# *V(is)owel*: An Interactive Vowel Chart to Understand What Makes Visual Pronunciation Effective in Second Language Learning

**Charlotte Kiesel**
yoder6@illinois.edu
Univ. of Illinois in Urbana-Champaign

**Dipayan Mukherjee**
dipayan2@illinois.edu
Univ. of Illinois in Urbana-Champaign

**Mark Hasegawa-Johnson**
jhasegaw@illinois.edu
Univ. of Illinois in Urbana-Champaign

**Karrie Karahalios**
kkarahal@illinois.edu
Univ. of Illinois in Urbana-Champaign

## Abstract

Visual feedback speeds up learners' improvement of pronunciation in a second language. The visual combined with audio allows speakers to see sounds and differences in pronunciation that they are unable to hear. Prior studies have tested different visual methods for improving pronunciation, however, we do not have conclusive understanding of what aspects of the visualizations contributed to improvements. Based on previous work, we created *V(is)owel*, an interactive vowel chart. Vowel charts provide actionable feedback by directly mapping physical tongue movement onto a chart. We compared *V(is)owel* with an auditory-only method to explore how learners parse visual and auditory feedback to understand how and why visual feedback is effective for pronunciation improvement. The findings suggest that designers should include explicit anatomical feedback that directly maps onto physical movement for phonetically untrained learners. Furthermore, visual feedback has the potential to motivate more practice since all eight of the participants cited using the visuals as a goal with *V(is)owel* versus relying on their own judgment with audio alone. Their statements are backed up by all participants practicing words with *V(is)owel* more than with audio-only. Our results indicate that *V(is)owel* is effective at providing actionable feedback, demonstrating the potential of visual feedback methods in second language learning.

## CCS Concepts

• **Human-centered computing** → **Empirical studies in visualization**; • **Applied computing** → *Sound and music computing*; **Computer-assisted instruction**.

## Keywords

Visual Feedback, Computer-Assisted Pronunciation Training, User Evaluation, Phonetics
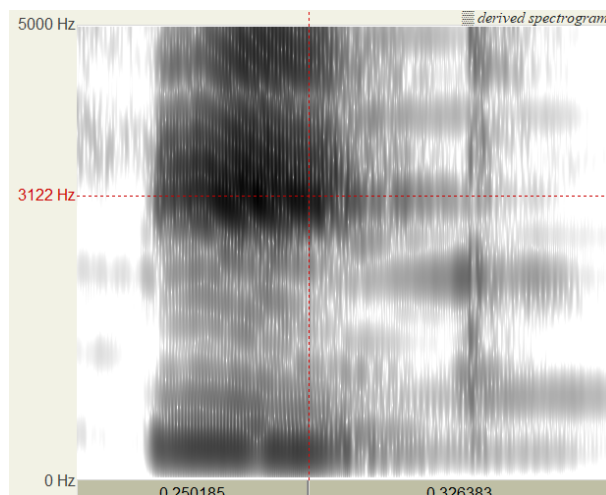
## 1 Introduction

Every speaker comes with their own language background and physical characteristics, thus, pronunciation feedback must be personalized for second language learners. Individualizing feedback requires more time from language teachers in a system already strained by the lack of instructors [20, 23, 24, 29]. Computer-Assisted Pronunciation Training (CAPT) has the potential to assist learners by providing feedback specific to their needs. Although CAPT introduces the potential for algorithmic error and requires access

to computing resources, it addresses the drawbacks of a teacher-dependent model. It does not require extra time from a teacher, is widely distributable, and provides access to practice in a self-paced environment. Among CAPT methods, visual feedback has arisen as an effective way to provide pronunciation feedback [9, 18, 25, 40].
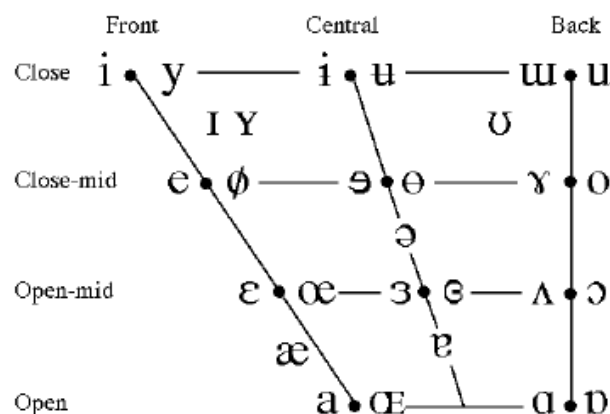
There are two main types of visual feedback for pronunciation: correctness indicators, such as color coding or percentages, and articulation-based representation of speech. Many well-known language learning apps, such as Duolingo, babbel, and Rosetta Stone, use correctness-based visual feedback, while previous research focuses on articulation-based visualizations (ABV) [9, 18, 25]. We choose an ABV because it satisfies all the criteria in Bliss et. al. [4] for effective feedback. Correctness-based feedback fails two of the criteria; namely, feedback must be (i) natural and logical and (ii) able to facilitate comparison [4]. Since color coding and percentages do not specify what about the production is incorrect, correctness-based feedback fails to be natural and logical. It fails the latter by forcing learners to use their ears to determine the differences instead of visually representing comparisons.

Three widely studied visual feedback methods are the spectrogram, animated face-cutaways, and the vowel chart. The spectrogram, which can be used for both vowels and consonants, is difficult to interpret without training and hard to modify pronunciation based on its feedback (Figure 1a). Animated face-cutaways and vowel charts, which provide feedback on a subset of sounds, have fewer parts to interpret and uses the principle of natural mapping to visualizes articulator movement (Figure 1b). The vowel chart follows more closely the principles of effective feedback and is effective in improving pronunciation [9, 30]. Previous implementation of vowel charts, however, do not fully integrate audio and visual. Users had to adjust pronunciation based either on visual or audio output. For this reason, we created an interactive version, *V(is)owel*. The chart displays the position of the tongue based on the first and second formants, prominent bands of energy based around a frequency [31]. *V(is)owel* provides broad applicability by accepting vowels within multiple contexts including rhotics, like the American /r/[1] and displays a vowel's change over time, allowing learners to train the difference between single vowels (monophthongs) and combined vowels (diphthongs). As far as we are aware, we are the first study to have vowels in multiple contexts, visuals for diphthongs, and **the first study in visual Computer-Assisted**

---

[1]The forward slash (/) represents the orthographic character in a language. Square brackets around a character indicate a phonetic description.
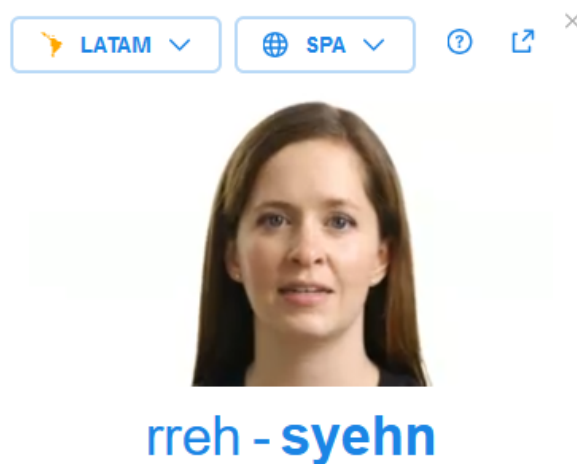
(a) A spectrogram from Praat of the English word "Heat".



(b) An example vowel chart [39]. The x-axis represents the second formant and frontedness or backedness of the tongue. The y-axis is the first formant and height of the tongue in the mouth.



(c) A screenshot of Duolingo's binary speaking feedback as of 2020 [3]



(d) A still frame from a video of a speaker's face on Spanish-Dict.com.

Figure 1: Four types of pronunciation visualizations.

**Pronunciation Training (CAPT) to capture how learners interpret visual results as they adjust their pronunciation**.

Using an iterative design approach, we created *V(is)owel* and an audio-only practice tool. We ran a within-subject study with 8 participants and elicited speakers' real-time thought processes as they used both tools.

*V(is)owel* encouraged participants to practice more because it gave them a goal to work toward or pushed them to reflect on the first language speaker's production and their own when visual results did not align with their expectations. Based on our findings, we suggest that future pronunciation feedback should include a visual representation of the closeness of the learner and target speaker.

Our work makes three main contributions.

(1) Learners' thought processes during interaction with visual pronunciation feedback.
(2) An interactive vowel chart, *V(is)owel*, that can be used in any word context.
(3) A technical implementation for extracting vowels from any context.

## 2  Related Work

Pronunciation in second language acquisition (SLA) is interconnected with the building blocks of language– grammar, vocabulary, and listening comprehension– as well as integration into a language community. Thus, improving pronunciation leads to improvement in other areas of SLA and vice versa [13, 33]. There are multiple ways that a first language (L1) of a speaker, usually a language spoken from a very young age, can interfere with acquiring sounds in a second language (L2). Sounds that distinguish meaning in an L1, called phonemes, may not distinguish meaning in an L2 or vice versa. As a result, interlocutors can be confused when communicating with an L2 speaker [29]. We discuss previous work that focuses on human- and computer- assisted methods for pronunciation improvement.

### 2.1  Human-Mediated Pronunciation Training

Little attention is given to pronunciation in textbooks and curricula other than at a collegiate level in the United States of America. This makes ad hoc reactive feedback for pronunciation common in foreign language classrooms from K-12 [11, 29, 38]. Learners who receive ad hoc reactive feedback do not improve significantly in comparison with intentional methods of training. As a result, we leave this out of the discussion despite it being a common approach in classrooms [11]. Beyond reactive feedback, three main methods of pronunciation instructor-dependent training are outlined below: reading out loud, mimicry of an L1 speaker, and minimal pair drills (drilling words that differ by one sound/phoneme[2]).

*2.1.1  Reading Out Loud*  A commonly studied method of pronunciation improvement is reading out loud [10, 37]. This can be achieved through paired practice, where peers listen and give feedback, mimicking practice, or group reading, where the whole classroom and teacher listen to the reader and provide feedback on pronunciation. The exercise improves oral reading fluency and reading comprehension [16], but is dependent on a teacher's input or practice with peers. Because feedback is given in real-time, L2 learners will not receive detailed pronunciation feedback.

*2.1.2  Mimicry of an L1 Speaker*  As mentioned in section 2.1.1, mimicry can be used in conjunction with reading out loud. It can also be used for shorter segments of pronunciation, such as words or sounds, which is the method we use in our study design. Research shows that mimicry is key in improving pronunciation [19, 35]. While mimicry is used in classrooms, it does not require an active listener. L2 learners can practice mimicry with any kind of medium that includes audio, such as movies or podcasts, though this means they have to rely on their own judgment to adjust pronunciation.

*2.1.3  Minimal Pair Drills*  Minimal pair drills (MPDs) rely on listening and mimicking an L1 speaker to bring out the contrast that different sounds have by putting them in identical contexts. MPDs are used in speech therapy and language learning to great effect though they limit the context that sounds can appear based on the L1 [1, 7, 15, 28]. Despite potential downsides of minimal pairs, we contend that visually contrasting vowels will benefit the L2 learner.

