

Hidden bawls, whispers, and yelps: can text be made to sound more than just its words?

Caluã de Lacerda Pataca, Paula Dornhofer Paro Costa

Abstract—Whether a word was bawled, whispered, or yelled, captions will typically represent it in the same way. If they are your only way to access what is being said, subjective nuances expressed in the voice will be lost. Since so much of communication is carried by these nuances, we posit that if captions are to be used as an accurate representation of speech, embedding visual representations of paralinguistic qualities into captions could help readers use them to better understand speech beyond its mere textual content. This paper presents a model for processing vocal prosody (its loudness, pitch, and duration) and mapping it into visual dimensions of typography (respectively, font-weight, baseline shift, and letter-spacing), creating a visual representation of these lost vocal subtleties that can be embedded directly into the typographical form of text. An evaluation was carried out where participants were exposed to this *speech-modulated typography* and asked to match it to its originating audio, presented between similar alternatives. Participants (n=117) were able to correctly identify the original audios with an average accuracy of 65%, with no significant difference when showing them modulations as animated or static text. Additionally, participants' comments showed their mental models of speech-modulated typography varied widely.

Index Terms—Affective computing, speech visualization, emotion representation, speech analysis.



1 INTRODUCTION

DESPITE its expressive richness, when speech is represented through captions it is typically reduced to its words, and its words only. Whatever nuance was originally conveyed by the ways in which the speaker modulated their voice — their mood, emotions, dispositions, etc — is lost in this flattened textual representation. This is particularly relevant when captions are used not as a complement to a readily available audio channel, but as its *replacement*. This can be true for deaf and hard of hearing (DHH) persons, but will potentially affect anyone, including hearing individuals facing a situational hearing impairment, e.g., someone affected by a situational hearing impairment such as watching a film on their mobile phone in a noisy environment [1]. If so much of communication is expressed in nuances not captured by written text, what is to be said of the experience of those who have no direct access to acoustic speech but only to its written forms?

Which is not to say that not capturing *all* of prosody¹ in speech is a shortcoming of written text. While both speech and text are of course tied together, one is not a simple variant of the other [2], and just as there are nuances of speech not captured by written text, the opposite is also frequently true, especially now that so much of internet-based communication centers around text [3]. In fact, it has been argued that, while at its origins written text's functions were that of supporting oral speech [4], as the medium developed its uses shifted — e.g., the invention of modern

punctuation and white spaces between words allowed for the emergence of *silent reading* [5], [6], a form detached from text's origins in sounded speech.

Still, there are contexts where this gap between speech and what is effectively captured by writing can be an obstacle, as has been shown when DHH individuals point out that captions are not a functional equivalent to hearing since they lack some meta-speech information such as 'speaker identification[,] punctuation, sentiment, [and] tone' [7]. This is not a new notion (see, for instance, [8]'s speculations from back in 1983 about how an enhanced captioning system could use 'special effects, color, and capital letters [to represent] the rich tonal information denied deaf children[.]'), but is becoming more and more relevant when we consider how real-time captioning has grown as a critical assistive technology. As the capabilities of automatic speech recognition (ASR) systems increase, captions are finding use in many novel contexts as a way of giving visual access to sounded speech for DHH individuals [9], [10].

As such, researchers have been exploring ways of enhancing the display of captions to include cues to inform both speaker behavior [11], [12] and readers' interpretations. Our work rests on this second case: creating a model capable of representing prosody over regular captions.

We propose a novel model of *Speech-Modulated Typography*, where acoustic features from speech are used to modulate the visual appearance of text. This could allow for a given utterance's transcription to not only represent words being said, but *how* they were said. With this, we hope to uncover typographic parameters that can be generally recognized as visual proxies for the prosodic features of amplitude, pitch, and duration. Our broad aim is to help further the development of real-time captioning technology,

• Caluã de Lacerda Pataca and Paula Dornhofer Paro Costa are with the University of Campinas. They can be contacted at calua.pataca@gmail.com and paulad@unicamp.br.

1. Prosody models the paralinguistic dimensions of *how* words are said, beyond their linguistic dimensions of *what* words are said.

for which the use of algorithmic systems for both processing of prosody and the subsequent modulation of typographic parameters is a requirement.

This approach can lead to better captioning systems that, just as traditional captions are already useful for DHH and hearing individuals [13], have potential benefits for anyone. While we do not aim to represent emotions directly, we argue that, in visually representing an acoustic correlate to emotion, our work is tied to applications and discussions of Affective computing, with additional insights from Visual design and Linguistics. Additionally, and even though our model was developed and evaluated in a Brazilian Portuguese context, our findings are plausibly applicable to any language that uses alphabets.

In this work, we also propose a method to empirically evaluate whether these visual modulations of typography are sufficiently clear to allow for their recognition and understanding by untrained participants. This, we posit, is an important condition if this technology is to eventually gain traction. This methodology informed the creation of a dataset of expressive vocal readings of poetry, which was used in an evaluation conducted with 117 participants aimed at answering two research questions:

- 1) Can readers of this speech-modulated typography use it to recognize prosodic features present in its originating audio but not in its textual content?
- 2) Are there differences in recognition between animated and static based versions of speech-modulated typography?

Lastly, in said experiment we also collected open-ended comments from participants about their impressions about the model and its possible uses. These aimed to better qualify the answers we obtained to our research questions, but also to uncover unforeseen potentials and limitations of our model not captured by our quantitative methods.

1.1 Outline of this paper

In Section 2 we present an overview of key concepts, challenges, and related work in fields such as linguistics, speech emotion recognition, and design and accessibility.

Following, in Section 3 we will explore our prosody-typography mapping model, discussing its choice of features, typographical mappings, and implementation details.

In Section 4, we will present an experiment we ran aimed at understanding if the choices we took with the model generated typographic results that were sufficiently clear to be understandable by a general audience.

The results of said experiment are presented in Section 5 and discussed in Section 6. Lastly, we present our broad conclusions, limitations, and proposed future work in Section 7.

2 RELATED WORK

This section is divided into three main areas. We first give an overview of how prosody has both linguistic and paralinguistic functions, and the latter grounds our approach. Following, we show how different authors have created visual representations for elements of sound that are typically missing in textual transcriptions. Lastly, we present works that, like our own, have used algorithmic approaches to modulate the visual form of typography to echo sound.

2.1 On the importance of prosody to shape meaning

Transcribing speech into captions is a lossy process. Latin scripts are able to codify phonemes but, beyond those, the means of representing expressive variations originally present in speech are restricted to just a handful of possibilities, such as typographical emphases (variants like *italic* and **bold**), special-purpose characters (!, ?, and more [14]), emojis and creative distortions of standard grammar and the use of punctuation [3], etc.

