

## Coda: Distant Listening and Resonance

Way back there when Hell was no bigger than Maitland, man found out something about the laws of sound. He had found out that sound could be assembled and manipulated and that such a collection of sound forms could become so definite and concrete as a war-ax or a food tool.

—Zora Neale Hurston (1999, 70–71)

I have described *distant listening* to literary texts as using computing to “distill the many-layered four-dimensional space of the text in performance (i.e., embodied within the performance network of interpretations with the listener in time and space) into a two-dimensional script called ‘code’” (Clement 2013). My definition strikes me as remarkably similar to Zora Neale Hurston’s description of folklore as “the boiled down juice of human living” (1999, 69). My definition is not as poetic or as entertaining as Hurston’s, of course, and when I wrote this definition, I was “listen[ing] in print” and referring to “the audio relationship between reader and text” (Furlonge 2011), but I believe the definition still holds when applied to audio files. In the definition, I think I could just as easily have used the word “data” and “algorithms” rather than “code.” Melanie Feinberg helpfully describes data as “just a description of a thing’s qualities”—it’s not the thing; it describes the thing—and algorithms as both mathematical operations and, humorously, mediation via “arcane spells” (2022, 3, 230, 3). To distill the many varied experiences of sound as described in this book into data operated on by algorithms, then, would be to filter the experience of human living down into a purified list (or matrix) of the signifying qualities of that sound and to deploy those via a mathematical operation intended to produce resonance or meaning. An updated definition might be that

*distant listening* is at root a technically complex matter of fitting a mathematical abstraction to collections of sounds at scale and trying to determine their resonances computationally as data.<sup>1</sup>

When this book was first reviewed for publication, one of the reviewers was dismayed to find out that the book did not make mention of distant listening. Why, the reviewer wondered, when I had spent so many years on big-team, funded grants studying the possibilities for distant listening and the use of cultural-analytic tools to develop knowledge about large caches of recorded audio, would I turn away from that prior work? In the past decade, I have been involved in multiple computational projects with archival audio collections through a project I lead called High Performance Sound Technologies for Access and Scholarship (HiPSTAS).<sup>2</sup> Some of the early HiPSTAS work was focused on using text-to-speech software to “listen” to literature (Clement et al. 2013), but much of the HiPSTAS research focused on using Adaptive Recognition with Layered Optimization (ARLO), a machine-learning application developed by my collaborator David Tcheng for analyzing large sound collections by extracting basic prosodic features such as pitch, rhythm, and timbre for clustering and classification (Clement et al., 2016). In May 2013 and May 2014, with funding from the National Endowment for the Humanities (NEH), HiPSTAS brought together twenty humanities junior and senior faculty and advanced graduate students as well as librarians and archivists from across the United States interested in developing and using ARLO to access and analyze spoken word audio collections.<sup>3</sup> The HiPSTAS Institute had two primary outcomes: (1) participants would produce new scholarship using audio collections with advanced technologies, such as classification, clustering, and visualizations; and (2) participants would engage in the scholarly work of digital infrastructure development by contributing to recommendations for the implementation of a suite of tools for cultural heritage institutions interested in supporting advanced digital scholarship in sound. Because (as the previous chapters show) access and scholarship with sound are entangled practices, we worked to develop ARLO in a second round of funding in 2014 to help users automate metadata description for undescribed sound collections (Clement 2018; Clement et al. 2018). Later, I led two more HiPSTAS projects through a project at the University of Texas at Austin called Good Systems, which concerned the ethics of artificial intelligence. The first project was focused on using machine learning to automatically generate metadata (Xu et al.

2020) and the second was concerned with training librarians and humanists on the ethics of data operations in the archives and libraries (Clement et al. 2022). It seems, the reviewer challenged, that given my experience, the book would provide a chance to comment productively on the debates around cultural analytics that are raging in literary studies right now.<sup>4</sup>

Perhaps. In my experience, machines cannot access or analyze the features of the recordings (silences, distortions, interferences, compressions, modes of reception) in the ways I found resonant for this book. That said, creating a computational model of sounds with the explicitness and consistency that computation requires can be generative. Computational models “force us to confront the radical difference between what we know and what we can specify computationally, leading to the epistemological question of *how we know what we know*” (McCarty 2004), and in this sense, my own distant listening projects have yielded productive inquiries into the cultural registry or the processes of archives, culture, data structures, information infrastructures, institutions, sound, culture, computation and other knowledge production systems.<sup>5</sup>

This coda is not a comment on *what* data to compute or *which* algorithms to use—instead, I have become more interested in the *so what?* of computational analysis with historical spoken word texts in the archives.<sup>6</sup> To the reviewer’s query, then, *Dissonant Records* as a whole is my response: interpretive practices with audio are difficult and subjective and always already entangled with the personal, cultural, sociopolitical, and technological context of listening as an agential process. Karen Barad writes that “Mattering is simultaneously a matter of substance and significance,” explaining that this is “why contemporary physics makes the inescapable entanglement of matters of being, knowing, and doing, of ontology, epistemology, and ethics, of fact and value, so tangible, so poignant” (2007, 3). Likewise, I find myself interested in distant listening in so far as its pursuit has resonances with larger questions about epistemologies, ontologies, and other areas of study concerned with better understanding dissonances in the cultural registry and imaginary. Distant listening is interesting to me to the extent that its processes hold traces of these entanglements. Thus, while I do not believe computers can do what I do in this book, I do believe it is essential for humanities scholars to better understand the limits and potentials of computational sound analysis for interpretive analysis, because these methods point to the limits and potentials of scholarly research with audio more

generally. Consequently, in this coda, I invite scholars to consider possible research questions about distant listening using the same principles I used to introduce resonance at the beginning of this book.