---

[2]For example, hat/hit, fit/feet, cat/bat are minimal pairs in English. Hat and catch are not since they differ by more than one sound

### 2.2  Computer-Assisted Pronunciation Training (CAPT)

CAPT has the advantage over traditional pronunciation feedback by its ability to be used by many and be automatically catered to the individual. CAPT feedback takes two different forms: audio-only and visual accompanied by audio. The latter outperforms traditional and audio-only feedback [18, 26]. Visual feedback falls into two main categories: spectrograms and articulatory-based visualizations.

*2.2.1  Spectrograms*  Spectrograms are useful for their flexibility in pronunciation feedback. They can be used to train vowels and consonants including the effects those sounds have on each other (Figure 1a). Studies have shown the positive impact on pronunciation when students reflect on their speech with the aid of a spectrogram. But a spectrogram can require hours of training to correctly interpret and previous studies use it as a static tool instead of an interactive one [14], likely because it is difficult to apply feedback. Using it as a static tool limits learners to reflect on their pronunciation without immediately adjusting their pronunciation.

*2.2.2  Articulatory-based Visualizations*  Articulatory-based visualizations, on the other hand, are used in phonetic training to teach learners about articulators, such as the tongue and lips, in many different languages. Previous work focuses on two ABV: face-cutaway animations and vowel charts.

Animated face-cutaway diagrams display how articulators move during speech. Learners can view correct movement based on the animation and attempt to mimic the example speaker [8]. However, this work does not display an animation of a learner's incorrect speech or how to change pronunciation from incorrect to be correct [8]. As a result, feedback targets mistakes but does not provide feedback to move from current pronunciation to the target pronunciation.

Vowel charts display a vowel's location on a trapezoid. There are two ways to represent vowels on a vowel chart – with IPA symbols (Figure 1b) or with formants (Figure 4). Previous studies show that vowel charts are effective within limited contexts in improving pronunciation [27, 30]. One method studied vowels by themselves while another limited it to vowel preceded by an 'h' and followed by a 'd' [27, 30]. Both restrictions limit the real-world use of the tool. Furthermore, how learners understand what they see on an interactive vowel chart is unclear since all the studies interview participants in a reflexive way, making it difficult to know how to improve the interaction experience to facilitate a better understanding of a vowel chart. Finally, these vowel charts do not integrate the audio with the visual, requiring learners to listen and interpret visuals separately [27, 30].

We designed a real-time interactive vowel chart system, *V(is)owel*, that would function for vowels in any context and provide integrated audio and visual feedback. To understand why the visualization is effective, we held think-aloud experiments and compared how learners think about visual feedback in contrast to audio-only feedback. By understanding how untrained learners interpret visual feedback, we can modify existing visualizations to promote improvement and minimize confusion.
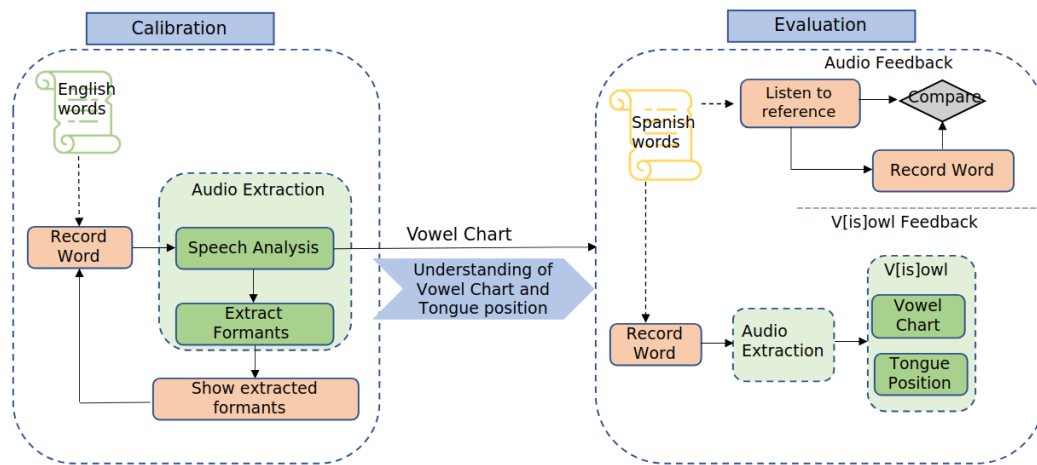
Figure 2: The study begins by calibrating a user's voice using extracted formants of four English vowels [i,u,a,ɒ]. After calibration, users interact with an audio-only feedback system and *V(is)owel*, which extracts formants and displays them on a vowel chart.

## 2.3 Research Questions

As we seek to understand what makes a visualization effective, we need to make sure we view it holistically. If a visualization is effective in improving the target sounds, such as vowels, but does not affect or negatively affects overall pronunciation, we should consider visualizations that communicate information for more sounds. In short, to understand what makes visualizations effective and if they limit speakers during practice, we seek to answer the following questions:

**RQ1:** How does interaction between a visual and audio-only feedback method differ?

**RQ2:** How do phonetically untrained users interpret vowel charts during interaction?

**RQ3:** How does a visualization that focuses on vowels affect users' perceptions of other aspects of pronunciation?
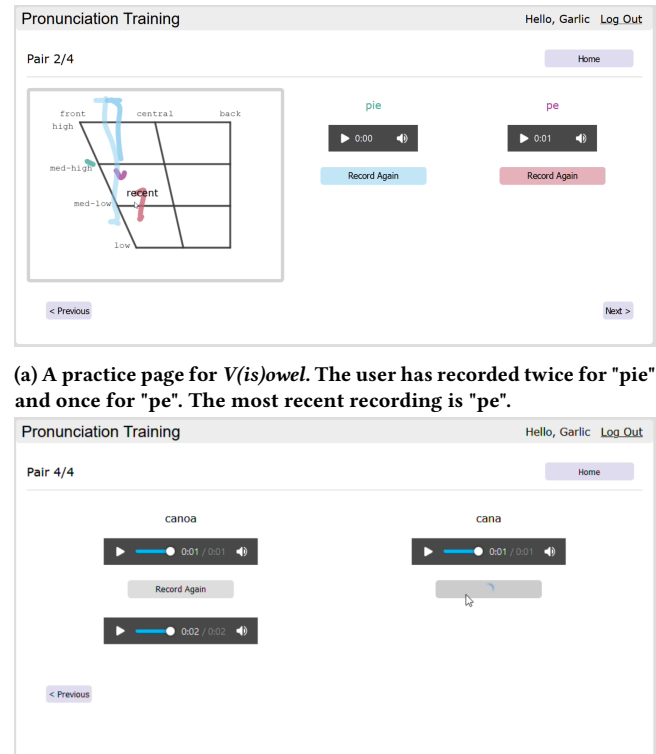
## 3 Design

In this section, we outline the design of our interactive vowel chart, *V(is)owel*. We discuss the interface of *V(is)owel* and its calibration and tutorial, then we present our design of *V(is)owel*'s signal processing backend.

## 3.1 User Interface

*V(is)owel* is comprised of a vowel chart, which displays the Spanish speaker's and user's vowels, recording buttons, and audio playback button for the Spanish speaker. When any vowel is clicked, including the Spanish speaker's, audio plays. In order for *V(is)owel* to display two speakers on one chart, the user must first have their voice calibrated. Once their voice is calibrated, users are able to complete a tutorial introducing them to *V(is)owel* where they practice with English words and Spanish words to begin to build mental models of how it works.

*3.1.1 V(is)owel* The final design of *V(is)owel* is trapezoidal shaped and has descriptive labels on the top and left sides, e.g. front, mid,



(a) A practice page for *V(is)owel*. The user has recorded twice for "pie" and once for "pe". The most recent recording is "pe".



(b) A practice page for audio-only. The user is waiting to record "cana", as shown by the semi-circle in the record button.

Figure 3: Example practice pages from audio-only and *V(is)owel*.

back, etc. Both of these design decisions were inspired by a traditional vowel chart (Figure 1b).
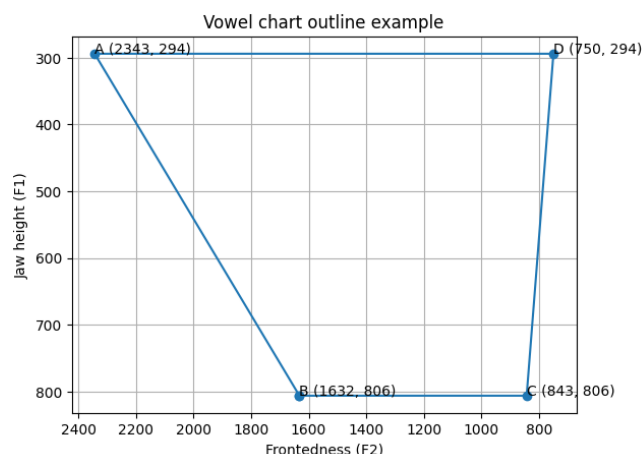
Figure 4: A vowel chart with first and second formant frequency markings.

**Vowel lines.** Vowels are plotted on the chart based on their change in frequencies over time. Clicking on a plotted vowel plays the associated audio for the extracted vowel then the whole word. If the vowel moved over time, an animation redraws the line starting from the beginning. Initially, we attached a simple arrow on the end of the line to indicate vowel change direction (Figure 6). The arrowhead helped, but due to changes in formants, it would sometimes point the wrong direction on the line. In the final version, we removed the arrowhead and added an animation of the line when it was clicked to communicate change over time.

The vowel trajectory is based on averaging neighboring formants. We began by directly plotting all extracted formant values but since extracted formants vary slightly, this created noisy lines. Because the difference in neighboring formants is due to slightly inaccurate formant extraction and not a speaker's tongue movement, we decreased the visual noise by averaging neighboring formants.

**Multiple recordings.** A user can record words as many times as they would like, however, only the most ten most recent recordings are displayed on the chart. Each time a new vowel is plotted on the chart, the previous vowels become more transparent. In the beginning phases of design, we did not limit the number of vowels on the chart nor the number of vowels extracted from the word (Figure 5). The mock-up tested well, but practically, plotting even one vowel proved visually overwhelming when the vowel was long, did not stay in the same place, or was recorded more than a few times. Many recordings created visual clutter on the chart, so we limited the total number of recordings to ten per word and use opacity to communicate the ordering of recordings.

**Grid lines and labels.** In our first design, we followed previous work and marked frequencies along the axes (Figure 4) [27, 30]. Discussion with the pretest group showed that people were unsure how to interpret vowels marked on the chart. Our research team decided to strip a vowel chart from all labels, then reintroduce labels one at a time to gauge understanding of the chart based on available labels. Frequency labels did not allow group members to understand what that meant with respect to their physical movements
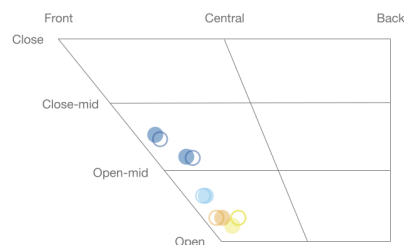


Figure 5: A mock-up of two vowels from one word plotted on a vowel chart. The first vowel's dot is colored in and the second vowel's dot is outlined.
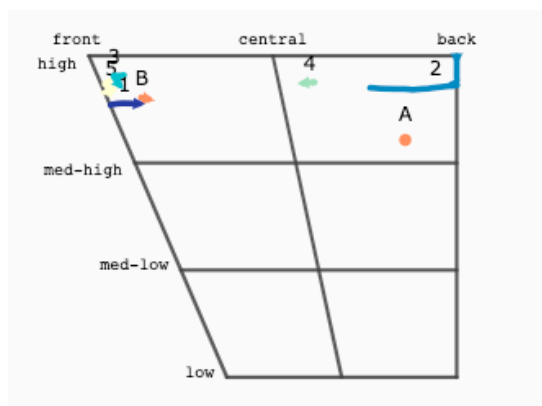


Figure 6: *V(is)owel* with Spanish vowels marked by A and B and user's vowels labeled numerically with the latest recording starting at 1.

while physical labels (front, back, high, low) were more helpful in interpretation.

