# Critical considerations for conducting web-based speech perception research

Rachel M. Theodore, Ph.D.
*Department of Speech, Language, and Hearing Sciences*
*Connecticut Institute for the Brain and Cognitive Sciences*
University of Connecticut

NIDCD    NSF    UCONN
UNIVERSITY OF CONNECTICUT

Hello, and welcome to the thematic session, "Changes in Space: Online Experimentation." I am excited to share with you:

**Critical considerations for conducting web-based speech perception research.**

*For more in-depth resources, including a copy of the slide deck and talk transcript:*

slaplab.uconn.edu

@rachelmtheodore

osf.io/8krg3/

So much so, that I've packed quite a bit into this 25-minute talk. I'd like to point you to some offline resources that you can find at our lab website, including a copy of this slide deck and a transcript of the talk.

At our OSF page, you can find data and analysis code for all of the examples that I'll share with you today, and many other web-based experiments.

## Overview

- Tools
- Successes
- Tips and tricks

In this talk, I'm going to share some tools, successes, and tips and tricks for promoting high quality web-based speech perception research. But first, I'd like to briefly talk about the why - as in, why would we ever consider abandoning the carefully controlled laboratory setting with our fancy headphones and sound booths and participants that we can actually see?

## Why?

- Web-based research promotes reproducibility

  *Larger samples, in-house replications, more diverse samples*

  *Efficient control and stimulus testing*

For me, there are myriad reasons to do so. My foray into web-based research was in direct consequence of me completely drinking the reproducibility Kool-Aid. As we moved to adopt emerging best practices for promoting reproducibility of research, we needed to find ways to collect data from larger sample sizes. And make in-house replication studies the norm. And not limit our samples to reflect the demographics of our university. And run more control experiments. And better vet our stimulus sets. And verify that our results aren't contingent on a single stimulus set. And so on.

## Why?

- Web-based research promotes reproducibility

  *Larger samples, in-house replications, more diverse samples*

  *Efficient control and stimulus testing*

- I like controlling the listening environment, but I don't (always) need to

Many of our studies don't actually require a laboratory level of control over the listening and response environment. We don't present auditory stimuli at threshold levels. Processing time effects that we're interested in exceed keyboard timing error. I like to think that the things we study (and thus the things we claim) might actually be relevant in a more natural listening environment...

## Why?

- Web-based research promotes reproducibility

  *Larger samples, in-house replications, more diverse samples*

  *Efficient control and stimulus testing*

- I like controlling the listening environment, but I don't (always) need it

- Productivity

It's hard to keep up with my extremely productive colleagues if we limit testing to our physical lab. For better or worse, science moves very quickly these days. Online data collection also facilitates productivity of our trainees, who are increasingly expected to have strong publication records.

## Why?

- Web-based research promotes reproducibility

  *Larger samples, in-house replications, more diverse samples*

  *Efficient control and stimulus testing*

- I like controlling the listening environment, but I don't (always) need it

- Productivity

- Technologies exist to provide high quality web-based data collection, even for speech perception studies

And finally, emerging technologies do exist to provide high quality web-based data collection, even for speech perception experiments.

- Software to build experiments + server to host web-based studies

- If you can dream it, Gorilla can build it

- Extensive features: *Collaboration, version control, open materials, support*

- Free to build experiments; payment model is based on token currency

If you're listening to this talk, you've probably got a few reasons of your own for considering web-based studies. Let's talk about some tools. I'm going to focus on two, Gorilla Experiment Builder and Prolific. Other tools exist, but are outside of my area of expertise. To be honest, I landed on these two because I just couldn't figure out how to make MTurk work for me.

Gorilla is a both a lovely piece of software to build experiments - think of it as an equivalent to E-Prime, or PsychoPy, OpenSesame, or SuperLab. In addition to building experiments, it's also a server that can host your web-based studies. In my experience, there's no design that Gorilla can't handle.

It has extensive features including collaboration, version control, a repository of open materials consisting of experiments that others have designed, in addition to an engaging and comprehensive, multimedia support. It's not an exaggeration to say that an hour of your time would be sufficient to get you going in Gorilla.

In Gorilla, it's completely free to build experiments; the payment model is based on a token currency such that you're charged one token for each participant who completes your experiment.

- Projects consist of experiments, tasks & questionnaires, and open materials
- Experiments are sequenced tasks & questionnaires
- Open materials are publicly available tasks and questionnaires



Materials in Gorilla are organized around projects. Projects consist of experiments, tasks & questionnaires, and open materials. Experiments are sequenced tasks and questionnaires; that is, experiments are formed by combining the smaller bits into a sequenced order.

Open materials are any aspects of your project that you make publicly available on Gorilla.

Here's an example of the interface for making a questionnaire; specifically, this is how we built a consent form. It's two images, one for each page of the consent form, along with a response option.

Here's an example of the task building interface; it's all super intuitive, does not require any programming knowledge, and is thoroughly described in the support documentation and videos.

And here's an example of the experiment interface. As you can see, the questionnaires, in green, and the tasks, in blue, have been ordered using logical branches, in orange. Here, people move from the start to the consent questionnaire, if they give consent, then they move to a headphone screen task, and if they pass, they then enter one of three experimental tasks. And so on —

## Tools: Gorilla Experiment Builder

- Seamless integration with Prolific, but can be used for any method of recruitment

As I mentioned, Gorilla is not only an experiment builder, but it's also a server that hosts your web-based study. Gorilla integrates seamlessly with Prolific, a tool for recruiting participants that I'll talk about next - but you can use Gorilla with any method of recruitment.

