

Western University

Scholarship@Western

Psychology Publications

Psychology Department

10-21-2023

Of words and whistles: Statistical learning operates similarly for identical sounds perceived as speech and non-speech
Of words and whistles: Statistical learning operates similarly for identical sounds perceived as speech and non-speech

Sierra J. Sweet

Stephen C. Van Hedger

Laura J. Batterink
lbatter@uwo.ca

Follow this and additional works at: <https://ir.lib.uwo.ca/psychologypub>

Citation of this paper:

Sweet, Sierra J.; Van Hedger, Stephen C.; and Batterink, Laura J., "Of words and whistles: Statistical learning operates similarly for identical sounds perceived as speech and non-speech" (2023). *Psychology Publications*. 249.
<https://ir.lib.uwo.ca/psychologypub/249>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

**Of words and whistles: Statistical learning operates similarly for identical sounds
perceived as speech and non-speech**

Sierra J. Sweet ¹ (ssweet4@uwo.ca)

Stephen C. Van Hedger ^{1 2 3} (svanhedg@uwo.ca)

Laura J. Batterink ^{1 2 *} (lbatter@uwo.ca)

¹ Department of Psychology, Western University, London, ON Canada

² Western Institute for Neuroscience, Western University, London, ON Canada

³ Department of Psychology, Huron University College, London, ON Canada

* Corresponding author: Laura J. Batterink (lbatter@uwo.ca); (519) 661-2111 x85409;

1151 Richmond St. WIRB 6124, London, ON Canada, N6A 5B7

20

Abstract

21 Statistical learning is an ability that allows individuals to effortlessly extract patterns from
22 the environment, such as sound patterns in speech. Some prior evidence suggests that
23 statistical learning operates more robustly for speech compared to non-speech stimuli,
24 supporting the idea that humans are predisposed to learn language. However, any
25 apparent statistical learning advantage for speech could be driven by signal acoustics,
26 rather than the subjective perception *per se* of sounds as speech. To resolve this issue,
27 the current study assessed whether there is a statistical learning advantage for
28 ambiguous sounds that are subjectively perceived as speech-like compared to the
29 same sounds perceived as non-speech, thereby controlling for acoustic features. We
30 first induced participants to perceive sine-wave speech (SWS)—a degraded form of
31 speech not immediately perceptible as speech—as either speech or non-speech. After
32 this induction phase, participants were exposed to a continuous stream of repeating
33 trisyllabic nonsense words, composed of SWS syllables, and then completed an explicit
34 familiarity rating task and an implicit target detection task to assess learning. Critically,
35 participants showed robust and equivalent performance on both measures, regardless
36 of their subjective speech perception. In contrast, participants who perceived the SWS
37 syllables as more speech-like showed better detection of individual syllables embedded
38 in speech streams. These results suggest that speech perception facilitates processing
39 of individual sounds, but not the ability to extract patterns across sounds. Our findings
40 suggest that statistical learning is not influenced by the degree of perceived linguistic
41 relevance of sounds, and that it may be conceptualized largely as an automatic,
42 stimulus-driven mechanism.

43 Keywords: statistical learning, speech, sine-wave speech, auditory perception

44 **1. Introduction**

45 Statistical learning, our ability to become sensitive to patterns in the environment,
46 has provided an important mechanistic explanation for language acquisition since its
47 initial documentation in the context of speech segmentation (Saffran et al., 1996a). In
48 this study, infants were presented with a continuous stream of trisyllabic nonsense
49 words, with no pauses or other acoustic cues to mark word boundaries. Thus, the
50 probabilities of syllables co-occurring with one another provided the only indication of
51 where individual words started and ended within the stream. After listening to the
52 stream, infants were able to successfully discriminate between words and foil items
53 through their looking time behaviour, providing evidence that they had extracted the
54 statistical information in the stream to discover the embedded words.

55 Since this seminal study, subsequent research has shown that statistical learning
56 is present across many domains outside of language (e.g., Conway & Christiansen,
57 2005; Fiser & Aslin, 2001; Saffran et al., 1999; Van Hedger et al., 2022). In one such
58 study, conducted by Saffran and colleagues (1999), participants were exposed to a
59 stream of six “tone words,” each of which consisted of a sequence of three pure tones.
60 On a subsequent two-alternative forced-choice recognition task, participants succeeded
61 in discriminating between tone words and foil sequences, providing a clear
62 demonstration that statistical learning also operates across non-linguistic auditory
63 stimuli – that is, auditory stimuli that lack a clear communicative purpose. Subsequent
64 research has found that listeners can also extract patterns embedded in non-linguistic
65 noises (Gebhart et al., 2009), everyday environmental sounds (Siegelman et al., 2018),

66 tactile sequences (Conway & Christiansen, 2005), visual stimuli (e.g., Bulf et al., 2011;
67 Kirkham et al., 2002; Fiser & Aslin, 2001), and multimodal contexts (Mitchel et al., 2014;
68 Seitz et al., 2007). Further, statistical learning is present not only in infants but also in
69 older children and adults (e.g., Moreau et al., 2022; Raviv & Arnon, 2018; Saffran et al.,
70 1996b, 1997), as well as in nonhuman animals, including dogs (Boros et al., 2021) and
71 cotton-top tamarins (Hauser et al., 2001). These observations have led to a general
72 consensus that statistical learning is not a “special” language-specific mechanism, but is
73 domain-general in that it is present across modalities, domains, and even species
74 (Aslin, 2017).

75 However, while statistical learning may be considered domain-general in that it is
76 present in many learning contexts, it shows important differences depending on
77 stimulus modality and learning domains, suggesting that it may not be a truly unitary
78 mechanism (Frost et al., 2015; Frost et al., 2019). For example, an early study found an
79 advantage for statistical learning of non-linguistic tones, as compared to tactile and
80 visual stimuli, which persisted even after controlling for low-level perceptual differences
81 between stimuli (Conway and Christiansen, 2005). Another study reported that changes
82 in presentation rate have opposite effects on auditory and visual statistical learning:
83 auditory statistical learning benefits from faster presentation rates, whereas visual
84 statistical learning benefits from slower rates (Emberson et al., 2011). In addition,
85 different types of statistical learning follow different developmental trajectories; statistical
86 learning for speech sounds is stable from childhood into adulthood; in contrast,
87 statistical learning improves with age for visual stimuli and non-linguistic tones (Arciuli &

88 Simpson, 2011; Moreau et al., 2022; Raviv & Arnon, 2018; Schlichting et al., 2017;
89 Shufaniya & Arnon, 2018; for review, Forest et al., 2023).

90 These findings, which indicate that statistical learning is not equivalent across
91 modalities, are not easily accommodated within frameworks that treat statistical learning
92 as a single unitary mechanism. Further evidence against a unitary view of statistical
93 learning comes from low interindividual correlations in statistical learning performance
94 across modalities and stimulus materials (Siegelman & Frost, 2015; Siegelman et al.,
95 2017). While an individual's statistical learning performance within a given domain is
96 relatively stable, as assessed by test-retest reliability, performance on one task does not
97 predict performance on a parallel tasks in a different domain (e.g. syllables to visual
98 shapes; Siegelman & Frost, 2015). Taken together, these results suggest that there are
99 nonoverlapping mechanisms supporting statistical learning abilities in different domains,
100 supporting a “pluralist” view of statistical learning (Frost et al., 2015; Frost et al., 2019).
101 According to this viewpoint, statistical learning is supported not only by domain-general
102 mechanisms (e.g. Schapiro et al., 2014; Covington et al., 2018; Conway, 2020;
103 Batterink et al., 2019), but also by modality-specific mechanisms that are united by
104 similar computational principles. These modality-specific mechanisms operate within
105 distinct networks and are governed by different constraints, depending on task domain
106 and modality (Frost et al., 2015, Frost et al., 2019; Conway, 2020).

107 **1.1. Is speech a privileged target for statistical learning?**

108 The consensus that there are important differences in statistical learning as a
109 function of learning domain raises a more specific question of whether statistical
110 learning operates differently—and perhaps more robustly—for speech than non-speech.

111 Human infants prefer to listen to speech compared to other auditory stimuli (Shultz &
112 Vouloumanos, 2010), and neuroimaging studies in adults have found greater activation
113 in left auditory cortex for speech compared to other sounds (Binder et al., 2000; Narain
114 et al., 2003; Parviainen et al., 2005; Scott et al., 2000; Vouloumanos et al., 2001).
115 These results are in line with the general idea that speech is “special,” engaging unique
116 neural and cognitive mechanisms not engaged by other auditory stimuli (Belin et al.,
117 2000; Liberman, 1982; Marno et al., 2015; Moore, 2000).

118 Infant studies of artificial grammar rule learning also support this notion,
119 suggesting that babies more readily extract simple grammar rules (e.g., “AAB” or “ABB”
120 rules) from speech than from non-speech auditory stimuli, such as tones or animal
121 sounds (Dawson & Gerken, 2009; Marcus et al., 2007). A number of theoretical
122 hypotheses (which are not mutually exclusive) have been proposed to account for this
123 speech advantage in rule learning, including that speech (1) better captures and holds
124 infants’ attention (Schultz & Vouloumanos, 2010; Vouloumanos & Werker, 2004), (2)
125 represents a communicative signal (Rabagliati et al., 2012; Ferguson & Lew-Williams,
126 2016), (3) is more familiar than other signals to infants, which facilitates learning
127 (Saffran et al., 2007; Thiessen, 2012), and/or (4) may be processed by specific
128 mechanisms that have been tuned to speech as humans evolved the capacity for
129 language (Rabagliati et al., 2012; Marcus & Rabagliati, 2008, as cited in Ferguson &
130 Lew-Williams, 2016). By extension, speech could also represent a privileged target for
131 the statistical learning of embedded units in continuous sound sequences, in infants and
132 adults alike.

133 Current evidence on whether there is indeed a statistical learning advantage for
134 speech sounds is conflicting. A recent study by Ordin and colleagues (2021) supports
135 the idea that there is a speech advantage in statistical learning. Participants were
136 presented with embedded triplet sequences that were fully linguistic in nature (made up
137 of natural syllables), semi-linguistic (made up of syllables that contained atypical
138 acoustic cues), and non-linguistic (made up of environmental sounds such as animal
139 noises and footsteps), and then asked to make old/new judgments for triplets from the
140 sequences and foils. Performance was highest in the syllable condition compared to the
141 semi-linguistic and non-linguistic conditions, providing support for a speech advantage
142 for statistical learning. This result also converges with rule learning studies in infants,
143 which have found a general advantage for speech stimuli over non-speech stimuli, as
144 described above (e.g., Dawson & Gerken, 2009; Marcus et al., 2007).

145 However, not all studies point to a clear linguistic advantage for statistical
146 learning. In the previously described “tone words” study by Saffran and colleagues
147 (1999), both age groups successfully segmented the tone stream, and no significant
148 differences were found between their performance on the tone version and the syllable
149 version of the task from a previous study (Saffran et al., 1996b). Similarly, another study
150 by Saffran (2002) presented adults and children with linguistic or non-linguistic auditory
151 “sentences,” made up of nonsense words for the linguistic group (e.g. kiff flor lum dupp)
152 and sequences of sounds such as bells, chimes, and drums for the non-linguistic group.
153 Both groups learned successfully and again, no significant differences were found
154 between conditions. Finally, a more recent study by Siegelman and colleagues (2018)
155 compared statistical learning of syllables and everyday environmental sounds. Overall

156 performance was similar between the two conditions, again suggesting that statistical
157 learning occurs with similar efficacy for speech and non-speech sounds.

158 Yet, even in situations where overall learning is comparable for linguistic and
159 non-linguistic items, there is evidence that linguistic items still might exhibit distinct
160 patterns of learning. For example, more nuanced analyses of the Siegelman and
161 colleagues (2018) data revealed that individual test items in the syllable condition
162 showed much lower internal consistency than in the sound condition. Additional
163 experiments indicated that participants' performance was influenced by the degree to
164 which test items corresponded to the phonotactics of their own native language of
165 Hebrew (see also Elazar et al., 2022). These results suggest that learners' prior
166 knowledge and expectations may critically impact statistical learning of linguistically-
167 relevant speech sounds, an effect that is less pronounced for non-linguistic sounds
168 (though see Van Hedger et al., 2022 for evidence of effects of prior knowledge on
169 statistical learning of instrument notes). Thus, even in the absence of overall
170 performance differences, there may be qualitative differences in how statistical learning
171 operates for speech versus non-speech sounds, particularly with respect to how
172 learning interacts with other cognitive factors.

173 **1.2. Differences between speech and non-speech sounds**

174 Part of the difficulty in assessing whether there may be a statistical learning
175 advantage for speech is that speech sounds and non-speech sounds, such as tones
176 and environmental noises, differ in many ways. Previous learning studies comparing
177 speech and non-speech have used different types of artificial languages, different
178 syllable inventories, and many different types of non-linguistic sounds (e.g. Marcus et

179 al., 2007; Ordin et al., 2021; Saffran et al., 1999; Saffran, 2002; Siegelman et al., 2018).
180 Thus, conflicting results across studies could—in principle—be at least partially
181 attributable to surface features of the learning materials. For example, speech sounds
182 and other natural auditory stimuli such as musical instruments and everyday object
183 sounds differ in fundamental frequency, timbre, aperiodicity, spectral variability, spectral
184 envelope, and temporal envelope (Ogg & Slevc, 2019). Any number of these low-level
185 acoustic features that differ between speech and non-linguistic stimuli may influence
186 perception, ease of encoding, and consequently statistical learning performance. In
187 other words, statistical learning differences between speech and non-speech—when
188 observed—could reflect signal-driven differences in lower-level processes, such as the
189 perception of individual items, rather than statistical learning *per se*.

