

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

Maryann Tan¹ & T. Florian Jaeger^{2,3}

¹ Centre for Research on Bilingualism, University of Stockholm

² Brain and Cognitive Sciences, University of Rochester

³ Computer Science, University of Rochester

Author Note

We are grateful to ### ommitted for review ###

Correspondence concerning this article should be addressed to Maryann Tan, Department of Bilingualism, Stockholm University, Sweden. E-mail: maryann.tan@biling.su.se

10 Abstract

11 YOUR ABSTRACT GOES HERE. All data and code for this study are shared via OSF,
12 including the R markdown document that this article is generated from, and an R library that
13 implements the models we present.

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

15 Word count: X

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

17 TO-DO

18 **0.1 Highest priority**

- 19 • MARYANN

20 **0.1.1 Lower Priority**

- 21 • Decide on PSE vs. category boundary
- 22 • standardize BE vs. AE spelling (categoriz/sation, label(1)ed, synthesiz/sed etc.)

23 **0.2 To do later**

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

1 Introduction

Human speech perception is now understood to be highly adaptive. Listeners' interpretation of acoustic input can change within minutes of exposure to an unfamiliar talker, improving recognition accuracy (Bradlow & Bent, 2008; Clarke & Garrett, 2004; Xie, Liu, & Jaeger, 2021; Xie et al., 2018). One of mechanisms thought to underlie this rapid adaptivity is distributional learning (Clayards, Tanenhaus, Aslin, & Jacobs, 2008; D. F. Kleinschmidt & Jaeger, 2015; idemaru-hold2011?; davis-sohoglu2020?). This hypothesis has gained considerable influence over the past decade, with findings that changes in listener perception are qualitatively predicted by statistics of exposure stimuli (Bejjanki, Clayards, Knill, & Aslin, 2011; Clayards et al., 2008; Nixon, Rij, Mok, Baayen, & Chen, 2016; Tan, Xie, & Jaeger, 2021; R. M. Theodore & Monto, 2019; idemaru2021?; kleinschmidt2012?; kleinschmidt-jaeger2015cogsci?; munson2011-thesis?; schertz-clare2019?; for important caveats, see harmon2018?).

We investigate an important constraints on this type of adaptivity that is suggested by recent findings. (kleinschmidt-jaeger2016?) exposed L1 US English listeners to over 200 recordings of /b/-/p/ minimal pair words like *beach* and *peach*. In English, the primary cue to this stop voicing contrast is voice onset timing (VOT), with /b/s having shorter VOTs (mean = XXX msec) than /p/s (mean = XXX msec). Kleinschmidt and Jaeger exposed separate groups of listeners to VOT distributions for which these category means had been shifted by XXX, XXX, ..., or XXX msec. In line with the distributional learning hypothesis, listeners' category boundary or point of subjective equality (PSEs)—i.e., the VOT for which listeners are equally likely to respond “d” and “t”—shifted in the same direction as the exposure distribution. Also in line with the distributional learning hypothesis, these shifts were larger the further the exposure distributions were shifted. However, Kleinschmidt and Jaeger also observed a previously undocumented property of these adaptive changes: shifts in the exposure distribution had less than proportional (sublinear) effect on shifts in PSE. While this finding—recently replicated in one more experiment (D. F. Kleinschmidt, 2020, Experiment 4)—is compatible with the hypothesis of distributional learning, it points to important not well-understood constraints on adaptive speech perception.

For example, the only distributional learning model that has been more extensively tested

against adaptive speech perception—incremental Bayesian belief-updating (D. F. Kleinschmidt & Jaeger, 2011)—predicts proportional, rather than sublinear, shifts (for proof, see SI ??). This model had previously been found to closely predict the cumulative effects of exposure in perceptual recalibration to audio-visually (D. F. Kleinschmidt & Jaeger, 2012; **kleinschmidt2011-jaeger?**) or lexically labeled speech (**cummings2023?**), as well as the type of exposure to unlabelled minimal pair words employed by Kleinschmidt and Jaeger (R. Theodore & Monto, 2019). However, all of these studies employed comparatively smaller changes in cue distributions, and lacked the design necessary to detect deviation from proportionality (we return to this point below). The findings presented in (**kleinschmidt-jaeger2016?**) would seem to reject this specific distributional learning model (though not necessarily the theory it is derived from, D. F. Kleinschmidt & Jaeger, 2015; for discussion of the relation between theory and model, see also D. F. Kleinschmidt, 2020; for a recent discussion of the importance of strongly predictive computational models, see **martin-XXX2021?**) Similarly, existing models of perceptual normalization—an alternative, but mutually compatible, hypothesis—also predict proportional changes in PSE (SI ??).