While words serve to encode concepts (e.g. the word *cat* can be deduced to mean the feline animal), only decoding them may not be enough to parse a statement's meaning — or, rather, to close in on a *plausible* interpretation among many possible alternatives. In a typical conversation, a listener will attempt to infer their understanding by using any available signals to close in on the most likely meaning among many credible alternatives. Prosody is one such signal [15]. In the *a cat was standing beside that car* statement, for instance, if one's voice emphasizes the word *was*, the cat might now be gone; if *standing*, maybe a contrast is being made with another cat that could have been lying down; if *that car*, maybe they are saying that they meant *that* car and not other plausible cars? And so on. It is only upon incorporating prosody in their interpretation that a listener will be able to settle between these concurring alternative meanings, hence its importance [16].

2.1.1 Paralinguistic prosody

One of prosody's functions is, then, to assign contrastive focus to one word in opposition to others, guiding us to sense what is more or less important at any given time. Yet, just as the voice generally has dimensions of meaning that go beyond language itself, prosody also conveys information about the speaker (e.g., age, gender, geographical origin, etc), their dispositions (are they tired, grumpy, sick, drugged, etc?), or their moods and emotions [17]. An example of this varied nature of prosody, [18] showed that pitch and rhythmic information in Brazilian Portuguese was able to encode both linguistic (whether an utterance was declarative, interrogative, or imperative) and paralinguistic dimensions (e.g., categorical emotions associated with an utterance) in speech.

From a perceptual point of view, [19] show how both linguistic and acoustic interpretations of emotion in speech can be present simultaneously. In their experiment, Swedish and Brazilian participants were presented with a set of speech samples in Brazilian Portuguese, which they had to label among a set of categorical emotions. The Brazilian group had a greater agreement on the labels chosen, but the Swedish — disconnected from the samples' linguistic but not paralinguistic dimensions — were still able to significantly agree on the labels.

An in depth discussion about *how* speech generally [20], [21], and prosody particularly [22], are able to encode emotion is not the focus of this paper. Suffice to say, this is important for areas such as speech emotion recognition, a diverse field with competing approaches that consider different acoustic features (e.g., pitch, formants, energy, timing, articulation, etc) to identify emotions, themselves understood and organized by different theories and frameworks [23]. These systems have varied applications in fields

ranging from customer satisfaction evaluation to depression diagnosis [24].

Still, given the complex nature of speech, even when using purposefully built data sets with state-of-the-art approaches, recognition rates can be generally low, both for machine-based or human classifiers [25]. This can be seen as a testament to the complexity of the task or, alternatively, as a sign of the limits to what can be inferred when one considers acoustic signals isolated from context.

Coming from a different paradigm, [26] proposes that we see emotion not as an *informational layer* embedded in speech (and therefore obtainable from the acoustic signal) but, rather, as a *culturally grounded phenomenon*, echoing [15]’s notion that prosodic signs must be understood not in isolation, but within a given cultural and social context.

For these authors, the same acoustic cues may have different meanings depending on their surrounding context, which imposes limits to how well an artificial intelligence system is able to perform considering only acoustic cues. While emotion is inevitably tied to the physiology of the bodies producing it, and even if it is methodologically convenient to be able to measure these physical correlates to emotion, the subjective experiencing of emotions is not constrained by only these signals [27].

Rather, we see emotion as something happening *between* subjects². An example of how this theoretical framework can inform affective computing is given by [28]. In it, a software takes snapshots of a person as they are writing an email, and those will accompany each paragraph of the message on the recipient’s end. The goal of this system is not to interpret the writer’s emotions as they compose the message but, instead, to provide additional information which the reader can use to inform their understanding. This is similar to how we conceptualized our own model.

2.1.2 Why represent prosody through typography?

We base our approach not in an attempt to directly represent *emotions* in speech but, rather, to represent *prosody*. If we see emotion as a *culturally grounded phenomenon* [26], it follows that it carries an inherently ambiguous dimension that depends on its originating context — who is involved (e.g., what we know and feel about someone can regulate our reactions)? What are their cultural and social backgrounds (e.g., certain groups might see some feelings as implausible in a certain setting, while for others it might seem natural)? How does their environment influence their conversation (e.g., shouting in a loud night club implies different meanings than shouting in a silent library, even if the acoustic phenomena coincide)?

As such, it is hard to ‘resolve’ emotions, as if there were culturally universal signs that could be inferred from physical signals carried by speech. Our approach is, then, not to try to detect emotions to then represent them but, rather, to represent the voice itself. In representing prosody, we hope to make it visible, and thus accessible for those unable to perceive sound. The viewer, we posit, will then be able to incorporate these raw acoustic elements within their given cultural and social context to inform their own understanding of speech.

2. These can be persons, but also a person and an object, for instance.

2.2 Visual depictions of speech, sound, and beyond

Different authors have tackled the issue of limits of what can be represented by the Latin alphabet, proposing different approaches to expand to it, either through changes in the shapes and typesetting parameters of letters, auxiliary graphical elements, or both. There is evidence that readers were able to successfully incorporate some of these systems to change their reading behavior.

[29] proposed the *Speechant* system of graphical elements that can be overlaid on a traditionally typeset text to give cues about its pronunciation, particularly pitch and rhythmic information, to help Portuguese speakers who were learning English as a second language. In their evaluation, students using this system were able to produce better-sounding readings of English utterances than a control group, as judged by a panel of trained phoneticists.

Similarly, [30] designed an extended version of the Times New Roman typeface, with special characters able to highlight and differentiate how some letters in the Dutch language can have different pronunciations depending on where they appear — something which, they argue, hinders how non-native speakers parse written text.

Attempting to tackle the known issue that children who read aloud in a monotonous, non-natural tone are more likely to become poorer readers later in life, [31] tested ways of directly representing prosody in typography. They were able to teach these visual cues to children, who were then able to use them to read aloud more expressively.

[32] developed a typeface that was able to have its characters’ shapes shift from positive valence (rounded strokes → smooth curves) to negative valence (harsh strokes → angled curves), along a continuous scale.

[33] worked with artists to create visual representations in closed captions of categorical emotions present in speech through animated closed captions. The recognition of these emotions was not better for the enhanced captions against a control group with traditional captions, but both hearing and DHH participants stated their preference for these new, enhanced captions.

Similarly, [34] describes the *Kinedit* system, later expanded by [35] to work for instant messaging, which allows inexperienced users to combine multiple animation behaviors to modulate typographical attributes such as font size, color, opacity, position, rotation, etc, which can be used to express prosody, emotions, the direction of attention, characters, etc.

Aiming at more specialized audiences, some authors developed more precise, albeit likewise hermetic, rule-based visualizations of speech to aid with, for instance, prosodic analysis [36], [37]. These mix traditional textual transcription with graphical elements that represent acoustic features of speech such as pitch, energy, and rhythm.

2.3 Algorithmic modulations of typographic form to represent speech

Some authors work exclusively with typography and rule-based manipulation of typographic parameters to echo acoustic features of speech. This approach allows for applications such as speech-modulated closed captions, with potential uses-cases not only in film media but also, as

we have seen, in the many scenarios where ASR-based captioning systems have been making inroads.