190 A study by Thiessen (2012) highlights the importance of considering acoustic
191 features when comparing statistical learning of speech versus non-speech sounds. The
192 authors of this study reasoned that speech contains more redundant cues to an abstract
193 rule than are typically available in non-linguistic stimuli, and that such redundancy may
194 facilitate rule learning. For example, a string such as “ga ti ga” instantiates the “ABA”
195 rule at multiple levels: at the syllable level, at the individual phoneme level (both the
196 initial consonant and final vowel differentiate the A and B elements) and at the level of
197 phonetic features (e.g., voicing). To test the importance of redundancy, the authors
198 presented infants with syllable sequences that contained reduced redundancy, in which
199 only the vowels, rather than both vowels and consonants, signaled the underlying rule
200 (e.g. “ba bi ba” rather than “ga ti ga”). When redundancy was reduced, infants’ rule
201 learning was impaired, suggesting that speech may allow for easier learning than non-

202 linguistic stimuli at least in part because of the redundant information in the acoustic
203 signal. These results underscore the importance of accounting for acoustic differences
204 in comparisons of statistical learning between speech and non-speech stimuli.

205 In addition to their acoustic differences, speech sounds also differ from non-
206 speech sounds in terms of their *subjective value* or *perceived relevance* to the listener.
207 In contrast to tones or environmental noises, speech sounds are a linguistically relevant
208 signal and serve a critical communicative purpose. This communicative value could in
209 part explain why speech captures infants' attention to a greater degree than non-speech
210 (e.g., Vouloumanos & Werker, 2004, 2007; Vouloumanos et al., 2010), or why auditory-
211 relevant regions within the left temporal lobe are more strongly activated for speech
212 than non-speech (Belin et al., 2000; Binder et al., 2000; Dick et al., 2007; Scott et al.,
213 2000), although here too acoustic differences cannot be ruled out. To our knowledge,
214 no previous studies have directly examined whether the communicative value of speech
215 *per se* may play a role in potential statistical learning differences between speech and
216 non-speech sounds.

217 In the current study, we tested the hypothesis that speech may serve as a
218 privileged target for statistical learning due to its *subjective value* as a communicative
219 signal, over and above any effects of acoustic differences between speech and non-
220 speech. To address this hypothesis, we leveraged "sine-wave speech" (SWS), a
221 manipulation that allows for comparing the processing of identical acoustic stimuli that
222 may be perceived from highly speech-like to un-speechlike. SWS is a degraded form of
223 natural speech consisting of time-varying sine waves modelling formant frequencies,
224 with fewer sine waves corresponding to greater degradation of the signal (Remez et al.,

1981). This degraded audio retains the phonetic properties of the original speech, but typically fails to be perceived as phonetic by naïve listeners, who may experience it as a sequence of whistles, chirps, and other types of “science fiction” sounds. SWS lacks many of the acoustic features that make speech sound natural, such as a fundamental frequency. However, it can still be perceived as speech if instructions to attend to the speech-like qualities of the stimuli, or information about its true nature, are given. For example, participants may suddenly perceive SWS as speech if they are played the intact, original audio immediately prior to the SWS version. Notably, once participants are induced into perceiving the SWS as speech, there is no known method to revert them back into hearing it as non-speech (Silva & Bellini-Leite, 2020). SWS thus provides a tool for manipulating listeners’ subjective, top-down perception of a signal as speech versus non-speech, while holding the physical stimuli constant. Essentially, this approach can be used to isolate speech-specific perceptual effects on statistical learning, independent of any acoustic differences.

1.3. The Current Study

The aim of the current experiment was to investigate whether statistical learning operates differently for sounds perceived as more speech-like compared to sounds perceived as non-speech in the absence of acoustic differences between stimuli. Participants initially completed an induction task, in which we attempted to induce them to perceive SWS syllables as either speech or non-speech sounds. They were then exposed to a continuous stream of repeating trisyllabic “words” composed of SWS syllables, and then completed two behavioural tasks to measure their statistical learning of the words: (1) an explicit familiarity rating task, in which participants rated their

270 A total of 200 participants were recruited from online participant recruitment
271 platforms Prolific (n = 65; Palan & Schitter, 2018) and Amazon Mechanical Turk through
272 CloudResearch (n = 135; Litman et al., 2017). Amazon Mechanical Turk participants
273 were initially recruited; however, because a substantial proportion failed the study's
274 attention check (as described in detail later), we recruited a second group of participants
275 from Prolific in hopes of obtaining participants who would perform better on this
276 attention check. All Amazon Mechanical Turk recruited participants were
277 CloudResearch-approved, indicating that they had been screened and shown proof that
278 they engage in tasks in an attentive manner. All Prolific participants had approval rates
279 between 90-100%, indicating that a high percentage of their submissions for other
280 research studies had been approved by the researchers. All participants reported
281 English as their primary language, were above 17 years old, and had normal or
282 corrected-to-normal hearing. Of the 200 participants, 100 were assigned to the speech
283 induction condition, while the remaining 100 were assigned to the non-speech induction.
284 Participants were financially compensated for their time.

285 Of the 200 participants, a total of 73 participants were excluded from analysis; 43
286 were excluded due to failing to pass both attention checks embedded in the exposure
287 stream (as described in greater detail later); 23 because their data failed to save to our
288 servers; and 3 due to making no responses during the target detection task. Finally, 1
289 participant was excluded due to not having normal or corrected-to-normal hearing, and
290 3 participants were excluded due to failing to meet the inclusion criteria of having
291 English as their primary language, based off their answers to the post-study survey.
292 Thus, final analyses comprise data from 71 participants in the speech induction (SI)

293 condition (mean age = 40.2 y; SD = 11.8 y; 37 men; 34 women), and 56 participants in
294 the non-speech induction (NSI) condition (mean age = 39.3 y; SD = 12.0 y; 29 men; 27
295 women).

296 **2.2. Stimuli**

297 The experimental stimuli consisted of 12 syllables recorded by a male native
298 English speaker, taken from Batterink and Paller (2019), in addition to 24 corresponding
299 SWS manipulated forms of these syllables, comprised of single-sine wave (highly
300 degraded) and three-sine wave (moderately degraded) versions of each of the original
301 syllables. Each syllable sound file was 300 ms. Manipulated forms of the syllables were
302 created in Praat (Boersma & Weenick, 2022) using a script by Darwin (2003). The
303 unmanipulated (original) forms and single-sine wave (highly degraded) forms of the
304 syllables were used only as primes in the induction task. The three-sine wave
305 (moderately degraded) forms comprised the key experimental stimuli that were used
306 throughout all statistical learning tasks, as well as the syllable transcription task.

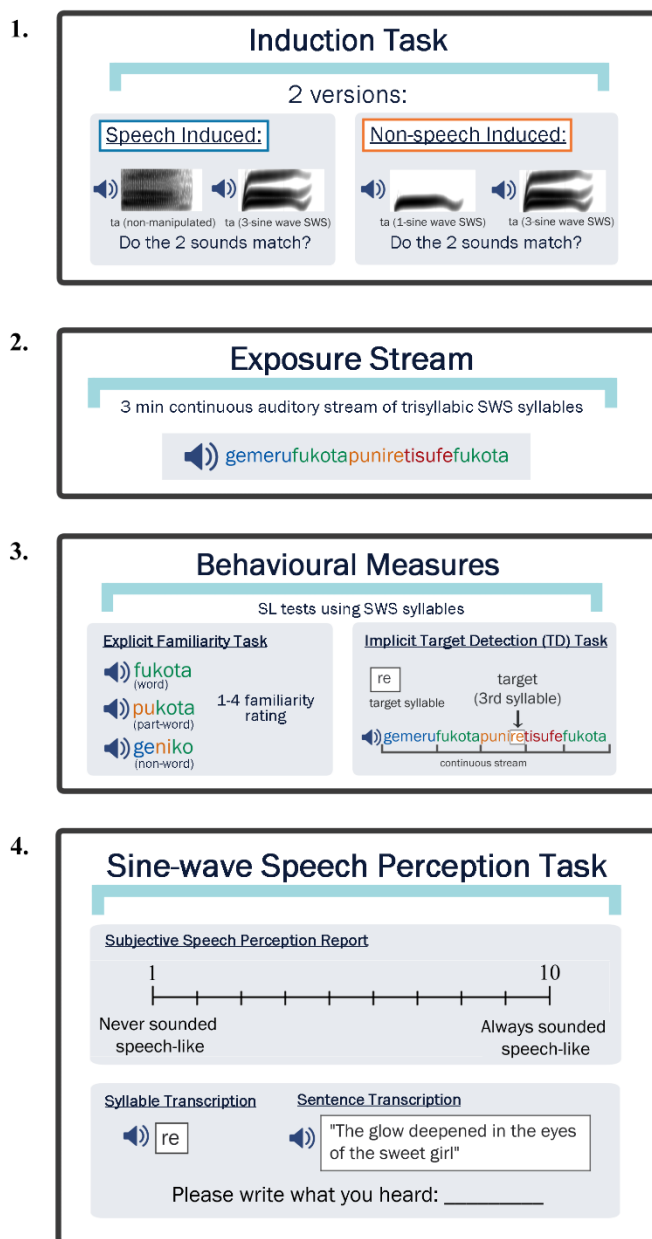
307 The 12 three-sine wave syllables were combined to create 4 trisyllabic nonsense
308 words (e.g. *tafuko*, *rigimi*, *rupuni*, *fitisu*). To form the continuous artificial speech stream,
309 these trisyllabic nonsense words were concatenated pseudorandomly, without pauses
310 between words, with the constraint that the same word never occurred consecutively.
311 Thus, the transitional probabilities of neighbouring syllables were 1.0 within a word, and
312 0.33 across word boundaries. The stream consisted of 600 syllables (200 words)
313 presented at a rate of 300 ms per syllable (i.e. 3.3 Hz), with each of the 4 words
314 repeated 50 times, for a total duration of 3 minutes. To control for potential syllable-
315 specific idiosyncrasies, the syllables in a given word were each assigned to the first,

316 second, and third position across three conditions, counterbalanced across participants
317 (Language A: *tafuko, rigimi, rupuni, fitisu*; Language B: *fukota, gimiri, puniru, tisufi*;
318 Language C: *kotafu, mirigi, nirupu, sufiti*). The experimental script was programmed in
319 jsPsych (de Leeuw et al., 2023).

320 **2.3. Procedure**

321 All tasks were performed online on the participants' own laptops or personal
322 computers. To minimize distractions during the study, participants were asked to
323 complete the tasks in a quiet listening environment and to use headphones for the
324 entire duration of the session. Each session began with a volume adjustment task
325 during which participants listened to a thirty-second noise and adjusted their sound
326 volume to a comfortable level.

327 The experimental procedure is summarized in Figure 1, and consisted of four
328 main phases, as described below. Participants completed one of two different versions
329 of the induction task depending on whether they were assigned to the SI or NSI
330 condition. All other tasks, including the key SWS stimuli, were identical between groups.



331

332 **Figure 1.** A summary of the experimental procedure. The induction task in the speech
 333 induced condition consisted of judging whether pairs of intact syllables and moderately
 334 degraded syllables matched. The induction task in the non-speech induced condition
 335 consisted of judging whether pairs of moderately degraded and heavily degraded
 336 syllables matched. Participants were exposed to 3 minutes of repeating nonsense
 337 words composed of the key SWS syllables. To measure learning, participants then
 338 completed a familiarity rating task, in which they rated the familiarity of words and foils,
 339 and a target detection task, in which they responded each time they detected a target
 340 syllable in a continuous stream consisting of the nonsense words. Finally, for each of

341 the 12 key SWS syllables, participants were asked to indicate how speech-like they
342 thought they were, and then transcribed the SWS syllables and sentences to the best of
343 their ability. Task order was identical for all participants.

344 **2.3.1. Induction Task**

345 This task was designed to induce participants to perceive the key SWS stimuli as
346 either speech (SI condition) or as non-speech (NSI condition). In this task, participants
347 were presented with “matched pairs” of syllables and instructed to intentionally learn the
348 syllable pairings. The SI participants were told that they would be listening to speech
349 syllables, and that each syllable would be followed by a distorted version of itself. They
350 were then presented with syllable pairs comprised of the intact, non-manipulated
351 version of each syllable (e.g. “fu”) followed by the target SWS version of the same
352 syllable (e.g. the three-sine-wave version of “fu”), in order to draw their attention to the
353 speech-like qualities of the SWS syllables. In contrast, the NSI participants were told
354 that they would be listening to robotic noises artificially generated by a computer. The
355 NSI participants were then presented with syllable pairs consisting of the highly
356 degraded version of each syllable (e.g. the single-sine wave version of “fu”) followed by
357 the target SWS version.

358 The task was made up of an initial training phase, followed by a test phase. In
359 the training phase, participants were simply presented with two repetitions of each of
360 the 12 pairs (24 total trials) and were instructed to pay careful attention as they would
361 be tested on the pairs later. Next, participants completed 40 test trials, comprised of 36
362 correctly paired syllables and 4 mismatched pairs. On each test trial, participants were
363 asked to judge whether the two sounds made up a correctly matched pair by pressing
364 one of two corresponding keys.