## Change Recruitment Policy

| Disable | Link | Email | Recruitment Service |

**Prolific**
Recruit your participants through Prolific.ac

**Sona Systems®**
Recruit your participants through Sona Systems®

**Amazon Mechanical Turk**
Recruit your participants through Amazon Mechanical Turk

**Cloud Research**
Recruit your participants through Cloud Research (formallly TurkPrime)

**Qualtrics**
Recruit your participants from Qualtrics

**Qualtrics Panel**
Recruit your participants from a Qualtrics Panel

**Kantar Profiles**
Recruit your participants through Kantar Profiles

**Research Now** BETA
Recruit your participants through Research Now

**Third Party**
Recruit your participants using a third-party recruitment company (e.g. marketing agency)

OK

For example, Gorilla has built-in integration with numerous participant platforms including MTurk, Sona, and Qualtrics.

Or you could simply generate a link to distribute anywhere you wish. We do this for in-house pilot testing. For example, we could build an experiment, generate the link, and then drop it in the lab Slack for whoever is free to take the experiment for a run.

## Tools: Gorilla Experiment Builder

- Seamless integration with Prolific, but can be used for any method of recruitment

- Real-time information on participants' progress

In addition to providing lots of ways to get participants to your study, Gorilla also provides great visualization of real-time progress while people are completing your study.

## Participant Progress

Public ID   ██████████

Status   **Complete**

Start Date   29/03/2020 17:41

End Date   29/03/2020 17:57

| | Time | |
|---|---|---|
| 🖥 **start** | 17:41 | |
| 📋 **questionnaire** questionnaire- | 17:41 | Name: Information-Sheet-20 Duration: 34 seconds |
| 🔀 **branch** branch-h3u2 | 17:41 | Branch: Yes |
| 📍 **task** task-ns6v | 17:41 | Name: Headphone-Check Duration: 57 seconds |
| 🔀 **branch** branch-rn36 | 17:42 | Branch: Pass |
| ✓ **checkpoint** checkpoint- | 17:42 | Name: Pass |
| 🎲 **randomiser** randomiser- | 17:42 | Branch: High |
| 📍 **task** task-64s8 | 17:42 | Name: PhLex-001 Duration: 6 min 10 sec |
| 📍 **task** task-xey1 | 17:48 | Name: VST Duration: 7 min 56 sec |
| 📋 **questionnaire** questionnaire- | 17:56 | Name: Demographics Duration: 20 seconds |
| 🖥 **finish** | 17:57 | |

For example, you can watch participants in real time in terms of where they are in your experiment tree. I find this part especially satisfying to monitor — so much so that my trainees joke that it is my favorite television show.

## Tools: Prolific

- Online participant pool with large, diverse sample
- Prolific uses numerous quality control methods to ensure high quality participants
- Prolific aims to provide a more ethical alternative to other platforms (e.g., minimum pay/hour)
- Prolific doesn't host experiments; they route participants to your experiment and handle incentive payments
- Prolific makes money by charging a 33% commission on participant payments

With Gorilla as the tool to build and host experiments, Prolific is a tool to find participants. This is an online pool with a large, diverse, sample. Anyone can sign up to join the pool!

Prolific does a host of things behind the scenes to promote high quality participants. They aim to provide a more ethical alternative to MTurk by setting a floor for participant incentives, among other researcher terms of service. Prolific doesn't host the experiment — they are the middle men between your online study and participants. They make money by charging a 33% commission on participant payments — so if you give the participant $3.00, you'll also give Prolific one dollar.

## Tools: Prolific

- Seamless integration with Gorilla, but can be used to distribute any web-based study
- Extensive participant filters
  - Age
  - Nationality/current residence
  - Language(s)
  - Previous studies
- System fosters efficiency in project administration and delivers *high quality participants*

Prolific integrates seamlessly with Gorilla, but you can use Prolific to distribute any web-based study.

When people join the Prolific pool, they first they do is complete a series of over 150 questions that researchers can then use to filter who is recruited for their study. These include things like age, nationality, residence, language experience, and your own previous studies. Prolific is very receptive to adding new things to the on-boarding form as researchers indicate that new criteria are needed. The interface really streamlines project administration, including submitting receipts for reconciliation. Most importantly, the system excels at delivering high quality participants.

Here's an example of the interface. I've indicated a study name; there's a place to provide an overview of the study that participants can see before they decide to do it; I've indicated some constraints, including that the study can't be completed on a tablet or mobile phone, and that there are audio stimuli.

You provide the link to your study — this is something that Gorilla generates for you.

And Prolific gives you a link to add to the end of your Gorilla study so that participants are automatically routed back to Prolific for their incentive payment.

You get a real time display of how many active users meet your filter constraints —

And the interface for applying participant filters is very easy to use.

