1        Unravelling the time-course of listener adaptation to an unfamiliar talker

2                        Maryann Tan[1,2] & T. Florian Jaeger[2,3]

3                [1] Centre for Research on Bilingualism, University of Stockholm

4                    [2] Brain and Cognitive Sciences, University of Rochester

5                    [3] Computer Science, University of Rochester

6                                        Author Note

8        Correspondence concerning this article should be addressed to Maryann Tan, Department

9   of Bilingualism, Stockholm University, Sweden. E-mail: maryann.tan@biling.su.se

Abstract

We investigate constraints on adaptive speech perception during the initial encounters with unfamiliar speech patterns. Such adaptive changes are now considered important to spoken language understanding, overcoming substantial cross-talker variability in the realization of speech categories. We present evidence from a novel incremental exposure-test paradigm to assess how previously experienced cross-talker variability guides (and thus constrains) listeners' adaptation. Specifically, we ask adaptation is constrained weakly—slower and sublinear, but continued adaptation with increasing exposure—or strongly—adaptation only up to a point, after which additional exposure has no benefits (at least not prior to, e.g., sleep). The results contribute to a proposed theoretical distinction between two hypotheses about the mechanisms underlying the intial moments of adaptation, model learning vs. model selection.

*Keywords:* speech perception; adaptation; incremental changes; distributional learning

Word count: X

Unravelling the time-course of listener adaptation to an unfamiliar talker

# 1  TO-DO

## 1.1  Highest priority

- MARYANN
- Please read this carefully.
- TIME TO STOP MESSY CODING. Let's have a zero-tolerance policy for that from now on in the main working branch (i.e., you can do what you'd like in branches that aren't the main branch, but you canNOT merge without cleaning up first). It is a real time-sink for everyone else and makes it near impossible for me to effectively help.

  - on the main working branch, functions should be in functions.R, in a clearly named section (see existing examples).

- Input data file:

  - There shouldn't be multiple data files that you're loading. E.g., I don't understand why there is an exposure trials data file in addition to the main data file. It's just confusing. Let's not do things like that.
  - Rename main data file to "experiment-results.csv"
  - Have a script in your other repo (for your thesis) that does all the data importing, variable and value formatting, etc. The input data file experiment-results.csv should already contain all the information you (and others might need) and be in the format that you'd like it to be. That's the only data file that will be in your paper repo.
    * Think carefully about how to name variables consistently and create all variants of variables you might need in the paper, e.g., Response, Item.ExpectedResponse, Response.Category, Item.ExpectedResponse.Category, Response.Voiced, Item.ExpectedResponse.Voiced (etc. if you indeed need all of those; we definitely need the first two pairs of these).
    * Also if you have to consistently rename levels for plotting, please just changed them once in the script that creates the file. E.g., there's various places in which

you deal with formatting the conditions and various names floating around (Shift0, 10, etc.; +0, +10, etc.; baseline, + 10 etc.). Pick one, do it at the top of the pipeline (i.e., in the input script). This will reduce the potential for error in your own coding, make your code in the main paper shorter, and it'll be much easier to read for others trying to follow your code (including me).

* Remove all data formatting code from the paper Rmd. There should only be a single load line.

* I've moved the code loading the chodroff data into the new pre-amble.R file. Consider doing the same for the experiment data. That way the data that we need throughout are available throughout.

- Clean up functions.R file:

  - PLEASE DO GET RID OF UNUSED FUNCTIONS. Search files for each function (cmd + shift + f). If it does not exist, remove it from functions.R
  - Use clearer function names. It often happens as a project develops that functions become ambiguous in their name. E.g., you have several functions that do similar things (like getting or plotting CIs from psychometric or IO models). Extend their names to be clear: e.g., compare get_CI to get_CI_from_ideal_observer; or make_CI to print_CI; or add_PSE_perception_median to add_PSE_median_to_plot (note how I also removed redundancy since PSEs are always about perception); etc. Rename the functions and use CMD + SHIFT + F to search and replace all mentions of those functions across all files.
  - Organize functions into sections with headings in functions.R

- Try to set local constants at top of chunk. e.g., Don't have stuff like empirical_means <- c(17, 62) in the middle of a chunk.

- It's best not to save unnecessary objects but if you do, remove them after they are no longer needed (e.g., the various excl.headphone, etc. in section 2: you could just have that code inline without ever storing them. But it's ok to do things the way you do. Just remove them after they have done their job).

## 1.2 Medium priority

- MARYANN

- FLORIAN

- think about table 1 and 2: how to change the wording on tables to actually refer to intercepts rather than PSEs or change the figures? Changing current representations of analyses to improve intuitive-ity.

- write overview of results

- restructure results presentation.

- write SI sections with proofs

### 1.2.1 Lower Priority

- MARYANN

- Combine data from exposure and test, use all together instead of coding block, code trial and code it as a smooth. That means using GAMM – that may require taking lapse (try it first without lapses because the GAMM takes care of the lapse. The RE will be expressed differently. It has to follow the GAMM syntax.) The primary thing we want to smooth over is "block", but could theoretically smooth over VOT and Block.

- Florian

- compare IBBU predictions over blocks with human behavioural data

## 1.3 To do later

- Everyone: Eat ice-cream and perhaps have a beer.

# 1   Introduction

Adaptivity is a hallmark of human speech perception, supporting faster and more accurate speech recognition. When exposed to an unfamiliar accent, the processing difficulty listeners might initially experience tends to alleviate with exposure (Bradlow, Bassard, & Paller, 2023; e.g., Bradlow & Bent, 2008; Clarke & Garrett, 2004; Sidaras, Alexander, & Nygaard, 2009; Xie et al., 2018; Xie, Liu, & Jaeger, 2021; for review, see **baeseberk2018?**). Research over the last few decades has made strides in identifying the conditions required for successful adaptation, its generalizability across talkers, and its longevity (for reviews, see Bent & Baese-Berk, 2021; Cummings & Theodore, 2023; Zheng & Samuel, 2020). It is now clear that listeners' categorization function—the mapping from acoustic or phonetic inputs to linguistic categories and, ultimately, word meanings—changes based on the phonetic properties of recent input (e.g., Bertelson, Vroomen, & De Gelder, 2003; Clayards, Tanenhaus, Aslin, & Jacobs, 2008; Eisner & McQueen, 2005; Idemaru & Holt, 2011; Kraljic & Samuel, 2005; McMurray & Jongman, 2011; Norris, McQueen, & Cutler, 2003; Reinisch & Holt, 2014; **cole2011?**; **kurumada2013?**; **xie2018jep?**; for review, Schertz & Clare, 2020; Xie, Jaeger, & Kurumada, 2023). This has led to the development of stronger theories and models of adaptive speech perception that explicitly link the distribution of phonetic properties in recent speech input to changes in subsequent speech recognition (e.g., Apfelbaum & McMurray, 2015; Assmann & Nearey, 2007; Harmon, Idemaru, & Kapatsinski, 2019; Johnson, 1997; Kleinschmidt & Jaeger, 2015; Lancia & Winter, 2013; Magnuson et al., 2020; Sohoglu & Davis, 2016; Xie et al., 2023).

As Cummings and Theodore (2023) point out, previous work has typically framed questions as an 'either-or'—adaptation is either observed or not—consistent with the focus on identifying the necessary conditions for adaptation and generalization. Recent reviews of the field instead emphasize the need to move towards stronger tests of existing theories, requiring the development of paradigms that support quantitative comparison and yield data that more strongly constrain the space of theoretical possibilities (Schertz & Clare, 2020; Xie et al., 2023; **baeseberk2018?**). This includes the need for data that characterize how adaptation develops *incrementally* as a function of both the *amount of exposure* and the *distribution of phonetic cues in the exposure input.* While existing theories differ in important aspects, they share critical predictions about

incremental adaptation that have remained largely untested: listeners' categorizations are predicted to change incrementally with exposure, and the direction and magnitude of that change should gradiently depend on (1) listeners' prior expectations based on previously experienced speech input from other talkers, and both (2a) the amount and (2b) distribution of phonetic evidence in the exposure input from the unfamiliar talker (for review, see Xie et al., 2023). We report initial results from a novel repeated exposure-test paradigm designed to test these predictions during the early moments of adaptation.