365 **2.3.2. Exposure Stream**

366 Next, participants were presented with the three-minute continuous stream of
367 nonsense words, made up of the same key SWS syllables for both induction groups.
368 They were instructed to pay attention to the stream, and were told they may be tested
369 on their knowledge of the stream later in the study. To ensure participant engagement in
370 the online testing environment, two attention checks were embedded within the
371 exposure stream, consisting of 4 s pauses inserted randomly at two of nine preselected
372 times in the stream. Prior to beginning the task, participants were instructed to listen for
373 pauses and to press the spacebar key within 4 s whenever they heard a pause. Failure
374 to detect both pauses resulted in participant exclusion from subsequent analyses.

375 **2.3.3. Statistical Learning Tasks**

376 Next, participants completed two behavioural tests of statistical learning, in the
377 order indicated below.

378 **2.3.3.1. Familiarity Rating Task**

379 This task is designed to assess explicit memory of the nonsense words (e.g.
380 Batterink & Paller, 2017, 2019). On each trial, participants listened to a syllable triplet
381 made of the key SWS syllables, and rated how familiar it sounded to them on a scale
382 from 1 (*very unfamiliar*) to 4 (*very familiar*). A total of 12 trials were presented, with 4
383 trials consisting of words from the exposure stream (e.g. *tafuko*), 4 trials consisting of
384 part-words (i.e. a syllable pair from a word in the exposure stream combined with an
385 additional syllable from a different word, e.g. *rufuko*), and 4 trials consisting of non-
386 words (syllables from the stream that had never occurred together, e.g. *rupufu*).

387 Evidence of explicit memory for the words would be provided by higher ratings to words,
388 followed by part-words, with non-words rated as least familiar.

389 **2.3.3.2. Target Detection Task**

390 This task measures participants' response times to target syllables embedded
391 within shortened versions of the speech stream, and can reveal statistical learning in the
392 form of prediction effects, in the absence of explicit memory or intentional retrieval of the
393 learned words (Batterink et al., 2015). On each trial, participants were presented with a
394 target SWS syllable; they were allowed to replay this target syllable as many times as
395 they wished. They then listened to a shortened version of the exposure stream (~14.5
396 s), containing the four trisyllabic nonsense words concatenated together four times each
397 in pseudorandom order (48 syllables total), in the same manner as the Exposure
398 stream. Participants were instructed to press the spacebar each time they heard that
399 target syllable as quickly and accurately as possible by pressing the spacebar.

400 Each of the 12 SWS syllables acted as a target three times overall, yielding a
401 total of 36 streams. Across all streams, this yielded a total of 144 targets, 48 within each
402 syllable position (1st, 2nd, 3rd). Stream order was randomized for every participant.

403 Successful learning of the speech stream would be reflected by faster reaction times to
404 target syllables that occurred in the medial or final position of a trisyllabic word relative
405 to syllables that occurred in the initial position, due to the opportunity to predict the
406 target (Batterink et al., 2015, 2019; Batterink & Paller, 2017).

407 **2.3.4. Speech Perception Task**

408 This task was designed to examine participants' perception and comprehension
409 of the key SWS stimuli, and contained three parts. As illustrated in Figure 1, this was
410 always the final task in the experiment, in order to avoid suggesting the communicative
411 nature of SWS to participants in the NSI group.

412 **2.3.4.1. Overall Subjective Speech Perception Rating.**

413 Participants were presented with an open-response textbox and asked to
414 describe the sounds that they had heard in the study. Using a slider, they were then
415 asked to rate the extent to which they had heard the SWS as speech-like, with the scale
416 ranging from 1 (*I never heard the sounds as speech*) to 10 (*I always heard the sounds*
417 *as speech*).

418 **2.3.4.2. Syllable Transcription**

419 Participants then listened to each of the 12 key SWS syllables one at a time and
420 were asked whether they thought it sounded like speech (yes/no response). If a
421 participant indicated that they heard a syllable as speech, they were then asked to
422 transcribe the syllable to the best of their ability by typing their response into an open-
423 response textbox.

424 **2.3.4.3. Sentence Transcription**

425 As a test of generalized SWS perception, participants listened to 10 SWS
426 sentences from the Harvard sentences database (IEEE, 1969) and transcribed each
427 one to the best of their ability. An example of one of the sentences is, "The glow
428 deepened in the eyes of the sweet girl." Participants were instructed to spell each word
429 as accurately as possible.

430 **2.3.5. Survey**

431 Finally, participants were redirected to a Qualtrics survey containing basic
432 demographic questions about age, gender identity, and language fluency.

433 **2.4. Statistical Analyses**

434 For all t-tests, the Student's t-test was utilized unless the assumption of equal
435 variances was violated. Welch's unequal variances t-tests were instead used whenever
436 Levene's Test was significant.

437 Bayes Factors were calculated for each test, using the default prior provided by
438 JASP. This prior uses a Cauchy distribution, centered around 0, with a width parameter
439 of 0.707. The reported Bayes Factors (BF_{10}) represent how likely the alternative
440 hypothesis is relative to the null hypothesis; values above 1 indicate evidence
441 supporting the alternative hypothesis, whereas values below 1 provide evidence
442 supporting the null hypothesis over the alternative hypothesis. As an example, a BF_{10} of
443 4 indicates that, given the data, the alternative hypothesis is four times likelier than the
444 null hypothesis. In contrast, a BF_{10} of 0.25 would indicate that the alternative hypothesis
445 is one-fourth as likely as the null hypothesis. Conventional means of interpreting the
446 relative strength of Bayes Factors regard $BF_{10} = 3-10$ as moderate evidence, such that a
447 BF_{10} of 4 suggests moderate evidence for the alternative hypothesis over the null
448 hypothesis (Schmalz et al., 2023). Bayes Factors can also be reported using BF_{01} , the
449 inverse of BF_{10} , which presents the likelihood of the null hypothesis relative to the
450 alternative hypothesis. Thus, a BF_{01} of 4 indicates that the null hypothesis is four times
451 likelier than the alternative hypothesis. BF_{10} values are reported for each test in this

452 study; however, for any tests that result in null findings, BF_{01} is also be reported for ease
453 of interpretation.

454 **2.4.1. Induction Task**

455 Each participant's accuracy on the matched pairs test was calculated.
456 Additionally, as there were many more "match" trials than "mismatch" trials, we also
457 computed d' scores as a bias-free measure of participants' sensitivity to the presence of
458 a match. D' was computed as the difference between the z-transforms of participants'
459 hit rate (i.e. the proportion of matched trials that they correctly identified as matching)
460 and false alarm rate (the proportion of mismatched trials that they incorrectly identified
461 as matching) in the task.

462 **2.4.2. Statistical Learning Tasks**

463 For all analyses of the statistical learning tasks, Greenhouse–Geisser corrections
464 were reported for factors with more than two levels.

465 **2.4.2.1. Familiarity Task**

466 Average familiarity ratings were computed for each word category (Word,
467 Partword, Nonword) and entered into a 2x3 mixed effects ANOVA with induction
468 condition (speech induced, non-speech induced) as a between-subjects factor and word
469 category (non-word, part-word, word) as a within-subjects factor.

470 Additionally, for subsequent correlational analyses, "familiarity rating scores"
471 (Batterink & Paller, 2017, 2019) were calculated by subtracting the average of a
472 participants' rating of partwords and nonwords from their average rating of a word.

473 Perfect sensitivity to words over foils on this measure would be a score of 3, with any
474 positive value suggestive of learning, as this would reflect higher scores for words
475 compared to both pseudo- and non-words.

476 **2.4.2.2. Target Detection Task**

477 Following the inclusion criteria of previous studies, responses that occurred
478 within 1200 ms following target onset were considered valid hits (Batterink & Paller,
479 2017, 2019). All other responses were considered false alarms.

480 **2.4.2.2.1. Detection Score**

481 For each participant, we first calculated the number of targets that were correctly
482 detected and the total number of false alarms. We then computed an overall “detection
483 score,” which represents a conservative estimate of a participant’s sensitivity to the
484 targets in the stream, computed as the overall number of hits divided by the overall
485 number of false alarms (Number of Hits/Number of False Alarms). Given that the “target
486 response” window (4 targets x 1200 ms = 4800 ms) for each stream was half the length
487 of the “false alarm” windows (total stream length of 14400 ms – “target response” length
488 of 4800 ms = 9600 ms), we reasoned that any score *greater* than 0.5 would provide
489 evidence of above-chance detection performance (with 0.5 indicating that hits occurred
490 half as frequently as false alarms, as would be expected if responses were distributed
491 randomly across the stream, without regard for the actual target locations). In other
492 words, a detection score of >0.5 would indicate that participant’s responses were more
493 likely to occur within a “target response” window than a “false alarm” window, providing
494 evidence of target detection at above-chance levels.

495 **2.4.2.2.2. Reaction Time**

496 In addition to already-reported exclusions (see section 2.1.), 3 additional
497 participants who only responded to initial targets were excluded from the RT analysis,
498 as their mean response times could not be computed for second and third position
499 targets. Furthermore, participants with a detection score of 0.5 or below were also
500 excluded from this analysis ($n = 32$). We reasoned that if a participant is unable to
501 detect the syllables at an above-chance level, any differences in their RTs cannot be
502 considered a valid measure of statistical learning. To summarize, 35 additional
503 participants were excluded from this analysis, yielding a final n of 92 participants. 52 of
504 these participants completed the speech induction (mean age = 40.0 y, SD = 11.5 y; 28
505 men; 24 women), and the remaining 40 were from the NSI group (mean age = 39.7 y,
506 SD = 13.1 y; 19 men; 21 women). For thorough reporting, a parallel analysis that also
507 includes data from participants who scored below chance on detection can be found in
508 Supplementary Materials ($n = 124$).

509 For each participant, mean RTs for detected targets were calculated for each
510 target position (initial, medial, final). Mean RTs were then entered into a 2 x 3 repeated-
511 measures ANOVA with induction group as the between-subject factor and target
512 position (initial, medial, final) as the within-subject factor. In addition, to quantify
513 statistical learning performance using a single metric while controlling for individual
514 differences in baseline response times, a “RT prediction score” was computed by
515 subtracting the average RT for the final syllable position from the average RT for the
516 initial syllable position and dividing it by the average RT for the initial syllable position
517 $[(RT_1 - RT_3) / RT_1]$; Batterink & Paller, 2019]. This calculation adjusts for potential

518 differences in baseline RTs between individuals, allowing us to measure statistical
519 learning across individuals with different RT baselines.

520 **2.4.3. Speech Perception Tasks**

521 **2.4.3.1. Syllable Transcription**

522 Scoring for this task was done by allocating 1 point for each syllable that was
523 fully correctly transcribed (with alternative spellings such as “mee” or “me” designated
524 as correct), and 0.5 points for each syllable that was partially correct, with either the
525 consonant or vowel transcribed correctly (e.g. typing “mee” when the SWS syllable
526 being played is “gee”). Average accuracy across the 12 total syllables in the task was
527 then computed for each participant.

528 **2.4.3.2. Sentence Transcription**

529 Each SWS sentence contained 5 keywords (e.g. in the sentence “Pluck the bright
530 rose without leaves” the keywords would be “pluck,” “bright,” “rose,” “without,” and
531 “leaves”). While participants wrote out the entire sentence, their scores were calculated
532 as the proportion of correctly transcribed keywords. Misspelled words were marked as
533 incorrect.

534 **3. Results**

535 We first report the results from the induction task. Following this, we then
536 characterize participants’ perception of the key SWS stimuli, as assessed through our
537 three speech perception tasks (Figure 1). Although these speech perception tasks were
538 completed at the end of the session, we report these results second, as they are

539 needed to understand the subsequent statistical learning analyses. We then turn to our
540 main set of results, which concerns performance on our two measures of statistical
541 learning—the familiarity rating and the target detection tasks—and how performance on
542 these tasks relates to perception of SWS stimuli.

543 **3.1. Induction Task**

544 Participants generally performed well on the matched pairs test, with an average
545 accuracy rate of 90.7% (SD = 8.1%). Not surprisingly, given that they were presented
546 with non-degraded syllable primes, speech induced (SI) participants outperformed non-
547 speech induced (NSI) participants on this task (SI: mean = 94.9%; SD = 5.2%; NSI:
548 mean = 85.4%; SD = 8.1%; $t(88.96) = -7.64$, $p < .001$, $d = -1.40$; $BF_{10} = 9.79 \times 10^9$).

549 The average d' was 2.33 (SD = 1.05), with SI participants also outperforming NSI
550 participants on this measure (SI: mean = 2.93; SD = 0.82; NSI: mean = 1.58; SD = 0.79;
551 $t(125) = -9.39$, $p < .001$, $d = -1.68$; $BF_{10} = 1.27 \times 10^{13}$).