Find the participants you need          42,899 participants      ✕

🔍  Search for screeners

| Demographics ① | ‹ Back |

Geographic

Languages                    **Current Country of Residence**

Custom Screener              Participants were asked the following question: **In what country do you currently reside?**

Work                         Please note that Prolific is currently only available for participants who live in OECD countries.
                             Read more about this
Education
                             Select the required responses or   select all
Health
                             Type to search...
Beliefs

Family & relationships       United Kingdom

Lifestyle and interests      **United States**                                    ✓

Technology and online        Ireland
behaviour
                             Germany

                             France

         Remove                                                          Apply

**RESEARCHER**

- New study
- Drafts
- Scheduled
- Active
- Completed

We've found **42,899** matching participants who have been active in the past 90 days

STUDY COST

How many participants are you looking to recruit?

20

How long will your study take to complete?                        ⓘ Max. time: 44 mins

Participants are paid according to your estimated study completion time. If the median completion time exceeds your estimate we will ask you to make additional payments. Read more about study completion time ☑

10 minutes

How much do you want to pay them?

$   1.67                                                        10.02/hr

*Hourly rate*

$6.50                          $10.02 *Good*                          $12.50+

**Total cost: $43.42**

Show cost breakdown ⌄

Save as draft        Preview        Publish ⌄

Everyone in a given study gets the same incentive payment, which is based on a good faith estimate of the average completion time.

The interface while a study is live, and also after it ends, is intuitive and informative.

## Tools: Headphone compliance

- Woods et al. (2017)
- Milne et al. (2020)

The last set of tools I'll tell you about are two tasks designed to assess headphone use in web-based studies. Both of these are extremely clever, quick, dichotic listening tasks —

**Tools: Headphone compliance (Woods et al., 2017)**

- Six-trial, loudness decision task; "pass" is defined as ≥ 5 correct responses

- On each trial, three tones with equal frequency and duration are presented

In phase    Out of phase    < Amplitude

*Amplitude*

*left*

*right*

*Time*

In the interest of time, I won't go into details, except to say that the Woods et al. task uses a dichotic phase cancelation manipulation to gauge headphone use from loudness judgments —

## Tools: Headphone compliance (Milne et al., 2020)

- Six-trial, tone detection task; "pass" is defined as 6 correct responses

- On each trial, three noise bursts are presented

- For one noise burst, noise is presented with a phase shift at 600 Hz

- Over headphones, listeners perceive the Huggins pitch

*Left noise*

*Right noise*

no pitch percept

frequency

$\Phi = 180°$

phase

frequency

phase

frequency

pitch percept

frequency

*Adapted from Figure 1 of Milne et al., 2020*

And the Milne et al. task is a Huggins pitch detection task —

## Tools: Headphone compliance

- The Huggins pitch task (Milne et al., 2020) shows more reliable detection than the loudness detection task (Woods et al., 2017)

- As reported in Milne et al. (2020), combining the two tasks **correctly identified 80%** of headphone users with a **false positive rate of 7%**

- If ear channel matters, be sure to supplement your headphone screens with a simple channel detection task…

Combining the tasks only adds 12 total trials to your study and results in reasonable sensitivity and specificity in detecting the use of stereo headphones.

## Successes

- Categorical perception/distributional learning
- Lexically guided perceptual learning
- Perceptual learning for noise-vocoded speech
- Talker adaptation
- Word familiarity ratings

In quick succession, I'm now going to share five successes we've had with web-based studies — selected to illustrate diversity in design and dependent measure. Preprints or postprints, along with data and analysis code, are available for all of these on our OSF repository.

**Success 1: Categorical perception/distributional learning**

**Block 1**

- 152 trials of phonetic ID for tokens drawn from a VOT continuum to form either short or long VOT input distributions

**Block 2**

- 152 trials of phonetic ID for tokens drawn from a VOT continuum to form either short or long VOT input distributions

Sample 1, Sample 2, Sample 3, Sample 4

p(k) — VOT (ms)

Input distributions
— Short
— Long

**To achieve sample (n = 320), we excluded n = 52 due to failure to perform the task and n = 27 due to failure to pass headphone screen; attrition = 20%.**

Success 1: Categorical perception and distributional learning.

In this study, four samples completed two blocks of a 2AFC category identification task. Across blocks, we manipulated the input distributions specifying the /g/ and /k/ categories. As you can see, all four samples showed the expected logistic relationship between category identification and VOT; critically, all four samples also yielded reliable evidence of distributional learning such that the identification function for the long VOT input distributions is displaced towards longer VOTs compared to the identification function for the short VOT input distributions.

Attrition due to failure to perform the task and failure to pass the headphone screen together yielded an attrition of 20%. As I go through these successes, you're going to see some variability in the attrition rate; in the tips and tricks section of this talk I'll share things we've learned to do to decrease our attrition rate.

## Success 2: Lexically guided perceptual learning

**Block: Exposure**

- 200 trials of a lexical decision task for word and nonword stimuli; critical ambiguous productions embedded in either /s/ or /ʃ/ biasing contexts

**Block: Test**

- 72 trials of phonetic ID for tokens drawn from an /asi/-/aʃi/ continuum

1A: Talker f1

1B: Talker m2

Bias
- SS
- SH

**To achieve sample (n = 560), we excluded n = 32 due to failure to perform the task and n = 112 due to failure to pass headphone screen; attrition = 20%.**

Success 2; lexically guided perceptual learning.