One of the challenges of such a line of work is that, until recently, typefaces’ digital files mostly worked as databases of fixed shapes, one for each glyph. While one can easily change compositional parameters such as font size, position, leading, etc, the drawing of each letter was immutable — to access a bold version of a font, for instance, you would necessarily need an additional file. This meant that, while acoustic features can be thought of as continuous, some typographic parameters would have to be changed in discrete steps. To bypass this constraint, some authors have developed custom type shaping engines, which allowed them greater visual freedom to modulate how each letter shape echoes an acoustic feature.

2.3.1 Projects working with custom type-shaping engines

[38] created a system where glyphs are composed by a set of primitive shapes, each of them subject to independent manipulation: size, mapped to loudness; pitch, mapped to vertical and horizontal stretching; rhythm, mapped by the speed at which words come into the screen. Additionally, the authors approximated the Latin alphabet with a phonetic one via the use of ligatures that merged glyphs when multiple letters would represent only one sound (e.g., *th*).

[39] created a dynamic font-shaping engine named *Voice-Driven Type Design*. It allows for the modulation of a custom font’s visual attributes (vertical and horizontal stroke thickness and letter width) to echo changes in prosody, which they explored for closed captions, text messaging, and expressive visualizations of poetry. This approach was later expanded in the *WaveFont* system, described in [40], which can be used in non-specialized software for the Arabic alphabet.

2.3.2 Projects working within traditional type-setting tools

While working within the constraints of already existing type shaping engines can impose limits to how one can expressively manipulate typography, it has two important advantages: first, it allows system designers access to hundreds of thousands of typefaces and their not easily replicable features, like for example how some fonts have an extended character set to include many languages, or others have their shapes built for specific purposes (e.g. helping dyslexic readers [41], working well in small, low-resolution screens [42], etc). Second, using off-the-shelf typographic engines allows for easier integration of custom typesetting algorithms into already existing workflows.

[43] mixed discrete (font-weight and letter repetition) and continuous (font-size and word spacing) typographic attributes to represent speech. They worked with traditional fonts, modified by scripts run in the Adobe InDesign software using the extracted prosody from audio files. As did [39], both approaches were empirically evaluated, with positive recognition outcomes — although the reliance on self-reporting of the former and exaggeratedly differentiated sound files of the latter casts some uncertainty over the results, which we hope to address in our own research.

[44] ran an evaluation of what typographical attributes participants would consider the most appropriate to represent two prosodic features: loudness and pitch. They tested

four typographic parameters: font weight, baseline shift, slant, and letter width. The first two were highly ranked as representations of, respectively, loudness and pitch. Of note, this approach uses variable fonts, a technology that allows for continuous variation of glyphs’ shapes along pre-determined axes (we will discuss this further in section 3.2.1).

3 A MODEL FOR SPEECH-MODULATED TYPOGRAPHY

In this section, we will present a model that extracts meaningful acoustic features from speech and, after they are processed, uses them to modulate parameters of typographic shape and composition. An overview of this process can be seen in figure 1.

3.1 Prosodic features

3.1.1 Extraction of prosodic features

To represent speech in typography, we modeled a set of acoustic dimensions that are related to how the voice’s expressiveness changes over time. We worked with three acoustic measures: *magnitude*, *pitch*, and *duration*.

One of the advantages of representing prosody and not other acoustic features, despite the good results in emotion recognition systems obtained with other metrics (e.g., [45]), is that with prosody the mapping between acoustic feature and written units is direct. By extracting features considering the unit of one syllable, we obtain values that can be then meaningfully mapped to this same syllable’s textual representation.

In figure 1, this process is represented by **A**, which receives as inputs the sound file **1**, from which the features are extracted, and the timed transcription **2**, which defines the timestamps to subdivide **1** into each of its syllables.

The first feature is related to perceived loudness. It is calculated by the root mean square (RMS) of all samples in a given region³. Given an array with k audio samples $A = [a_1, a_2, \dots, a_k]$, RMS is calculated by:

$$A_{rms} = \sqrt{\frac{1}{k} \sum_{j=1}^k a_j^2}$$

Pitch is related to the perceived melody of the voice and was calculated through an auto-correlation algorithm, as described by [46], applied at the level of the syllable, limited to between 50 and 350 Hz (typical for a masculine voice, as was the case for the dataset used in our evaluation).

Lastly, *duration* relates to rhythm and is the simple temporal measurement of each syllable.

3.1.2 Processing of prosodic features

While magnitude, pitch and duration are measured absolutely, they are perceived relatively. This means that each syllable is perceived not by its absolute values, but rather by how it is related to its neighboring syllables. Prosody creates contrast: *this* is different than *that* [16].

3. While the two other prosodic features were calculated using audio segments that corresponded to whole syllables, for magnitude we only considered excerpts that matched vowels, excluding consonants.

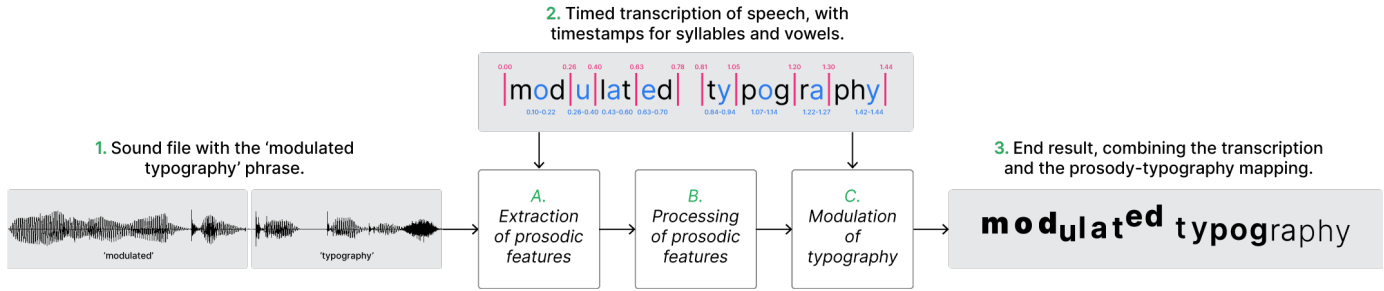


Fig. 1. Diagram of how the model combines the sound of a given spoken utterance (identified by the green 1) with its timed transcription (2), generating text that has visual cues for prosody (3). This processing is done in two parts: first, prosodic features are extracted from the audio and processed (A and B, discussed, respectively, in sections 3.1.1 and 3.1.2). Following, these values are used to modulate typographic parameters (C, discussed in section 3.2).