Although ***V(is)owel*** has bounding grid lines, vowels are not constrained to them. When we forced vowels within the trapezoid, crucial information was lost regarding the current tongue position with respect to previous recordings (Figure 6). For example, the vowel "ea" ([i][3]) in the English word "beat", could be plotted as slightly above or further forward than the previous recordings, which would put it outside the chart. Removing this information would inaccurately represent the tongue's position as the same between two recordings even though it had moved. As a result, we let vowels be plotted outside the chart. We kept vowels within a wider box around the grid lines, so that even if our extraction algorithm failed to pick out the correct frequencies for the formant, users would have visual feedback.

**Calibration** User voices can vary significantly depending on the individual. This variety will change the shape of the vowel space per individual, as shown by the left image in Figure 7. We transform the individual vowel space into a fixed trapezoid, as shown by the right image in Figure 7, which removes idiosyncratic speech characteristics and retains language-dependent features. We collect

---

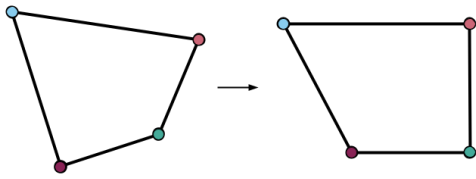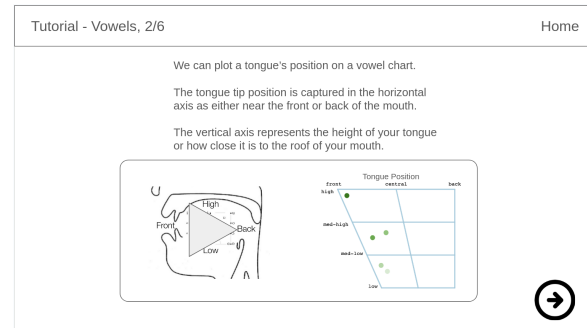[3]Square brackets around a character indicate a phonetic description.

**Figure 7: On the left is an example frequency shaped trapezoid using corner vowels. We use a projective transformation to map the observed formant frequencies and to a consistent vowel chart representation.**

the user's specific vowel space using their elicitation of the English "corner" vowels, [i,u,æ,ɒ(ɕ)] extracted from the words beet, boot, bat, and bought. We use projective transformation (homography) to map the user-specific frequencies to our vowel chart. Since there is no set ratio for vowel chart trapezoids, we chose the trapezoid ratio for height, top width, and bottom width as 1, 1.33, and 0.833 respectively. The height represents the spread of the first formant ($f_1$) frequency, and the width represents the spread of the second format at different $f_1$. This shape is consistent with the known spread of the vowel chart. The consistent shape of the vowel chart enhances the aesthetics of the chart and the interpretability of results across users.
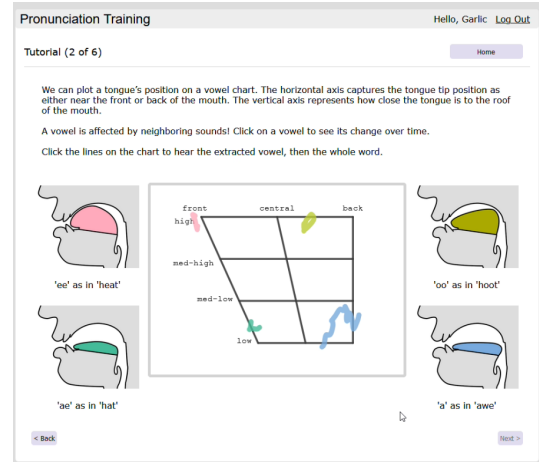
**Tutorial** The tutorial for ***V(is)owel*** leads learners through understanding a vowel chart. It introduces the tongue's position for producing vowels, how to understand the axes of the vowel chart, and finally, how to mimic a speaker based on their plotted vowel. We settled on this flow based on pretest user interaction with paper and interactive prototypes (details in Appendix F). We begin with informing users of the tongue's role in producing a vowel, specifically focusing on closeness to the roof of the mouth, tongue height, and where it is in the mouth, close to the teeth or away. We started with the tongue's position, since pretest users were not sure how the tongue's position related to vowel production. After an overview of the chart in relation to tongue position, we give them opportunities to focus on pairs of English vowels that differ in tongue position and height. By starting with known vowels, we were able to bake interaction into the examples and avoid a lot of written instruction, since pretest users often skipped reading.

## 3.2 Technical Backend

We discuss the design of the signal processing algorithm and its accuracy based on pre-study results. The algorithm extracts the vowel and vowel formants, then transforms the formants into the html viewbox according to a speaker's transformation (Figure 11). An audio file with a single word is passed as a signal to the vowel extractor. The vowel extraction uses the autocorrelation of the signal to determine the onset and offset times of the vowel. A ".wav" file is created using these timestamps. The new file is passed into formant extraction which tests different numbers of formants and calculates a goodness metric. The algorithm extracts $f_1$ and $f_2$ values using the formant with the best metric. The lists of $f_1$ and $f_2$ are transformed into the viewbox then displayed on the vowel chart.



**(a) The second page of the paper prototype. One face cutaway is on the left and on the right a vowel chart has multiple different vowels plotted on it.**



**(b) Final interface for the tutorial's second page. The vowels on the chart share the same color as the tongue in their respective face cutaway.**

**Figure 8: Design cycle of page 2**

*3.2.1 Vowel Extraction* We designed a vowel extraction algorithm based on the autocorrelation of the signal after testing two existing algorithms that did not have the precision needed, Dynamic Time Warping (DTW) [2] and PocketSphinx Phoneme Aligner (PS) [21], as shown in Table 1. We tested the algorithms on 5 voices from our lab, 3 female. Even though DTW generally had a lower absolute difference, it tended to mark the beginning of the vowel during the consonant, as represented by the negative. We determined the precision needed through trial and error. Starting 20 milliseconds after the expected beginning of the vowel will not miss diphthong information but starting more than 10ms before the start of the vowel will pick up on too much of the preceding consonant vice versa for vowel offset.

Our novel signal processing algorithm extracts the beginning and end of the vowel by calculating the degree of autocorrelation within a sliding window, where autocorrelation is computed as follows:
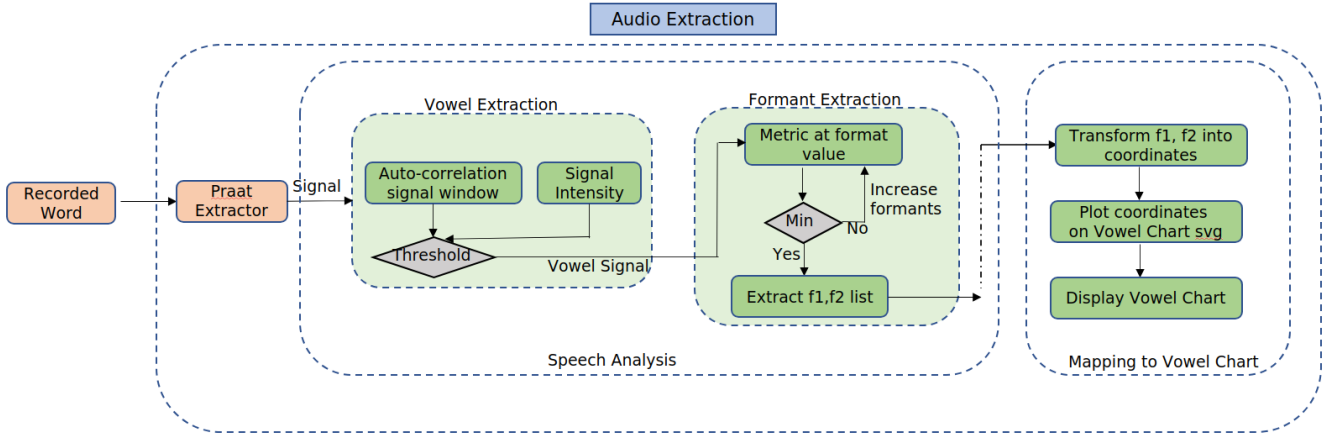
Figure 9: A flowchart representing the algorithmic process of extracting vowels and formants.
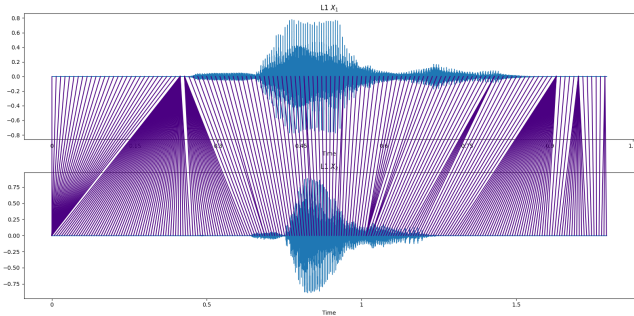


Figure 10: An example graph we made demonstrating the Dynamic Time Warping (DTW) algorithm [2] which matches similar parts of two different signals. The purple lines represent the connections DTW made.

$$c_k = \sum_n v_{n+k} \cdot \bar{v}_n$$

where $v$ is the audio signal and $\bar{v}_n$ is the complex conjugate. Once the autocorrelation goes over a predefined threshold and then back under, the code marks the timestamps (Figure 11). We set the threshold by experimenting on different words and speakers to find the degree of autocorrelation needed to be considered a vowel. This algorithm produced the most accurate timestamps for the beginning of the vowel though it still produced inaccurate endings based on whether the following consonant was also fairly resonant.

| Algorithm | Average diff. to Manual Calculation |
|-----------|-------------------------------------|
| DTW       | -8.2ms                              |
| sphinx    | 25.4ms                              |
| Ours      | 13.4ms                              |

Table 1: The average difference in onset times for three vowel extraction algorithms. We ended with ours because it was better to lag than lead the beginning of the vowel



Figure 11: An example degree of autocorrelation for the word 'bata' of a signal with the threshold marked horizontally in red. The two peaks are the vowels.