One possib

Xie and colleagues (**xie2018?**) distinguish between two types of mechanisms that might underlie representational changes, *model learning* and *model selection*. The former refers to the learning of a new category representations—for example, learning a new generative model for the talker (D. F. Kleinschmidt & Jaeger, 2015, pt. II) or storage of new talker-specific exemplars (Sumner, 2011). (**xie2018?**) hypothesize that this process might be much slower than is often assumed in the literature, potentially requiring multiple days of exposure and memory consolidation during sleep (see also **fenn2013?**; **tamminen2012?**; **xie2018sleep?**). Rapid adaptation that occurs within minutes of exposure might instead be achieved by selecting between *existing* talker-specific representations that were learned from previous speech input—e.g., previously learned talker-specific generative models (see mixture model in D. F. Kleinschmidt & Jaeger, 2015, pp. 180–181) or previously stored exemplars from other talkers (**johnson1997?**). Model learning and model selection both offer explanations for the sublinear effects observed in (**kleinschmidt-jaeger2016?**). But they suggest different predictions for the

evolution of this effect over the course of exposure.

Under the hypothesis of model learning, sublinear shifts in PSEs can be explained by assuming a hierarchical prior over talker-specific generative models ($p(\Theta)$ in D. F. Kleinschmidt & Jaeger, 2015, p. 180). This prior would ‘shrink’ adaptation towards listeners’ priors—similar to the effect of random by-subject or by-item effects in generalized linear mixed-effect models, which shrink group-level effect estimates towards the population mean of the data (**bates?**). Critically, as long as these priors attribute non-zero probability to even extreme shifts (e.g., the type of Gaussian prior used in mixed-effects models), this predicts listeners’ PSEs will continue to change with increasing exposure until they have converged against the PSE that is ideal for the exposure statistics. In contrast, the hypothesis of model selection predicts that rapid adaptation is more strongly constrained by previous experience: listeners can only adapt their categorisation functions up to a point that corresponds to (a mixture of) previously experienced talker-specific generative models.

Contrastive tests against alternative hypotheses remain lacking (**xie2023?**). This is at least in part due to often informal and vague

- **THE AIM OF THIS STUDY-** The study we report here builds on the pioneering work of Clayards et al. (2008) and D. F. Kleinschmidt and Jaeger (2016) with the aim to shed more light on how listeners’ initial interpretation of cues from a novel talker incrementally change as they receive progressively more informative input of her cue-to-category mappings.

POINTS-TO-MAKE

- The strength of these beliefs has bearing on listener propensity to adapt to a new talker – the stronger the prior beliefs the longer it takes to adapt. Listeners’ strengths in prior beliefs about the means and variances are represented by parameters in the computational model. Listener behaviour observed collectively, thus far which speaks to this framework of thinking should by now be able to indicate roughly what those parameter values are. But it looks like those parameters are biased by the length of exposure and the outcome during experiments. No one has confronted this issue of very quick but limited adaptation which can’t be solved by giving more exposure trials.

- How do we distinguish the results from normalization accounts which can also explain adaptation but is not usually regarded as learning? + will discuss constrain under other hypotheses

A secondary aim of the present study was to *begin* to address possible concerns about ecological validity in research on distributional learning. The pioneering works that inspired the present study employed highly unnatural sounding stimuli that were clearly identifiable as robotic speech (Clayards et al., 2008; **kleinschmidt-jaeger2016?**). These studies also followed the majority of research on distributional learning in language (e.g., **maye2003?**; **pajak2012?**) and *designed* rather than *sampled* the exposure distributions. As a consequence, exposure distributions in these experiments tend to be symmetrically balanced around the category means—unlike in everyday speech input. Indeed, all of the works we follow here further used categories with *identical* variances (e.g., identical variance along VOT for /b/ and /p/, Clayards et al., 2008; **kleinschmidt-jaeger2016?**; or /g/ and /k/, R. Theodore & Monto, 2019). This, too, is highly atypical for everyday speech input (Chodroff & Wilson, 2017; **lisker-abrahamson1964?**). We take modest steps to improve the ecological validity of our stimuli (building on Nixon et al., 2016; R. Theodore & Monto, 2019), and exposure distributions.