Figure 1 illustrates our approach. The experiment builds on computational and behavioral findings from separate lines of research on unsupervised distributional learning during speech perception (DL, Clayards et al., 2008; Kleinschmidt, 2020; Theodore & Monto, 2019), lexically- or visually-guided perceptual learning (LGPL, Cummings & Theodore, 2023; VGPL, Kleinschmidt & Jaeger, 2012; Vroomen, Linden, De Gelder, & Bertelson, 2007), and accent adaptation (Hitczenko & Feldman, 2016; Tan, Xie, & Jaeger, 2021). These paradigms and findings have complementing strengths that we seek to combine and extend. Following previous work on distributional learning in speech perception, we expose different groups of listeners to phonetic distributions that are shifted to different degrees (Bejjanki, Beck, Lu, & Pouget, 2011; Clayards et al., 2008; Kleinschmidt, Raizada, & Jaeger, 2015; Munson, 2011; Nixon, Rij, Mok, Baayen, & Chen, 2016; Theodore & Monto, 2019). Unlike this work, we incrementally assess changes in listeners' categorization from pre-exposure onward.
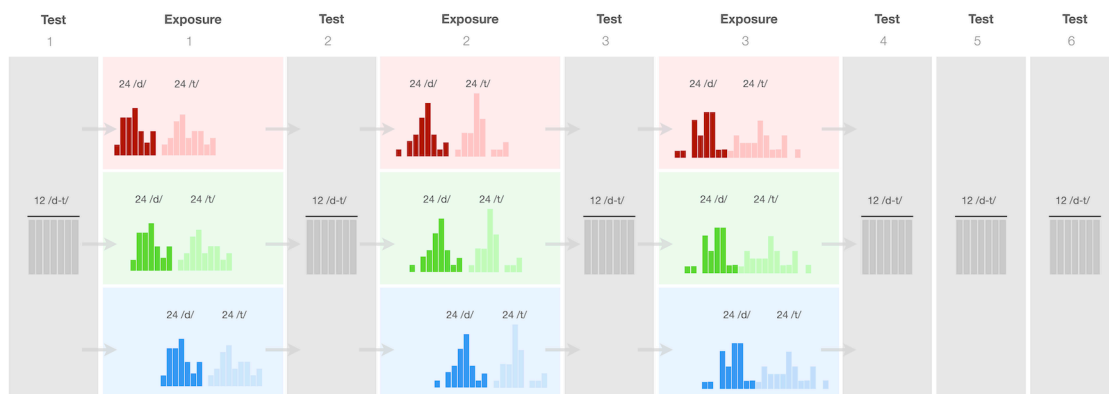


*Figure 1.* Exposure-test design of the experiment. Exposure conditions (rows) differed in the distribution of voice onset time (VOT), the primary phonetic cue to word-initial /d/ and /t/ in English (e.g., "dip" vs. "tip"). Test blocks presented identical VOT stimuli within and across conditions.

146     Following previous DL studies, we use of phonetically manipulated stimuli. This gives

147 researchers control over the distribution of acoustic-phonetic properties that listeners experience

148 during exposure and test (unlike AA, LGPL, and VGPL paradigms). This control is an

149 important prerequisite for stronger tests of predictions (1) and (2a,b). For example, recent

150 findings from LGPL and VGPL provide evidence in support of prediction (2a)—that the amount

151 of phonetic evidence during exposure gradiently affects the magnitude of subsequent changes in

152 listeners' categorization response (Cummings & Theodore, 2023; see also Liu & Jaeger, 2018,

153 2019). This includes some initial evidence that these changes accumulate incrementally

154 (Kleinschmidt & Jaeger, 2012; Vroomen et al., 2007), in ways consistent with models of adaptive

155 speech perception. LGPL and VGPL paradigms—at least as used traditionally—do, however,

156 limit experimenters' control over the phonetic properties of the exposure stimuli: shifted sound

157 instances are selected to be perceptually ambiguous (e.g., between "s" and "sh"), not to exhibit

158 specific phonetic distributions. To the extent that LGPL and VGPL research has assessed the

159 effects of phonetic properties on the degree of boundary shift following exposure, this has been

160 limited to qualitative post-hoc analyses (Drouin, Theodore, & Myers, 2016; Kraljic & Samuel,

161 2007; Tzeng, Nygaard, & Theodore, 2021?). This makes it difficult to test predictions (1) and

162 (2b) about the effects of phonetic distributions in prior and recent experience.

163     Support for prediction (2b) has thus primarily come from research in DL paradigms. In an

164 important early study, Clayards et al. (2008) exposed two different groups of US English listeners

165 to instances of "b" and "p" that differed in their distribution along the voice onset time

166 continuum (VOT). VOT is the primary phonetic cue to word-initial /b/-/p/, /d/-/t/, /g/-/k/ in

167 US English: the voiced category (e.g. /b/) is produced with lower VOT than the voiceless

168 category (e.g., /p/). Clayards and colleagues held the VOT means of /b/ and /p/ constant

169 between the two exposure groups, but manipulated whether both /b/ and /p/ had wide or

170 narrow variance along VOT. Exposure was unlabeled: on any trial, listeners saw pictures of, e.g.,

171 bees and peas on the screen while hearing a synthesized recording along the "bees"-"peas"

172 continuum (obtained by manipulating VOT). Listeners' task was to click on the picture

173 corresponding to the word they heard. If listeners adapt by learning the VOT distributions of /b/

174 and /p/, listeners in the wide variance group were predicted to exhibit a more shallow

175  categorization function than the narrow variance group. This is precisely what Clayards and

176  colleagues found (see also Nixon et al., 2016; Theodore & Monto, 2019). Together with more

177  recent findings from adaptation to natural accents (Hitczenko & Feldman, 2016; Tan et al., 2021;

178  Xie, Buxó-Lugo, & Kurumada, 2021), this important finding suggests that the *outcome* of

179  adaptation qualitatively follows the predictions of distributional learning models (e.g., exemplar

180  theory, Johnson, 1997; ideal adaptors, Kleinschmidt & Jaeger, 2015). However, the findings in

181  this line of work relied on tests that either averaged over, or followed, hundreds of trials of

182  exposure. This leaves open how adaptation proceeds from the earliest moments of exposure—i.e.,

183  whether listeners categorization behavior indeed changes in the way predicted by models of

184  adaptive speech perception, developing from expectations based on previously experienced

185  phonetic distributions to increasing integration of the phonetic distributions observed during

186  exposure to the unfamiliar talker. It also leaves open whether potential constraints on the extent

187  to which listeners' behavior changes with exposure (for initial evidence and discussion, see

188  Cummings & Theodore, 2023; Kleinschmidt, 2020; Kleinschmidt & Jaeger, 2016) reflect hard

189  limits on adaptivity or simply the incremental learning outcome—'how far the learner has

190  gotten'—at the only point at which adaptation is assessed (i.e., following exposure).

191      The repeated exposure-test paradigm in Figure 1 aims to address these knowledge gaps.

192  The experiment starts with a test block that assesses listeners' state prior to informative

193  exposure—often assumed, but not tested, to be identical across exposure conditions. Additional

194  intermittent tests—opaque to participants—then assess incremental changes up to the first 144

195  informative exposure trials. By employing physically identical test trials both across block within

196  exposure conditions and across exposure conditions, we aim to facilitate assumption-free

197  comparison of cumulative exposure effects (we additionally also measure adaptation during

198  exposure). As we detail under Methods, the use of repeated testing deviates from previous work

199  (Clayards et al., 2008; Harmon et al., 2019; Idemaru & Holt, 2011, 2020; Kleinschmidt, 2020;

200  Kleinschmidt & Jaeger, 2016; Munson, 2011; Nixon et al., 2016; Theodore & Monto, 2019), and is

201  not without challenges.

202      Finally, we took several modest steps towards addressing concerns about ecological validity

203  that have been argued to limit the generalizability of DL results. This includes concerns about

the ecological validity of both the stimuli and their distributions in the experiment (see discussion in **baseberk2018?**). For example, previous distributional learning studies have used highly unnatural, 'robotic'-sounding, speech (but see Theodore & Monto, 2019). Beyond raising questions about what types of expectations listeners apply to such speech, these stimuli also failed to exhibit naturally occurring covariation between phonetic cues that listeners are known to expect (see, e.g., Idemaru & Holt, 2011; Schertz, Cho, Lotto, & Warner, 2016). We instead developed stimuli that both sound natural and exhibit the type of phonetic covariation that listeners expect from everyday speech perception. We return to these and additional steps we took to increase the ecological validity of the phonetic *distributions* under Methods.

All data and code for this article can be downloaded from https://osf.io/hxcy4/. The article is written in R markdown, allowing readers to replicate our analyses with the press of a button using freely available software (R, R Core Team, 2022; RStudio Team, 2020), while changing any of the parameters of our models (see SI, **??**).

## 2   Methods

### 2.1   Participants

We recruited 126 participants from the Prolific crowdsourcing platform. We used Prolific's pre-screening to limit the experiment to participants (1) of US nationality, (2) who reported to be English speaking monolinguals, and (3) had not previously participated in any experiment from our lab on Prolific. Prior to the start of the experiment, participants had to confirm that they (4) had spent the first 10 years of their life in the US, (5) were in a quiet place and free from distractions, and (6) wore in-ear or over-the-ears headphones that cost at least $15. An additional 115 participants loaded the experiment but did not start or complete it.[1]

Participants took an average of 31.6 minutes to complete the experiment (SD = 20 minutes) and were remunerated $8.00/hour. An optional post-experiment survey recorded

---

[1] Unlike in lab-based experiments, for which participants' right to stop the experiment at any point is costly (both in terms of physical effort and perceived social cost), exercising this right in web-based experiments is essentially cost free—in particular, if exercised early in the experiment.

participant demographics using NIH prescribed categories, including participant sex (59 = female, 60 = male, 3 = NA), age (mean = NA years; 95% quantiles = 20-62.1 years), race (6 = Black, 31 = White, 85 = NA), and ethnicity (6 = Hispanic, 113 = Non-Hispanic, 3 = NA).