Because of this, to modulate typography we must use not the original prosodic measurements, but rather their relative values. In figure 1, this transformation of the original acoustic features extracted in A is represented by B. These values are obtained by normalizing each prosodic measurement considering maximum and minimum values at both an utterance and local level. For the former, we used:

$$x'_i = \frac{x_i - x'_i \min}{x'_i \max - x'_i \min}$$

where x_i is the unprocessed obtained value for each feature, and where $x'_i \min$ and $x'_i \max$ are, respectively, the minimum and maximum values for that same feature considering the whole utterance. To calculate the *local* normalization, we used:

$$x''_i = \frac{x_i - x''_i \min}{x''_i \max - x''_i \min}$$

where x''_i is a normalized value for the x_i feature that considers an unequally-spaced window of 15 syllables around it, as such:

$$x''_i \max = \max\{x_j\}_{j=i-10}^{j=i+5} \text{ and } x''_i \min = \min\{x_j\}_{j=i-10}^{j=i+5}$$

The final value was an arithmetic mean of x'_i and x''_i . An example of how the two variables produce different normalizations is seen in figure 2.

A camel cannot see its own hump

Whole utterance normalization (x')

A camel cannot see its own hump

15-word window normalization (x'')

A camel cannot see its own hump

Arithmetic mean between x' and x''

Fig. 2. The three normalizations. From top to bottom, the one that considers the whole utterance, the one that considers a window of 15 positions around the current one, and the arithmetic mean of the two. To accentuate the effect for illustrative purposes, in this figure we have used random numbers applied to each *letter* — rather than to each syllable, as was used in our model.

3.2 Modulation of typography

With the prosodic features processed and ready for use, the next step was to map them as transformations of typographic parameters. In figure 1, these features come from B

(processing of prosodic features) and, with the transcription coming from 2, are processed by C (modulation of typography) to generate the end-result in 3.

As mentioned in Section 2.3.2, in [44] participants ranked how well different visual parameters were able to represent the prosodic features of magnitude and pitch, for which a clear preference arose for the use of, respectively, font-weight and baseline shift. We followed this recommendation in our current work.

Font-weight is a parameter that sets the thickness of each letter, e.g., higher values make letters thicker, while *baseline shift* controls the vertical displacement of letters. While these changes can be discrete (e.g., most fonts will have separate files for a regular and bold version), we favored variable fonts with an axis for font-weight, allowing for continuous changes. In the example shown in figure 3, we used the *Inter* typeface; in figures 2 and 5 (and in the evaluation itself), we used the *Recursive* typeface. Both are available on the Google Fonts website [47].

Baseline shift is a compositional parameter. This means that it is not related to the shape glyphs themselves but, rather, to how they are placed on the line. As such, it is independent of the typeface used.⁴

We were still left with having to find a typographic modulation to represent the duration of syllables. While [39] and [31] explored using letter width to echo speech's rhythmic patterns, [44] points that it, along with slant, are not ideal modulation candidates — they were rarely chosen over font-weight and baseline shift and, when they were, they served more as a way to represent a voice that was flat, i.e., inexpressive, than prosodic variations per se.

Inspired by how [48] represented the duration of pauses through changes in spacing between words, we used the compositional parameter of *letter-spacing* (the horizontal spaces between each letter) to represent each syllable's duration. Following [49], who shows that decreasing letter-spacing can reduce reading speed, we worked only with positive values, avoiding *squeezing* letters together. Thus, a faster-than-average syllable was displayed with normal letter-spacing, while a slower than average one would have wider spacing.

4. It is worth mentioning, however, that fonts with long ascenders and descenders (the parts of the letters that extend beyond its main body, such as the tail in a 'j' or the hook in an 'f') will allow for smaller variations in baseline shift before lines start overlapping.

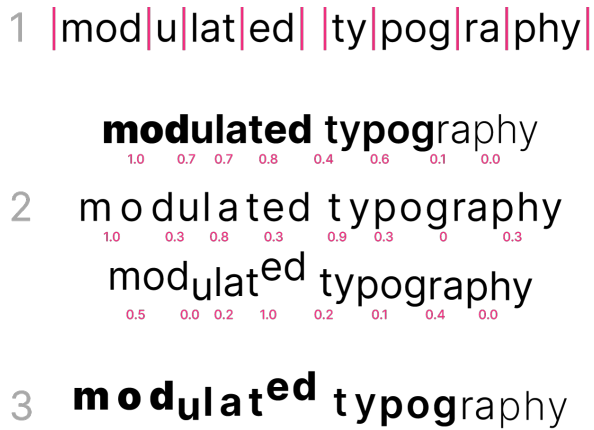


Fig. 3. The words ‘modulated typography’ are shown with a decomposition of how the model will modulate each syllable using the three typographic parameters used in our model. On the row identified by the number 1, the words are in their default state with the syllabic boundaries highlighted (for this example we are using the *Inter* typeface). On row 2, each of the three typographic modulations is represented, one in each line. They are, respectively, font-weight, letter-spacing, and baseline shift. Underneath each syllable, in pink, we marked the relative values each parameter received. On row 3, the finished product, with the three typographic modulations applied at the same time.

The complete model, with the three prosodic-typographic mappings, can be seen in action in figure 3. From the image, we can deduce that the *modulated* was spoken at a louder volume than *typography*, and that its *mod* syllable was the slowest and that its *ed* syllable the one with the highest pitch. Some intuition of these patterns can be seen in the sound wave representation seen in figure 1, which is from the recording used to generate this example.

3.2.1 Speech-modulated typography and variable fonts

To actually render the modulated typography, we decided to work with variable fonts because of the flexibility it allowed. Introduced in 2016 OpenType’s 1.8 specification, variable fonts are a format where each glyph’s shape is defined not as a collection of points with fixed x and y positions but, rather, by shifts in these x and y positions along *axes of variation* [50]. An example of this process is shown in figure 4. Note how some points travel long distances while others remain static — a change in the axis’ value can mean a big shift in position for some points, but little to none for others.

While we only used one variation axis in our model (font-weight), variable fonts allow for an arbitrary number of axes to be defined. Since each axis defines the changes in position for each point (a set of deltas) when more than one axis is present their combined values can be calculated by summing these differences, which allows for independent control of each axis. In the future this could help expand and modify the currently proposed model.

4 EVALUATION

Our evaluation consisted of using our model to create the visual representations of specific speech utterances and then showing participants these representations along with two different audio files, one the original recording, another with the same text spoken with different prosodic patterns. The

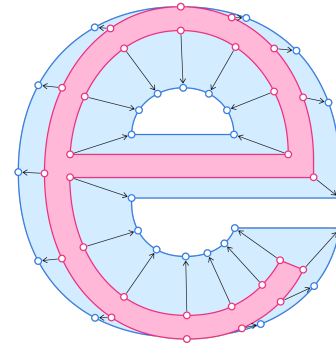


Fig. 4. Example of how a single glyph changes shape between two extremes of an axis of variation. In this case, we are showing how the font-weight axis changes the letter *e* from the *Inter* typeface, going from the lighter value of 200 (in pink) to the bolder value of 900 (in blue).