All data and code for this article can be downloaded from XXX. 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, 2021; RStudio Team, 2020), while changing any of the parameters of our models (see SI, ??).

2 Experiment

We aimed to design our experiment to provide high statistical power to detect effects of exposure, both incrementally within each exposure condition, and cumulatively across exposure conditions. To this end, we employed the repeated exposure-test design shown in Figure 1. The use of test blocks that repeated same stimuli across blocks and exposure conditions deviates from previous work (Clayards et al., 2008; D. F. Kleinschmidt, 2020; **kleinschmidt-jaeger2016?**). This design feature allowed us to assess how increasing exposure affects listeners’ perception without making

strong assumptions about the nature of these changes (e.g., linear changes across trials). Since previous work has found that repeated testing over uniform test continua can reduce or undo the effects of informative exposure (Liu & Jaeger, 2018, 2019; **cummings202X?**), we kept test blocks short, each consisting of only 12 trials. The final test blocks were intended to ameliorate the potential risks of this novel design: in case adaptation remains stable despite repeated testing, those additional test blocks were meant to provide additional statistical power to detect the effects of cumulative exposure. Finally, as we detail below, our design also allowed us to measure adaptation during exposure.

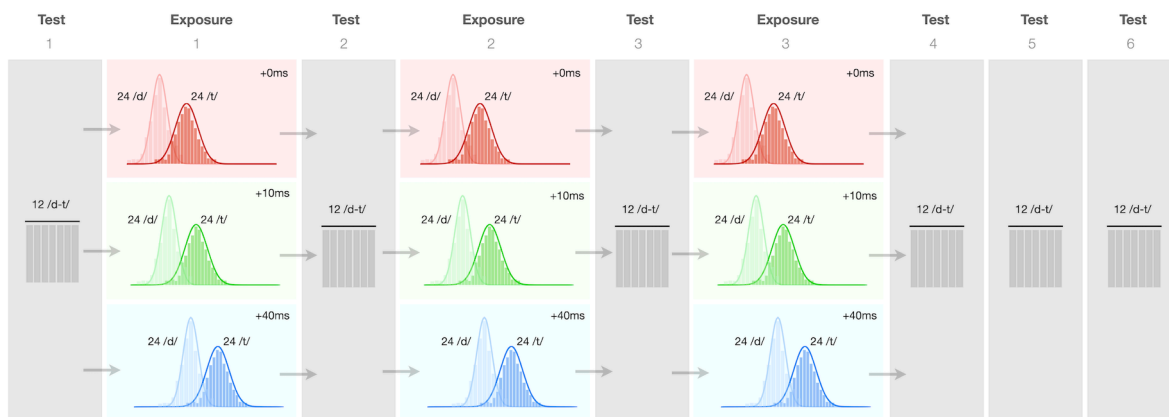


Figure 1. Exposure-test design of the experiment. Test blocks presented identical stimuli within and across conditions

2.1 Methods

2.1.1 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.

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 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?**).

2.1.2 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. These recordings were used to create four natural-sounding minimal pair VOT continua using a script (Winn, 2020) in Praat (**praat?**). The VOTs generated for each continuum ranged from -100 to +130 msec in 5 msec steps.¹ 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, ??).

This approach resulted in continuum steps that sound natural (unlike the highly robotic-sounding stimuli employed in Clayards et al., 2008; D. F. 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 (in contrast to other experiments we have conducted with robotic-sounding stimuli like those of **clayards?**). In addition to the critical minimal pair

¹ For simplicity's sake, we follow previous work (D. F. Kleinschmidt, 2020; **OTHERS?**) and refer to prevoicing as negative VOTs though we note that prevoicing is perhaps better conceived of as a separate phonetic feature (for discussion, see **REF?**). In L1-US English, the occurrence of prevoicing varies between study 20% - 48% of word-initial voiced stops and 0% of voiceless stops (**lisker-abramson1967?**; **smith1978?**).

continua we also recorded three words that 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.