Participants' responses were collected via Javascript developed by the Human Language Processing Lab at the University of Rochester (**JSEXP?**) and stored via Proliferate developed at, and hosted by, the ALPs lab at Stanford University (Schuster, S, 2020).

## 2.2   Materials

We recorded 8 tokens each of four minimal word pairs (*dill/till*, *dim/tim*, *din/tin*, and *dip/tip*) from a 23-year-old, female L1-US English talker from New Hampshire, judged to have a "general American" accent. In addition to these critical minimal pairs we also recorded three words that did not did not contain any stop consonant sounds ("flare", "share", and "rare"). These word recordings were used for catch trials. Stimulus intensity was normalized to 70 dB sound pressure level for all recordings.

The critical minimal pair recordings were used to create four VOT continua using a script (Winn, 2020) in Praat (Boersma & Weenink, 2022). This approach resulted in continuum steps that sound natural (unlike the highly robotic-sounding stimuli employed in Clayards et al., 2008; Kleinschmidt & Jaeger, 2016). A post-experiment survey asked participants: "*Did you notice anything in particular about how the speaker pronounced the different words (e.g. till, dill, etc.)?*" No participant reported that the stimuli sounded unnatural. The procedure also maintained the natural correlations between the most important cues to word-initial stop-voicing in L1-US English (VOT, F0, and vowel duration). Specifically, the F0 at vowel onset of each stimulus was set to respect the linear relation with VOT observed in the original recordings of the talker. The duration of the vowel was set to follow the natural trade-off relation with VOT (Allen & Miller, 1999). Further details on the recording and resynthesis procedure are provided in the supplementary information (SI, **??**). We note that our effort to use more human-sounding stimuli had a substantial effect on participant attention especially in the beginning of the experiment. As we will elaborate in the results section, participant lapse rates (the portion of trials where participants were not paying attention) were found to be very low even from the earliest moments

in the experiment.

The VOTs generated for each continuum ranged from -100 to +130 ms in 5 ms steps.[2] A
norming experiment (N = 24 participants) reported in the SI (**??**) was used to select the three
minimal pair continua that elicited the most similar categorization responses (*dill-till*, *din-tin*, and
*dip-tip*). These three continua were used to create the exposure conditions shown in Figure 1.

## 2.3   Procedure

At the start of the experiment, participants acknowledged that they met all requirements and
provided consent, as per the Research Subjects Review Board of the University of Rochester.
Participants also had to pass a headphone test (Woods, Siegel, Traer, & McDermott, 2017), and
were instructed to not change the volume throughout the experiment. Following instructions,
participants completed 234 two-alternative forced-choice categorization trials (Figure 2).
Participants were instructed that they would hear a female talker say a single word on each trial,
and were asked to select which word they heard. Participants were asked to listen carefully and
answer as quickly and as accurately as possible. They were also alerted to the fact that the
recordings were subtly different and therefore may sound repetitive.

Unbeknownst to participants, the 234 trials were split into exposure (54 trials each) and
test blocks (12 trials each). Participants were given the opportunity to take breaks after every 60
trials, which was always during an exposure block. Finally, participants completed an exit survey
and an optional demographics survey.

*Test blocks.*   The experiment started with a test block. Test blocks were identical within
and across conditions, always including 12 minimal pair trials assessing participants'
categorization at 12 different VOTs (-5, 5, 15, 25, 30, 35, 40, 45, 50, 55, 65, 70 ms). A test block

―――――――

[2] We follow previous work (Kleinschmidt, 2020; Lisker & Abramson, 1964) and refer to pre-voicing as negative
VOTs though we note that pre-voicing is perhaps better conceived of as a separate phonetic feature (for discussion,
see **REF?**). Estimates of the proportion of voiced stops produced with pre-voicing in L1-US English vary
substantially between studies (between 20% and 57%) (Dmitrieva, Llanos, Shultz, & Francis, 2015; e.g. Lisker &
Abramson, 1967; Smith, 1978; Westbury, 1979). Because pre-voicing is not regarded as a phonemic determinant of
English, some studies either discard such data or ignore them altogether (e.g. Zue (1976); Klatt (1975); Chodroff
and Wilson (2017)). In some studies that do report pre-voicing, the majority of the tokens were attributed to a
minority of talkers (Flege & Brown Jr, 1982; e.g. Lisker & Abramson, 1967). Although speakers tend to prefer one
type of production over the other they do not typically use one type exclusively (Docherty, 2011).
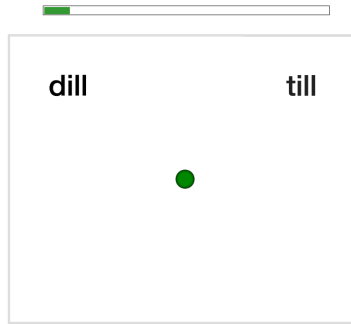
*Figure 2.* Example trial display. When the green button turned bright green, participants had to click on it to play the recording.

followed each exposure block to assess the effects of cumulative exposure. As noted earlier while repeated testing is one of the core innovations in the present design it does come with challenges which informed the decision to keep testing short. First, listeners' attention span is limited. Second, previous work has found that repeated testing over uniform test continua can reduce or undo the effects of informative exposure (Cummings & Theodore, 2023; Liu & Jaeger, 2018, 2019; Tzeng et al., 2021). Third, holding the distribution of test stimuli constant across exposure condition inevitably means that the relative unexpectedness of these test stimuli differs between the exposure conditions. By keeping tests short relative to exposure, we aimed to minimize the influence of test trials on adaptation while still being able to estimate changes in listeners categorization function.

A uniform distribution over VOTs was chosen to maximize the statistical power to determine participants' categorization function. The assignment of VOTs to minimal pair continua was randomized for each participant, while counter-balancing it within and across test blocks. Each minimal pair appear equally often within each test block (four times), and each minimal pair appear with each VOT equally often (twice) across all six test blocks (and no more than once per test block).

Each trial started with a dark-shaded green fixation dot being displayed. At 500ms from trial onset, two minimal pair words appeared on the screen, as shown in Figure 2. At 1000ms from trial onset, the fixation dot would turn bright green and participants had to click on the dot to play the recording. This was meant to reduce trial-to-trial correlations by resetting the mouse pointer to the center of the screen at the start of each trial. Participants responded by clicking on

the word they heard and the next trial would begin.

*Exposure blocks.*   Each exposure block consisted of 24 /d/ and 24 /t/ trials, as well as 6 catch trials that served as a check on participant attention throughout the experiment (2 instances for each of three combinations of the three catch recordings). With a total of 144 trials, exposure was substantially shorter than in similar previous experiments (cf. 228 trials in Clayards et al., 2008; 222 trials in Kleinschmidt, 2020; 2 x 236 trials, Theodore & Monto, 2019; 456 trials, Nixon et al., 2016).

The distribution of VOTs across the 48 /d/-/t/ trials depended on the exposure condition. Specifically, we first created a *baseline* condition. Although not critical to the purpose of the experiment, we aimed for the VOT distribution in this condition to closely resemble participants' prior expectations for a 'typical' female talker of L1-US English (for details, see SI, **??**). The mean and standard deviations for /d/ along VOT were set at 5 ms and 8.9 ms, respectively. The mean and standard deviations for /t/ were set at 50 ms and 16 ms, respectively. To create more realistic VOT distributions, we *sampled* from the intended VOT distribution (top row of Figure 3). This creates distributions that more closely resemble the type of distributional input listeners experience in everyday speech perception, deviating from previous work, which exposed listeners to highly unnatural fully symmetric samples (Clayards et al., 2008; Kleinschmidt, 2020; Kleinschmidt & Jaeger, 2016).

Half of the /d/ and half of the /t/ trials were labeled, the other half was unlabeled. Earlier distributional learning studies have mostly used fully unlabeled exposure (Bejjanki et al., 2011; Clayards et al., 2008; Nixon et al., 2016). This contrasts with visually- or lexically-guided perceptual learning studies, which use labeled exposure (Bertelson et al., 2003; Kraljic & Samuel, 2005; Norris et al., 2003; Vroomen et al., 2007). Such labeling is known to facilitate adaptation (**burchill2018?**; **burchill2023?**; but see Kleinschmidt et al., 2015)—indeed, if shifted pronunciations are embedded in minimal pair or nonce-word context, listeners do no longer shift their categorization boundary (Norris et al., 2003; **REF-theodore?**; **babel?**). While lexical contexts often disambiguate sounds in everyday speech, that is not *always* the case: especially, when confronted with unfamiliar accents, listeners often have uncertainty about the word they are hearing, and must either use contextual information to label the input or adapt from unlabeled

<sup></sup>input. Here, we thus aimed to strike a compromise between always and never labeling the input

(paralleling one of the conditions in Kleinschmidt et al., 2015).