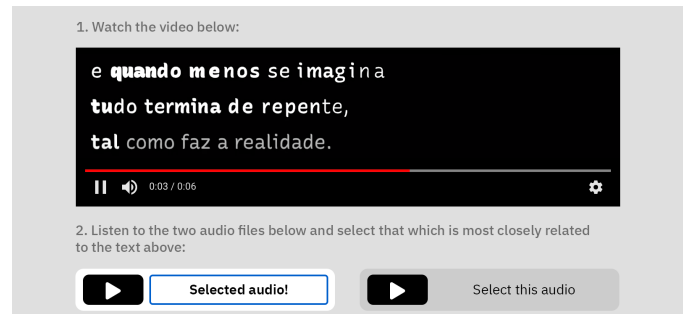


Fig. 5. A translated screenshot of one round of the test. Participants had to match one of the audio files below with the video above, which was muted. Both audio files had a reading of the same text as presented in the video, but with different prosodic emphases. The words on the video say, in Portuguese, ‘and when one least expects it / everything ends suddenly, / as does reality.’ [51, p. 75]

goal was to measure how successful participants would be in matching the speech-modulated typography to its originating audio, which was only subtly different from the other audio it was presented alongside with.

We proposed this experimental method to evaluate whether our model allowed for the *recognition of vocal prosody*: a greater-than-chance rate of success for this matching between the typography and its corresponding, correct audio, would indicate that the model’s visual representation of prosody was sufficiently clear to allow participants to decode its correspondence with audio. (The evaluation was reviewed and approved by an IRB.)

4.1 Expressive prosody dataset

The sound recordings used in the test were created for the purposes of our evaluation. Like [43], we decided to base our evaluation on readings of poetry, since they can accommodate a greater variation in vocal expression without sounding *artificial* — even if the poem’s author had a specific set of emphases in mind when writing a poem, readers are free to have their own different interpretations, a freedom we took advantage of. However, because we need two different readings for each poetic stanza, we realized we would need to create our own dataset.

We instructed a hired actor to repeatedly record a reading of the same poem, at each round emphasizing different

aspects of their voice: loudness, pitch, rhythm, a mix of the three, a monotonic reading, a naturalistic reading. With these different versions, we sliced the audio files into word-sized units, which we recombined into the two final versions of each poem, making sure that the three prosodic features had heterogeneous patterns between the readings.⁵

4.2 Rendering speech-modulated typography

To display the modulated typography generated after processing the audio files, we created a script that ran in the browser, receiving as input the video and a subtitle file that, beyond the transcription itself, contained cues about how to change the typography of each syllable, how long these changes should take (i.e., the duration of the syllable), and at which moment (i.e., the moment that syllable was said).

Since variable fonts are already supported in all major browsers [53], by working in a web environment we were able to take advantage of certain built-in facilities, such as how a native HTML video player can fire syncing events related to when each subtitle text block is supposed to enter or exit the screen once the video is played, or how CSS allows for typographic parameters to change with an eased, non-jarring animation. However, while we were able to run the animations in real-time in our own setups, to accommodate for participants in lower-end systems such as older mobile phones, we decided to screen capture the animations and make them available as YouTube videos.

4.3 Test platform

In the experiment, participants would be randomly assigned to match either animated or static instances of speech-modulated typography. The test itself consisted of 15 rounds of sound-and-image matching. Participants would be given no instructions about the prosodic-typographic mappings used by our model, which they would have to figure out on their own. Figure 5 shows a screenshot of the test in action.

We divided our participants into two groups: the first received a *static-image* based version of the test, i.e., each of the 15 stanzas was displayed as an image to which all typographic modulations were already applied. This is similar, for instance, to how the typography is shown in the third row of figure 3.

The second group of participants did an animated closed-caption version of the test. This consisted of the same stanzas and audio files, but each version of speech-modulated typography was provided as an embedded, muted YouTube video. This video started with a version of the stanza in its default state (i.e., neutral typographic parameters), and as each syllable sounded its typographic parameters animated towards their end state (see figure 6 for a frame-by-frame example).

By showing changes synchronously with audio, these animations emphasize the rhythmic aspect of the model, but it can also make it noisier to parse, possibly adding a cognitive cost for readers. Having both this animated and the static version of the test running in parallel aimed to answer our second research question.

5. A pilot version of our experiment showed us that these versions had to have roughly the same duration otherwise participants would try to sync video and audio to find the correct match, an effect that [52] had already found in her own similar experiment.

4.4 Open-ended comments

At the end of the evaluation, participants were prompted to leave us comments about ‘the test, speech-modulated typography, its possible applications, or whatever else they wanted.’ This step was optional, but nevertheless 43 participants took it. We analyzed these messages using the thematic analysis approach, outlined in [54]: after familiarizing ourselves with the messages, we generated an initial batch of codes, which were then reviewed and collated having our research questions in mind, but also our additional goal of helping to uncover limitations of our model and methodological approach to evaluate it. Lastly, these codes were summarized into two themes, which we will present in section 5.3, and discuss in section 6.1.

5 RESULTS

5.1 Participants

We ran the experiment from November to December 2020. Participants were recruited through social media and academic mailing lists. Of the 117 participants that concluded the test, 47% were female, 51% male, and 2% non-binary; 90% had undergraduate degrees or higher; 25% had between 18 and 24 years, 49% between 25 and 39, 23% between 40 and 59, and 3% had 60 or more years. 52% of the participants were randomly assigned the static-image-based test, while 48% received the animated closed captions variant.

5.2 Rates of success

The average rate of correct answers for the 61 participants who took the static-image-based test was 67% (95% CI [.64, .70]), slightly higher than the 63% (95% CI [.59, .67]) obtained by the 56 participants who took the animated closed captions variant, as shown in figure 7. This 4% difference between the two distributions was found to be not significant ($t(115) = 1.54, p > 0.05$).

5.2.1 Effect size

To simulate a control group we assumed that an ineffective prosodic-typographic model would have no effect on participants, generating results indistinguishable to those of a random distribution. We compared this simulated data to our actual results to measure the effect size of each of our two approaches⁶. Both the static-image-based test ($t(59) = 4.18, p < 0.05, d = 0.75$) and the animated closed caption-based variant ($t(54) = 2.91, p < 0.05, d = 0.55$) were significantly different from the random distribution, with effect sizes of *medium* magnitudes (Cohen’s d between 0.5 and 0.8).

5.3 Open-ended comments

Sending us comments at the end of the evaluation was optional. Nevertheless, 43 out of the 117 participants left their thoughts. In the quotes below, participants will be marked by the letter P along with an identifying number. We have translated the comments, which were originally in Portuguese, and, when necessary, edited for clarity. We summarized the data using two themes, which name the following sections.

6. Results reported averaged over 1,000 runs of the statistical tests, after which fluctuations of the randomly generated values stabilized.