552 **3.2. Speech Perception Tasks**

553 **3.2.1. Overall Subjective Speech Perception Rating**

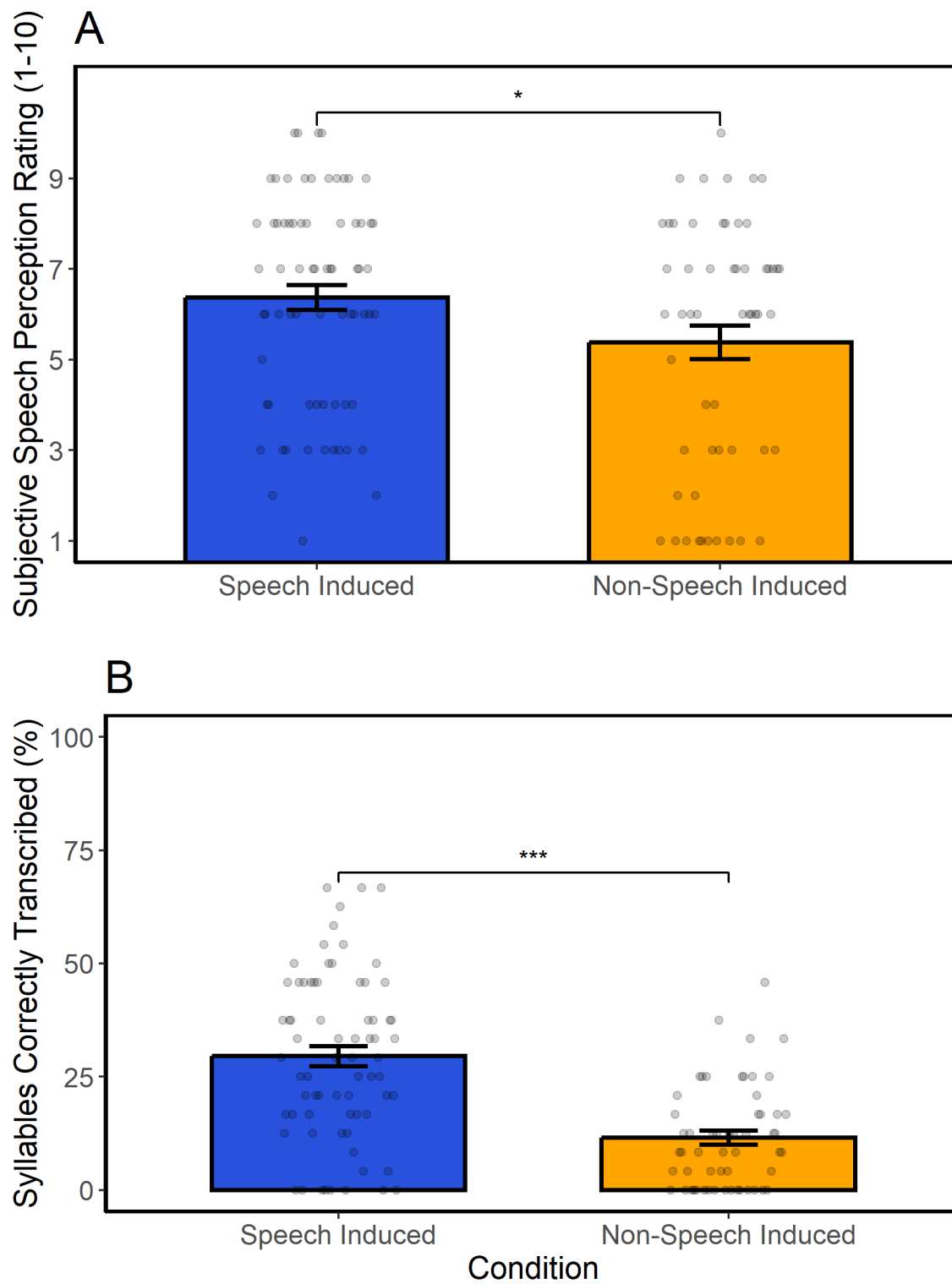
554 Reponses on the scale, ranging from 1 to 10, showed that SI participants (M =
555 6.37, SD = 2.32) rated the SWS as sounding significantly more speech-like overall than
556 the NSI participants (M = 5.38, SD = 2.79), $t(106.71) = -2.14$, $p = .035$, $d = -0.39$; $BF_{10} =$
557 1.62. Nonetheless, there was considerable overlap in the scores, such that some NSI
558 participants perceived the stimuli to sound more speech-like, while some SI participants
559 perceived the stimuli to not sound very speech-like. The distribution of participant
560 responses on the scale are presented in Figure 2A.

561 **3.2.2. Syllable Transcription**

562 As expected, participants in the SI group (M = 53.8%, SD = 28.3%) judged a
563 significantly higher percentage of SWS syllables to be speech-like compared to the NSI
564 participants (M = 35.7%, SD = 29.2%), $t(125) = -3.52$, $p < .001$, $d = -0.63$; $BF_{10} = 44.24$.
565 Additionally, SI participants (M = 29.5%, SD = 18.3%) also correctly transcribed a
566 significantly larger proportion of the 12 SWS syllables than the NSI participants (M =
567 11.5%, SD = 11.2%), $t(118.36) = -6.82$, $p < .001$, $d = -1.19$; $BF_{10} = 4.34 \times 10^6$ (see
568 Figure 2B).

569

570



571

572 **Figure 2.** (A) The distribution of participant responses on the subjective speech
 573 perception scale. The error bars represent the standard error of the mean. (B)

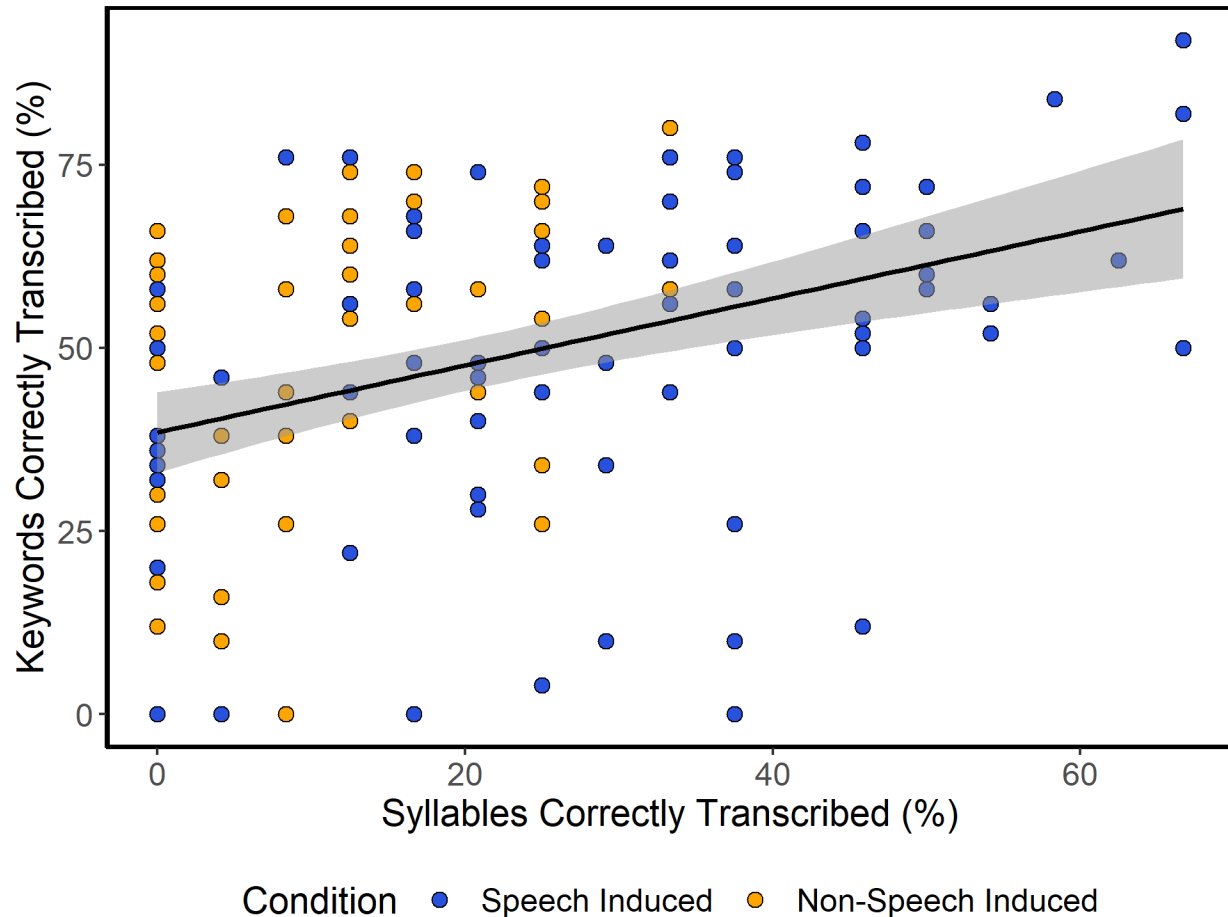
574 Participants' accuracy in syllable transcription task. The error bars represent the
575 standard error of the mean. * $p < .05$; *** $p < .001$.

576 **3.2.3. Sentence Transcription**

577 Participants correctly transcribed 48.4% of the keywords in total (SD = 21.4%).

578 Somewhat unexpectedly, there was no significant difference in the keyword
579 transcription accuracy between SI participants (M = 49.4%, SD = 22.2%) and NSI
580 participants (M = 47.0%, SD = 20.5%), $t(125) = -0.64$, $p = .521$, $d = -0.12$; $BF_{10} = 0.23$
581 [$BF_{01} = 4.35$]. This suggests that the speech induction training on individual syllables did
582 not generalize to novel sentences. However, across all participants, there was a
583 significant positive correlation between accuracy on the syllable transcription task and
584 sentence transcription task, $r(125) = 0.38$, $p < .001$; $BF_{10} = 1867.35$ (see Figure 3),
585 suggesting that performance on these two tasks reflects a common ability.

586



587
 588 **Figure 3.** The correlation between the percentage of SWS sentences and the key SWS
 589 syllables that participants transcribed accurately ($r = 0.38, p < .001$).

590 3.3. Statistical Learning Tasks

591 As just described, while the two induction groups showed significant differences
 592 on self-reported subjective speech perception and on SWS syllable transcription
 593 accuracy, there was considerable overlap between the groups on these measures. In
 594 addition, there were no group differences on the sentence transcription task. These
 595 results indicate that our speech perception manipulation only partially altered
 596 participants' perception of the key SWS syllables, rather than producing a dramatic
 597 transformation of participants' percepts. Thus, as a further test of the relationship
 598 between statistical learning and speech perception, we examined correlations between

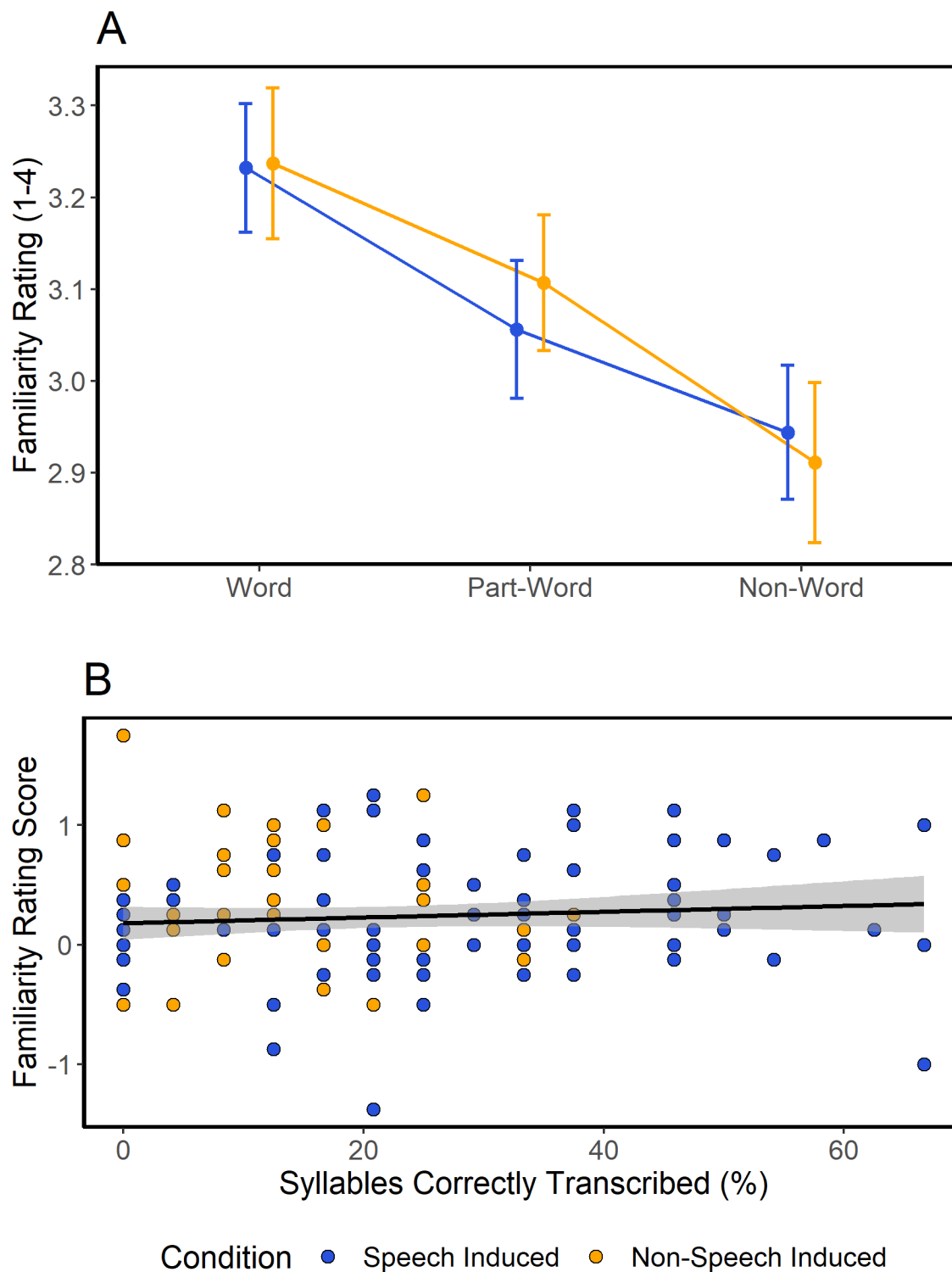
599 participants' accuracy on the SWS syllable transcription task—taking this as a measure
600 of speech perception—and their statistical learning performance. Hence, in the following
601 section, for both our measures of statistical learning, we report (1) differences in
602 performance between our two *a priori* defined groups and (2) correlations between
603 accuracy on the syllable transcription task and statistical learning performance.