A norming experiment ($N = 24$ participants) reported in the SI (??) was used to select the three minimal pairs that elicited the most similar categorization responses (dill-till, din-tin, and dip-tip). These three continua were used to create the three exposure conditions shown in Figure 1.

2.1.3 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 (REF?), and were instructed to not change the volume throughout the experiment. Following instructions, participants completed 234 two-alternative forced-choice categorisation trials (Figure ??). 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.

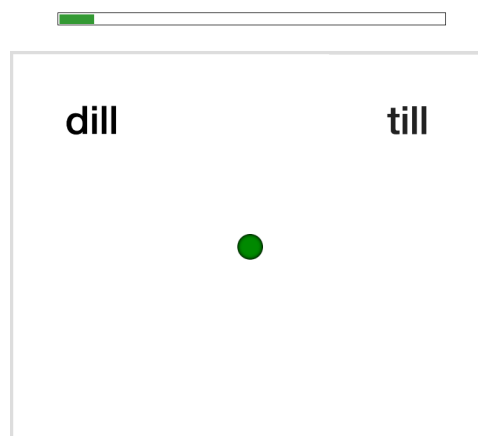


Figure 2. Example trial display. When the green button turned bright green, participants had to click on it to play the recording.

Unbeknownst to participants, the 234 trials were split into exposure blocks (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 msec). A uniform distribution over VOTs was chosen to maximize the statistical power to determine participants' categorisation 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 ???. 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 D. F. Kleinschmidt, 2020; 2 x 236 trials, R. 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 5 msec and 50 msec, respectively.

The mean and standard deviations for /t/ were set 80 msec and 270 msec, 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; D. F. Kleinschmidt, 2020; **kleinschmidt-jaeger2016?**).

Half of the /d/ and half of the /t/ trials were labeled, the other half was unlabeled (paralleling one of the conditions in D. Kleinschmidt, Raizada, & Jaeger, 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. 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 msec (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.

2.1.4 Exclusions

```
## Warning: There were 42 warnings in `mutate()`.
## The first warning was:
## i In argument: `CategorizationModel = map(...)`
## i In group 2: `ParticipantID = 119`, `Experiment = AE-DLVOT`, `Condition.Exposure = Shift0`
## Caused by warning:
## ! glm.fit: fitted probabilities numerically 0 or 1 occurred
## i Run `dplyr::last_dplyr_warnings()` to see the 41 remaining warnings.
```