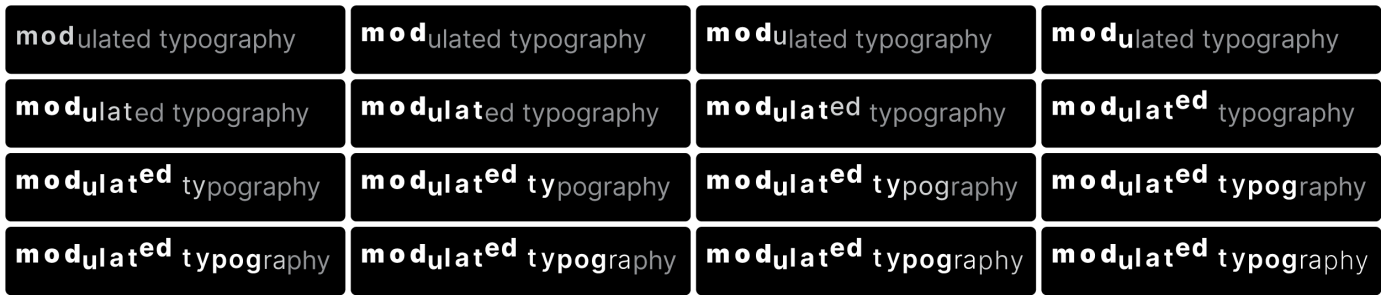


Fig. 6. Example of how each syllable’s typographic parameters change in time, as was shown to participants that did the animated closed-caption based test. The image is meant to be read from left to right and, subsequently, from top-to-bottom.

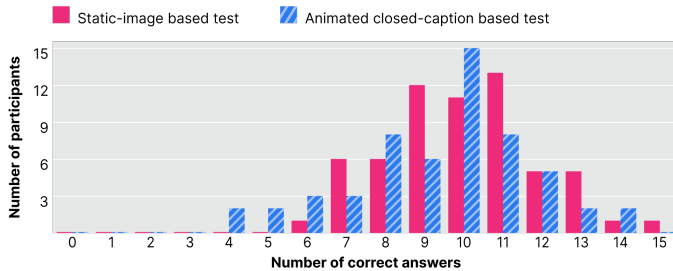


Fig. 7. Distribution of how many participants (in the y axis) got how many correct answers (in the x axis), considering both versions of the test.

5.3.1 Theme 1: Divergent interpretations

Regardless of how well they did towards correctly matching audio and typography, we received comments both complimenting how easy the model was to understand or criticizing how obscure it was. Among the discontents, P1 commented that baseline shift was not an intuitive modulation and that some more subtle differences in intonation became ambiguous, which echoed P2, that among all modulations found those the hardest to associate with the audio. P3 tried, but could not find a clear logic behind the modulations, something also present in P4’s comment that ‘each round seemed to have a unique pattern,’ and in P5’s that ‘the meaning behind how letters moved wasn’t clear.’

Some participants found the model easier to grasp. P6 understood that thicker letters represented louder sounds, while P7, P8, and P9 all managed to correctly match the three typographic modulations (font-weight, baseline shift, letter-spacing) to their corresponding prosodic features (respectively, magnitude, pitch, and duration). P9, particularly, said that:

When seeing the video without the audio and then comparing the two [audio] options it was perfectly possible to interpret, using only the sound, when the actor paused, elongated words, shouted, made their voice thicker or thinner. It was not possible, however, to extract from the text whether his voice was calm, commanding, irritated, etc.

5.3.2 Theme 2: Expansions to the model

Some participants felt that the model should be changed and/or expanded to include additional typographic modulations. P10 suggested that we could have tested also using uppercase words to represent emphases. P1 suggested that

we should have explored a visual language closer to that of comic books. Similarly, P7 thought that to be able to represent voice qualities such as a raspy voice, the fonts should have pointy corners, while a soft voice could have been represented by smoother corners (similar, then, to [32]’s proposed font). They go on to say that

[it should be] more cartoonish, more expressive. How would these captions represent different voices that are feminine, elderly, young, stuttery, twangy, sleepy, non-native, [etc]?

There were comments that made reference to how our experiment’s format made the text hard to understand. P11 told us that at each step they felt the need to know how the actor had read aloud the previous stanzas, since understanding ‘if a clause is subordinate or coordinate (...) has great impact in how the voice is modulated.’

We received some comments about the model being hard to read. P5, for instance, suggested we ‘make the text bigger and more legible.’ P12 mentioned being dyslexic, and that they had a hard time reading the texts,

since the movement of the letters and changes to their typography disrupted me, and I had to watch the videos many times before I was able to understand them.

P13 wrote that, when words had their letter-spacing increased, it became hard to differentiate one word from the other:

For example: at first glance, I considered that ‘A D O R⁷’ was a single word, waiting for the rest [that would complete the word], but eventually I realized that there were two words, ‘A D O R⁸.’ I had to read the captions again to make sense of them.

6 DISCUSSION

We asked whether *visual modifications in the shapes of letters composing a text could allow for the recognition of vocal prosody*. Our experiment was able to measure a strong-enough effect to affirm that yes: participants were able to intuitively grasp our prosodic-typographic mappings in such a way as to be able to *reverse-engineer* the images and find their originating audio files. This reinforces [44]’s proposed mapping of font-weight and baseline shift to, respectively, magnitude and

7. This is not a word in Portuguese.

8. ‘The pain,’ in Portuguese.

pitch, as it does to our addition of letter-spacing mapping changes in duration.

Curiously, and in answering our questions of whether there would be *differences in recognition between animated and static versions of speech-modulated typography*, the animated typography did not seem to make a difference in terms of increasing (or decreasing) recognition performance. While this is an indication that static-image-based speech-modulated typography could find useful applications on its own, it is also surprising considering that in the animated tests there were additional signals helping tie images to sound.

6.1 What did participants tell us?

While our quantitative results point to our model's success in conveying speech through typography, participants' answers paint a more nuanced picture. Many of them were able to correctly identify each of the model's mappings between prosody and typography, while others found it at times inconsistent, with some specifically pointing out problems with our use of baseline shift. These divergent interpretations show that, even if the model as a whole was effective, more work is needed to understand how well each typographic parameter is able to represent its corresponding prosodic feature, and whether there are perceptual and semantic changes when they are used in tandem.

Moreover, and perhaps surprisingly, some participants felt the model should be able to embody more nuanced and complex representations of speech, which it should do with a more varied and expressive visual vocabulary. While this is not a simple task, it is nevertheless encouraging to imagine how different applications of speech-modulated typography could branch out as this new field develops.

Care must be taken, however, to whether the modulations can hinder legibility. The use of animations, and the change of positioning and spacing parameters caused discomfort to some participants. Future work is needed to understand how models such as the one we propose here could make text harder to read, and applications of these models should consider giving the user options to either tone down the modulations or turn them off completely.

7 CONCLUSION

We have proposed a model for transforming acoustic cues present in speech into visual modulations of typography, allowing for a transcription of not only words but also the paralinguistic dimensions of an utterance. This model was evaluated with hearing individuals that, having had no previous training, were able to use this speech-modulated typography to correctly match speech and typography. The model — which mapped magnitude, pitch, and rhythm to, respectively, font-weight, baseline shift, and letter-spacing — worked equally well when its modulations were presented as animations or static images.