In this study, listeners completed two experimental blocks, an exposure phase and then a test phase. During exposure we manipulated the biasing lexical context for an ambiguous fricative. At test, listeners completed a 2AFC identification task for an *ashi* to *asi* continuum. Robust perceptual learning was observed for both tasks, with more *asi* responses at test for those biased to perceive the ambiguity as /s/ during exposure compared to those who were biased to perceive it as /ʃ/.

## Success 3: Perceptual learning for vocoded speech

**Block: Pre-test**

- 30 trials of a transcription task for vocoded sentences w/o feedback

**Block: Training**

- 150 trials with vocoded sentences
  - *Control*: Sentence transcription w/o feedback
  - *Lexical:* Sentence transcription w/ feedback
  - *Talker*: Talker ID w/ feedback

**Block: Post-test**

- 30 trials of a transcription task for vocoded sentences w/o feedback

**To achieve sample (n = 108), we excluded n = 2 due to failure to perform the task and n = 12 due to failure to pass headphone screen; attrition = 11%.**

Success 3; perceptual learning of noise-vocoded speech.

In this study, listeners completed pre-test, training, and post-test blocks. The task at pre- and post-test was free transcription of noise-vocoded sentences. Robust perceptual learning was observed, with transcription accuracy improved following training. Not shown here is a one-week follow-up test; web-based studies have truly opened logistical doors for us in terms of longitudinal experimental designs.

**Success 4: Talker adaptation**

- Four blocks (64 trials/block) of a speeded word ID task
- Blocks crossed talker variability and phonemic ambiguity
- Dependent measure was reaction time
- Can effects < 100 ms be reliably detected in web-based protocols?

Low ambiguity: /i/ - /o/
57 ± 98
p < 0.001

High ambiguity: /o/ - /u/
106 ± 92
p < 0.001

RT (ms)

Variability

To achieve sample (n = 320), we excluded n = 30 due to failure to meet accuracy criterion and n = 38 due to failure to pass headphone screen; attrition = 17%.

Success 4; talker adaptation.

In this study, listeners completed a speeded 2AFC word identification task for two blocks of stimuli, a single talker block, low variability, and a mixed talker block, high variability. This was our first foray into using RT as a dependent measure for a web-based design. As you can see, we had no challenges in reliably detecting variability effects under 100 ms in this sample —

## Success 4: Talker adaptation

- Four blocks (64 trials/block) of a speeded word ID task
- Blocks crossed talker variability and phonemic ambiguity
- Dependent measure was reaction time
- Can effects < 100 ms be reliably detected in web-based protocols?

Low ambiguity: /i/ - /o/

50 ± 84
$p < 0.001$

High ambiguity: /o/ - /u/

64 ± 99
$p < 0.001$

RT (ms)

1500
1000
500

Low    High    Low    High

Variability

**To achieve sample (n = 320), we excluded n = 30 due to failure to meet accuracy criterion and n = 38 due to failure to pass headphone screen; attrition = 17%.**

And in this sample.

# Success 4: Talker adaptation

- Four blocks (64 trials/block) of a speeded word ID task

- Blocks crossed talker variability and phonemic ambiguity

- Dependent measure was reaction time

- Can effects < 100 ms be reliably detected in web-based protocols?
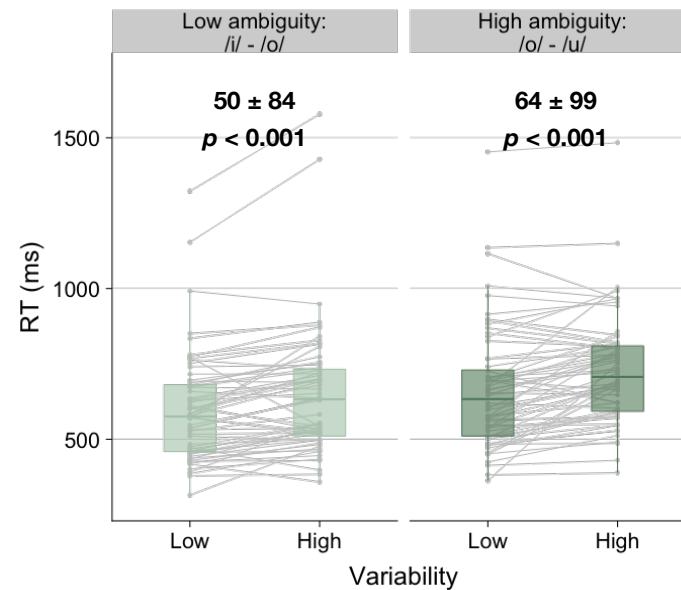


**To achieve sample (n = 320), we excluded n = 30 due to failure to meet accuracy criterion and n = 38 due to failure to pass headphone screen; attrition = 17%.**

And in this sample.

## Success 4: Talker adaptation

- Of 81,920 trials, the audio lag ranged between 0 and 177 ms

- 88% of trials had a lag < 2 ms

- 98% of trials had a lag < 5 ms

- Because Gorilla reports lag time, RTs can be adjusted relative to audio onset

**To achieve sample (n = 320), we excluded n = 30 due to failure to meet accuracy criterion and n = 38 due to failure to pass headphone screen; attrition = 17%.**

Gorilla does a lot of things in the background to optimize stimulus presentation and response timing. In your data file, you get not only the timing of a button response, but also the lag between when the audio stimulus was set to play and when it actually did, which might vary based on a participant's particular system. Because of this, you can correct your RTs to reflect the actual onset of stimulus presentation. We analyzed the lag across all trials for 320 participants in this study and it was exquisite; 98% of the trials had a lag less than 5 milliseconds.