604 **3.3.1. Familiarity Task**

605 As expected, across both induction groups, words were rated as the most
606 familiar, followed by part-words, with non-words rated as the least familiar, leading to a
607 significant effect of word type, $F(1.98,248.18) = 18.00$, $p < .001$, $\eta^2p = 0.13$; $BF_{10} = 2.48$
608 $\times 10^5$ (see Figure 4A).

609 Supporting the hypothesis that statistical learning operates in a similar manner
610 across stimuli perceived as linguistically-relevant and irrelevant, performance on the
611 familiarity rating task was not significantly different between the two induction groups
612 (Main Effect of Induction: $F(1,125) = 6.45 \times 10^{-3}$, $p = .936$, $\eta^2p = 5.16 \times 10^{-5}$; $BF_{10} = 0.22$
613 [$BF_{01} = 4.54$]; Word Type \times Induction: $F(1.98,248.18) = 0.34$, $p = .714$, $\eta^2p = 0.0027$;
614 $BF_{10} = 0.07$ [$BF_{01} = 14.3$]).

615 Further, there was no significant correlation between participants' syllable
616 transcription accuracy and their familiarity rating scores, $r(125) = 0.09$, $p = .34$, with the
617 Bayes Factor indicating moderate evidence (Schmalz et al., 2023) for the null
618 hypothesis of no relation between these two measures ($BF_{10} = 0.18$ [$BF_{01} = 5.55$]; see
619 Figure 4B). This result indicates that more accurate perception of the stimuli as syllables
620 did not lead to better performance on the familiarity measure.



621

622 **Figure 4.** (A) Participants' ratings of triplet familiarity from the familiarity rating task. The
 623 error bars represent the standard error of the mean. (B) The correlation between

624 participants' familiarity rating score and the percentage of key SWS syllables that they
625 transcribed accurately ($r = 0.09$, $p = .34$).

626 **3.3.2. Target Detection Task**

627 **3.3.2.1. Overall Detection Rate**

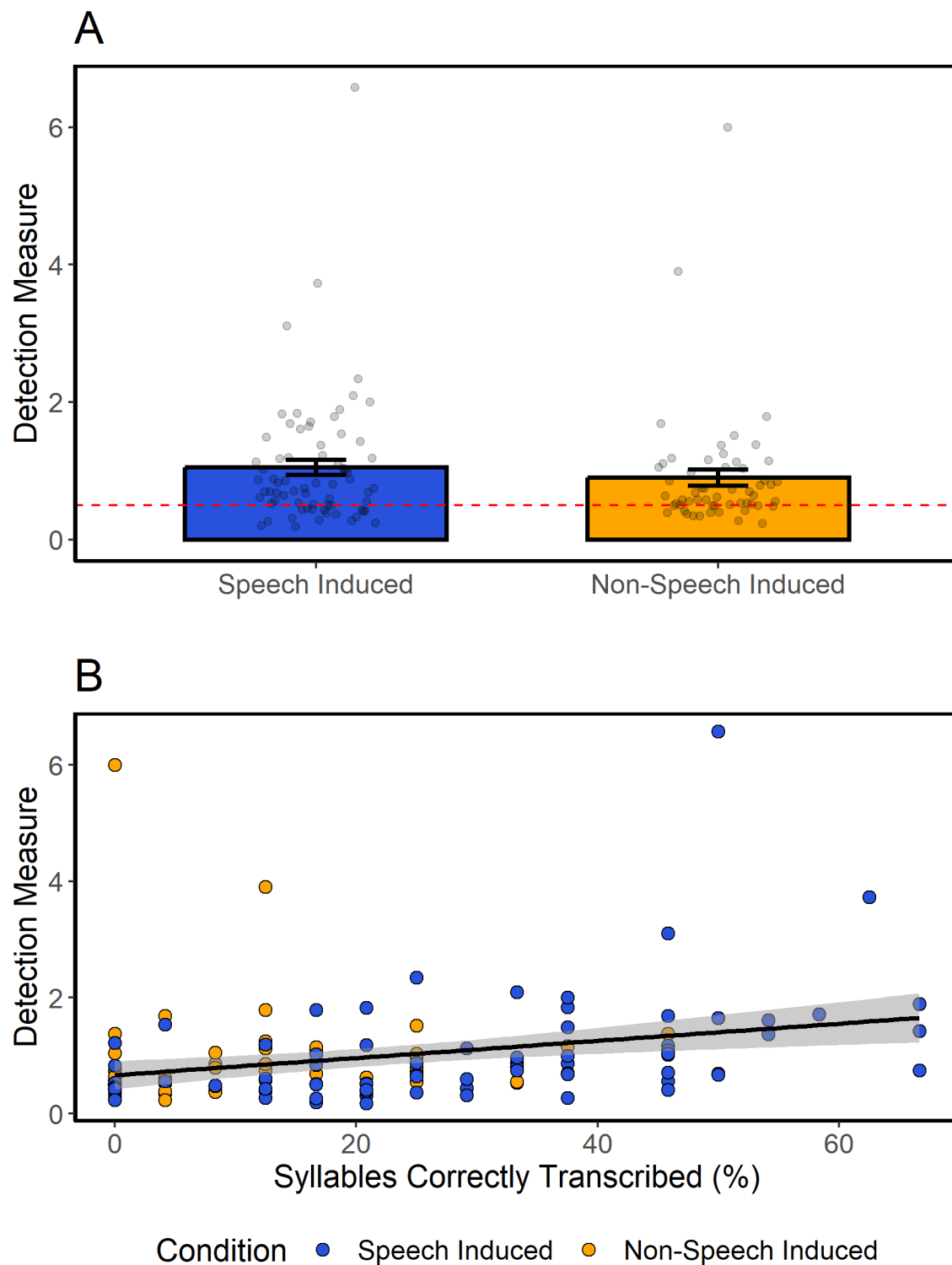
628 Participants correctly responded to 67.4% (SD = 20.0%) of the targets on
629 average and made an average of 148.7 false alarms total (SD = 101.2). Accuracy rate
630 was relatively low and false alarms were relatively high compared to previous versions
631 of this task (e.g. Batterink et al., 2015; Batterink & Paller, 2017, 2019). This relatively
632 poor performance may be attributed to the manipulated nature of the syllables, which
633 made them more difficult to identify. Nonetheless, participants performed significantly
634 above chance, as assessed by the detection score ($M = 0.98$, $SD = 0.92$; $t(126) = 5.91$,
635 $p < .001$, $d = 0.52$; chance is 0.5 on this measure), with no significant difference in
636 performance between the SI participants ($M = 1.05$, $SD = 0.95$) and NSI participants (M
637 $= 0.90$, $SD = 0.89$), $t(125) = -0.93$, $p = .355$, $d = -0.17$; $BF_{10} = 0.28$ [$BF_{01} = 3.57$].

638 Interestingly, there was a significant positive correlation between the Detection
639 Measure values and syllable transcription accuracy, $r(125) = 0.29$, $p = .001$; $BF_{10} =$
640 22.39 , as presented in Figure 5. This result indicates that participants who more
641 accurately perceived the stimuli as syllables were also better able to detect them in the
642 continuous speech sequences.

643

644

645



646

647 **Figure 5.** (A) Participants' detection score values on the target detection task (chance is
 648 0.5). The error bars represent the standard error of the mean. (B) The correlation

649 between participants' Detection Measure values on the target detection task and the
650 percentage of key SWS syllables that they transcribed accurately ($r = 0.29$, $p = .001$).

651 **3.3.2.2. Reaction Time**

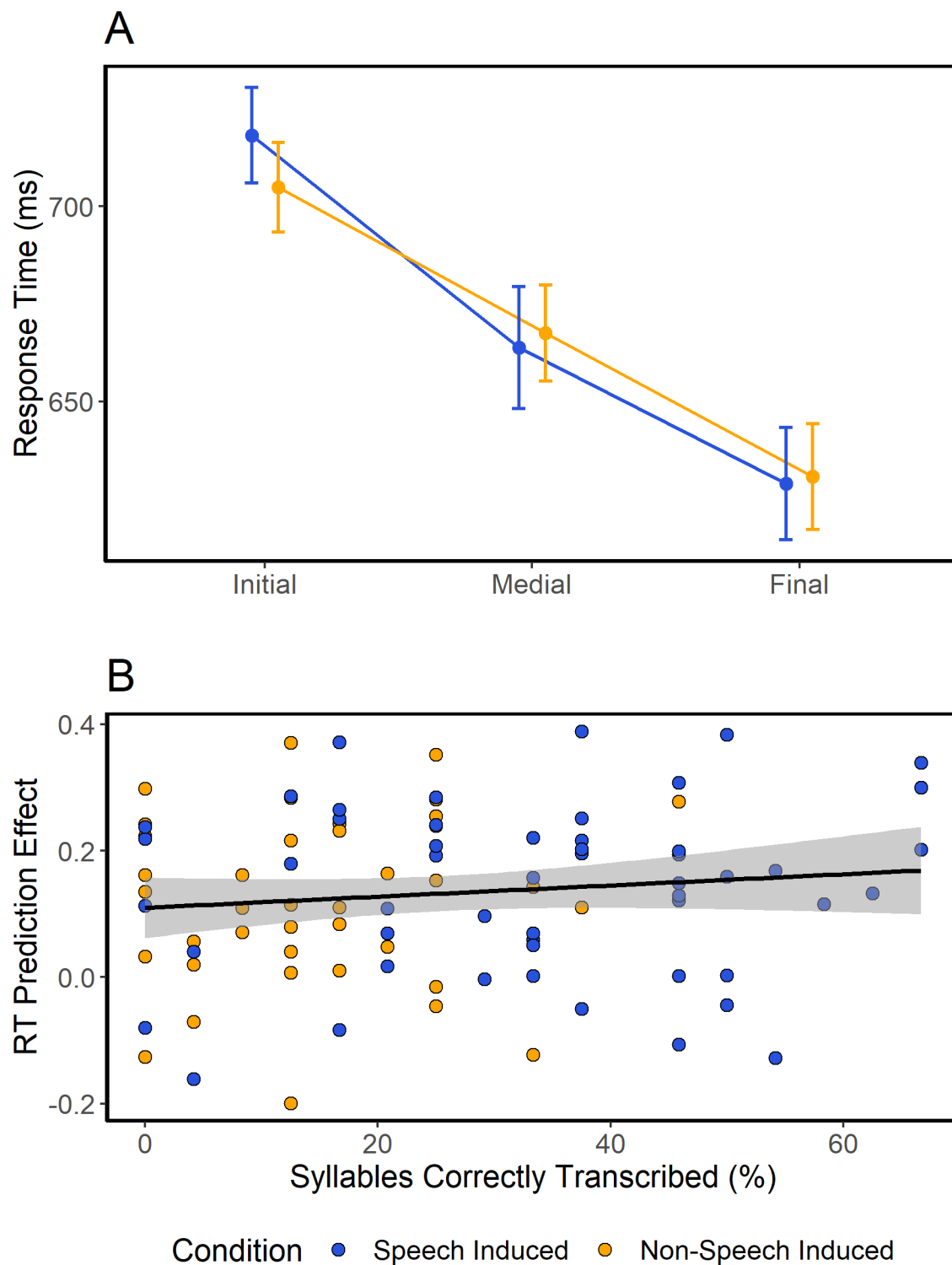
652 As expected, across both groups, RTs were the fastest for final-position
653 syllables, second fastest for medial-position syllables, and slowest for initial-position
654 syllables, as shown in Figure 6A, leading to a significant effect of syllable position,
655 $F(1.68, 150.76) = 61.69$, $p < .001$, $\eta^2p = 0.41$; $BF_{10} = 4.17 \times 10^{18}$. Notably, there was no
656 significant difference in the RTs between induction groups, either overall or as a
657 function of syllable position (Main Effect of Induction: $F(1, 90) = 0.06$, $p = .802$, $\eta^2p =$
658 7.00×10^{-4} ; $BF_{10} = 0.24$ [$BF_{01} = 4.17$]; Position x Induction: $F(1.68, 150.76) = 1.14$, p
659 $= .315$, $\eta^2p = 0.01$; $BF_{10} = 0.19$ [$BF_{01} = 5.26$]).

660 Additionally, there was no significant correlation between RT prediction effect
661 and syllable transcription accuracy, $r(90) = 0.12$, $p = .253$; $BF_{10} = 0.25$ [$BF_{01} = 4.00$], as
662 shown in Figure 6B. This suggests more accurately perceiving the SWS stimuli as
663 syllables did not lead to an enhanced ability to predict final position syllables. For a
664 summary of the Bayes Factors for the study's statistical learning measures, see Table
665 1.

666 While the above analysis excludes participants who failed to detect syllables at
667 above-chance levels, we also report results from the full sample (see Supplementary
668 Materials). We note that the overall pattern of findings is largely similar between the two
669 analyses.

670

671



672

673 **Figure 6.** (A) Participants' average reaction times for each of the syllable positions in
 674 the target detection task. The error bars represent the standard error of the mean. (B)

675 The correlation between participants' RT prediction effect and the percentage of key
676 SWS syllables that they transcribed accurately ($r = 0.12, p = .253$).