Although we envision that our approach could be beneficial to DHH individuals, the fact that we conducted an evaluation with only hearing participants limits how generalizable our results are, particularly considering DHH persons — it is plausible, for instance, that a firsthand based intuition of speech might be needed for one to make sense of our model.

Future experiments should attempt to investigate how well the model performs with DHH persons, and what adjustments could be made to it considering their perceptions of speech-modulated typography. Considering how ASR-based captioning systems have made in-roads in many contexts, understanding if and how our approach could bring benefits to specific settings could also be fruitful.

The claims by some participants about how the model is at times difficult to read should be investigated.

Lastly, and in addition to simply measuring the model's *recognition rate*, an important avenue for future exploration should also include attempting to understand what *effects* these speech-modulated closed captions could have on viewers' subjective experience, e.g., immersion in film-media, quality of communication in online meetings.

REFERENCES

- [1] P. Biswas, P. Langdon, J. Umadikar, S. Kittusami, and S. Prashant, "How interface adaptation for physical impairment can help able bodied users in situational impairment," in *Inclusive Designing*. Springer, 2014, pp. 49–58.
- [2] M. Seidenberg, *Language at the Speed of Sight: How We Read, Why So Many Can't, and What Can Be Done About It*, 1st ed. New York: Basic Books, 2017, kindle version.
- [3] G. McCulloch, *Because Internet: Understanding the New Rules of Language*, 1st ed. New York: Riverhead Books, 2019, kindle version.
- [4] R. Nünlist, "Users of literature," in *A companion to Greek Literature*, M. Hose and D. Schenker, Eds. West Sussex: John Wiley & Sons, 2016, ch. 19, pp. 296–297.
- [5] M. W. Küster, "Writing beyond the letter," *Tijdschrift voor Mediageschiedenis*, vol. 19, no. 2, 12 2016.
- [6] R. W. McCutcheon, "Silent reading in antiquity and the future history of the book," *Book History*, vol. 18, no. 1, pp. 1–32, 2015. [Online]. Available: <https://doi.org/10.1353/bh.2015.0011>
- [7] R. S. Kushalnagar and C. Vogler, "Teleconference accessibility and guidelines for deaf and hard of hearing users," in *The 22nd International ACM SIGACCESS Conference on Computers and Accessibility*, 2020, pp. 1–6.
- [8] V. Murphy-Berman and L. Whobrey, "The impact of captions on hearing-impaired children's affective reactions to television," *The Journal of Special Education*, vol. 17, no. 1, pp. 47–62, 1983.
- [9] F. Loizides, S. Basson, D. Kanevsky, O. Prilepova, S. Savla, and S. Zaraysky, "Breaking boundaries with live transcribe: Expanding use cases beyond standard captioning scenarios," in *The 22nd International ACM SIGACCESS Conference on Computers and Accessibility*, 2020, pp. 1–6.
- [10] J. R. Mallory, M. Stinson, L. Elliot, and D. Easton, "Personal perspectives on using automatic speech recognition to facilitate communication between deaf students and hearing customers," in *Proceedings of the 19th International ACM SIGACCESS Conference on Computers and Accessibility*, 2017, pp. 419–421.
- [11] M. Seita, K. Albusays, S. Kafle, M. Stinson, and M. Huenerfauth, "Behavioral changes in speakers who are automatically captioned in meetings with deaf or hard-of-hearing peers," in *Proceedings of the 20th International ACM SIGACCESS Conference on Computers and Accessibility*, 2018, pp. 68–80.
- [12] E. J. McDonnell, P. Liu, S. M. Goodman, R. Kushalnagar, J. E. Froehlich, and L. Findlater, "Social, environmental, and technical: Factors at play in the current use and future design of small-group captioning," *Proceedings of the ACM on Human-Computer Interaction*, vol. 5, no. CSCW2, pp. 1–25, 2021.
- [13] M. A. Gernsbacher, "Video captions benefit everyone," *Policy insights from the behavioral and brain sciences*, vol. 2, no. 1, pp. 195–202, 2015.
- [14] M. A. Buchanan, "@ face value//expanding our typographic repertoire," *Communication Design*, vol. 3, no. 1, pp. 27–50, 2015.
- [15] D. Wilson and T. Wharton, "Relevance and prosody," *Journal of pragmatics*, vol. 38, no. 10, pp. 1559–1579, 2006.
- [16] P. A. Barbosa, "Conhecendo melhor a prosódia: aspectos teóricos e metodológicos daquilo que molda nossa enunciação," *Revista de Estudos da Linguagem*, vol. 20, no. 1, pp. 11–27, 2012.