Unlabeled trials were identical to test trials except that the distribution of VOTs across those trials was bimodal (rather than uniform), and determined by the exposure condition.[3] Labeled trials instead presented two response options with identical stop onsets (e.g., *din* and *dill*). This effectively labeled the input as belonging to the intended category (e.g., /d/).

Next, we created the two additional exposure conditions by shifting these VOT distributions by +10 or +40 ms (see Figure 3). This approach exposes participants to heterogeneous approximations of normally distributed VOTs for /d/ and /t/ that varied across blocks, while holding all aspects of the input constant across conditions except for the shift in VOT. The order of trials was randomized within each block and participant, with the constraint that no more than two catch trials would occur in a row. Participants were randomly assigned to one of 3 (exposure condition) x 3 (block order) x 2 (placement of response options) lists.

### 2.3.1 Exclusions

Due to data transfer errors 4 participants' data were not stored and therefore excluded from analysis. We further excluded from analysis participants who committed more than 3 errors out of the 18 catch trials (<83% accuracy, N = 1), participants who committed more than 4 errors out of the 72 labelled trials (<94% accuracy, N = 0), participants with an average reaction time more than three standard deviations from the mean of the by-participant means (N = 0), participants who had atypical categorization functions at the start of the experiment (N = 2, see SI, **??** for details), and participants who reported not to have used headphones (N = 0). This left for analysis 17,136 exposure and 8,568 test observations from 119 participants (94% of total), evenly split across the three exposure conditions.

---

[3] Previous studies have estimated changes in participants' categorization responses by analyzing responses on unlabeled exposure trials (e.g., Clayards et al., 2008; Kleinschmidt & Jaeger, 2016; Theodore & Monto, 2019). This approach compares responses across different values of acoustic-phonetic cues (since the exposure inputs differed by exposure condition), so that assumptions baked into the analysis approach (e.g., linearity along the acoustic-phonetic continuum) can potentially bias the results. Here we avoid this issue by holding test stimuli constant (see also Kleinschmidt, 2020, Experiment 4).
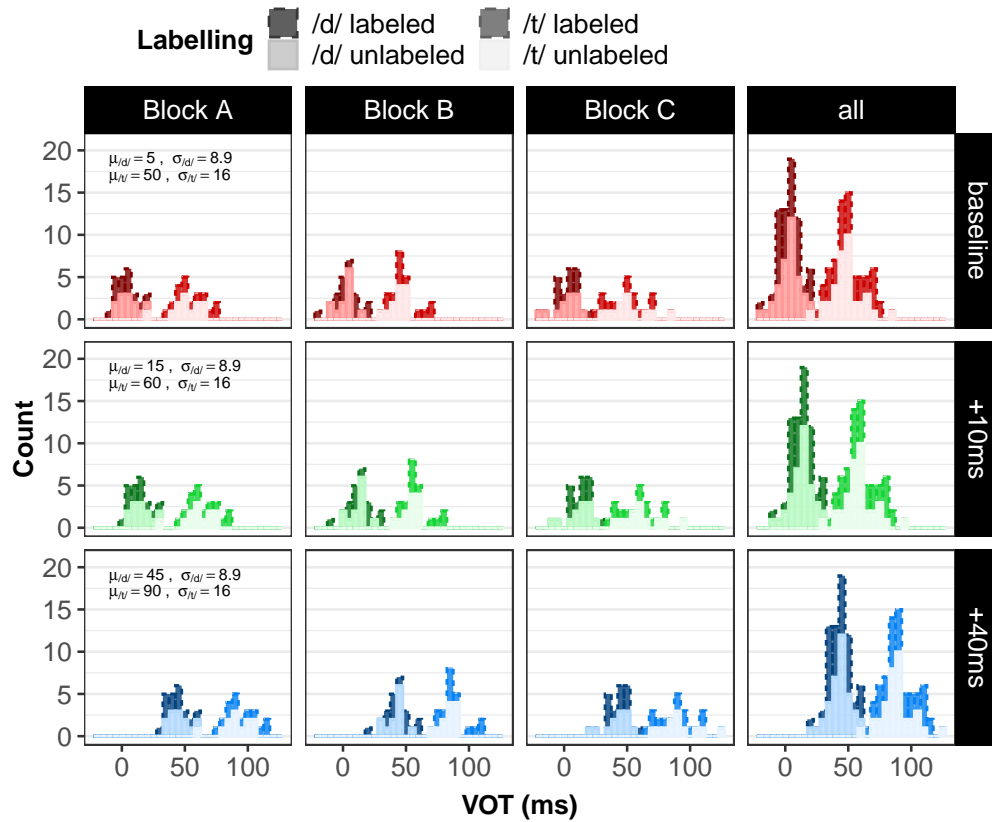
*Figure 3.* Histogram of voice onset times (VOTs) for each of the three exposure blocks A-C by trial type (/d/ or /t/, labeled or unlabeled) and exposure condition (baseline vs. +10 vs. +40). Each exposure block contained 12 labeled /d/, 12 labeled /t/, 12 unlabeled /d/, and 12 unlabeled /t/ trials, as well as 6 catch trials (not shown). Except for the shift in VOTs (+0, 10 or 40 ms VOT to each trial), the VOT distribution of these trials was identical across exposure conditions. The order of exposure blocks A-C was counter-balanced across participants using a Latin-square design.

## 2.4 Results

We analyzed participants' categorization responses during exposure and test blocks in two separate Bayesian mixed-effects psychometric models, using brms (Bürkner, 2017) in R (R Core Team, 2022; RStudio Team, 2020, for details, see SI, **??**). Psychometric models account for attentional lapses while estimating participants' categorization functions. Failing to account for attentional lapses—while commonplace in research on speech perception (but see Clayards et al., 2008; Kleinschmidt & Jaeger, 2016)—can lead to biased estimates of categorization boundaries (Prins, 2011; Wichmann & Hill, 2001). For the present experiment, however, lapse rates were negligible (0.8%, 95%-CI: 0.4 to 1.5%), and all results replicate in simple mixed-effects logistic

regressions (Jaeger, 2008). The estimated lapse rate compares favourably against the lapse rates reported in prior work (Clayards et al., 2008; Kleinschmidt, 2020; Kleinschmidt & Jaeger, 2016). Particularly, Kleinschmidt (2020) reported that the lapse rates in the first sixth of trials was estimated at 12% before it stabilised at an average of 5% for the rest of the experiment. The disparity in lapse rates may be indicative of the qualitative difference between the stimuli employed in both experiments.

Each psychometric model regressed participants' categorization responses against the full factorial interaction of VOT, exposure condition, and block, while including the maximal random effect structure (see SI, **??**. Figure 4 summarizes the results that we describe in more detail next. Panels A and B show participants' categorization responses during exposure and test blocks, along with the categorization function estimated from those responses via the mixed-effects psychometric models. These panels facilitate comparison between exposure conditions within each block. Panels C and D show the slope and point of subject equality (PSE)—i.e., the point at which participants are equally likely to respond "d" and "t"—of the categorization function across blocks and conditions. These panels facilitate comparison across blocks within each exposure condition. Here we focus on the test blocks, which were identical within and across exposure conditions. Analyses of the exposure blocks are reported in the SI (**??**), and replicate all effects found in the test blocks.

