens in each context into 10 evenly-sized bins, and smoothed the counts in each bin (adding
 1142 a count of 1e-5 to each bin). We then used these counts to arrive at a probability distribution over the 10 bins, which were used to
 1143 calculate KL-divergence. As KL-divergence is not a symmetric measure, for each pair of contexts, we summed the KL-divergence of
 1144 Context 1 against Context 2 and the KL-divergence of Context 2 against Context 1 to arrive at a final measure, that is plotted in Figure
 1145 S2. We find that the expected qualitative pattern still emerges across all three test cases even using a different metric. It is, however,
 1146 interesting that the difference in magnitude across the test cases changes: here, Dutch ADS and IDS show the greatest KL divergences,
 1147 while Japanese/French showed the greatest Earthmover’s distances. This suggests that this effect is largely insensitive to the particular
 1148 metric used, but which metric is used could make a difference in specific predictions the account makes.

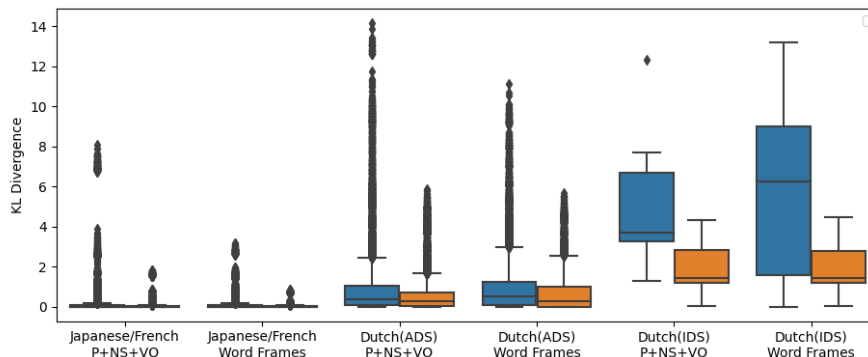


Fig. S2. Distribution of KL divergences by test case. The same qualitative pattern holds, though the difference in magnitude changes across the three test cases. P+NS+VQ = prosodic position + neighboring sounds + vowel quality.