- [17] —, *Prosódia*. Parábola Editorial, 2019.
- [18] J. A. de Moraes and A. Rilliard, “Prosody and emotion in brazilian portuguese,” in *International Grammar in Ibero-Romance*. John Benjamins, 2016, pp. 135–152.
- [19] W. d. Silva, P. A. Barbosa, and Á. Abelin, “Cross-cultural and cross-linguistic perception of authentic emotions through speech: An acoustic-phonetic study with brazilian and swedish listeners,” *DELTA: Documentação de Estudos em Linguística Teórica e Aplicada*, vol. 32, no. 2, pp. 449–480, 2016.
- [20] S. G. Koolagudi and K. S. Rao, “Emotion recognition from speech: a review,” *International journal of speech technology*, vol. 15, no. 2, pp. 99–117, 2012.
- [21] M. El Ayadi, M. S. Kamel, and F. Karray, “Survey on speech emotion recognition: Features, classification schemes, and databases,” *Pattern recognition*, vol. 44, no. 3, pp. 572–587, 2011.
- [22] K. S. Rao, R. Reddy, S. Maity, and S. G. Koolagudi, “Characterization of emotions using the dynamics of prosodic features,” in *Speech Prosody 2010-Fifth International Conference*, 2010.
- [23] L. Stark and J. Hoey, “The ethics of emotion in artificial intelligence systems,” in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 2021, pp. 782–793.
- [24] P. A. Perez-Toro, J. C. Vasquez-Correa, T. Bocklet, E. Noth, and J. R. Orozco-Arroyave, “User state modeling based on the arousal-valence plane: Applications in customer satisfaction and healthcare,” *IEEE Transactions on Affective Computing*, 2021.
- [25] D. Ong, Z. Wu, Z.-X. Tan, M. Reddan, I. Kahhale, A. Mattek, and J. Zaki, “Modeling emotion in complex stories: the stanford emotional narratives dataset,” *IEEE Transactions on Affective Computing*, 2019.
- [26] K. Boehner, R. DePaula, P. Dourish, and P. Sengers, “Affect,” in *Proceedings of the 4th decennial conference on Critical computing between sense and sensibility - CC '05*. ACM Press, 2005. [Online]. Available: <https://doi.org/10.1145/1094562.1094570>
- [27] —, “How emotion is made and measured,” *International Journal of Human-Computer Studies*, vol. 65, no. 4, pp. 275–291, 2007.
- [28] J. Ängeslevä, C. Reynolds, and S. O’Modhrain, “Emotemail,” in *ACM SIGGRAPH 2004 Posters*, 2004, p. 9.
- [29] J. dos Reis and V. Hazan, “Speechant: a vowel notation system to teach english pronunciation,” *ELT Journal*, vol. 66, no. 2, pp. 156–165, Jun. 2011. [Online]. Available: <https://doi.org/10.1093/elt/ccr019>
- [30] W. Verbaenen, “Phonotype. the visual identity of a language according to its phonology,” Master’s thesis, PXL-MAD, 2019.
- [31] A. Bessemans, M. Renckens, K. Bormans, E. Nuyts, and K. Larson, “Visual prosody supports reading aloud expressively,” *Visible Language*, vol. 53, no. 3, pp. 28–49, 2019.
- [32] S. Promphan, “Emotional type: Emotional expression in text message,” Master’s thesis, Basel School of Design, Switzerland, 2017.
- [33] R. Rashid, Q. Vy, R. Hunt, and D. I. Fels, “Dancing with words: Using animated text for captioning,” *Intl. Journal of Human-Computer Interaction*, vol. 24, no. 5, pp. 505–519, 2008.
- [34] J. Forlizzi, J. Lee, and S. Hudson, “The kinedit system: affective messages using dynamic texts,” in *Proceedings of the SIGCHI conference on Human factors in computing systems*, 2003, pp. 377–384.
- [35] K. Bodine and M. Pignol, “Kinetic typography-based instant messaging,” in *CHI’03 extended abstracts on Human factors in computing systems*, 2003, pp. 914–915.
- [36] A. Albert, F. Cangemi, and M. Grice, “Using periodic energy to enrich acoustic representations of pitch in speech: A demonstration,” in *Proceedings Speech Prosody*, vol. 9, 2018, pp. 13–16.
- [37] A. Öktem, M. Farrús, and L. Wanner, “Prosograph: a tool for prosody visualisation of large speech corpora,” in *Proceedings of the 18th Annual Conference of the International Speech Communication Association (INTERSPEECH 2017); 2017 Aug. 20-24; Stockholm, Sweden. Baixas: ISCA; 2017. p. 809-10*. International Speech Communication Association (ISCA), 2017.
- [38] T. Rosenberger-Shankar and R. L. MacNeil, “Prosodic font: translating speech into graphics,” in *CHI’99 Extended Abstracts on Human Factors in Computing Systems*, 1999, pp. 252–253.
- [39] M. Wölfel, T. Schlippe, and A. Stütz, “Voice driven type design,” in *2015 international conference on speech technology and human-computer dialogue (SpeD)*. IEEE, 2015, pp. 1–9.
- [40] T. Schlippe, S. Alessai, G. El-Taweel, M. Wölfel, and W. Zaghouni, “Visualizing voice characteristics with type design in closed captions for arabic,” in *2020 International Conference on Cyberworlds (CW)*. IEEE, 2020, pp. 196–203.
- [41] X. Zhu, K. Kageura, and S. Satoh, “Analysis of typefaces designed for readers with developmental dyslexia,” in *International Workshop on Document Analysis Systems*. Springer, 2020, pp. 529–543.
- [42] J. Benson, K. Olewiler, and N. Broden, “Typography for mobile phone devices: The design of the qualcomm sans font family.” New York, NY, USA: AIGA: American Institute of Graphic Arts, 2005.
- [43] J. C. Castro *et al.*, “Máquina de ouver—representação do discurso oral pela tipografia,” Master’s thesis, Universidade de Coimbra, 2019, license: Creative Commons BY-NC-ND 4.0.
- [44] C. de Lacerda Pataca and P. D. P. Costa, “Speech modulated typography: towards an affective representation model,” in *Proceedings of the 25th International Conference on Intelligent User Interfaces*, 2020, pp. 139–143.
- [45] F. Eyben, K. R. Scherer, B. W. Schuller, J. Sundberg, E. André, C. Busso, L. Y. Devillers, J. Epps, P. Laukka, S. S. Narayanan *et al.*, “The geneva minimalistic acoustic parameter set (gemaps) for voice research and affective computing,” *IEEE transactions on affective computing*, vol. 7, no. 2, pp. 190–202, 2015.
- [46] P. Boersma *et al.*, “Accurate short-term analysis of the fundamental frequency and the harmonics-to-noise ratio of a sampled sound,” in *Proceedings of the institute of phonetic sciences*, vol. 17, no. 1193. Citeseer, 1993, pp. 97–110.
- [47] R. Andersson and S. Nixon, “Recursive and inter typefaces,” accessed on January 5, 2022. [Online]. Available: <https://fonts.google.com/share?selection.family=Inter%7CRecursive>
- [48] J. C. e. Castro, P. Martins, A. Boavida, and P. Machado, “Máquina de ouver—from sound to type: Finding the visual representation of speech by mapping sound features to typographic variables,” in *Proceedings of the 9th International Conference on Digital and Interactive Arts*, 2019, pp. 1–8.
- [49] S. T. Chung, “The effect of letter spacing on reading speed in central and peripheral vision,” *Investigative Ophthalmology & Visual Science*, vol. 43, no. 4, pp. 1270–1276, 2002.
- [50] M. Jacobs and P. Constable, “Opentype specification version 1.8,” <https://docs.microsoft.com/en-us/typography/opentype/otspec18/>, 2018, accessed on January 5, 2021.
- [51] P. Britto, *Macau*. São Paulo, Brazil: Companhia das Letras, 2003.
- [52] T. Rosenberger-Shankar, “Prosodic font: The space between the spoken and the written,” Master’s thesis, Massachusetts Institute of Technology, 1998.
- [53] A. Deveria, “Can i use: Variable fonts?” <https://caniuse.com/#feat=variable-fonts>, 12 2021, accessed December 20, 2021.
- [54] V. Braun and V. Clarke, “Using thematic analysis in psychology,” *Qualitative research in psychology*, vol. 3, no. 2, pp. 77–101, 2006.



Caluá de Lacerda Pataca received his MSc degree in Computer Engineering from the University of Campinas, Brazil, in 2021, and is currently working towards a Ph.D. degree at the Computing and Information Sciences department at the Rochester Institute of Technology, USA. His research interests include speech accessibility, visual design, and human-computer interaction.



Paula Dornhofer Paro Costa received her Ph.D. degree in Computer Engineering from the University of Campinas (Unicamp), Brazil, in 2015. In 2016, she joined the Dept. of Computer Engineering and Automation (DCA) of the School of Electrical and Computer Engineering, at Unicamp, as a research scientist and assistant professor. Her research interests focus on affective computing and multimodal artificial intelligence.