677

678 **Table 1**

679 *Summary of Bayes Factor Results for Statistical Learning Performance*

Task	BF ₀₁	Strength of evidence in favour of null
Familiarity Task		
Main Effect of Induction	4.54	Moderate
Word Type x Induction	14.29	Strong
Correlation	5.55	Moderate
Target Detection Task		
Main Effect of Induction	4.17	Moderate
Position x Induction	5.26	Moderate
Correlation	4.00	Moderate

680 *Note.* Moderate evidence: BF₀₁ = 3-10. Strong evidence: BF₀₁ = 10-30. The null
681 hypothesis here indicates no impact of speech perception on statistical learning
682 performance.

683

4. Discussion

684 In the current study, we examined whether statistical learning occurs more
685 robustly for sounds subjectively perceived as speech relative to those perceived as non-
686 speech, independently of stimulus acoustics. The key novel aspect of the current study
687 was the use of SWS to eliminate acoustic differences between stimuli perceived
688 linguistically versus non-linguistically. Overall, we found that statistical learning operates
689 similarly for stimuli, regardless of the degree to which they are perceived as
690 linguistically-relevant. Participants who were induced into hearing syllables as speech-
691 like did not show any significant differences in performance on our two statistical
692 learning measures compared to participants induced into hearing the syllables as non-

693 linguistic sounds. In addition, participants' ability to linguistically label individual SWS
694 syllables did not predict their statistical learning performance. Taken together, these
695 results provide no strong evidence of a statistical learning advantage for sounds
696 perceived as more speech-like, instead suggesting that statistical learning occurs
697 indiscriminately across auditory stimuli, regardless of their linguistic relevance.

698 More specifically, on the familiarity rating task, we observed no significant
699 difference in ratings between the speech induced and non-speech induced group, as
700 well as no significant correlation between participants' accuracy in transcribing the SWS
701 syllables and their familiarity rating score. Similarly, on the target detection task, there
702 was no significant difference in the RTs between the induction groups, nor was there a
703 significant correlation between participants' SWS syllable transcription accuracy and the
704 magnitude of their RT prediction effect. Thus, taken together, our results suggest that
705 statistical learning operates largely similarly across physically identical auditory stimuli,
706 regardless of participants' perception of the stimuli as more or less speech-like.

707 Importantly, we found that the speech induced (SI) group was better at identifying
708 the SWS syllables by their linguistic labels than the non-speech induced (NSI) group, as
709 demonstrated by significantly higher accuracy on the syllable transcription task (30%
710 accuracy for the SI group versus 12% for the NSI). We also found that participants in
711 the SI group rated the syllables as subjectively more speech-like than participants in the
712 NSI group, although the difference in subjective ratings were small. These findings
713 provide a key manipulation check and indicate that our induction task did produce
714 differences in the subjective perception of SWS syllables between the two groups.
715 However, we note that our induction task did not produce a dramatic perceptual

716 transformation of the syllables, as can be found when sentences are used as stimuli
717 (Davis & Johnsrude, 2007; Remez et al., 1981), and was also limited in its
718 generalizability, with no effect on participants' ability to transcribe full sentences. We
719 return to this general point in the Limitations section.

720 Previous findings in the literature have suggested that statistical learning shows
721 important differences across domains and may be governed by modality- and domain-
722 specific constraints (e.g., Siegelman & Frost, 2015; Siegelman et al., 2017; Frost et al.,
723 2015; Conway et al., 2020; Van Hedger et al., 2022). For example, several findings
724 point to the idea that statistical learning is influenced by the shared resemblance
725 between novel words in the speech stream and existing words in learners' native
726 language, with words that share native language phonotactic patterns being more easily
727 segmented and/or subsequently recognized (Siegelman et al., 2018; Elazar et al., 2022;
728 Finn & Hudson Kam, 2008). Our results provide initial evidence that domain-specific
729 constraints for statistical learning are at least partially attributable to sensory-level
730 processes, and not necessarily to higher-level cognitive mechanisms related to the
731 conceptual categorization of incoming stimuli. For example, networks in auditory cortex
732 may be better equipped to process and encode incoming novel words that have high
733 acoustic overlap with existing words in the learner's lexicon, which in turn could facilitate
734 binding between syllables and lead to observed "linguistic entrenchment" effects
735 (Siegelman et al., 2018). In contrast, the judged linguistic relevance of an ambiguous
736 signal may be a later-occurring, downstream process that does not directly impact
737 statistical learning.

738 Our approach differed from several previous statistical learning studies in that we
739 did not directly compare learning of speech versus non-speech stimuli (cf. Hoch et al.,
740 2013; Marcus et al., 2007; Ordin et al., 2021; Saffran, 2002; Saffran et al., 1999;
741 Siegelman et al., 2018), which differ in both low-level acoustic features and in
742 communicative relevance. Instead, we assessed the statistical learning of acoustically
743 identical ambiguous stimuli that differed in the degree to which they were subjectively
744 perceived as speech, allowing us to address the more specific question of whether the
745 subjective linguistic value (Berent et al., 2021; Rabagliati et al., 2018) of auditory
746 stimuli—in and of itself—influences statistical learning. To our knowledge, no previous
747 study has directly examined this question in adults. However, there is some relevant
748 prior work in infants, which has examined whether the meaningfulness or
749 communicative relevance of stimuli increases infants' success in learning abstract
750 repetition rules (such as AAB or ABA). Ferguson and Lew-Williams (2016) presented
751 infants with a video prime in which tones were embedded in a natural conversation
752 between two actors, thereby inducing the infants to believe that tones are a
753 communicative signal. In a subsequent rule learning phase, infants who were
754 communicatively primed successfully learned abstract rules from tones, whereas
755 unprimed infants failed to show learning. This finding suggests that infants learn better
756 from stimuli that are communicatively relevant. Supporting this conclusion, a recent
757 meta-analysis of 20 papers (Rabagliati et al., 2018) found that infants are better able to
758 learn abstract repetition rules from stimuli that are communicatively or ecologically
759 meaningful—such as spoken syllables, communicatively primed tones, or natural
760 categories such as dogs or faces—than meaningless stimuli such as geometric shapes

761 or tones. In a follow-up experiment designed to directly test this idea, Rabagliati and
762 colleagues (2018) had infants view either a prime video that portrayed gestures as
763 communicative and meaningful, or a control video, and then exposed them to
764 sequences of gestures following an ABB or ABA pattern. Again, as in Ferguson and
765 Lew-Williams (2016), only infants primed to view gestures as a communicative signal
766 displayed evidence of rule learning. Altogether, these studies suggest that the
767 communicative status of a stimulus enhances abstract rule learning in infants.

768 In contrast to this general finding in infants, the present results fail to support the
769 idea that the perceived linguistic relevance of auditory stimuli influences or enhances
770 statistical learning in adults. This divergence could potentially be attributed to any
771 number of factors that differ between prior work in infants and the current study,
772 including the population under investigation (adults versus infants), the type of learning
773 (abstract grammatical rule learning versus statistical learning of embedded words in
774 continuous speech), and/or the experimental manipulation used to bias the linguistic
775 relevance of the stimuli. For example, it may be the case that infants show larger
776 differences in learning between communicative and noncommunicative signals
777 compared to adults, in line with the idea that infancy represents a critical period for
778 language acquisition, during which the brain is highly tuned to speech and other
779 communicative signals (Vouloumanos et al., 2010; Vouloumanos & Werker, 2004, 2007;
780 Werker & Hensch, 2015). Another possibility is that findings from abstract grammatical
781 rule learning (e.g., learning of rules such as AAB or ABA) are not directly generalizable
782 to the type of statistical learning under investigation in the current study. Rule learning
783 involves extracting an abstract rule and generalizing to novel instances, whereas

784 statistical learning involves extracting repeating, item-based regularities from
785 unsegmented input, without a generalization component. While these two types of
786 learning appear to be closely related in certain ways (Aslin & Newport, 2012), they may
787 be influenced by different factors and operate under different sets of constraints
788 (Endress & Bonatti, 2007; Endress & Mehler 2009; Peña et al., 2002; Thiessen, 2017).

789 Finally, we must also consider the possibility that our SWS manipulation did not
790 produce sufficiently diverse percepts of the identical stimuli across individual
791 participants to produce robust differences in statistical learning. Most prior work
792 investigating the processing and intelligibility of SWS have used meaningful sentences
793 (Corcoran et al., 2023; Khoshkhoo et al., 2018; Remez et al., 1981). In contrast, we
794 applied the sine-wave manipulation to isolated syllables, such that participants'
795 perception of the SWS stimuli could not benefit from top-down prediction provided by
796 semantic context. Thus, it is conceivable that even participants who achieved high
797 scores on syllable transcription accuracy may not have experienced a clear speech
798 percept for each syllable. However, a critical point arguing against this possibility is that
799 we did find a significant and highly robust correlation between participants' individual
800 syllable transcription accuracy and overall detection performance for individual syllables
801 in the target detection task. Based on this result, we can conclude that participants
802 experienced real, meaningful variability in their perceptions of the SWS stimuli that was,
803 at minimum, sufficient to robustly predict performance on a separate task. That we did
804 not find similar robust correlations between syllable identification and statistical learning
805 performance suggests that any speech-perception advantage in statistical learning—if it
806 exists at all—is likely to be very small.

807 The finding that syllable comprehension accuracy predicted overall syllable
808 detection performance in the target detection task is also interesting in and of itself. This
809 result suggests that ability to perceive ambiguous auditory stimuli as more speech-like
810 and the ability to correctly assign linguistic labels to those stimuli facilitate the online
811 identification of the ambiguous stimuli under challenging circumstances, i.e., when the
812 target stimulus is embedded within a continuous stream of similar-sounding sounds. An
813 analogous finding has been reported in the visual domain using a visual search
814 paradigm (Lupyan & Spivye, 2008; Klemfuss et al., 2012). Participants in these studies
815 were presented with arrays of rotated numbers (“2” and “5”), and were asked to indicate
816 for each trial whether the display was homogenous or contained an oddball.
817 Interestingly, participants who were given the linguistic labels or who spontaneously
818 noticed that the shapes were rotated numbers were faster to respond to the arrays
819 compared to participants who were told that the stimuli were abstract shapes. One
820 proposed explanation for this result is that the top-down effects of a linguistic cue may
821 sharpen visual feature detectors, with feedback connections from linguistic
822 representations providing a mechanism for biasing or amplifying activity in perceptual
823 detectors associated with those representations (Lupyan & Spivye, 2008). An
824 alternative explanation is that the benefit of linguistic cues on stimulus identification may
825 occur because language provides a “ready form of efficient coding,” thereby reducing
826 the burden on working memory (Klemfuss et al., 2012). Similar mechanisms operating
827 at both the perceptual and post-perceptual level could also explain the current findings.
828 The ability to perceptually transform a degraded, ambiguous target stimulus into a
829 verbalizable syllable (e.g. “ba”) may have sharpened auditory feature detectors for that

830 sound signal, and may also have facilitated the maintenance of the target stimulus in
831 working memory during the subsequent stream presentation.

832 **4.1. Limitations**

833 As previously alluded to, a limitation in this study was that the speech induction
834 task had only a moderate impact on participants' overall subjective speech perception.
835 As shown in Figure 2A, the speech induction manipulation did not cleanly divide
836 participants into two groups, as some speech-induced participants indicated that they
837 perceived the sounds as relatively un-speechlike, and vice-versa for the non-speech
838 induced participants. In addition, the speech induced group's transcription accuracy of
839 the SWS syllables—while better than the non-speech induced group's—was still fairly
840 low (approximately 30% accuracy). An ideal induction manipulation would have led all
841 the speech-induced participants to accurately perceive the SWS stimuli as speech, and
842 the non-speech induced participants to report hearing the stimuli as non-speech, as was
843 our original intention. This would have allowed for a cleaner comparison between
844 participants speech-induced and non-speech-induced participants, capitalizing on the
845 benefits of an experimental design using random assignment. Because our induction
846 did not result in a clear division between groups, and to account for the continuous, non-
847 binary nature of speech perception, we adopted a complementary approach that tested
848 whether an individual's syllable transcription accuracy predicted their statistical learning
849 performance. However, with this approach there is a possibility that any correlations
850 between transcription performance and statistical learning performance (should they be
851 observed) could be inflated by unintended third variables, such as an individual's
852 general motivation or interest in the experimental tasks. Ultimately, we believe it would

853 be challenging to design a perfectly effective speech induction task when using isolated
854 syllables as SWS stimuli, given their processing cannot benefit from top-down lexical
855 information, which plays an important role in the perceptual learning of distorted speech
856 (Davis et al., 2005). To further probe the role of linguistic relevance in statistical
857 learning, future work could leverage other types of experimental manipulations, such as
858 using priming videos to induce participants into believing that neutral stimuli are a
859 communicative signal (e.g., Ferguson & Lew-Williams, 2016; Rabagliati et al., 2018).

860 Finally, while the current study demonstrates that overall statistical learning
861 performance is similar as a function of listeners' subjective speech perception, our study
862 design does not allow us determine whether this equivalent performance is supported
863 by a common underlying mechanism or set of mechanisms, or by different mechanisms
864 that depend on speech perception. For example, it is possible that triplets perceived as
865 nonspeech may be segmented and learned as holistic or gestalt-like units, whereas
866 triplets perceived as speech may be learned by extracting sequential syllable patterns—
867 pairs and then triplets—unfolding over time. The theoretical possibility of different
868 mechanisms varying by stimulus material is supported by findings by Siegelman and
869 colleagues (2018), as previously mentioned in the Introduction. This study
870 demonstrated similar overall levels of statistical performance for auditory non-verbal
871 stimuli (everyday sounds) and syllables, which nonetheless belied important differences
872 in the internal consistency of test items between conditions, reflecting different
873 influences on performance that vary by domain. Although we would consider that the
874 possibility of different mechanisms that are equally effective to not necessarily represent
875 the most parsimonious explanation for the current data, the present study design cannot

876 rule it out. Future studies could leverage approaches such as EEG or neuroimaging to
877 examine this possibility directly.