*3.2.2 Formant Extraction.* We use Praat's [5] format extraction library to retrieve the first two formants. Praat's method requires a frequency ceiling and the number of formants. We set the frequency ceiling at 5500 Hz. However, since the number of formants in a signal changes based on the word, not the speaker, we designed a metric for a signal, described below, to measure the best formant for the signal. We need the first two formants to be smooth because when the algorithm incorrectly chooses a formant, the formant tends to jump multiple frequencies between points.

| ID | Word List | 1st Tool | 2nd Tool | Sex |
|----|-----------|----------|----------|-----|
| P2 | first | Control | Vowel | Male |
| P3 | second | Vowel | Control | Female |
| P4 | second | Control | Vowel | Female |
| P5 | first | Control | Vowel | Male |
| P6 | second | Control | Vowel | Male |
| P7 | first | Vowel | Control | Female |
| P8 | second | Vowel | Control | Male |
| P9 | second | Vowel | Control | Female |

**Table 2: Participant information**

$$smoothness(sig) = stddev(sig)$$

$$metric(sig, fmt\_cnt) = smoothness(sig_{f1}) + smoothness(sig_{f2})$$

$$+ \sum_{i=3}^{fmt\_cnt} 0.8^{i-2} smoothness(sig_{fi})$$

Based on this metric, we choose the formant number which has the lowest metric value. It ensures that the first two formants are consistent for the signal, and the higher formants are also somewhat stable. The stability of the formant ensures that our algorithm has found the correct frequency for the first and second formant.

## 4 User Study Methodology

We held a within-subject study with 8 participants at a large Midwestern university in the United States of America during December of 2024 (Figure 12). In this section, we describe the process of recruitment and human-subject evaluation for *V(is)owel* and audio-only pronunciation practice. We ran two pilot sessions before our study, thus participant numbers start from 2.

### 4.1 Recruitment

We recruited participants through the campus e-newsletter, and physical posters throughout campus buildings (See Table 2). The participants were between the ages of 18 and 24, had no history of speech or auditory impairment, and spoke American English as their first language. We limited participants to those who had not spent more than a semester in college learning Spanish, a week studying Spanish in a Spanish-speaking country, or two months in a Spanish-speaking country. We asked participants to choose the sex of the Spanish speaker with which they wished to practice.

### 4.2 Interface Evaluation

To understand how participants process their pronunciation during practice with and without a visualization, we conducted one-hour think-aloud experiments in-person on campus. During the practice portion of the experiment, we often prompted them to express why they decided to record a second time. We randomized the order in which participants saw the two practices, *V(is)owel* and audio-only (Figure 3). Since we wanted to test the robustness of our signal processing algorithm and visualization, participants practiced with words that included vowels in varied linguistic contexts. These lists were also counterbalanced to appear with both visualizations.

### 4.3 Experimental Setup

Figure 12 shows the flow of our study. We collected consent, gave a brief introduction to the project, and had participants record English words for a baseline task load (NASA-TLX). Depending on which condition they saw first, the next step was interaction with *V(is)owel* or audio-only. Audio-only's interface was identical to *V(is)owel* without a vowel chart (Figure 3b). The tutorial for audio-only was a spoken introduction to its mimicry setup. Other than an initial calibration step, interactions with the tools followed the same steps: a tutorial introducing the practice, practice with the tool with four different Spanish minimal pairs, followed by a NASA-TLX and System Usability Score (SUS) surveys. After interacting with both tools, we held a semi-structured exit interview.

We ran the tools locally on a Lenovo laptop with a Bluetooth microphone recording at 48kHz. The microphone was set on a 6″ stand. Participants were compensated for their time and effort by a $20 Amazon gift card. They viewed the interface either on a monitor or laptop screen and used a mouse to navigate the interface.

After explaining and signing the consent form, we briefly introduced the goal of the project. We used NASA TLX to measure workload [17], which measures the workload based on mental, physical, and temporal demand as well as performance, effort, and frustration while doing a task. We took a baseline workload assessment to mitigate potential effects of participant, meeting time, location, and screen difference. Participants were asked to separately record three English words as the baseline assessment. We used the process as a baseline since it would show how difficult participants found it to record and speak on our platform without the mingled effect of practicing pronunciation. Then, we introduced the tools to the participants and asked that they let us know if they had any questions before they began. The audio-only introduction was verbal, while the introduction to *V(is)owel* began with calibration and a 6-step tutorial. After the tutorial, they took a 3-question vowel-chart literacy check. We used the check to allow participants to apply what they'd learned and ask clarification questions about the visualization.

We wrapped up the experiment with a semi-structured survey about their experience using both pronunciation practice tools.

### 4.4 Data Analysis

One person on the research team transcribed the interviews and coded the discussions. In the first round of coding, they marked statements as either relating to pronunciation, visual information, auditory feedback, or reaction to the system. During the second round, they broke up the main codes into sub-codes, which are presented in results. We discovered two main themes during analysis of the coded interview data: the role of **auditory feedback** in understanding speech with and without visual feedback and how **visual feedback** can help and confuse learners in their goal to improve pronunciation.

## 5 Results

To answer our research questions, we must look at (i) the difference in thought processes between *V(is)owel* and audio-only practices, (ii) how and what users considered during interpretation of the
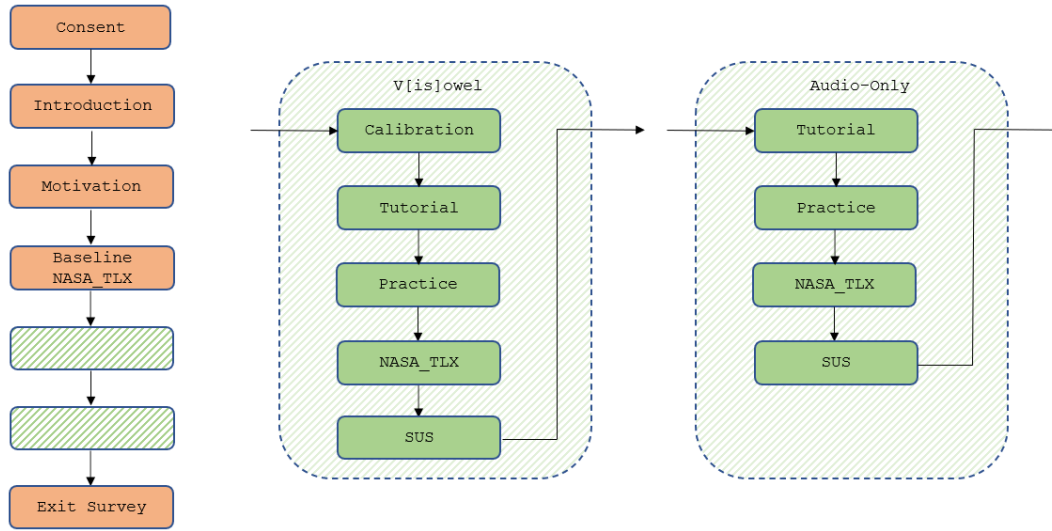
**Figure 12: A flowchart representing the progress of the study. The orange boxes represent the fixed portions of the study, the two green striped boxes are the counterbalanced pronunciation practice tools, *V(is)owel* and audio-only.**

vowel chart, and finally, (iii) the effects on participants' focus on pronunciation as a whole during *V(is)owel* interaction.

With this in mind, we present our results in two parts: first, with the themes drawn from the qualitative analysis of interaction and interview data then present quantitative findings; second, the technical evaluation of our signal processing algorithm. Quotes from participants are edited for clarity.

Our qualitative results start with participants' reactions to audio feedback since it was present in both practices. The differences and similarities in statements between the two practices answer (i) RQ1. After audio feedback comments, we present results related to visual feedback. The interactions answer all three RQs by showing how participants process visual pronunciation data previously unknown to them.

The front end results conclude with quantitative findings: the number of times participants recorded with each tool, NASA-TLX, and SUS. Each of the findings help complete the picture painted with qualitative results.

In general, participants enjoy vowel charts but often do not understand how to interpret visual results or successfully change their pronunciation to affect the visualization. They talked about their vowels more, but did mention other aspects of pronunciation, causing results for RQ2 to be inconclusive. In general, participants engaged with *V(is)owel* more and appreciated the guidance it gave including the change of a vowel over time in comparison with the audio-only practice due to the reliance on self-perception of audio.

### 5.1 User Study: During Practice

*5.1.1 Auditory Feedback* We found four main categories of auditory feedback that participants considered while practicing pronunciation. They brought up similarities to English in words or sounds and connected the sound to how they felt or believed their mouth moved. Linguistic categories, vowels, consonants, and pitch, also impacted participants' understanding of their speech. Finally,

they brought up uncertainty regarding the success of pronunciation based on auditory feedback with or without visual feedback.

**Similarities to English (3/8)**   Two of the participants used the similarity of Spanish words to English words as a reference for practicing Spanish. P2 focused on the sounds in Spanish in relation to English sounds. "'Loa' wasn't too difficult either. Just like boa constrictor, I think." They did not bring up the relation to English during *V(is)owel* interaction. P3 and P4 mentioned their pronunciation in relation to English three times during both practices. In both the *V(is)owel* and audio-only practice they mentioned that they felt some of their pronunciations had a strong American accent. These comments begin to answer RQ1 by suggesting that learners consider L1 speech patterns during interaction regardless of represented visuals of those parts of speech.

**Mouth Movement (4/8)**   The connection between mouth movement and subsequent pronunciation was consistent for both conditions. Four of the eight participants brought up how their mouth affected what they were hearing. Half saw the audio-only condition first (N=2). Three participants connected the sound to their mouth movement during interaction with *V(is)owel* and two mentioned it during interaction with the audio-only feedback, one who saw the audio-only first. Those who mentioned it during interaction with *V(is)owel* said that they were trying to connect the sound that caused the line to show up on the chart to the tongue position in the mouth (N=3).

> *I'm trying to replicate the position on the chart and think about how the sound is going to be connected to, you know, the tongue position and everything.* (P9)

Based on these responses, it is possible that seeing a visualization that provides feedback on tongue position brings attention to articulators like the tongue than feedback that does not include that information.

**Aspects of Sound (3/8)**    The participants also discussed different sounds in their speech: vowels, consonants, pitch, and emphasis. Two participants mentioned thinking about consonants, one during interaction with audio-only (P9) and the other during *V(is)owel* (P7).

> *It kind of sounds like the speaker is saying* [ðɐ] *instead of* [dɐ] *and I was saying* [dɐ] *instead of* [ðɐ]. (P7)

Without visual feedback, participants considered overall pronunciation such as pitch and emphasis during practice (N=4) while the vowel-specific visualization caused them to focus on the difference vowels (N=3). Only one participant, who saw *V(is)owel* first, mentioned the vowels during the audio-only portion. Participants compared their pronunciations to that of the Spanish speaker's (N=3) and reflected on modifying their pronunciation based on the Spanish speaker's pronunciation (N=1) with *V(is)owel*.

> *There's more of an* [ɛ] *to the speaker's sound and there's more of an* [a] *to mine.* (P8)

These results begin to answer RQ1 and RQ2 by indicating that a visualization focused on vowels will focus learners on but does not seem to limit their attention to them. On the other hand, participants without visual feedback to guide them did not consider segmentals, other than one who may have been primed by experiencing *V(is)owel* first.

**Uncertainty During Interaction (4/8)**    Four participants mentioned uncertainty during practice with audio-only and/or *V(is)owel*. While interacting with *V(is)owel*, participants mentioned feeling uncertain about some aspect of their pronunciation without being able to articulate what was causing their uncertainty (N=3). Most of the comments expressed doubt regarding how the vowel was plotted on the chart in comparison to what they heard.