We begin by presenting the overall effects, averaging across all test blocks. This part of our analysis matches previous work, which has focused on the overall effect of exposure across the entire experiment ('batch tests,' e.g., Clayards et al., 2008; Kleinschmidt & Jaeger, 2016; Nixon et al., 2016; Theodore & Monto, 2019) and/or during a single post-exposure test block (e.g., Kleinschmidt, 2020). Then we turn to the results that address the empirical gaps concerning incremental changes in participants' categorization responses as a function of 1) the amount and 2) distribution of phonetic evidence from the exposure input.
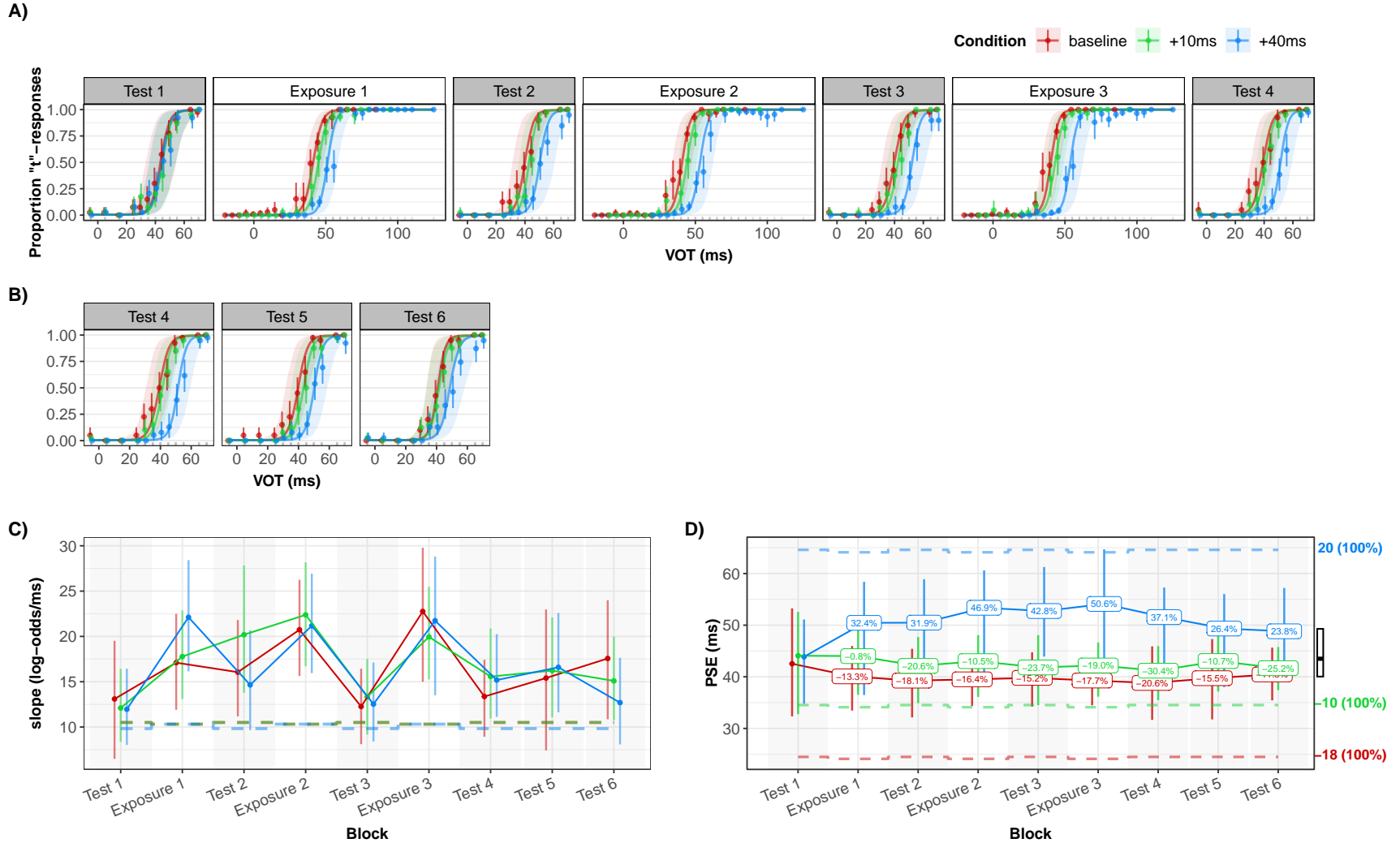
*Figure 4.* Summary of results. **Panel A:** Changes in listeners psychometric categorization functions as a function of exposure, from Test 1 to Test 4 with all intervening exposure blocks (only unlabeled trials were included in the analysis of exposure blocks since labeled trials provide no information about listeners' categorization function). Point ranges indicate the mean proportion of "t"-responses and their 95% bootstrapped CI. Lines and shaded intervals show the MAP predictions and 95% posterior CIs of a Bayesian mixed-effects psychometric model fit to participants' responses. **Panel B:** Same as Panel A but for the final three test blocks without intervening exposure. Test 4 is shown as part of both Panels A and B. **Panels C & D:** Changes across blocks in the slope and boundary (point-of-subjective-equality, PSE) of the categorization functions shown in Panels A-B. Point ranges represent the posterior medians and their 95% CI. Dashed reference lines show the intercepts and PSEs that naive (non-rational) learner would be expected to converge against after sufficient exposure (an ideal observer model that knows the exposure distributions). Percentage labels indicate the amount of shift

### 2.4.1 Does exposure affect participants' categorizations (averaging across all blocks)?

We first used the psychometric mixed-effects model to assess whether the exposure conditions had the expected effects across all test blocks *relative to each other*. Unsurprisingly, participants were more likely to respond "t" the larger the VOT $(\hat{\beta} = 15.09, \ 90\%-\text{CI} = [12.377, 17.625], \ BF = Inf, \ p_{posterior} = 1)$. Critically, exposure affects participants' categorization responses in the expected direction. Marginalizing across all blocks, participants in the +40 condition were less likely to respond "t" than participants in the +10 condition $(\hat{\beta} = -2.26, \ 90\%-\text{CI} = [-3.258, -1.228], \ BF = 162.3, \ p_{posterior} = 0.994)$ or the baseline condition $(\hat{\beta} = -3.08, \ 90\%-\text{CI} = [-4.403, -1.669], \ BF = 215.2, \ p_{posterior} = 0.995)$. There was also evidence—albeit less decisive—that participants in the +10 condition were less likely to respond "t" than participants in the baseline condition $(\hat{\beta} = -0.82, \ 90\%-\text{CI} = [-1.887, 0.282], \ BF = 8.9, \ p_{posterior} = 0.899)$. That is, the +10 and +40 conditions resulted in categorization functions that were shifted rightwards compared to the baseline condition, as also visible in Figures 4.

This replicates previous findings that exposure to changed VOT distributions changes listeners' categorization responses (for /b/-/p/: Clayards et al., 2008; Kleinschmidt, 2020; Kleinschmidt & Jaeger, 2016; for /g/-/k/, Theodore & Monto, 2019). Having established that exposure affected categorization, we turn to the questions of primary interest. Incremental changes in participants' categorization responses can be assessed from three mutually complementing perspectives. First, we compare how exposure affects listeners' categorization responses relative to other exposure conditions. This tests how early in the experiment differences between exposure conditions began to emerge. Second, we compare how exposure affects listeners' categorization responses within each condition relative to listeners' responses prior to any exposure. Third and finally, we compare changes in listeners' responses to those expected from an ideal observer that has fully learned the exposure distributions. This investigates the degree of boundary shift listeners make at each test block relative to their expectations before informative exposure.

### 2.4.2   Comparing across exposure conditions: How quickly does exposure begin to affect participants' responses?

Figure 4A suggests that differences between exposure conditions emerged early in the experiment: already in Test 2, listener's categorization functions seem to be shifted rightwards (larger PSEs) in the +40 condition compared to the +10 condition, and in the +10 condition compared to the baseline condition. This is confirmed by the Bayesian hypothesis tests summarized in Table 1. Prior to any exposure, during Test 1, participants' responses did not differ across exposure condition (all BFs > 3). After exposure to only 24 /d/ and 24 /t/ stimuli, during Test 2, participants' responses differed between exposure conditions (BFs > 13.7). The difference between the +40 condition and the +10 or baseline condition kept increasing with exposure up to Test 4. Additional hypothesis tests in Table 2 show that the change from Test 1 to 2 was largest (BF = 57.82), followed by the change from Test 2 to 3 (BF = 10), with only minimal changes from Test 3 to 4 (BF = 1.68). Qualitatively paralleling the changes across blocks for the +40 condition, the change in the difference between the +10 and baseline conditions was largest from Test 1 to 2 (BF = 5.42), and then somewhat decreased from Test 2 to Test 4 (BFs < 1). The comparison across exposure conditions thus suggests that changes in listeners' categorization responses emerged quickly—indeed, they were present already *during* the first exposure block (see SI, **??**)—but then leveled off. The comparison across exposure conditions also yields one result that is, at first blush, surprising: while the difference between the +10 and the baseline condition emerged already after the first exposure block, this difference *de*creased, rather than increased, with additional exposure from Test 2 to 3 (see second row of Table 2). We return to this effect below.

    Tables 1 and 2 also reveal the consequences of repeated testing. The difference between exposure conditions decreased from Test 4 to 6 (BFs > 4.3; see also Figure 4B & D). On the final test block, the +10 condition did not differ any longer from the baseline condition. Only the differences between the +40 condition relative to the +10 and baseline conditions persisted, albeit substantially reduced compared to Test 4. This pattern of results replicates previous findings that repeated testing over uniform test continua can undo the effects of exposure (Cummings & Theodore, 2023; Liu & Jaeger, 2018, 2019), and extends them from perceptual recalibration paradigms to distributional learning paradigms (see also Kleinschmidt, 2020). One important

methodological consequence of these findings is that longer test phases do not necessarily increase

the statistical power to detect effects of adaptation (unless analyses take the effects of repeated

testing into account, as in the approach developed in Liu & Jaeger, 2018). Analyses that average

across all test tokens—as remains the norm—are bound to systematically underestimate the

adaptivity of human speech perception.