878 **4.2. Conclusions**

879 In summary, our results provide evidence that statistical learning operates largely
880 indiscriminately across auditory stimuli, regardless of the degree to which they are
881 perceived linguistically. In contrast, linguistic perception robustly improves the
882 identification of individual target stimuli embedded in a continuous auditory sequence.
883 These results generally support previous findings of similar statistical learning
884 performance for speech stimuli and non-speech stimuli (Saffran et al., 1999; Saffran,
885 2002; Siegelman et al., 2018), and raise the possibility that previous demonstrations of
886 the statistical learning advantage for verbal materials (e.g., Hoch et al., 2013; Ordin et
887 al., 2021) may mainly be driven by acoustic differences between the classes of stimuli.
888 These results contribute to the literature on domain-specific versus domain-general
889 contributions to statistical learning, suggesting that statistical learning may be
890 conceptualized as a largely bottom-up mechanism that undiscerningly captures
891 regularities in input regardless of higher-level context.

892

893

894

895

896

897 **Data Availability**

898 All data associated with this manuscript are available on Open Science
899 Framework (https://osf.io/jqmx/?view_only=d7a7d891d2e54a05ad15fe2277dfeb05).

900

901 **Competing interests**

902 The authors declare no competing interests.

903

904 **Acknowledgements**

905 This research was supported by a Natural Sciences and Engineering Research
906 Council (NSERC) Discovery Grant (2019-05132) to Laura Batterink.

907

908 **CRedit Author Statement**

909 **Sierra Sweet:** writing – original draft, formal analysis, visualization, resources,
910 investigation; **Stephen Van Hedger:** visualization, writing – review and editing,
911 resources, investigation, conceptualization; **Laura Batterink:** conceptualization, writing
912 – review and editing, resources, methodology, supervision, funding acquisition

913

914

915

References

- 916
917 Arciuli, J., & Simpson, I. C. (2011). Statistical learning in typically developing children:
918 The role of age and speed of stimulus presentation. *Developmental Science*,
919 14(3), 464–473. <https://doi.org/10.1111/j.1467-7687.2009.00937.x>
- 920 Aslin, R. N. (2017). Statistical learning: A powerful mechanism that operates by mere
921 exposure. *Wiley Interdisciplinary Reviews: Cognitive Science*, 8(1–2),
922 10.1002/wcs.1373. <https://doi.org/10.1002/wcs.1373>
- 923 Aslin, R. N., & Newport, E. L. (2012). Statistical learning: From acquiring specific items
924 to forming general rules. *Current Directions in Psychological Science*, 21(3),
925 170–176. <https://doi.org/10.1177/0963721412436806>
- 926 Batterink, L. J., & Paller, K. A. (2017). Online neural monitoring of statistical learning.
927 *Cortex*, 90, 31–45. <https://doi.org/10.1016/j.cortex.2017.02.004>
- 928 Batterink, L. J., & Paller, K. A. (2019). Statistical learning of speech regularities can
929 occur outside the focus of attention. *Cortex*, 115, 56–71.
930 <https://doi.org/10.1016/j.cortex.2019.01.013>
- 931 Batterink, L. J., Paller, K. A., & Reber, P. J. (2019). Understanding the Neural Bases of
932 Implicit and Statistical Learning. *Topics in Cognitive Science*, 11(3), 482–503.
933 <https://doi.org/10.1111/tops.12420>
- 934 Batterink, L. J., Reber, P. J., Neville, H. J., & Paller, K. A. (2015). Implicit and explicit
935 contributions to statistical learning. *Journal of Memory and Language*, 83, 62–78.
936 <https://doi.org/10.1016/j.jml.2015.04.004>

- 937 Belin, P., Zatorre, R. J., Lafaille, P., Ahad, P., & Pike, B. (2000). Voice-selective areas in
938 human auditory cortex. *Nature*, *403*, 309-312. <https://doi.org/10.1038/35002078>
- 939 Berent, I., de la Cruz-Pavía, I., Brentari, D., & Gervain, J. (2021). Infants differentially
940 extract rules from language. *Scientific Reports*, *11*, Article 20001.
941 <https://doi.org/10.1038/s41598-021-99539-8>
- 942 Binder, J. R., Frost, J. A., Hammeke, T. A., Bellgowan, P. S. F., Springer, J. A.,
943 Kaufman, J. N., & Possing, E. T. (2000). Human Temporal Lobe Activation by
944 Speech and Nonspeech Sounds. *Cerebral Cortex*, *10*(5), 512–528.
945 <https://doi.org/10.1093/cercor/10.5.512>
- 946 Boersma, Paul & Weenink, David. (2022). *Praat: doing phonetics by computer* (Version
947 6.2.08). <http://www.praat.org/>
- 948 Boros, M., Magyari, L., Török, D., Bozsik, A., Deme, A., & Andics, A. (2021). Neural
949 processes underlying statistical learning for speech segmentation in dogs.
950 *Current Biology*, *31*(24), 5512-5521.e5. <https://doi.org/10.1016/j.cub.2021.10.017>
- 951 Bulf, H., Johnson, S. P., & Valenza, E. (2011). Visual statistical learning in the newborn
952 infant. *Cognition*, *121*(1), 127-132.
953 <https://doi.org/10.1016/j.cognition.2011.06.010>
- 954 Conway, C. M. (2020). How does the brain learn environmental structure? Ten core
955 principles for understanding the neurocognitive mechanisms of statistical
956 learning. *Neuroscience and Biobehavioral Reviews*, *112*, 279–299.
957 <https://doi.org/10.1016/j.neubiorev.2020.01.032>

- 958 Conway, C. M., & Christiansen, M. H. (2005). Modality-constrained statistical learning of
959 tactile, visual, and auditory sequences. *Journal of Experimental Psychology:*
960 *Learning, Memory, and Cognition*, 31(1), 24-39. [https://doi.org/10.1037/0278-](https://doi.org/10.1037/0278-7393.31.1.24)
961 7393.31.1.24
- 962 Corcoran, A. W., Perera, R., Koroma, M., Kouider, S., Hohwy, J., & Andrillon, T. (2023).
963 Expectations boost the reconstruction of auditory features from
964 electrophysiological responses to noisy speech. *Cerebral Cortex*, 33(3), 691–
965 708. <https://doi.org/10.1093/cercor/bhac094>
- 966 Covington, N. V., Brown-Schmidt, S., & Duff, M. C. (2018). The Necessity of the
967 Hippocampus for Statistical Learning. *Journal of Cognitive Neuroscience*, 30(5),
968 680–697. https://doi.org/10.1162/jocn_a_01228
- 969 Darwin, C. (2003). *SWS produced automatically using a script for the PRAAT program*.
970 University of Sussex School of Life Sciences.
971 http://www.lifesci.sussex.ac.uk/home/Chris_Darwin/SWS/
- 972 Davis, M. H., & Johnsrude, I. S. (2007). Hearing speech sounds: Top-down influences
973 on the interface between audition and speech perception. *Hearing Research*,
974 229, 132–147. <https://doi.org/10.1016/j.heares.2007.01.014>
- 975 Davis, M. H., Johnsrude, I. S., Hervais-Adelman, A., Taylor, K., & McGettigan, C.
976 (2005). Lexical Information Drives Perceptual Learning of Distorted Speech:
977 Evidence From the Comprehension of Noise-Vocoded Sentences. *Journal of*
978 *Experimental Psychology: General*, 134(2), 222–241.
979 <https://doi.org/10.1037/0096-3445.134.2.222>

- 980 Dawson, C., & Gerken, L. (2009). From Domain-Generality to Domain-Sensitivity: 4-
981 Month-Olds Learn an Abstract Repetition Rule in Music That 7-Month-Olds Do
982 Not. *Cognition*, 111(3), 378–382. <https://doi.org/10.1016/j.cognition.2009.02.010>
- 983 de Leeuw, J. R., Gilbert, R. A., & Luchterhandt, B. (2023). jsPsych: Enabling an Open-
984 Source Collaborative Ecosystem of Behavioral Experiments. *Journal of Open
985 Source Software*, 8(85), 5351. <https://doi.org/10.21105/joss.05351>
- 986 Dick, F., Saygin, A. P., Galati, G., Pitzalis, S., Bentrovato, S., D'Amico, S., Wilson, S.,
987 Bates, E., & Pizzamiglio, L. (2007). What is Involved and What is Necessary for
988 Complex Linguistic and Nonlinguistic Auditory Processing: Evidence from
989 Functional Magnetic Resonance Imaging and Lesion Data. *Journal of Cognitive
990 Neuroscience*, 19(5), 799–816. <https://doi.org/10.1162/jocn.2007.19.5.799>
- 991 Elazar, A., Alhama, R. G., Bogaerts, L., Siegelman, N., Baus, C., & Frost, R. (2022).
992 When the “Tabula” is Anything but “Rasa:” What Determines Performance in the
993 Auditory Statistical Learning Task? *Cognitive Science*, 46(2), e13102.
994 <https://doi.org/10.1111/cogs.13102>
- 995 Emberson, L. L., Conway, C. M., & Christiansen, M. H. (2011). Timing is everything:
996 Changes in presentation rate have opposite effects on auditory and visual implicit
997 statistical learning. *Quarterly Journal of Experimental Psychology*, 64(5), 1021–
998 1040. <https://doi.org/10.1080/17470218.2010.538972>
- 999 Endress, A., & Bonatti, L. (2007). Rapid learning of syllable classes from a perceptually
1000 continuous speech stream. *Cognition*, 105, 247–299.
1001 <https://doi.org/10.1016/j.cognition.2006.09.010>

- 1002 Endress, A. D., & Mehler, J. (2009). Primitive computations in speech processing.
1003 *Quarterly Journal of Experimental Psychology*, 62(11), 2187–2209.
1004 <https://doi.org/10.1080/17470210902783646>
- 1005 Ferguson, B., & Lew-Williams, C. (2016). Communicative signals support abstract rule
1006 learning by 7-month-old infants. *Scientific Reports*, 6(1), 25434.
1007 <https://doi.org/10.1038/srep25434>
- 1008 Finn, A. S., & Hudson Kam, C. L. (2008). The curse of knowledge: First language
1009 knowledge impairs adult learners' use of novel statistics for word segmentation.
1010 *Cognition*, 108(2), 477–499. <https://doi.org/10.1016/j.cognition.2008.04.002>
- 1011 Fiser, J., & Aslin, R. N. (2001). Unsupervised Statistical Learning of Higher-Order
1012 Spatial Structures from Visual Scenes. *Psychological Science*, 12(6), 499–504.
1013 <https://doi.org/10.1111/1467-9280.00392>
- 1014 Forest, T. A., Schlichting, M. L., Duncan, K. D., & Finn, A. S. (2023). Changes in
1015 statistical learning across development. *Nature Reviews Psychology*, 2(4), Article
1016 4. <https://doi.org/10.1038/s44159-023-00157-0>
- 1017 Frost, R., Armstrong, B. C., Siegelman, N., & Christiansen, M. H. (2015). Domain
1018 generality versus modality specificity: The paradox of statistical learning. *Trends*
1019 *in Cognitive Sciences*, 19(3), 117-125. <https://doi.org/10.1016/j.tics.2014.12.010>
- 1020 Frost, R., Armstrong, B. C., & Christiansen, M. H. (2019). Statistical learning research:
1021 A critical review and possible new directions. *Psychological Bulletin*, 145(12),
1022 1128–1153. <https://doi.org/10.1037/bul0000210>

- 1023 Gebhart, A. L., Newport, E. L., & Aslin, R. N. (2009). Statistical learning of adjacent and
1024 nonadjacent dependencies among nonlinguistic sounds. *Psychonomic Bulletin &*
1025 *Review*, 16(3), 486–490. <https://doi.org/10.3758/PBR.16.3.486>
- 1026 Hauser, M. D., Newport, E. L., & Aslin, R. N. (2001). Segmentation of the speech
1027 stream in a non-human primate: Statistical learning in cotton-top tamarins.
1028 *Cognition*, 78(3), B53–B64. [https://doi.org/10.1016/S0010-0277\(00\)00132-3](https://doi.org/10.1016/S0010-0277(00)00132-3)
- 1029 Hoch, L., Tyler, M. D., & Tillmann, B. (2013). Regularity of unit length boosts statistical
1030 learning in verbal and nonverbal artificial languages. *Psychonomic Bulletin &*
1031 *Review*, 20(1), 142–147. <https://doi.org/10.3758/s13423-012-0309-8>
- 1032 IEEE. (1969). IEEE Recommended Practice for Speech Quality measurements. *Institute*
1033 *of Electronic Engineers*, New York.
1034 <https://doi.org/10.1109/IEEESTD.1969.7405210>
- 1035 Khoshkhoo, S., Leonard, M. K., Mesgarani, N., & Chang, E. F. (2018). Neural correlates
1036 of sine-wave speech intelligibility in human frontal and temporal cortex. *Brain and*
1037 *Language*, 187, 83–91. <https://doi.org/10.1016/j.bandl.2018.01.007>
- 1038 Kirkham, N. Z., Slemmer, J. A., & Johnson, S. P. (2002). Visual statistical learning in
1039 infancy: Evidence for a domain general learning mechanism. *Cognition*, 83(2),
1040 B35–B42. [https://doi.org/10.1016/S0010-0277\(02\)00004-5](https://doi.org/10.1016/S0010-0277(02)00004-5)
- 1041 Klemfuss, N., Prinzmetal, B., & Ivry, R. (2012). How Does Language Change
1042 Perception: A Cautionary Note. *Frontiers in Psychology*, 3.
1043 <https://www.frontiersin.org/articles/10.3389/fpsyg.2012.00078>