> *I recorded again partially because I want to match my vowel... partially because I want to see how it interprets what I'm saying and see if I agree with that and feel like I'm seeing that.* (P8)

All four participants mentioned uncertainty when interacting with the audio-only practice. Each of the comments directly related to being unsure why they sounded different to the Spanish recording. In terms of RQ1, it appears that participants can feel equally uncertain about what makes them sound different from a Spanish speaker with and without visuals. Visual information adds another layer of uncertainty– do they match up with what participants think they heard.

*5.1.2    Visual Feedback* All comments regarding visual feedback were made with *V(is)owel*. Participants expressed opinions on their closeness to the Spanish recording lines and related what they saw appear on the chart with what they felt or heard when they said the word. The results answer the second research question by revealing how participants interpreted results and what they found difficult during interpretation.

**Using the Spanish Speaker's Line as a Goal (8/8)**    *V(is)owel* provided participants with a tangible goal to work toward (N=8), which is in stark contrast with all participants relying on their perception of auditory signals to practice with the audio-only condition (N=8). Three expressed their satisfaction with their speech as getting closer to the Spanish speaker based on the proximity of the Spanish speaker's line.

> *That was cool... I didn't expect it to be that much of a linear progression. It really did go from furthest to closer there.* (P2)

Five mentioned that they were surprised by where the vowel showed up on the chart, either because they felt they had not said the vowel well and it showed up closer to the Spanish speaker than expected or because they felt like it had sounded correct and it was not close. In general, they tended to trust the visualization over what they heard.

> *Oh wait... I feel like the "a" was like- it was very much like an* [æ] *sound.* (P4, when a vowel they expected to show up far from the target showed up closer.)

Learners enjoyed experimenting with the chart when it matched with what they heard and when it seemed like they were improving. Otherwise, they felt discouraged with interaction if they couldn't affect the location of the line or if it was far away from the spa line. They also tended to over rely on the visualization– if a first recording was close to the Spanish speaker's line, they would not record another time.

**Connecting Visual and Auditory Feedback (5/8)**    Difficulties in interpretation of *V(is)owel* fell into two main categories: the length of the vowel line and the area in which the vowel landed. While some participants ignored the length and shape of lines, others reported uncertainty interpreting them (N=4). Dissatisfaction with the length of the line caused two participants to record again. None of the participants considered how tongue movement affected the lines during speech.

> *I wasn't able to make a clear connection between the length of the line drawn and what I was saying. I think it was related to, like, how long I spent on a vowel, but it was hard for me to actually determine how long I was spending.* (P2)

The location of lines also confused participants, though part of this can be attributed to inaccuracy in formant extraction (N=5). However, participants expressed confusion about location when the algorithm correctly found formants, with one in particular expressing confusion by vowels showing up outside the lines (N=4).

> *In the chart, it says that my tongue is, like, high front. So now, I tried to lower it, but that doesn't seem to be working.* (P4)

A few participants attempted to adjust their pronunciation in order to move the line (N=3), but only two expressed success in achieving their goal. Both of the participants explicitly used the English corner vowels introduced during calibration to move their tongue in that direction. The success of these participants suggests that including more known vowels during the tutorial phase might help learners associate tongue positions with visuals.

**Creating a Mental Model (4/8)**    Participants came up with different reasons for why their vowels showed up in different locations. Some thought that the loudness of their speech affected the vowel

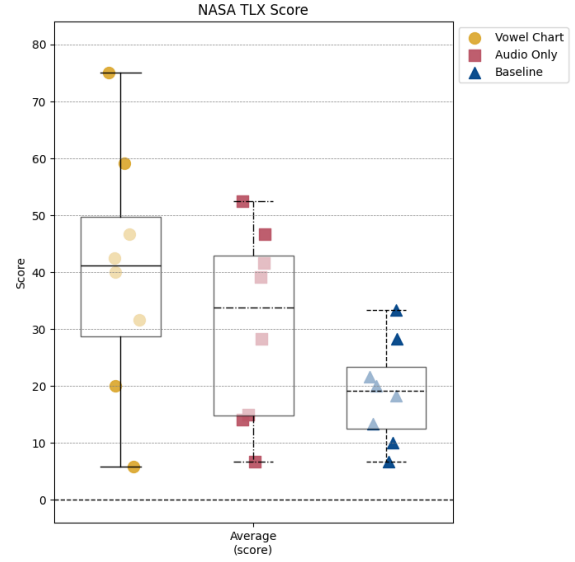| ID | *V(is)owel* | Audio-Only |
|---|---|---|
| P2 | 2.25 | 1.25 |
| P3 | 2.625 | 1.625 |
| P4 | 1.4 | 1.25 |
| P5 | 1.375 | 1.25 |
| P6 | 2.375 | 1.25 |
| P7 | 1.875 | 2.25 |
| P8 | 2 | 1.5 |
| P9 | 4.25 | 2.75 |
| Total Avg. | 2.26875 | 1.640625 |

**Table 3: Average number of recordings per word separated by practice tool across the four words within each practice.**

extractor (N=2). Others thought that their pitch might affect where the vowel showed up with a higher pitch making the line end higher and a lower pitch making the vowel show up lower (N=3). P8 suggested that the calibration or microphone might affect why a vowel showed up in a different place from expected. Finally, P6 thought that speaking English might make vowels show up higher on the chart. It seems that peoples' mental models of speech displayed on graphs may tend toward interpreting increasing and decreasing lines as associated to pitch and amplitude.
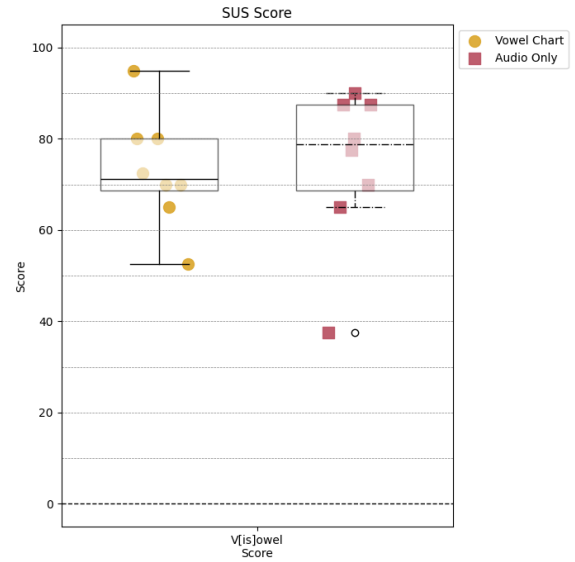
**Familiarity** One of the challenges that participants faced was lack of familiarity with the tongue's role in producing speech (N=2). Pt expressed how this made it more difficult to reach their goals. "I don't think about how I move my tongue, like, ever. So when I'm given instructions ... I don't have the muscle memory of that, like, ideology or specifically how to get there." Although none of the participants said that they had seen a vowel chart before, a few had heard some of the practice words (N=2) and another had seen face-cutaway photos in linguistic content online (N=1). This suggests that interpretation of visuals may not be the greatest barrier in adjusting pronunciation if users do not have previous experience adjusting articulators in accordance with their goals.

*5.1.3 Engagement* We define engagement as the number of times participants recorded practice words during *V(is)owel* and audio-only practice, disregarding a recording if our system did not find a vowel in *V(is)owel* practice. Because participants stated that they chose to record again when our system did not find a vowel during audio-only practice, we include those recordings. Table 3 shows the average number of recordings per word. Most participants recorded with *V(is)owel* more than with AOF. This is consistent with qualitative results. Every participant found *V(is)owel* to provide them with actionable feedback that did not depend on their individual judgement.

*5.1.4 Task Load* The results for the task load can be seen in Figure 13a. The majority of difference in scores is between the pronunciation practice with the tools and baseline task. It appears that *V(is)owel* was more of a burden which is complemented by some participants saying that they thought harder to situate the visual with the auditory feedback for their recording and the Spanish speaker's (N=2).



**(a) A box and whisker plot of the average NASA-TLX scores for *V(is)owel*, audio-only, and baseline. The y-axis ranges from 0 to 80.**



**(b) A box and whisker plot of the average SUS for *V(is)owel* and audio-only. The y-axis goes from 0 to 100.**

**Figure 13: Average scores from surveys.**

*5.1.5 System Usability* All but P2 rated both practice systems with an above 65 for system usability (Figure 13b). Again, audio-only seems to be easier to use, likely because even if the algorithm did not detect a vowel, an playback box would appear. In *V(is)owel*, nothing would appear because there were no formants to extract.

## 5.2 Interview: Reflection on Practice

*5.2.1 Audio-Only* In general, participants appreciated the ability to record and listen to themselves as much as they wanted in the audio-only practice, but some wanted to have some sort of feedback instead of completely relying on their own judgment (N=2). P7 noted that the amount of feedback depended on the word they were practicing, stating that a more complex word would be better with *V(is)owel*.

The simplicity afforded participants the ability to listen to themselves and the Spanish speaker as many times as they liked (N=6). P2 felt that the lack of feedback didn't ultimately help them differentiate between pronunciation that was "correct" or "incorrect". On the other hand, P3 preferred the simpler feedback because they felt more confident recording with it.

> *For some reason, I felt like this one I could listen to it more. And so, like, I could focus on listening to it, if that makes sense.* (P3)

Participants mentioned using the Spanish speaker's audio and their own recently recorded audio to help them hear differences between what they pronounced (N=4). P2 used a trial and error approach to changing their pronunciation based on the difference in emphasis, while P6 and P7 spoke generally about hearing what was wrong in their recording based on what they heard in the Spanish recording. P9 reflected on how the audio helped them move away from how English letters did not have the same pronunciation of Spanish letters.

> *V(is)owel help me go forward and understand how exactly to pronounce different vowels in different situations.. It'll help me with the pronunciation in the future. Just thinking more carefully about that.* (P9)

Most of the participants felt confident about their pronunciation (N=5), though some mentioned that their confidence could be misplaced since it was based on their own judgment (N=2). While P4 expressed it as a leap of faith, P3 appreciated the lack of visual feedback because, "I wasn't seeing how I was getting farther and farther away... I feel like I could match myself to the speaker more, if that makes sense." Two participants said that they were confident, but some words were just hard for them to pronounce because they did not have a Spanish accent. The rest of the participants did not feel confident interacting with only auditory feedback (N=3). Two stated this was because they were just guessing. P9 explained further that, "having an untrained ear is just kind of hard to tell exactly what I did."

*5.2.2 V(is)owel* Participants found *V(is)owel* interesting, though not always accurate (N=3). P8 assumed that something was off with the algorithm while P3 and P7 attributed differences in their pronunciation to lack of experience with the chart and Spanish respectively. On the other hand, three participants preferred the visual feedback to audio-only because it gave them tangible information to act on, either as a visual goal (N=2), or feedback on the position of the tongue (N=1). Participants appreciated having feedback that directly related to a physical space since that gave them a direct path to change how they were speaking (N=2), while others expressed appreciation for a visualization that showed the

progression of vowel over time since a vowel sound could start off in the right space but end somewhere different from the goal (N=3).