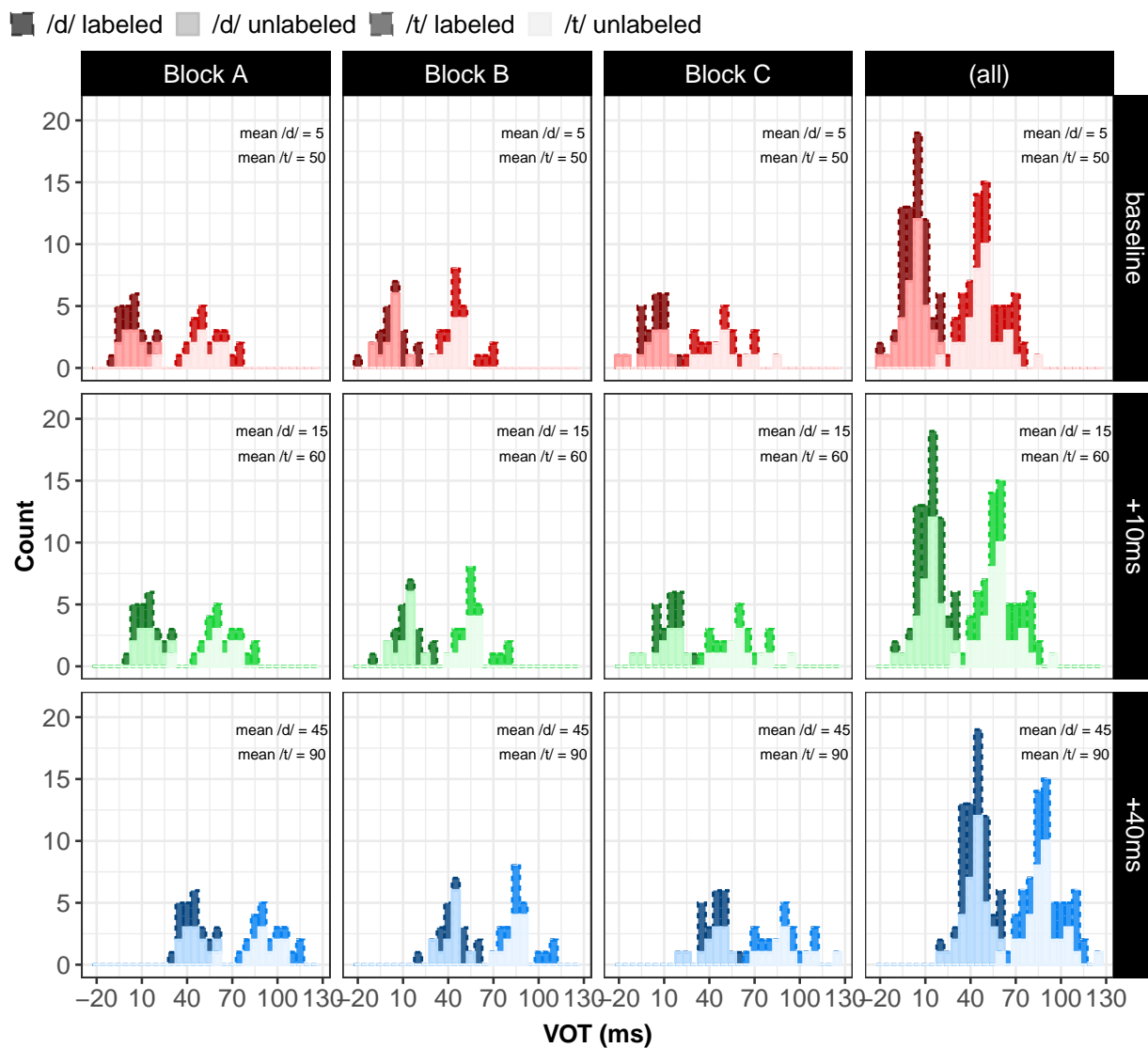


Figure 3. Histogram of VOTs across the 48 trials of all three exposure blocks by exposure condition. The dashed gray line shows the theoretical (Normal) distribution that the baseline condition was sampled from. The order of blocks was counter-balanced across participants.

```

253 ## Warning: Using one column matrices in `filter()` was deprecated in dplyr 1.1.0.
254 ## i Please use one dimensional logical vectors instead.
255 ## This warning is displayed once every 8 hours.
256 ## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was generated.

```

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

266 2.2 Results

267 2.2.1 Research questions and hypotheses

- 268 1. Do listeners change their categorization behaviour in the direction predicted by their
 269 respective exposure distributions?
- 270 2. At what stage in the experiment did the behavioural change first emerge?
- 271 3. Are the shifts in categorisation behaviour proportional to the differences between the
 272 exposure conditions?
- 273 4. Do the differences between exposure conditions diminish with repeated testing and without
 274 intermittent exposure?

275 [MORE HERE]

2.2.2 Analysis approach

Figures 4A-B summarize participants' categorisation responses during exposure and test blocks, depending on the exposure condition and VOT. We analyzed participants' categorisation responses during exposure and test blocks in two separate Bayesian mixed-effects psychometric models, fit using brms (Bürkner, 2017) in R (R Core Team, 2021; RStudio Team, 2020, for details, see SI, ??). These models account for attentional lapses while estimating participants' categorisation functions. Failing to account for attentional lapses—while commonplace in research on speech perception (but see Clayards et al., 2008; D. F. Kleinschmidt & Jaeger, 2016)—can lead to biased estimates of categorization boundaries (Prins, 2012; Wichmann & Hill, 2001). For the present experiment, however, lapse rates were negligible (0.9%, 95%-CI: 0.4 to 1.5%), and all results replicate in simple mixed-effects logistic regressions (Jaeger, 2008).