Overall for this study, the magnitude of effects, standard deviations of effects, and proportion of RT outliers were incredibly similar to in-lab work with similar paradigms.

Success 5; word familiarity ratings.

One thing that I think is especially frightening when moving to web-based studies, especially if you're using the Prolific pool, is that you can't see or interact with your participants. As a consequence, researchers often fear that they aren't who they say they are.

To try and develop a tool that might help researchers verify some aspect of language competence, such as, are they a native English speaker as they say they are, we ported a paper-and-pencil vocabulary assessment to Gorilla. This is the word familiarity test developed by David Pisoni and colleagues. On each of 150 trials, participants see a word and are asked to indicate their familiarity with this word.

Success 5: Word familiarity ratings

• Mean ratings by frequency category for the Prolific sample were very similar to existing norms, both by subjects and by items

Sample
Norms
Prolific

Mean rating

Frequency bin

Low Medium High

To achieve sample (n = 100), we excluded n = 2 due to failure to perform the task and n = 0 due to failure to pass headphone screen; attrition = 2%.

The 150 items represent 50 items in each of three frequency bins. David had normative data for this assessment so we were able to compare the mean ratings for the Prolific sample to the existing norms of the in-lab Hoosier participants. And look at that — mean ratings were *incredibly* similar between the two samples.

## Success 5: Word familiarity ratings

- Mean ratings by frequency category for the Prolific sample were very similar to existing norms, both by subjects and by items



$r = 0.80$

Rating (Prolific) — Rating (Norms)

**To achieve sample (n = 100), we excluded n = 2 due to failure to perform the task and n = 0 due to failure to pass headphone screen; attrition = 2%.**

Not only were norms similar across samples by subjects, but they also tracked closely by items.

Here are individual subject functions; all subjects show the expected frequency effect. One also sees robust, and sensible, individual variation. For example, E2.019 shows overall higher ratings than E2.020, but both show the expected frequency effect.

Here are the other 50 people in the sample; lovely individual patterns as well.

## Success 5: Word familiarity ratings

- Mean ratings by frequency category for the Prolific sample were very similar to existing norms, both by subjects and by items
- Ratings showed high split-half reliability



$r = 0.80$

Rating (B items) vs Rating (A items)

To achieve sample (n = 100), we excluded n = 2 due to failure to perform the task and n = 0 due to failure to pass headphone screen; attrition = 2%.
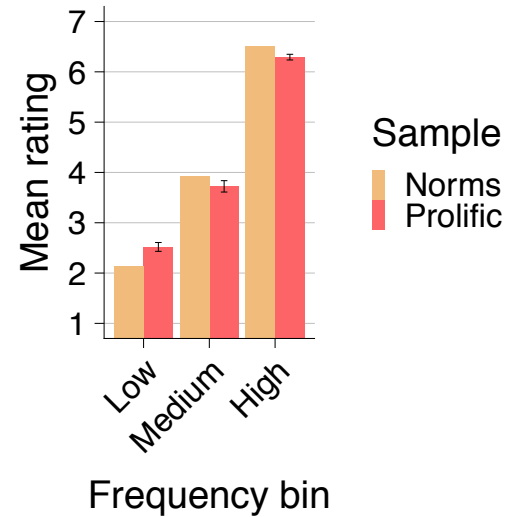
In this first sample of 100 participants, we observed incredibly high split-half reliability, which we used as motivation to try and develop an even briefer assessment.

## Success 5: Word familiarity ratings

- A second experiment was conducted to examine test-retest reliability
- 100 participants were tested in session 1; 85 returned for session 2
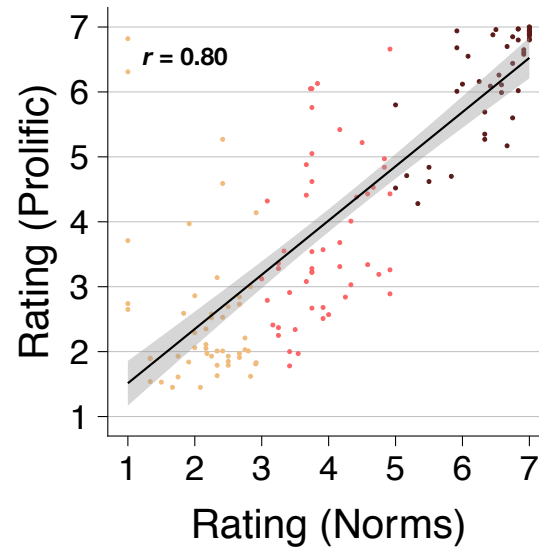- Mean completion time for the brief WordFAM versions was very quick



**To achieve sample (n = 85), we excluded n = 1 due to failure to perform the task and n = 0 due to failure to pass headphone screen; attrition = 2%.**

Specifically, a second experiment was conducted that included 85 participants who completed two brief versions of the WordFam test, separated by about two weeks in time. Mean completion time was around four minutes —