> *I thought it was interesting to kind of see, especially, the transitions between certain parts of sounds in the vowel.* (P8)

Two participants noted that they liked seeing what they might be doing wrong in the context of thinking that they had done a good job (N=2).

While participants appreciated the ability to work toward a goal, they were not always sure how to interpret what they were seeing (N=2). P6 ignored the length of the line altogether and focused on getting their lines to be in the same region of the Spanish speaker. P5 noted that although it would take longer to understand *V(is)owel*, they found it more useful for improving pronunciation. P9 noted that *V(is)owel* not only helped them practice their phonetics, but also helped them focus on which part of the word was stressed by the Spanish speaker.

All but one participant felt confident about adjusting their pronunciation with *V(is)owel* (N=7), but attributed their confidence to different aspects. P5 said that they were comfortable adjusting their pronunciation, but noted that they had to think about the visual and auditory information so it took longer. Four said they were confident adjusting their pronunciation based on the feedback because they could see their progress, while one stated that they did not think the chart always accurately displayed what was wrong.

> *I felt confident in my ability to change the line a little bit, but not in a way that I would get an answer that was necessarily more accurate.* (P8)

In general, participants expressed hesitancy regarding the accuracy of the vowel chart (N=6). Half explicitly stated that they thought it might not being giving accurate feedback, while the rest attributed the feedback to a difficulty on their part understanding how to move their tongue to get the lines to change on the chart.

## 5.3 Technical Evaluation

For our tool to work long term in the real world, the formant extraction must choose the correct frequencies for $f_1$ and $f_2$ and extract vowel onset and offset times within a set tolerance.

Any major mistakes during formant extraction will result in inaccurate results for the position of the tongue as displayed on *V(is)owel*, creating more confusion for the user and potentially negatively impacting their pronunciation. Slight variation for a formant is considered a small error. We report the mean absolute error (MAE) in such cases. If the algorithm picks the wrong frequency for a formant, that is categorized as a large error. We report the probability of such an error occurring.

We calculate the mean absolute error (MAE) for vowel onset and offsets. We use the MAE to compare our algorithmic results with manually extracted formant values by Praat. The Praat algorithm computes the formant values within the whole signal by the Burg algorithm [5]. By selecting the vowel, we outputted a list of formants and timestamps for the vowel.

*5.3.1 Vowel Boundary Accuracy* We randomly selected audio files to hand calculate the onset and offset times for the vowel. MAE for vowel boundaries was generally within tolerance (<10ms) for

| | Word | MAE Onset (ms) | MAE Offset (ms) |
|---|---|---|---|
| P2 | pie | 5 | 15 |
| | tara | 5 | **25** |
| P3 | canoa | **25** | **265** |
| | duda | **15** | 5 |
| P4 | polo | **15** | **15** |
| | talla | 5 | 5 |
| P5 | pe | 5 | 6 |
| | cana | 5 | **135** |
| P6 | oreo | 5 | **155** |
| | pie | 5 | 5 |
| P7 | talla | 5 | 5 |
| | poleo | 5 | **205** |
| P8 | oro | 5 | **125** |
| | duda | 5 | 5 |
| P9 | oro | 5 | **85** |
| | cana | 5 | **145** |

**Table 4: MAE for a randomly sampled set of words. Bold indicates an offset that is outside allowed tolerance.**

| ID | MAE $f_1$ (Hz) | MAE $f_2$ (Hz) |
|---|---|---|
| P2 | 43.44 | 83.82 |
| P3 | 41.26 | 94.15 |
| P4 | 47.19 | **109.69** |
| P5 | 65.60 | 95.51 |
| P6 | 80.08 | **122.24** |
| P7 | 53.99 | **125.74** |
| P8 | 84.19 | **182.78** |
| avg. MAE | 59.39 | 116.2757143 |

**Table 5: The average Mean Absolute Error (MAE) of $f_1$ and $f_2$ for each participant. Bold indicates an offset that is outside allowed tolerance.**

onset times but was completely out of tolerance for offset (Table 4), which increases the likelihood that participants saw accurate formants for the beginning of a vowel and inaccurate for the end.

*5.3.2 Formant Extraction Accuracy* On average, MAE of the formant frequency was within the threshold for acceptable difference of 100Hz for $f_1$ (Table 5). The formant extraction for $f_2$ was not as accurate, only being under the threshold for a few participants (N=3), which is expected given the low accuracy for vowel offsets. The likelihood of a large error varied significantly between participants, ranging from 5% to 26% (Table 6). Because of the high likelihood of $f_2$ inaccuracies, participants often saw visually inconsistent results between recordings. One recording would show up in the appropriate quadrant of the chart but a similar production's recording could show up on the opposite side.

## 5.4 RQ1: How does interaction between a visual and audio-only feedback method differ?

Despite mistakes made by our formant extractor and confusion in interpreting what the lines on the chart meant, participants still chose to practice with *V(is)owel* more than audio-only feedback

| ID | # of Large Errors | Total Recordings | % Incorrect |
|---|---|---|---|
| P2 | 4 | 28 | 14.29 |
| P3 | 3 | 32 | 9.38 |
| P4 | 3 | 29 | 10.34 |
| P5 | 1 | 21 | 4.76 |
| P6 | 7 | 27 | 25.93 |
| P7 | 8 | 32 | 25.00 |
| P8 | 4 | 24 | 16.67 |

**Table 6: The percentage of large formant errors for each participant, where large is getting the formant completely wrong.**

(AOF). While some of this could be due to the novelty effect, based on verbal feedback during interaction, participants were also motivated by producing speech that would get the vowel just a little bit closer to the Spanish speaker's. They chose to practice with the vowel chart because **they had guidance on what to change**, the visual giving them a basis to adjust their pronunciation beyond auditory perception and a way to evaluate after recording. Despite showing a desire to engage with *V(is)owel*, participants expressed a degree of uncertainty that did not exist with AOF. Many participants said that they felt confident practicing with audio-only feedback, while recognizing that the confidence was based on their own perception. Confidence was lower during interaction with *V(is)owel* because of the conflict between visual and auditory interpretation. The conflicting information could contribute to the higher average recordings made with *V(is)owel* since participants were trying to bring the audio and visual feedback into agreement.

Neither method of feedback affected whether participants thought about other aspects of speech, which indicates that despite placing an emphasis on visual feedback, *V(is)owel* may not cause a blinder effect. Just under half of the participants noted that they liked seeing the progression of the vowel over time (N=3), suggesting that our novel approach to vowel visualization is a desired feature. In short, *V(is)owel* was treated with some distrust, but that did not cause participants to dislike it or to practice with it less.

## 5.5 RQ2: How do phonetically untrained users interpret vowel charts during interaction?

Participants liked the target Spanish lines on the chart because it gave them a tangible goal to work toward. We noticed that all participants were not sure how to interpret the length of the lines, how to affect the length, or what kinds of interactions were available. We believe that there are three aspects that led to participants' confusion based on the qualitative and quantitative data. First, participants were not used to thinking about how their tongues moved, which would hinder recognition of the tongue's movement during the vowel. Second, the algorithm picked up too much of the surrounding consonants. Finally, although the *V(is)owel* tutorial included a description that discussed what the line length meant, it did not allow participants to experiment with line length.

Many participants thought of where their tongue's position based on visual feedback, despite being unaccustomed. Not everyone felt successful adjusting their tongue position, but those

who were used known sounds to move their tongue in that direction. As a result, we believe that having known target vowels is important for success during initial interaction with *V(is)owel* since the tongue will know what position to take for the known vowel sound, giving it a direction to move.

More practice in the tutorial would provide learners with more known words on the chart, leading to a better understanding of how to move the tongue to get the desired change in location and reduce the possibility of misunderstanding which axes map to different parts of the tongue's position. For example, a few participants mentioned the calibration words to remind them where known words landed. But they were not always correct, indicating a conceptual error [6]. The strength of the vowel chart lies in its mapping to the physical realm. If participants are confused by what that mapping is, that conflicts with their ability to implement feedback. Including examples in the tutorial that provide learners with more points of reference, they will have a better chance of knowing how to get their line to move across the chart.

Participants often thought that pitch might affect what they saw on the chart. We did not bring up any of these in our introduction to *V(is)owel*, which implies that people might naturally attach pitch to lines that vary in height on a chart. One participant incorrectly concluded that the origin of language would affect the chart, which is not the case because of our method of calibration. A short exercise in the tutorial focusing on pitch should alleviate misunderstandings.

Due to our algorithm, participants saw mistakenly plotted vowels. They put greater emphasis on the visualization than on their perception of pronunciation, which lines up with statements made about how discerning if their pronunciation was correct based only on audio feedback sometimes was just guessing. Due to the trust put in the visualization, it is important that there are minimal mistakes made by the algorithm.

In conclusion, users found the target based method of a vowel chart to help them practice pronunciation. They struggled to understand what made vowel lines long or short and attributed changes in the angle of the line to a change in their pitch during production.

## 5.6 RQ3: How does a visualization that focuses on vowels affect users' perceptions of other aspects of pronunciation?

A concern with practicing on a specific sound in a language would be that it creates a blinder effect toward other important aspects of pronunciation. Based on comments made during interaction, it is not clear whether this happens or not. Multiple participants mentioned a desire to improve their pitch, But only two participants thought about consonants, one during interaction with *V(is)owel*, the other during audio-only after interaction with *V(is)owel*. It appears that participants continue to consider multiple aspects of pronunciation based on their comments, but more research is needed to confirm these preliminary findings.

## 6 Discussion

**Real-world use.** We note that *V(is)owel* is a vowel-specific interactive tool. It still struggles with certain consonants which can limit its real-world use. Future work may address this limitation by using additional linguistic properties in conjunction with our

methods. Additionally, users could also be given greater freedom on which vowel should be targeted for practice.

Any incorrect feedback will lead users astray, which is likely to be detrimental over an extended period of time. Based on offset times, *V(is)owel* currently gives faulty feedback. Another possible effect of using *V(is)owel* for a longer time is that it could put undue emphasis on vowels. While initial interactions are encouraging since participants continued to think about pronunciation as a whole, there were many more comments regarding vowels.

**Longer tutorial.** We noticed a tendency for participants to trust the visual feedback over their own ears. Because there is no statistical way to capture just the vowel and exclude all of a consonant, introducing learners to the concept of sonorant consonants should help reduce confusion regarding longer lines on the chart. To further reduce excess visual information, we could simplify vowel lines, potentially allowing toggle between an average location and the whole vowel.

Given two participants success with anchoring vowels based on English vowels, we believe that including more information during practice explaining how to interact with *V(is)owel* could help. This might include more English words with different vowels to provide more reference points, and feedback to emphasize the physical connection to the chart by giving relational feedback, e.g. "Your vowel is lower than the goal, try bringing your tongue up to get closer".

One of the takeaways from participant responses to *V(is)owel* is the need for greater transparency on how the system extracts their vowels. We received multiple comments regarding the confusion of what affected the position on the chart, one being that participants thought the pitch of their voice might make the line show up in different places. The tutorial could have practice that contrasts loud production with soft and high with low pitched speech.