2.2.3 Does exposure affect participants' categorisations?

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. Unsurprisingly, participants were more likely to respond “t” the larger the VOT ($\hat{\beta} = 15.68$, 90%-CI = [13.149, 18.4], $BF = 7999$, $p_{posterior} = 1$). Critically, exposure affects participants' categorisation 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 ($\hat{\beta} = -2.43$, 90%-CI = [-3.541, -1.363], $BF = 443.4$, $p_{posterior} = 0.998$) or the baseline condition ($\hat{\beta} = -3.39$, 90%-CI = [-4.969, -1.93], $BF = 332.3$, $p_{posterior} = 0.997$). 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 ($\hat{\beta} = -0.97$, 90%-CI = [-2.241, 0.298], $BF = 9.2$, $p_{posterior} = 0.902$). That is, the +10 and +40 conditions resulted in categorisation 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; D. F. Kleinschmidt, 2020;

kleinschmidt-jaeger2016?; for /g/-/k/, theodore-monto2018?). Having established that exposure affected categorization, we turn to the questions of primary interest.

2.2.4 Incremental changes in listeners' categorisation with increasing exposure (Test 1 to 4)

As already visible in Figure 4A, effects of exposure emerged early in the experiment. Table ?? summarizes the simple effects of exposure condition during each of the first four test blocks. Prior to any exposure, during Test 1, participants' responses did not differ across exposure condition. After exposure to only 24 /d/ and 24 /t/ stimuli, during Test 2, participants' responses already differed between exposure conditions. 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 ?? show that the change from Test 1 to 2 was largest ($BF = 27.8$), followed by the change from Test 2 to 3 ($BF = 19.2$), with only minimal changes from Test 3 to 4 ($BF = 1.7$). 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 ($BF = 13.5$), and then somewhat decreased from Test 2 to Test 4 ($BFs < 4$).

This pattern of changes is also evident in Figure 4D, which shows how participants' point of subject equality (PSE)—i.e., the point at which “d” and “t” responses are equally likely—changes with increasing exposure. 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 +10 and +40 conditions were indeed shifted rightwards compared to the baseline condition (relatively larger PSEs), both the +10 and the baseline condition actually shift leftwards *relative* to their pre-exposure starting point in Test 1. Second, the reason for the slight decrease in the difference between the +10 and baseline conditions observed in Tables ?? and ?? (visible in Figure 4D as the decreasing difference between 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 stops changing after Test 2, which brings the +10 condition increasingly closer to the baseline condition. To understand this pattern, it is necessary to relate our exposure conditions to the distribution of VOT in listeners' prior experience.

2.2.5 Relating incremental changes in categorisation to listeners' prior experience (Test 1 to 4)

Figure ?? shows the mean and covariance of our exposure conditions relative to the distribution of VOT by talkers of L1-US English (based on Chodroff & Wilson, 2018). This comparison offers an explanation as to why the baseline condition (and to some extent the +10 condition) shift leftwards with increasing exposure, whereas the +40 condition shifts rightwards: relative to listeners' prior experience our baseline condition actually presented lower-than-expected category means; of our three exposure conditions, only the +40 condition presented larger-than-expected category means. That is, once we take into account how our exposure conditions relate to listeners' prior experience, both the direction of changes from Test 1 to 4 *within* each exposure condition, and the direction of differences *between* exposure conditions receive an explanation.

2.2.6 Constraints on cumulative changes

2.2.7 Effects of repeating testing

Finally, we turn the consequences of repeated testing. As evident in Panel B and D of Figure 4, repeated testing without additional exposure resulting in partial undoing of the effects described so far. Bayesian hypothesis tests confirmed that the difference in the PSE decreased from Test 4 to 6, both for the +40 compared to the +10 condition ($\hat{\beta} = 1.98$, 90%—CI = $[-0.418, 4.338]$, $BF = 12.2$, $p_{posterior} = 0.924$) and the +10 compared to the baseline condition ($\hat{\beta} = 0.93$, 90%—CI = $[-0.921, 2.908]$, $BF = 4.3$, $p_{posterior} = 0.811$).

This replicates previous findings that repeated testing over uniform test continua can undo the effects of exposure (Liu & Jaeger, 2018, 2019; **cummings?**; **others?**), and extends them from perceptual recalibration paradigms to distributional learning paradigms. One important methodological consequence of this 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.

```
358 ## Warning: Using `size` aesthetic for lines was deprecated in ggplot2 3.4.0.  
359 ## i Please use `linewidth` instead.  
360 ## This warning is displayed once every 8 hours.  
361 ## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was generated.
```