## Success 5: Word familiarity ratings

- A second experiment was conducted to examine test-retest reliability

- 100 participants were tested in session 1; 85 returned for session 2

- Mean completion time for the brief WordFAM versions was very quick

- Test-retest reliability was very high in the aggregate and by category
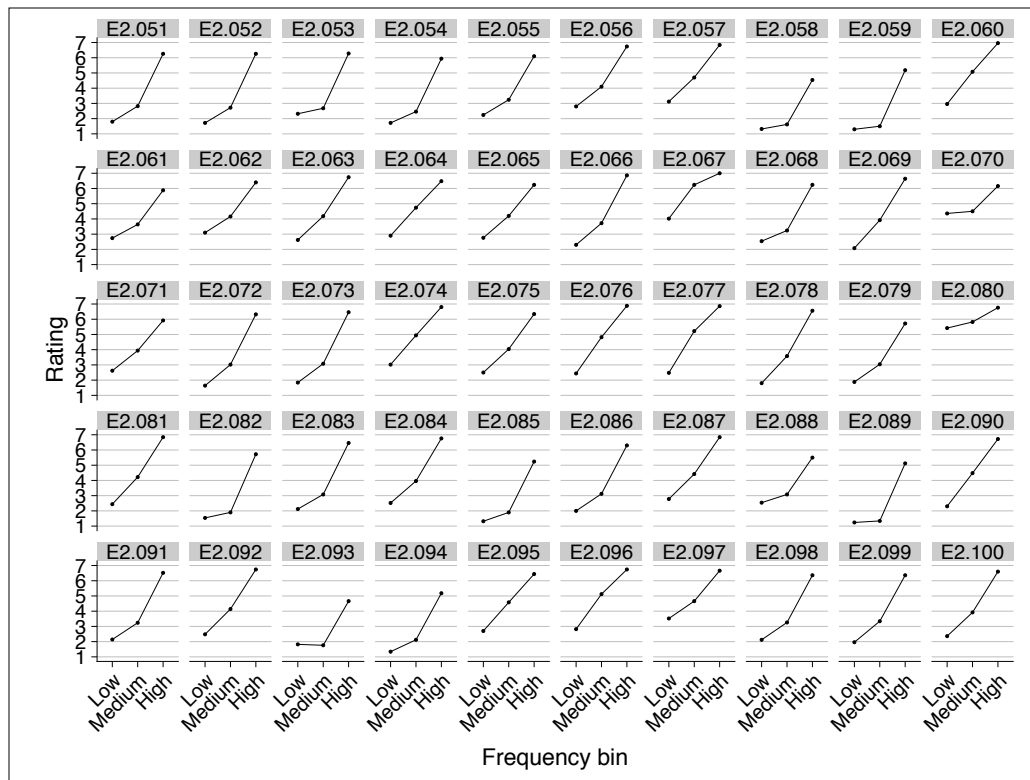


$r = 0.82$

WordFAM (Session 2) vs WordFAM (Session 1)

**To achieve sample (n = 85), we excluded n = 1 due to failure to perform the task and n = 0 due to failure to pass headphone screen; attrition = 2%.**

Test-retest reliability was incredibly high in the aggregate —

**Success 5: Word familiarity ratings**

| Low | Medium | High |

$r = 0.73$  $r = 0.80$  $r = 0.72$

WordFAM (Session 2)

WordFAM (Session 1)

And within each frequency category. These assessments, along with a second measure of vocabulary knowledge, will soon be available as Open Materials on Gorilla.

## Challenges

- With the methods I've described, you can't see your participants and (usually) can't answer questions in real time

- You have less control over the technology

- You have less control over the listening environment

These five success are just a sample of what's possible — what a time to be alive. That being said, it's also true that web-based speech perception studies are not without challenges given the loss of some control over both participant's specific hardware and the listening environment.

Challenges

- With the method                                    e

- 

- Y
  e

BuzzFeed

News    Buzz    Life    Entertainment    Quizzes    Videos

TOP POST
3,192,591 VIEWS

**12 Things That Will Absolutely Fix Almost Everything That's Wrong with Remote Experiments**

To offset these challenges, here's 12 tips and tricks that we've found useful for promoting high quality data.

## Tips and tricks

1. Be *exceptionally clear* with your participants in terms of technology requirements and study instructions

Be very clear with your participants.

**LDTN-005-d**

Hosted by *Rachel Theodore*

$1.67 • 10 minutes • $10.02/hr • 33 places remaining

The purpose of this study is to examine how listeners comprehend speech. You will be asked to listen to words and sounds and make decisions about what you hear. Then, you will be asked to fill out demographic information.

This study needs to be completed on a desktop or laptop while hearing headphones. **Wearing headphones is really important for this task. Any headphones or earbuds are fine so long as they deliver a stereo signal, meaning that different sounds can go to the left and right ears.** Participants who do not meet these requirements will be asked to return their submissions.

Auto-play for sound files must be enabled in your browser for the study to run.

Let them know what they need to do your study well. Define jargon, like stereo headphones. Tweaking our instructions to provide this definition, really helped to decrease our attrition rate.

You will see a central arrow appear on the screen.

** Press **a** as in *apple* if it is pointing left. **

** Press **l** as in *lemon* if it is pointing right. **

Ignore the arrows on either side, and just pay attention to the central arrow.