Table 1
*When did exposure begin to affect participants' categorization responses? When, if ever, were these changes undone with repeated testing? This table summarizes the simple effects of the exposure conditions for each test block.*

| Hypothesis | Estimate | SE | 90%-CI | BF | $p_{posterior}$ |
|---|---|---|---|---|---|
| **Test block 1 (pre-exposure)** | | | | | |
| +10 vs. baseline = 0 | -0.34 | 0.75 | [-2.025, 1.437] | 3.3 | 0.77 |
| +40 vs. +10 = 0 | 0.25 | 0.73 | [-1.338, 1.903] | 3.7 | 0.79 |
| +40 vs. baseline = 0 | -0.08 | 0.91 | [-2.124, 2.082] | 4.8 | 0.83 |
| **Test block 2** | | | | | |
| +10 vs. baseline | -1.45 | 0.88 | [-2.933, 0.181] | 13.7 | 0.93 |
| +40 vs. +10 | -2.08 | 0.99 | [-3.824, -0.173] | 24.3 | 0.96 |
| +40 vs. baseline | -3.49 | 1.24 | [-5.635, -1.072] | 54.2 | 0.98 |
| **Test block 3** | | | | | |
| +10 vs. baseline | -0.78 | 0.62 | [-1.888, 0.364] | 7.9 | 0.89 |
| +40 vs. +10 | -2.80 | 0.82 | [-4.188, -1.113] | 86.0 | 0.99 |
| +40 vs. baseline | -3.56 | 0.97 | [-5.202, -1.582] | 110.1 | 0.99 |
| **Test block 4** | | | | | |
| +10 vs. baseline | -0.88 | 0.85 | [-2.36, 0.847] | 4.8 | 0.83 |
| +40 vs. +10 | -3.32 | 0.89 | [-4.883, -1.636] | 128.0 | 0.99 |
| +40 vs. baseline | -4.16 | 1.21 | [-6.275, -1.882] | 122.1 | 0.99 |
| **Test block 5 (no additional exposure)** | | | | | |
| +10 vs. baseline | -1.33 | 0.71 | [-2.556, -0.003] | 19.1 | 0.95 |
| +40 vs. +10 | -2.38 | 0.86 | [-3.893, -0.796] | 65.1 | 0.98 |
| +40 vs. baseline | -3.25 | 1.24 | [-5.307, -0.923] | 53.0 | 0.98 |
| **Test block 6 (no additional exposure)** | | | | | |
| +10 vs. baseline | -0.22 | 0.72 | [-1.485, 1.114] | 1.6 | 0.62 |
| +40 vs. +10 | -1.70 | 0.79 | [-3.078, -0.171] | 25.0 | 0.96 |
| +40 vs. baseline | -2.57 | 1.22 | [-4.58, -0.191] | 24.0 | 0.96 |

Table 2

*Was there incremental change from test block 1 to 4? Did these changes dissipate with repeated testing from block 4 to 6? This table summarizes the interactions between exposure condition and block, whether the differences between exposure conditions changed from test block to test block.*

| Hypothesis | Estimate | SE | 90%-CI | BF | $p_{posterior}$ |
|---|---|---|---|---|---|
| **Difference in +10 vs. baseline** | | | | | |
| Block 1 to 2: increased $\Delta_{PSE}$ | -0.85 | 0.78 | [-2.166, 0.632] | 5.42 | 0.84 |
| Block 2 to 3: increased $\Delta_{PSE}$ | 0.34 | 0.77 | [-1.144, 1.761] | 0.48 | 0.32 |
| Block 3 to 4: increased $\Delta_{PSE}$ | 0.06 | 0.77 | [-1.382, 1.532] | 0.89 | 0.47 |
| *Block 1 to 4: increased $\Delta_{PSE}$* | -0.42 | 1.26 | [-2.759, 1.963] | 1.70 | 0.63 |
| Block 4 to 5: decreased $\Delta_{PSE}$ | -0.33 | 0.60 | [-1.43, 0.785] | 0.41 | 0.29 |
| Block 5 to 6: decreased $\Delta_{PSE}$ | 1.03 | 0.65 | [-0.234, 2.164] | 11.95 | 0.92 |
| *Block 4 to 6: decreased $\Delta_{PSE}$* | 0.70 | 0.82 | [-0.896, 2.177] | 3.83 | 0.79 |
| **Difference in +40 vs. +10** | | | | | |
| Block 1 to 2: increased $\Delta_{PSE}$ | -2.36 | 0.89 | [-3.811, -0.754] | 57.82 | 0.98 |
| Block 2 to 3: increased $\Delta_{PSE}$ | -1.16 | 0.83 | [-2.592, 0.312] | 10.00 | 0.91 |
| Block 3 to 4: increased $\Delta_{PSE}$ | -0.27 | 0.82 | [-1.694, 1.162] | 1.68 | 0.63 |
| *Block 1 to 4: increased $\Delta_{PSE}$* | -3.78 | 1.22 | [-5.865, -1.447] | 84.11 | 0.99 |
| Block 4 to 5: decreased $\Delta_{PSE}$ | 1.14 | 0.77 | [-0.244, 2.514] | 11.38 | 0.92 |
| Block 5 to 6: decreased $\Delta_{PSE}$ | 0.45 | 0.77 | [-0.985, 1.787] | 2.58 | 0.72 |
| *Block 4 to 6: decreased $\Delta_{PSE}$* | 1.59 | 1.00 | [-0.3, 3.323] | 12.68 | 0.93 |
| **Difference in +40 vs. baseline** | | | | | |
| Block 1 to 2: increased $\Delta_{PSE}$ | -3.16 | 1.02 | [-4.958, -1.185] | 79.00 | 0.99 |
| Block 2 to 3: increased $\Delta_{PSE}$ | -0.82 | 1.08 | [-2.749, 1.145] | 3.39 | 0.77 |
| Block 3 to 4: increased $\Delta_{PSE}$ | -0.20 | 1.08 | [-2.146, 1.741] | 1.34 | 0.57 |
| *Block 1 to 4: increased $\Delta_{PSE}$* | -4.19 | 1.71 | [-7.219, -0.93] | 45.78 | 0.98 |
| Block 4 to 5: decreased $\Delta_{PSE}$ | 0.80 | 0.92 | [-0.971, 2.493] | 4.16 | 0.81 |
| Block 5 to 6: decreased $\Delta_{PSE}$ | 1.48 | 0.94 | [-0.36, 3.117] | 10.85 | 0.92 |
| *Block 4 to 6: decreased $\Delta_{PSE}$* | 2.27 | 1.27 | [-0.12, 4.442] | 16.47 | 0.94 |

### 2.4.3   Comparing within exposure conditions: How quickly does exposure begin to affect participants' responses?

Next, we compared how exposure affects listeners' categorization responses within each condition relative to listeners' responses prior to any exposure. These changes are summarized for the slope and PSE in Figure 4C & D, respectively. This visualization makes apparent two aspects of participants' behavior that were not readily apparent in the statistical comparisons we have summarized so far. First, while the PSEs for the +40 and +10 conditions were shifted rightwards compared to the baseline condition, both the +10 and the baseline condition actually shift leftwards relative to their pre-exposure starting point in Test 1. This is confirmed by Bayesian

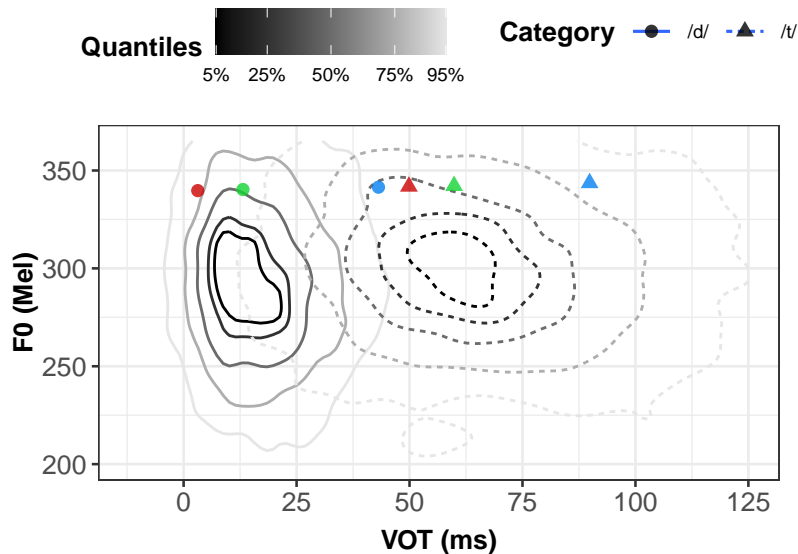456  hypothesis tests summarized in Table **??**.



*Figure 5.* Placement of exposure stimuli relative to an estimate of typical phonetic distributions for 4212 word-initial /d/ and /t/ productions in L1-US English (based on 65 female talkers in Chodroff & Wilson, 2018). The outermost contour of each category shows the 95% density quantile. Points show the category means of the exposure condition.