**Longitudinal study.** Our goal is to improve intelligibility in a second language, not attain native-likeness. Our visualization, though an improvement over frequency markings and addressing previous limitations, does not provide visual markers sufficient for beginner learners. Based on the feedback, and to de-emphasize native-like pronunciation, incorporating an intelligibility threshold would be helpful for a second language learners.

Another limitation to our findings is the novelty effect; none of the participants had seen the visualization before. However, since participants mentioned using the Spanish speaker's visual as a goal better explain why participants greater engagement with *V(is)owel*. A study over a longer period of time would allow us to understand how much of the engagement was due to novelty.

**Applicability for other languages.** *V(is)owel* can be extended to other languages, provided that they primarily differ from the first language in tongue location. We could add to the current visuals with additional language specific contexts, such as differences in lip rounding or nasalization. For future work, we could create a catalog of language specific idiosyncrasies, which could be loaded based on user needs.

## 7 Conclusion

State-of-the-art CAPT tools do not provide personalized audio and video feedback for pronunciation improvement. We find that personalized feedback can have a significant impact on a learners' pronunciation in a second language. To that effect, we introduced *V(is)owel*, an interactive vowel chart that provides visual and auditory feedback based on a learner's speech patterns. Our system creates a personalized vowel chart by extracting their L1 vowel space in a calibration phase. During pronunciation practice, our algorithm extracts the vowels from the recorded words, and plots the formants in the vowel chart, along with the associated audio for feedback, allowing for learners to interpret their pronunciation based on both visual and audio feedback. Findings from our user study suggest that our system supports engaging and personalized learning experiences. Future work includes expanding to more attributes of vowels and other languages.

## Acknowledgments

## References

[1] Jessica Barlow and Judith Gierut. 2002. Minimal Pair Approaches to Phonological Remediation. *Seminars in speech and language* 23 (03 2002), 57–68. https://doi.org/10.1055/s-2002-24969

[2] R. Bellman and R. Kalaba. 1959. On adaptive control processes. *IRE Transactions on Automatic Control* 4, 2 (1959), 1–9. https://doi.org/10.1109/TAC.1959.1104847

[3] Cindy Blanco and Ari Moline. 2025. Covering all the bases: Duolingo's approach to listening skills. https://blog.duolingo.com/covering-all-the-bases-duolingos-approach-to-listening-skills/

[4] Heather Bliss, Jennifer Abel, and Bryan Gick. 2018. Computer-Assisted Visual Articulation Feedback in L2 Pronunciation Instruction: A Review. *Journal of Second Language Pronunciation* 4 (05 2018), 129–153. https://doi.org/10.1075/jslp.00006.bli

[5] Paul Boersma and David Weenink. 1992-2022. Praat: doing phonetics by computer. https://www.fon.hum.uva.nl/praat/ [Computer program].

[6] Paul A. Booth. 1991. Errors and theory in human-computer interaction. *Acta Psychologica* 78, 1 (1991), 69–96. https://doi.org/10.1016/0001-6918(91)90005-K

[7] Adam Brown. 1995. Minimal pairs: minimal importance? *ELT Journal* 49, 2 (04 1995), 169–175. https://doi.org/10.1093/elt/49.2.169 arXiv:https://academic.oup.com/eltj/article-pdf/49/2/169/9791152/169.pdf

[8] Yaohua Bu, Tianyi Ma, Weijun Li, Hang Zhou, Jia Jia, Shengqi Chen, Kaiyuan Xu, Dachuan Shi, Haozhe Wu, Zhihan Yang, et al. 2021. PTeacher: a Computer-Aided Personalized Pronunciation Training System with Exaggerated Audio-Visual Corrective Feedback. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–14.

[9] Michael Carey. 2004. CALL Visual Feedback for Pronunciation of Vowels: Kay Sona-Match. *CALICO Journal* 21, 3 (2004), 571–601. https://www.jstor.org/stable/24149798 Publisher: Equinox Publishing Ltd..

[10] Qiaoyi Chen, Siyu Liu, Kaihui Huang, Xingbo Wang, Xiaojuan Ma, Junkai Zhu, and Zhenhui Peng. 2024. RetAssist: Facilitating Vocabulary Learners with Generative Images in Story Retelling Practices. In *Proceedings of the 2024 ACM Designing Interactive Systems Conference* (Copenhagen, Denmark) *(DIS '24)*. Association for Computing Machinery, New York, NY, USA, 2019–2036. https://doi.org/10.1145/3643834.3661581

[11] Isabelle Darcy, Brian Rocca, and Zoie Hancock. 2021. A Window into the Classroom: How Teachers Integrate Pronunciation Instruction. *RELC Journal* 52, 1 (2021), 110–127. https://doi.org/doi/10.1177/0033688220964269

[12] Keya Rani Das and AHMR Imon. 2016. A brief review of tests for normality. *American Journal of Theoretical and Applied Statistics* 5, 1 (2016), 5–12.

[13] Tracey Derwing. 2003. What Do ESL Students say about Their Accents? *Canadian Modern Language Review-revue Canadienne Des Langues Vivantes - CAN MOD LANG REV* 59 (06 2003), 547–567. https://doi.org/10.3138/cmlr.59.4.547

[14] Beth G. Greene, David B. Pisoni, and Thomas D. Carrell. 1984. Recognition of speech spectrograms. *The Journal of the Acoustical Society of America* 76, 1 (July 1984), 32–43. https://doi.org/10.1121/1.391035

[15] Mohd Hilmi Hamzah and Abdullah Bawodood. 2019. TEACHING ENGLISH SOUNDS VIA MINIMAL PAIRS: THE CASE OF YEMENI EFL LEARNERS. *Journal of English Language and Literature* 6 (10 2019), 97–102. https://doi.org/10.33329/joell.63.97

[16] Feifei Han. 2016. Integrating pronunciation instruction with passage-level reading instruction. *Pronunciation in the classroom: The overlooked essential* (2016), 143–151.

[17] Sandra G Hart. 1986. NASA task load index (TLX). (1986).

[18] Soon-Hin Hew and Mitsuru Ohki. 2004. Effect of Animated Graphic Annotations and Immediate Visual Feedback in Aiding Japanese Pronunciation Learning: A Comparative Study. *CALICO Journal* 21, 2 (2004), 397–419. https://www.jstor.org/stable/24149403 Publisher: Equinox Publishing Ltd..

[19] Martin Hinton. 2013. An Aptitude for Speech: The Importance of Mimicry Ability in Foreign Language Pronunciation. In *Teaching and Researching English Accents in Native and Non-native Speakers*, Ewa Waniek-Klimczak and Linda R. Shockey (Eds.). Springer, 103–111. https://doi.org/10.1007/978-3-642-24019-5_8

[20] Murat Hismanoglu and Sibel Hismanoglu. 2010. Language teachers' preferences of pronunciation teaching techniques: traditional or modern? *Procedia-Social and Behavioral Sciences* 2, 2 (2010), 983–989.

[21] David Huggins-Daines, Mohit Kumar, Arthur Chan, Alan W Black, Mosur Ravishankar, and Alexander I Rudnicky. 2006. Pocketsphinx: A free, real-time continuous speech recognition system for hand-held devices. In *2006 IEEE international conference on acoustics speech and signal processing proceedings*, Vol. 1. IEEE, I–I.

[22] Yannick Jadoul, Bill Thompson, and Bart De Boer. 2018. Introducing parselmouth: A python interface to praat. *Journal of Phonetics* 71 (2018), 1–15.

[23] Kelly Moser and Tianlan Wei. 2024. COVID-19 and the pre-existing language teacher supply crisis. *Language Teaching Research* 28, 5 (2024), 1940–1975.

[24] Murray J Munro, Tracey M Derwing, and Ron I Thomson. 2015. Setting segmental priorities for English learners: Evidence from a longitudinal study. *International Review of Applied Linguistics in Language Teaching* 53, 1 (2015), 39–60.

[25] Heather M. Offerman and Daniel J. Olson. 2016. Visual feedback and second language segmental production: The generalizability of pronunciation gains. *System* 59 (July 2016), 45–60. https://doi.org/10.1016/j.system.2016.03.003

[26] Daniel J Olson. 2014. Benefits of visual feedback on segmental production in the L2 classroom. (2014).

[27] Annu Paganus, Vesa-Petteri Mikkonen, Tomi Mäntylä, Sami Nuuttila, Jouni Isoaho, Olli Aaltonen, and Tapio Salakoski. 2006. The vowel game: continuous real-time visualization for pronunciation learning with vowel charts. In *Advances in Natural Language Processing: 5th International Conference on NLP, FinTAL 2006 Turku, Finland, August 23-25, 2006 Proceedings*. Springer, 696–703.

[28] Indah Rahman and Isna Nur. 2017. THE USE OF MINIMAL PAIR TECHNIQUE IN TEACHING PRONUNCIATION AT THE SECOND YEAR STUDENTS OF SMAN 4 BANTIMURUNG. 4 (2017), 276. https://doi.org/10.24252/Eternal.V42.2018.A11

[29] Rajiv Rao. 2019. *Key issues in the teaching of Spanish pronunciation*. Routledge London, United Kingdom.

[30] Ivana Rehman. 2021. *Real-time formant extraction for second language vowel production training*. Ph. D. Dissertation. Iowa State University.

[31] Ronald W. Schafer and Lawrence R. Rabiner. 1970. System for Automatic Formant Analysis of Voiced Speech. *The Journal of the Acoustical Society of America* 47, 2B (02 1970), 634–648. https://doi.org/10.1121/1.1911939 arXiv:https://pubs.aip.org/asa/jasa/article-pdf/47/2B/634/18766538/634_1_online.pdf

[32] Michael R Sheldon, Michael J Fillyaw, and W Douglas Thompson. 1996. The use and interpretation of the Friedman test in the analysis of ordinal-scale data in repeated measures designs. *Physiotherapy Research International* 1, 4 (1996), 221–228.

[33] Laura Sicola and Isabelle Darcy. 2015. Integrating pronunciation into the language classroom. *The handbook of English pronunciation* (2015), 471–487.

[34] Lars St, Svante Wold, et al. 1989. Analysis of variance (ANOVA). *Chemometrics and intelligent laboratory systems* 6, 4 (1989), 259–272.

[35] Magdalena Szyszka. 2015. Good English Pronunciation Users and Their Pronunciation Learning Strategies. *Research in Language* 13 (03 2015), 93–106. https://doi.org/10.1515/rela-2015-0017

[36] Jeff Tennant and RC Gardner. 2004. The computerized mini-AMTB. *Calico Journal* (2004), 245–263.

[37] R. I. Thomson and T. M. Derwing. 2014. The Effectiveness of L2 Pronunciation Instruction: A Narrative Review. 36, 3 (12 2014), 326–344. https://doi.org/10.1093/applin/amu076

[38] Michael Wei. 2006. A Literature Review on Strategies for Teaching Pronunciation. *Online submission* (2006).

[39] Martijn Wieling, Eliza Margaretha, and John Nerbonne. 2011. Inducing phonetic distances from dialect variation. *Computational Linguistics in the Netherlands Journal* 1 (2011), 109–118.