- 1044 Liberman, A. M. (1982). On finding that speech is special. *American Psychologist*,
1045 37(2), 148–167. <https://doi.org/10.1037/0003-066X.37.2.148>
- 1046 Litman, L., Robinson, J., & Abberbock, T. (2017). TurkPrime.com: A versatile
1047 crowdsourcing data acquisition platform for the behavioral sciences. *Behavior*
1048 *Research Methods*, 49(2), 433–442. <https://doi.org/10.3758/s13428-016-0727-z>
- 1049 Lupyan, G., & Spivey, M. J. (2008). Perceptual processing is facilitated by ascribing
1050 meaning to novel stimuli. *Current Biology*, 18(10), R410–R412.
1051 <https://doi.org/10.1016/j.cub.2008.02.073>
- 1052 Marcus, G. F., Fernandes, K. J., & Johnson, S. P. (2007). Infant Rule Learning
1053 Facilitated by Speech. *Psychological Science*, 17(5), 387-391.
1054 <https://doi.org/10.1111/j.1467-9280.2007.01910.x>
- 1055 Marcus, G. F., & Rabagliati, H. (2008). In J. Colombo, P. McCardle & L. Freund (Eds.),
1056 *Infant Pathways to Language: Methods, Models and Research Directions*.
1057 Lawrence Erlbaum Associates.
- 1058 Mitchel, A. D., Christiansen, M. H., & Weiss, D. J. (2014). Multimodal integration in
1059 statistical learning: Evidence from the McGurk illusion. *Frontiers in Psychology*,
1060 5. <https://doi.org/10.3389/fpsyg.2014.00407>
- 1061 Marno, H., Farroni, T., Vidal Dos Santos, Y., Ekramnia, M., Nespors, M., & Mehler, J.
1062 (2015). Can you see what I am talking about? Human speech triggers referential
1063 expectation in four-month-old infants. *Scientific Reports*, 5.
1064 <https://doi.org/10.1038/srep13594>

- 1065 Moore, D. R. (2000). Auditory neuroscience: Is speech special? *Current Biology*, *10*(10),
1066 R362–R364. [https://doi.org/10.1016/S0960-9822\(00\)00479-6](https://doi.org/10.1016/S0960-9822(00)00479-6)
- 1067 Moreau, C. N., Joanisse, M. F., Mulgrew, J., & Batterink, L. J. (2022). No statistical
1068 learning advantage in children over adults: Evidence from behaviour and neural
1069 entrainment. *Developmental Cognitive Neuroscience*, *57*, 101154.
1070 <https://doi.org/10.1016/j.dcn.2022.101154>
- 1071 Narain, C., Scott, S. K., Wise, R. J. S., Rosen, S., Leff, A., Iversen, S. D., & Matthews,
1072 P. M. (2003). Defining a Left-lateralized Response Specific to Intelligible Speech
1073 Using fMRI. *Cerebral Cortex*, *13*(12), 1362–1368.
1074 <https://doi.org/10.1093/cercor/bhg083>
- 1075 Ogg, M., & Slevc, L. R. (2019). Acoustic Correlates of Auditory Object and Event
1076 Perception: Speakers, Musical Timbres, and Environmental Sounds. *Frontiers in*
1077 *Psychology*, *10*. <https://www.frontiersin.org/articles/10.3389/fpsyg.2019.01594>
- 1078 Ordin, M., Polyanskaya, L., & Samuel, A. (2021). An evolutionary account of
1079 intermodality differences in statistical learning. *Annals of New York Academy of*
1080 *Science*, *1486*(1), 76–89. <https://doi.org/10.1111/nyas.14502>
- 1081 Palan, S., & Schitter, C. (2018). Prolific.ac—A subject pool for online experiments.
1082 *Journal of Behavioral and Experimental Finance*, *17*, 22–27.
1083 <https://doi.org/10.1016/j.jbef.2017.12.004>
- 1084 Parviainen, T., Helenius, P., & Salmelin, R. (2005). Cortical differentiation of speech and
1085 nonspeech sounds at 100 ms: Implications for dyslexia. *Cerebral Cortex*, *15*(7),
1086 1054–1063. <https://doi.org/10.1093/cercor/bhh206>

- 1087 Peña, M., Bonatti, L. L., Nespor, M., & Mehler, J. (2002). Signal-Driven Computations in
1088 Speech Processing. *Science*, 298(5593), 604–607.
1089 <https://doi.org/10.1126/science.1072901>
- 1090 Rabagliati, H., Ferguson, B., & Lew-Williams, C. (2018). The profile of abstract rule
1091 learning in infancy: Meta-analytic and experimental evidence. *Developmental*
1092 *Science*, 22(1). <https://doi.org/10.1111/desc.12704>
- 1093 Rabagliati, H., Senghas, A., Johnson, S., & Marcus, G. F. (2012). Infant Rule Learning:
1094 Advantage Language, or Advantage Speech? *PLOS ONE*, 7(7), e40517.
1095 <https://doi.org/10.1371/journal.pone.0040517>
- 1096 Raviv, L., & Arnon, I. (2018). The developmental trajectory of children’s auditory and
1097 visual statistical learning abilities: Modality-based differences in the effect of age.
1098 *Developmental Science*, 21(4). <https://doi.org/10.1111/desc.12593>
- 1099 Remez, R. E., Rubin, P. E., Pisoni, D. B., & Carrell, T. D. (1981). Speech Perception
1100 without Traditional Speech Cues. *Science*, 212(4497), 947–950.
1101 <https://doi.org/10.1126/science.7233191>
- 1102 Saffran, J. R. (2002). Constraints on Statistical Language Learning. *Journal of Memory*
1103 *and Language*, 47(1), 172-196. <https://doi.org/10.1006/jmla.2001.2839>
- 1104 Saffran, J. R., Aslin, R. N., & Newport, E. L. (1996a). Statistical learning by 8-month-old
1105 infants. *Science*, 274(5294), 1926-1928.

- 1106 Saffran, J. R., Newport, E. L., & Aslin, R. N. (1996b). Word Segmentation: The Role of
1107 Distributional Cues. *Journal of Memory and Language*, 35(4), 606–621.
1108 <https://doi.org/10.1006/jmla.1996.0032>
- 1109 Saffran, J. R., Newport, E. L., Aslin, R. N., Tunick, R. A., & Barrueco, S. (1997).
1110 Incidental Language Learning: Listening (and Learning) Out of the Corner of
1111 Your Ear. *Psychological Science*, 8(2), 101–105. [https://doi.org/10.1111/j.1467-](https://doi.org/10.1111/j.1467-9280.1997.tb00690.x)
1112 [9280.1997.tb00690.x](https://doi.org/10.1111/j.1467-9280.1997.tb00690.x)
- 1113 Saffran, J. R., Johnson, E. K., Aslin, R. N., & Newport, E. L. (1999). Statistical learning
1114 of tone sequences by human infants and adults. *Cognition*, 70(1), 27–52.
1115 [https://doi.org/10.1016/S0010-0277\(98\)00075-4](https://doi.org/10.1016/S0010-0277(98)00075-4)
- 1116 Saffran, J. R., Pollak, S. D., Seibel, R. L., & Shkolnik, A. (2007). Dog is a dog is a dog:
1117 Infant rule learning is not specific to language. *Cognition*, 105(3), 669–680.
1118 <https://doi.org/10.1016/j.cognition.2006.11.004>
- 1119 Schapiro, A. C., Gregory, E., Landau, B., McCloskey, M., & Turk-Browne, N. B. (2014).
1120 The Necessity of the Medial Temporal Lobe for Statistical Learning. *Journal of*
1121 *Cognitive Neuroscience*, 26(8), 1736–1747.
1122 https://doi.org/10.1162/jocn_a_00578
- 1123 Schlichting, M. L., Guarino, K. F., Schapiro, A. C., Turk-Browne, N. B., & Preston, A. R.
1124 (2017). Hippocampal Structure Predicts Statistical Learning and Associative
1125 Inference Abilities during Development. *Journal of Cognitive Neuroscience*,
1126 29(1), 37–51. https://doi.org/10.1162/jocn_a_01028

- 1127 Schmalz, X., Biurrun Manresa, J., & Zhang, L. (2023). What is a Bayes factor?
1128 *Psychological Methods*, 28(3), 705-718. <https://doi.org/10.1037/met0000421>
- 1129 Scott, S. K., Blank, C. C., Rosen, S., & Wise, R. J. S. (2000). Identification of a pathway
1130 for intelligible speech in the left temporal lobe. *Brain*, 123(12), 2400–2406.
1131 <https://doi.org/10.1093/brain/123.12.2400>
- 1132 Seitz, A. R., Kim, R., Van Wassenhove, V., & Shams, L. (2007). Simultaneous and
1133 Independent Acquisition of Multisensory and Unisensory Associations.
1134 *Perception*, 36(10), 1445–1453. <https://doi.org/10.1068/p5843>
- 1135 Shufaniya, A., & Arnon, I. (2018). Statistical Learning Is Not Age-Invariant During
1136 Childhood: Performance Improves With Age Across Modality. *Cognitive Science*,
1137 42(8), 3100–3115. <https://doi.org/10.1111/cogs.12692>
- 1138 Shultz, S., & Vouloumanos, A. (2010). Three-Month-Olds Prefer Speech to Other
1139 Naturally Occurring Signals. *Language Learning and Development*, 6(4), 241–
1140 257. <https://doi.org/10.1080/15475440903507830>
- 1141 Siegelman, N., Bogaerts, L., Christiansen, M. H., & Frost, R. (2017). Towards a theory
1142 of individual differences in statistical learning. *Philosophical Transactions of the*
1143 *Royal Society B: Biological Sciences*, 372(1711).
1144 <https://doi.org/10.1098/rstb.2016.0059>
- 1145 Siegelman, N., Bogaerts, L., Elazar, A., Arciuli, J., & Frost, R. (2018). Linguistic
1146 entrenchment: Prior knowledge impacts statistical learning performance.
1147 *Cognition*, 177, 198-213. <https://doi.org/10.1016/j.cognition.2018.04.011>

- 1148 Siegelman, N., & Frost, R. (2015). Statistical learning as an individual ability: Theoretical
1149 perspectives and empirical evidence. *Journal of Memory and Language*, *81*,
1150 105–120. <https://doi.org/10.1016/j.jml.2015.02.001>
- 1151 Silva, D. M. R., & Bellini-Leite, S. C. (2020). Cross-modal correspondences in sine
1152 wave: Speech versus non-speech modes. *Attention, Perception, &*
1153 *Psychophysics*, *82*(3), 944–953. <https://doi.org/10.3758/s13414-019-01835-z>
- 1154 Thiessen, E. D., (2012). Effects of Inter- and Intra-modal Redundancy on Infants' Rule
1155 Learning. *Language Learning and Development*, *8*(3), 197-214.
1156 <https://doi.org/10.1080/15475441.2011.583610>
- 1157 Thiessen, E. D. (2017). What's statistical about learning? Insights from modelling
1158 statistical learning as a set of memory processes. *Philosophical Transactions of*
1159 *the Royal Society B: Biological Sciences*, *372*(1711), 20160056.
1160 <https://doi.org/10.1098/rstb.2016.0056>
- 1161 Van Hedger, S. C., Johnsrude, I. S., & Batterink, L. J. (2022). Musical instrument
1162 familiarity affects statistical learning of tone sequences. *Cognition*, *218*, 104949.
1163 <https://doi.org/10.1016/j.cognition.2021.104949>
- 1164 Vouloumanos, A., Hauser, M. D., Werker, J. F., & Martin, A. (2010). The Tuning of
1165 Human Neonates' Preference for Speech. *Child Development*, *81*(2), 517–527.
1166 <https://doi.org/10.1111/j.1467-8624.2009.01412.x>
- 1167 Vouloumanos, A., Kiehl, K. A., Werker, J. F., & Liddle, P. F. (2001). Detection of Sounds
1168 in the Auditory Stream: Event-Related fMRI Evidence for Differential Activation to

- 1169 Speech and Nonspeech. *Journal of Cognitive Neuroscience*, 13(7), 994–1005.
1170 <https://doi.org/10.1162/089892901753165890>
- 1171 Vouloumanos, A., & Werker, J. F. (2004). Tuned to the signal: The privileged status of
1172 speech for young infants. *Developmental Science*, 7(3), 270-276.
1173 <https://doi.org/10.1111/j.1467-7687.2004.00345.x>
- 1174 Vouloumanos, A., & Werker, J. F. (2007). Listening to language at birth: Evidence for a
1175 bias for speech in neonates. *Developmental Science*, 10(2), 159–164.
1176 <https://doi.org/10.1111/j.1467-7687.2007.00549.x>
- 1177 Werker, J. F., & Hensch, T. K. (2015). Critical Periods in Speech Perception: New
1178 Directions. *Annual Review of Psychology*, 66(1), 173–196.
1179 <https://doi.org/10.1146/annurev-psych-010814-015104>