457       To understand this pattern, it is helpful to relate our exposure conditions to the

458  distribution of VOT in listeners' prior experience. Figure 5 shows the category means of our

459  exposure conditions relative to the distribution of VOT by talkers of L1-US English (based on

460  Chodroff & Wilson, 2018). This comparison offers an explanation as to why the baseline

461  condition (and to some extent the +10 condition) shift leftwards with increasing exposure,

462  whereas the +40 condition shifts rightwards: relative to listeners' prior experience our baseline

463  condition actually presented lower-than-expected category means; of our three exposure

464  conditions, only the +40 condition presented larger-than-expected category means. That is, once

465  we take into account how our exposure conditions relate to listeners' prior experience, both the

466  direction of changes from Test 1 to 4 *within* each exposure condition, and the direction of

467  differences *between* exposure conditions receive an explanation.

468       Second, the reason for the slight decrease in the difference between the +10 and baseline

469  conditions observed in Tables 1 and 2 (visible in Figure 4D as the decreasing difference between

470  the green and red line) is *not* due to a reversal of the effects in the +10 condition. Rather, both

conditions are changing in the same direction but the baseline condition did not move much after Test 2 which reduces the difference between the +10 and baseline conditions (see Table 1). The relative distances between the baseline and +10 condition will become clearer when we assess them with ideal observers.

The comparison across blocks leaves us with mixed impressions. Firstly, across all conditions participants' responses initially changed rapidly with exposure. The pattern that follows after this initial change becomes less clear with increasing exposure, and depends on the direction the exposure condition was shifted relative to participants' initial expectations. In the rightward-shifted +40 condition incremental shifting was observed in Test 3 albeit at a smaller increase compared to Test 2. By Test 4 participants appear to have retracted their boundaries. Taking the general trajectory across test blocks into account, it is possible that listeners reached a limit to the amount they were willing to shift after the end of 144 exposure trials although the evidence for a plateau is not strong given the very wide range of posterior estimates. Participants in the leftwards-shifted baseline condition did not show clear evidence of incremental shifting after Test 2, and instead moved their boundaries within a tight band. In the +10 condition, also leftward-shifted, we see similar boundary movements along a narrow range although notably the shifts up to Test 4 did progressively increase.

### 2.4.4   Constraints on cumulative changes

Finally, Figures 4C & D also compare participants' responses against those of an ideal observer that has fully learned the exposure distributions. The dashed lines represent the respective optimal boundaries of each condition while the labels indicate the amount of shift made at each block as a proportion of the distance between the ideal PSE and the PSE at Test 1. Notably, shifts were always in the right direction but none of the groups converged on the ideal boundary. We also see that while the +10 condition fell short of the ideal boundary changes in PSEs consistently and incrementally moved towards the target up to test 4. Even so, the magnitude of shift was relatively low with the group achieving at most 30% of the maximal shift. What is most striking from the figure is the asymmetry in listener behavior between the leftward-shifted and rightward-shifted groups: when the exposure distribution is rightward shifted listeners showed a

greater propensity to move their category boundaries further from initial expectations. When the exposure distribution is leftward shifted, listeners are far more conservative with their shifts and appear to be under greater constraints. This is most obvious between the baseline and +40 condition; the baseline condition is almost a mirror opposite in shift (-18ms from the PSE at Test 1) compared to the +40 condition (+20ms from the PSE at Test 1) but the maximum shift achieved by the former was just over 20% compared to 43% in the latter.

## 3   General discussion

- discuss consequences of findings for other accounts (decision-making; normalization)

- discuss fact that test stimuli deviate from exposure stimuli to different extent. on the one hand, it's just 1/4 of all trials. on the other hand, we do see relatively systematic changes in slopes each time we test. so there is evidence that even these 12 trials can affect categorisation slopes (though it is worth keeping in mind that this is a comparison across different sets of stimuli). could this explain shrinkage? unlikely since it wasn't the case in kleinschmidt and jaeger. could it explain the constraint on adaptation? that's less clear. we can, however, compare the relative mean of exposure and test. future studies could rerun the exact same paradigm but only test at position x (i.e., a between-subject version of our design)

- could some form of moving window with historical decay explain the findings? On the one hand if the moving window is very small, that would not explain why we do see some *cumulative* changes across blocks (window must be at least $48 + 12 = 60$ trials). on the other hand, the qualitative changes in the PSEs and slopes suggest that 12 trials can be enough to change some aspects of the categorisation function. it's thus *possible* that something that ways recent input much more strongly but also considers less recent input beyond 48 trials might explain the overall pattern.

- discuss potential that observed adaptation maximizes accuracy under the choice rule. use psychometric function fit during unlabeled exposure trials to calculate *accuracy* (not

525    likelihood) on labeled trials under criterion and under proportional matching decision rules.

526    compare against accuracy if ideal observers categorization functions are used instead.

# 4    References

528    Allen, J. S., & Miller, J. L. (1999). Effects of syllable-initial voicing and speaking rate on
529        the temporal characteristics of monosyllabic words. *The Journal of the Acoustical*
530        *Society of America*, *106*(4), 2031–2039.

531    Apfelbaum, K. S., & McMurray, B. (2015). Relative cue encoding in the context of
532        sophisticated models of categorization: Separating information from categorization.
533        *Psychonomic Bulletin & Review*, *22*, 916–943.

534    Assmann, P. F., & Nearey, T. M. (2007). Relationship between fundamental and formant
535        frequencies in voice preference. *The Journal of the Acoustical Society of America*,
536        *122*(2), EL35–EL43.

537    Bejjanki, V. R., Beck, J. M., Lu, Z.-L., & Pouget, A. (2011). Perceptual learning as
538        improved probabilistic inference in early sensory areas. *Nature Neuroscience*, *14*(5),
539        642–648.

540    Bent, T., & Baese-Berk, M. (2021). Perceptual learning of accented speech. *The Handbook*
541        *of Speech Perception*, 428–464.

542    Bertelson, P., Vroomen, J., & De Gelder, B. (2003). Visual recalibration of auditory
543        speech identification: A McGurk aftereffect. *Psychological Science*, *14*(6), 592–597.

544    Boersma, P., & Weenink, D. (2022). *Praat: Doing phonetics by computer. Version 6.2. 12.*

545    Bradlow, A. R., Bassard, A. M., & Paller, K. A. (2023). Generalized perceptual
546        adaptation to second-language speech: Variability, similarity, and intelligibility. *The*
547        *Journal of the Acoustical Society of America*, *154*(3), 1601–1613.

548    Bradlow, A. R., & Bent, T. (2008). Perceptual adaptation to non-native speech.
549        *Cognition*, *106*(2), 707–729.

550    Bürkner, P.-C. (2017). brms: An R package for Bayesian multilevel models using Stan.
551        *Journal of Statistical Software*, *80*(1), 1–28. https://doi.org/10.18637/jss.v080.i01

552    Chodroff, E., & Wilson, C. (2017). Structure in talker-specific phonetic realization:

Covariation of stop consonant VOT in american english. *Journal of Phonetics*, *61*, 30–47.

Chodroff, E., & Wilson, C. (2018). Predictability of stop consonant phonetics across talkers: Between-category and within-category dependencies among cues for place and voice. *Linguistics Vanguard*, *4*(s2).

Clarke, C. M., & Garrett, M. F. (2004). Rapid adaptation to foreign-accented english. *The Journal of the Acoustical Society of America*, *116*(6), 3647–3658.

Clayards, M., Tanenhaus, M. K., Aslin, R. N., & Jacobs, R. A. (2008). Perception of speech reflects optimal use of probabilistic speech cues. *Cognition*, *108*(3), 804–809.

Cummings, S. N., & Theodore, R. M. (2023). Hearing is believing: Lexically guided perceptual learning is graded to reflect the quantity of evidence in speech input. *Cognition*, *235*, 105404.

Dmitrieva, O., Llanos, F., Shultz, A. A., & Francis, A. L. (2015). Phonological status, not voice onset time, determines the acoustic realization of onset f0 as a secondary voicing cue in spanish and english. *Journal of Phonetics*, *49*, 77–95.

Docherty, G. J. (2011). The timing of voicing in british english obstruents. In *The timing of voicing in british english obstruents*. De Gruyter Mouton.

Drouin, J. R., Theodore, R. M., & Myers, E. B. (2016). Lexically guided perceptual tuning of internal phonetic category structure. *The Journal of the Acoustical Society of America*, *140*(4), EL307–EL313.

Eisner, F., & McQueen, J. M. (2005). The specificity of perceptual learning in speech processing. *Perception & Psychophysics*, *67*(2), 224–238.

Flege, J. E., & Brown Jr, W. S. (1982). The voicing contrast between english/p/and/b/as a function of stress and position-in-utterance. *Journal of Phonetics*, *10*(4), 335–345.