Please respond as quickly and accurately as possible. Keep your index fingers on top of the **a** and **l** keys as shown in the figure below to help make fast responses.

Press "Next" to see an example.

Next

Give guidance; be very clear in your instructions.

**This is the central arrow. It is pointing right, so you should press the "I" key. Press the "I" key now to continue.**

< < > < <

If your task is tricky, let them practice, like we did with the flanker task.

In this part, you will hear two tone sequences on each trial. Your job is decide if the two sequences are the same or if they are different.

Let me hear an example.

Or like we did for a sound discrimination task.

Here's an example where the two tone sequences are the same. You can listen to this example a few times .

▶ Play

Let me hear another example.

In this example, the two tone sequences are different. You can listen to this example a few times to hear the difference.

▶ Play

I'm ready to begin.

## Tips and tricks

1. Be *exceptionally clear* with your participants in terms of technology requirements and study instructions

2. Give people multiple chances to pass the headphone screen, *along with reminders* of the headphones requirement

Give people multiple chances to pass the headphone screen, with a reminder of the study requirements.

It's easy to set up a branch for this in Gorilla —

Stereo headphones were not detected.

As stated in the study description, you must be wearing headphones that deliver a stereo signal to do this study. This means that your headphones need to be able to send different sounds to your left and right ears.

If you have started this study by accident, then you are welcome to return your submission on Prolific without penalty.

My headphones are connected and I want to try again

And when we introduced this, our attrition due to lack of headphone compliance went down more than half.

1. Be *exceptionally clear* with your participants in terms of technology requirements and study instructions

2. Give people multiple chances to pass the headphone screen, *along with reminders* of the headphones requirement

3. Make sure any constraints set in Prolific and Gorilla are *mirrored* across systems

Be sure to mirror requirements across all systems you're using.

For example, if you set a constraint for computer only participation in Gorilla, but don't do that in Prolific too — then Prolific is going to send people to Gorilla only for Gorilla to reject them. This will lead to frustrated participants and a gazillion messages for you to respond to…
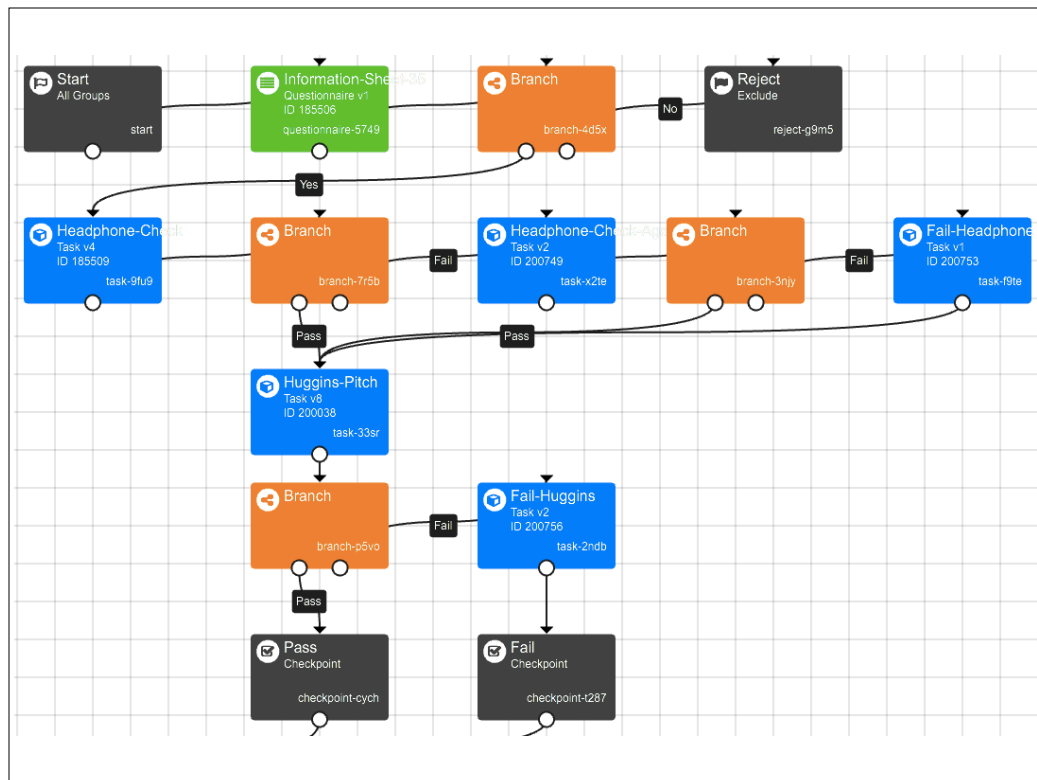
## Tips and tricks

1. Be *exceptionally clear* with your participants in terms of technology requirements and study instructions

2. Give people multiple chances to pass the headphone screen, *along with reminders* of the headphones requirement

3. Make sure any constraints set in Prolific and Gorilla are *mirrored* across systems

4. To decrease variability in reaction times, use within-subjects designs and provide a visual cue for hand placement

For reaction time studies, use within-subjects designs when you can, because the largest source of variability is going to come from different hardware/software set-ups across participants.