[40] Ka-Wa Yuen, Wai-Kim Leung, Peng-fei Liu, Ka-Ho Wong, Xiao-jun Qian, Wai-Kit Lo, and Helen Meng. 2011. Enunciate: An internet-accessible computer-aided pronunciation training system and related user evaluations. In *2011 International Conference on Speech Database and Assessments (Oriental COCOSDA)*. 85–90. https://doi.org/10.1109/ICSDA.2011.6085985

# A Statistical Analysis

## A.1 Engagement

A Q-Q visual test revealed that the number of times participants recorded is not normally distributed [12]. Since the samples are not independent, we used the Friedman test to determine statistical significance [32]. We found that participants recorded more with *V(is)owel* practice than with audio-only (p=0.0339). Because the visualization in *V(is)owel* depended on speaker, we tested for a difference between participants who practiced with a female voice (N=4) versus a male voice and did not find a statistically significant difference (p=1).

## A.2 NASA-TLX

In order to determine the suitability of performing an ANOVA on average NASA scores [34], we used visual Q-Q graphs to see how much our distribution varied from a normal distribution. We compared the visual results with a Shapiro-Wilk numerical test and determined that the data could be considered to be pulled from a normal distribution (Figure 13a). We performed a two-factor ANOVA without Replication on the average NASA TLX scores collected after the baseline check, audio-only, and *V(is)owel* practices (p=0.008). We excluded the baseline scores in a second ANOVA and found that there was no statistical difference between the load of *V(is)owel* and audio-only practices (p=0.102), which is expected since the baseline condition does not require thinking about pronunciation.

## A.3 System Usability Scores

A test for normalcy in the distributions showed that the *V(is)owel* scores were close to a normal distribution, but the audio-only was not. As a result, we used the Friedman test to determine if the difference in usability between the tools was statistically significant. There was no statistical difference in usability (p=0.480).

# B Motivation

To better understand learning goals, we had participants complete a motivation survey, modeled after the mini-AMTB by Tennant and Gardner [36]. The survey gauges the motivation behind learning Spanish using five point Likert scales.

The overall motivation of a participant as defined by Tennant and Gardner is:

$$MOTIV = MI + D + AL$$

where MI is the numerical response to the motivation question, D is the response to the desire question, and AL is the attitude toward learning the target language (see Appendix B.1 for a list of questions). To return to a 5-pt Likert scale, we divide the result by 3.

## B.1 Questions

**Integrative Orientation (IO)**    If I were to rate my feelings about learning Spanish in order to interact with Hispanic or Latino speakers, I would have to say they are, "weak ... strong"

**Attitude toward Hispanic or Latino Americans (AFA)**    My attitude toward Hispanic or Latino speakers is, "unfavorable ... favorable"

| ID | D | MI | ALSpa | Average MOTIV |
|----|---|----|-------|---------------|
| P2 | 4 | 3 | 5 | 4.0 |
| P3 | 5 | 3 | 5 | 4.0 |
| P4 | 5 | 3 | 5 | 3.7 |
| P5 | 5 | 5 | 5 | 5.0 |
| P6 | 4 | 3 | 5 | 3.7 |
| P7 | 4 | 4 | 5 | 4.0 |
| P8 | 4 | 2 | 5 | 3.3 |
| P9 | 4 | 2 | 5 | 3.3 |

Table 7: Responses to the motivation questionnaire, where average motivation is calculated as ( Desire + Motivational Intensity + Attitude toward Learning Spanish ) / 3

**Interest in Foreign Languages (IFL)**    My interest in languages other than Spanish and English is, "very low ... very high"

**Desire to learn Spanish (D)**    My desire to learn Spanish is: "weak ... strong"

**Attitude toward learning Spanish (AL)**    My attitude toward learning Spanish is: "unfavorable ... favorable"

**Instrumental Orientation (IO)**    If I were to rate my feelings about learning Spanish for practical purposes such as to improve my occupational opportunities, I would have to say they are "weak ... strong"

**Motivational Intensity (MI)**    I would characterize how hard I work at learning Spanish as "very little ... very much"

# C Results

The participants all had a greater than 3 average motivation to learn Spanish (Table 7). The motivation to engage with Spanish practice is reinforced by the thoughtful comments participants made during practice (Sections 5.1.1 and 5.1.2 ).

# D Exit Interview Questions

Each question was asked for both practices.

- How did you like the feedback?
- How was the feedback useful?
  - What feedback could you act on?
- How trustworthy was the feedback?
- What did you find confusing or distracting?

# E Design Decisions

## E.1 Color Palettes

We considered colorblind friendly colors for differentiating the words and speakers, green and purple for Spanish, and blue and red for the learner (Figure 3a).

## E.2 Layout

We evaluated three potential layouts for the minimal pair recording buttons and *V(is)owel* (Figure 15). The layouts represented the potential orderings that were possible given the chart and buttons. We ended with a horizontal alignment of the buttons with the chart
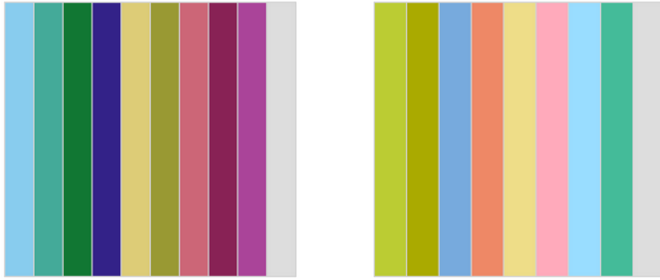
**Figure 14: Two colorblind friendly palettes. From left to right: tol_muted and tol_light.**

because the vowel chart has a straight side which makes the best use of the available space.

## F Tutorial

The tutorial begins with face-cutaways that show the position of the tongue during the production of four different vowels in extreme positions (two front vowels, high and low, and two back vowels, high and low. See Figure 17a). The second page we introduce the vowel chart alongside a face-cutaway to show correlation between the position of the tongue and the points on the chart. The third and fourth steps break down the horizontal and vertical axes by bringing attention to the position of the tongue to a user and letting them record English words that differ in the respective axis. After allowing users to interact with English words, we presented them with the idea of getting close to a Spanish recording. The final step was a review of the tongue's mapping to the vowel chart. Based on pretests, we modified the tutorial to include the same four vowels to cement the association from the first step (Figure 8b). We created an interactive prototype and ran another pretest. We learned that people were unaware that the face cutaways on the first page were interactive. When clicked, they played the associated vowel. To bring attention to the interaction, we show an outline when the mouse hovers over any of the pictures (Figure 17b). Originally, all the vowels on the second page were the same color. Based on feedback, we colored the tongues to match the colors of the vowels on the chart to create a visual association. An overall modification we made based on pretest information was to put the instructions above the vowel chart to encourage users to read the instructions before plowing ahead. We also pared down the instructions to be more concise.

## G Vowel Boundary Algorithms

We outline brief descriptions of the algorithms we tested for vowel boundary extraction and why we did not use them.

We started using a DTW algorithm [2], which maps similar points between two audio samples regardless of speed. There were two flaws with DTW. First, if an English speaker didn't pronounce the consonants like the Spanish reference, then the mapping would include multiple points near the beginning of the vowel, potentially throwing off the beginning and ending boundaries for the vowel (Figure 10). Second, even when DTW was accurate, the window

size could not be made small enough (<10ms) to exclude enough of the consonant in the extraction.

Next, we tested extracting timestamps by using pocketSphinx (PS) [21]. PS is a speaker-independent continuous speech-recognition algorithm begun by Huggins-Daines et. al. It takes in an audio file and returns the phones and silences with associated timestamps detected in the file. The timestamps had a similar degree of inaccuracy as DTW.

We decided to test Praat's vowel extractor suitability for extraction using the parselmouth library, a python interface for Praat [22]. The original Praat script written by Hugo Quené only extracts a portion of the vowel so we modified the script to output the timestamps for the whole vowel [5]. The boundaries using this method were also inaccurate.

## H Timing Considerations

We wanted to keep the timing between recording and visualization low. Some parts of the server side code have a set amount of time: speech recognition, vowel boundary and formant extraction. While the algorithm ran these processes, we displayed "Processing" on the record button and disabled interaction with it. For the formants to be displayed on the chart, they must be transformed from the frequency to svg domain. Initially, we transformed each pair of formants, $f_1$ and $f_2$, one by one. We optimized the process by transforming formants in a batch.
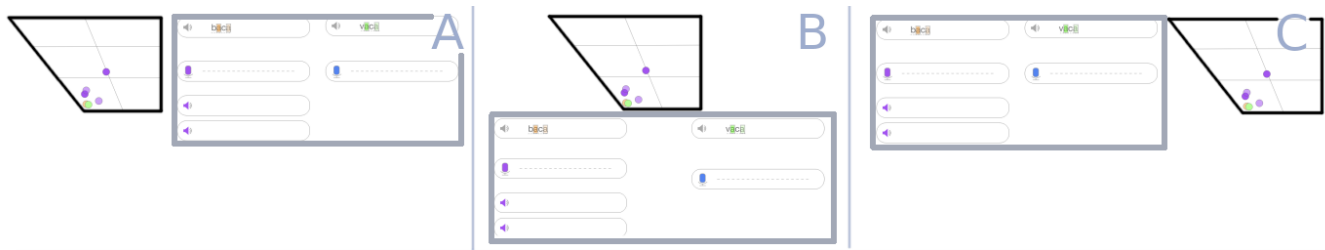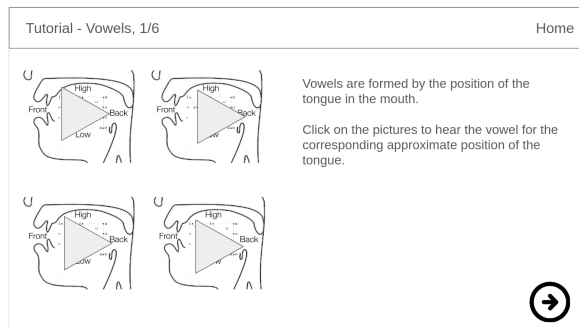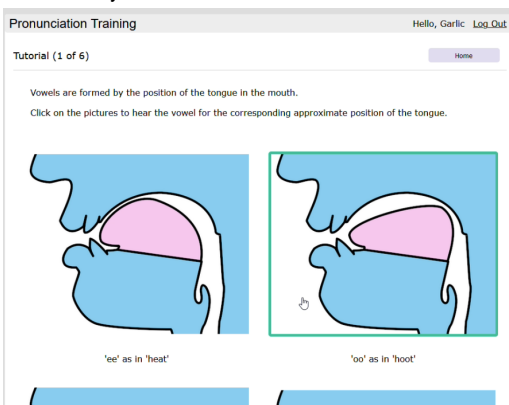
Figure 15: Three potential layouts for *V(is)owel* practice. The vowel chart is outlined in black and the recording buttons and recordings are outlined in grey. Audio-only practice is just the grey portion.



(a) The paper prototype of the first page of the tutorial. A grid of four face cutaways is on the left. Triangles on each of the pictures indicated that they were short videos.



(b) Final interface for the tutorial's first page. The face cutaways are static due technical limitations. Hovering over a picture lines it with green indicating that a user can interact with it.

Figure 17: Design cycle of page 1