Harmon, Z., Idemaru, K., & Kapatsinski, V. (2019). Learning mechanisms in cue reweighting. *Cognition*, *189*, 76–88.

Hitczenko, K., & Feldman, N. H. (2016). Modeling adaptation to a novel accent. *Proceedings of the Annual Conference of the Cognitive Science Society*.

Idemaru, K., & Holt, L. L. (2011). Word recognition reflects dimension-based statistical learning. *Journal of Experimental Psychology: Human Perception and Performance*,

583     *37*(6), 1939.

584     Idemaru, K., & Holt, L. L. (2020). Generalization of dimension-based statistical learning.

585     *Attention, Perception, & Psychophysics*, *82*, 1744–1762.

586     Jaeger, T. F. (2008). Categorical data analysis: Away from ANOVAs (transformation or

587     not) and towards logit mixed models. *Journal of Memory and Language*, *59*(4),

588     434–446.

589     Johnson, K. (1997). Speech perception without speaker normalization. In K. Johnson &

590     J. W. Mullennix (Eds.), *Talker variability in speech processing* (pp. 145–146). San

591     Diego: Academic Press.

592     Klatt, D. H. (1975). Voice onset time, frication, and aspiration in word-initial consonant

593     clusters. *Journal of Speech and Hearing Research*, *18*(4), 686–706.

594     Kleinschmidt, D. (2020). *What constrains distributional learning in adults?*

595     Kleinschmidt, D., & Jaeger, T. F. (2012). A continuum of phonetic adaptation:

596     Evaluating an incremental belief-updating model of recalibration and selective

597     adaptation. *Proceedings of the Annual Meeting of the Cognitive Science Society*, *34*.

598     Kleinschmidt, D., & Jaeger, T. F. (2015). Robust speech perception: Recognize the

599     familiar, generalize to the similar, and adapt to the novel. *Psychological Review*,

600     *122*(2), 148.

601     Kleinschmidt, D., & Jaeger, T. F. (2016). What do you expect from an unfamiliar talker?

602     *CogSci*.

603     Kleinschmidt, D., Raizada, R. D., & Jaeger, T. F. (2015). Supervised and unsupervised

604     learning in phonetic adaptation. *CogSci*.

605     Kraljic, T., & Samuel, A. G. (2005). Perceptual learning for speech: Is there a return to

606     normal? *Cognitive Psychology*, *51*(2), 141–178.

607     Kraljic, T., & Samuel, A. G. (2007). Perceptual adjustments to multiple speakers. *Journal

608     of Memory and Language*, *56*(1), 1–15.

609     Lancia, L., & Winter, B. (2013). The interaction between competition, learning, and

610     habituation dynamics in speech perception. *Laboratory Phonology*, *4*(1), 221–257.

611     Lisker, L., & Abramson, A. S. (1964). A cross-language study of voicing in initial stops:

612     Acoustical measurements. *Word*, *20*(3), 384–422.

Lisker, L., & Abramson, A. S. (1967). Some effects of context on voice onset time in english stops. *Language and Speech*, *10*(1), 1–28.

Liu, L., & Jaeger, T. F. (2018). Inferring causes during speech perception. *Cognition*, *174*, 55–70.

Liu, L., & Jaeger, T. F. (2019). Talker-specific pronunciation or speech error? Discounting (or not) atypical pronunciations during speech perception. *Journal of Experimental Psychology: Human Perception and Performance*, *45*(12), 1562.

Magnuson, J. S., You, H., Luthra, S., Li, M., Nam, H., Escabi, M., et al.others. (2020). EARSHOT: A minimal neural network model of incremental human speech recognition. *Cognitive Science*, *44*(4), e12823.

McMurray, B., & Jongman, A. (2011). What information is necessary for speech categorization? Harnessing variability in the speech signal by integrating cues computed relative to expectations. *Psychological Review*, *118*(2), 219.

Munson, C. M. (2011). *Perceptual learning in speech reveals pathways of processing* ({PhD} dissertation). The University of Iowa.

Nixon, J. S., Rij, J. van, Mok, P., Baayen, R. H., & Chen, Y. (2016). The temporal dynamics of perceptual uncertainty: Eye movement evidence from cantonese segment and tone perception. *Journal of Memory and Language*, *90*, 103–125.

Norris, D., McQueen, J. M., & Cutler, A. (2003). Perceptual learning in speech. *Cognitive Psychology*, *47*(2), 204–238.

Prins, N. (2011). The psychometric function: Why we should not, and need not, estimate the lapse rate. *Journal of Vision*, *11*(11), 1175–1175.

R Core Team. (2022). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing. Retrieved from https://www.R-project.org/

Reinisch, E., & Holt, L. L. (2014). Lexically guided phonetic retuning of foreign-accented speech and its generalization. *Journal of Experimental Psychology: Human Perception and Performance*, *40*(2), 539.

RStudio Team. (2020). *RStudio: Integrated development environment for r*. Boston, MA: RStudio, PBC. Retrieved from http://www.rstudio.com/

Schertz, J., Cho, T., Lotto, A., & Warner, N. (2016). Individual differences in perceptual adaptability of foreign sound categories. *Attention, Perception, & Psychophysics*, *78*, 355–367.

Schertz, J., & Clare, E. J. (2020). Phonetic cue weighting in perception and production. *Wiley Interdisciplinary Reviews: Cognitive Science*, *11*(2), e1521.

Schuster, S. (2020). *Praat: Doing phonetics by computer [computer program]*. Stanford, CA: Interactive Language Processing Lab Stanford. Retrieved from https://docs.proliferate.alps.science/en/latest/contents.html

Sidaras, S. K., Alexander, J. E., & Nygaard, L. C. (2009). Perceptual learning of systematic variation in spanish-accented speech. *The Journal of the Acoustical Society of America*, *125*(5), 3306–3316.

Smith, B. L. (1978). Effects of place of articulation and vowel environment on 'voiced' stop consonant production. *Glossa*, *12*, 163–175.

Sohoglu, E., & Davis, M. H. (2016). Perceptual learning of degraded speech by minimizing prediction error. *Proceedings of the National Academy of Sciences*, *113*(12), E1747–E1756.

Tan, M., Xie, X., & Jaeger, T. F. (2021). Using rational models to interpret the results of experiments on accent adaptation. *Frontiers in Psychology*, 4523.

Theodore, R. M., & Monto, N. R. (2019). Distributional learning for speech reflects cumulative exposure to a talker's phonetic distributions. *Psychonomic Bulletin & Review*, *26*, 985–992.

Tzeng, C. Y., Nygaard, L. C., & Theodore, R. M. (2021). A second chance for a first impression: Sensitivity to cumulative input statistics for lexically guided perceptual learning. *Psychonomic Bulletin & Review*, *28*, 1003–1014.

Vroomen, J., Linden, S. van, De Gelder, B., & Bertelson, P. (2007). Visual recalibration and selective adaptation in auditory–visual speech perception: Contrasting build-up courses. *Neuropsychologia*, *45*(3), 572–577.

Westbury, J. R. (1979). Aspects of the temporal control of voicing in consonant clusters in english. *Texas Linguistic Forum Austin, Tex*, 1–304.

Wichmann, F. A., & Hill, N. J. (2001). The psychometric function: I. Fitting, sampling,

and goodness of fit. *Perception & Psychophysics*, *63*(8), 1293–1313.

Winn, M. B. (2020). Manipulation of voice onset time in speech stimuli: A tutorial and flexible praat script. *The Journal of the Acoustical Society of America*, *147*(2), 852–866.

Woods, K. J., Siegel, M. H., Traer, J., & McDermott, J. H. (2017). Headphone screening to facilitate web-based auditory experiments. *Attention, Perception, & Psychophysics*, *79*, 2064–2072.

Xie, X., Buxó-Lugo, A., & Kurumada, C. (2021). Encoding and decoding of meaning through structured variability in intonational speech prosody. *Cognition*, *211*, 104619.

Xie, X., Jaeger, T. F., & Kurumada, C. (2023). What we do (not) know about the mechanisms underlying adaptive speech perception: A computational framework and review. *Cortex*.

Xie, X., Liu, L., & Jaeger, T. F. (2021). Cross-talker generalization in the perception of nonnative speech: A large-scale replication. *Journal of Experimental Psychology: General*, *150*(11), e22.

Xie, X., Weatherholtz, K., Bainton, L., Rowe, E., Burchill, Z., Liu, L., & Jaeger, T. F. (2018). Rapid adaptation to foreign-accented speech and its transfer to an unfamiliar talker. *The Journal of the Acoustical Society of America*, *143*(4), 2013–2031.

Zheng, Y., & Samuel, A. G. (2020). The relationship between phonemic category boundary changes and perceptual adjustments to natural accents. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *46*(7), 1270.

Zue, V. W. (1976). *Acoustic characteristics of stop consonants: A controlled study.* MASSACHUSETTS INST OF TECH LEXINGTON LINCOLN LAB.