Showing this display for finger placement in RT studies not only led to faster RTs overall, but also less variable RTs.

## Tips and tricks

5. Sign up to be a participant on Prolific

Tip 5; join the Prolific pool as a participant yourself!

I've learned so much about how to make better studies that lead to a better participant experience by being a participant myself.

Tips and tricks

5. Sign up to be a participant on Prolific

6. Monitor and contribute to the Prolific subreddit:
   https://www.reddit.com/r/ProlificAc/new/

Related, monitor and contribute to the Prolific subreddit; I can't stress enough how much we can learn from our participants! You'll quickly tune in to what drives them crazy, what they enjoy — and then you can use this information to optimize your designs.

## Tips and tricks

5. Sign up to be a participant on Prolific

6. Monitor and contribute to the Prolific subreddit: https://www.reddit.com/r/ProlificAc/new/

7. Stay on top of your Prolific messages in real time

Prolific has a great messaging system for researchers and participants. Stay on top of those messages!

Just wanted to let you know that the audio did not work when I used Google Chrome, but it worked on another browser (Mozilla). Not sure if that was an issue just related to me/my computer, but I figured I would let you know, in case it's a problem for other people as well. Everything went well on Mozilla and was able to complete it.

11 Apr 2020, 12:02

Thanks so much for letting me know! We've vetted in on multiple browsers (including Chrome), but it could be a version issue. Or it could be that autoplay is not set-up in your Chrome browser but is in your Mozilla browser. Either way, I'm so glad that you could complete it, and am very grateful that you took the time to reach out to me. Thanks again for participating in our research! -Rachel

11 Apr 2020, 12:06

In my experience, participants are quick to report when something goes wrong.

That test made me question my hearing lol. I swear it sounded like goat was being said for most of the test. I was trying to ignore the effect the statement leading up to the word was. Naturally whenever it said something like "For it's safety, I caged the...." my brain wanted to automatically assume goat since caging a goat doesn't make any sense. The same with "I ironed the...". You'd naturally assumed you ironed a coat and not a goat.

2 May 2021, 14:59

What you experienced is exactly what we're trying to learn more about in this study! We're studing how listeners integrate the meaning of a person's message with how words are pronounced. Some participants get the meaning before the target word (e.g., For it's safety, she caged the ---) , and others get the meaning after the target word (e.g., The --- was caged for it's safety). Our prediction is that the meaning part will be more important than the actual pronunciation when the meaning comes before instead of after the target word.

Thanks for reaching out - and thanks so much for participating in our study. We couldn't do our research without you!

Rachel

AND — messaging with your participants is an excellent forum for science communication.

5. Sign up to be a participant on Prolific

6. Monitor and contribute to the Prolific subreddit: https://www.reddit.com/r/ProlificAc/new/

7. Stay on top of your Prolific messages in real time

8. Run a small sample through your experiment and *check everything* before running your full sample

9. Keep your tasks *as quick and as engaging* as you can; I highly recommend the **simr** package in R for power analyses

A few more tips and tricks — run a small sample through before you run the full sample. The only downside of being able to collect data from 300 participants in an hour is that one is also able to make a fatal mistake that affects 300 participants in an hour…

Keep your tasks as quick as you can while also ensuring adequate power. Gorilla has just released a game builder feature that I'm really excited about as a means to make our boring psychophysical tasks more engaging for participants, which will only benefit data quality.

## Tips and tricks

5. Sign up to be a participant on Prolific

6. Monitor and contribute to the Prolific subreddit: https://www.reddit.com/r/ProlificAc/new/

7. Stay on top of your Prolific messages in real time

8. Run a small sample through your experiment and *check everything* before running your full sample

9. Keep your tasks *as quick and as engaging* as you can; I highly recommend the **simr** package in R for power analyses

10. Use MP3 format instead of WAV for sound files

Tip 10; use MP3 format instead of WAV for sound files.

MP3 is native to browsers and you will run into glitches with some participants if you use WAV files. I know, I know, we've all been trained to avoid lossy formats. BUT, conversion algorithms are very good these days and I've yet to find anyone who can detect important missing information in our MP3 conversions. You can try yourself with the examples on our website. All of the studies I showed in this presentation used MP3 audio stimuli.

Tips and tricks

11. Calibrate expectations; technological glitches will occur, people will fail your headphone screen, you will get a low effort participant

12. Apply everything else you know about running great experiments to web-based testing; in-lab and web-based methods are more similar than different

Last two tips: Calibrate your expectations; you're going to have glitches, you're going to have a low effort participant; these things happen even in the laboratory. Look for them; and design tasks that make it easy to detect low effort responses.

And finally — don't forget to apply everything else you know about running great experiments; in-lab and web-based methods are more similar than different.
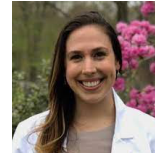
## Acknowledgements



Nikole Giovannone    Lee Drown    Julia Drouin    Nick Monto

David Pisoni    Lynne Nygaard    Christian Stilp    Christina Tzeng

LabPhon18

slaplab.uconn.edu

@rachelmtheodore

osf.io/8krg3/

*For more in-depth resources and a transcript, visit:*
https://slaplab.uconn.edu

Additional resources are available on our website and OSF page; and please don't hesitate to reach out to me offline if you have questions. Thank you very much.