### The Material Particularities of the Apparatus Matter

What, exactly, is the matter of audio data? What gets filtered in the process of digitizing sound? Modeling sound for computational methods is an algorithmic process of discretization. Sound is air pressure variation over time. Ears and hands can turn the pressure differences into neural activations while microphones create digital sound by translating these variances into voltage differences. An audio signal is a sequence of mathematical abstractions that map voltage (or pressure) over time in a wave, and frequency is the number of times per second that a sound pressure wave repeats itself (McFee 2020a). To make things more complicated, audio signal processing methods do not work with continuous signals. Instead, before being processed by a computer, sound engineers and computational analysts discretize the sound pressure wave through sampling and quantization, using a mathematical representation of the continuous signal through samples that indicate the whole sound without capturing it fully (McFee 2020b). The sampling process is considered more or less precise when more or fewer discrete samples are used to represent a signal across a period of time, but all the information (all the qualities) is never represented at once. Imagine a sound event as a pebble dropped in a still lake. Digitization would be concerned with the rings that fan around the pebble, in a mathematical matrix of variables that represents the physical dimension of the plop, including information about the water depth, the air temperature, and whatever else impacts the shape, size, and speed of expansion of the rings. Indeed, reconstructing a sound event from a Fourier transform, which describes a sound's frequency mathematically,<sup>7</sup> is entirely possible. In reality, however, while the digitized transformation of the plop has a large enough matrix of its qualities to replay the plop, it is not the plop.

Sampling and quantization imply absence: I can imagine adding *sampling* as a cultural keyword to my list of sound technologies. As *Dissonant Records* has indicated, technical operations such as amplification, distortion, interference, compression, and reception reveal absences, historical absences based on media type and genre, information infrastructures, professional

protocols, institutional values and social mores, and personal experience. The process of sampling is, by definition, similarly lossy. I get the sense that when Hurston refers to the “boiled down juice of human living,” that juice comprises the traces of life that resonate throughout this book: the violences and offenses, but also the reconciliations, the passion, and the love. Adding water to these qualities of life might make a muddy mess. Likewise, sound can be reconstructed from its frequency samples by computational algorithms, but something is gone. Something new is generated. A scholar might ask: What present absences does sampling reveal?

### Meaning Making Is Dialogic

Distant listening research requires a dialogue between humanities scholars, information professionals, signal processing engineers, and machine learning scientists who are all interested in making matter and making meaning with audio differently. On HiPSTAS, I have worked with anthropologists, archivists, data analysts, historians, librarians, literary scholars, machine-learning scientists, musicologists, ornithologists, and poets, among others. When humanists listen, they abstract from sonic event features to consider how events like speech or music influence an understanding of cultural phenomena (Bernstein 2011; Clement and McLaughlin 2016; Francis et al. 2016, 354; MacArthur, Zellou, and Miller 2018; Mustazza 2014, 2016, 2018). When signal processing scientists listen, they consider damping ratios, gain, frequencies, spectra, and pitch energy, and talk about how these features influence sound fidelity. Machine-learning scientists are listening for feature selection, models, clustering and classification, correlation and probability, validation, and optimization (Clement et al. 2014). Librarians and archivists want to consider how large datasets of sound features could potentially be used for basic tasks that facilitate discovery with audio collections, including creating transcriptions, metadata generation (event detection, keyword extraction, speaker disambiguation or diarization, and speaker recognition), and quality analysis (Dunn et al. 2018; Gref, Köhler, and Leh 2018; Harrington 2019; Oard 2012).

Digital humanities projects like HiPSTAS and the Digging into Data Challenge (2009–2019) provide a sample of collections, methods, and outcomes in distant listening research with historical audio collections. Funded by agencies across the Western Hemisphere,<sup>8</sup> most of the projects in the Digging into Data Challenge analyzed image and text; a few provided new

methods for discovery with audio files, such as the Structural Analysis of Large Amounts of Music (SALAMI), which included approximately 50,000 hours of a wide variety of music genres from different venues, including live concert recordings, folk, jazz, orchestral, and twentieth-century avant-garde music from across the world. The SALAMI researchers attempted to use machine-learning algorithms to discern musical genres, musical similarity (to mark repeated themes), function (e.g., interlude) and lead instrument (e.g., vocal or guitar) (Bay et al. 2009). Analyzing natural language usage, the Mining a Year of Speech project used approximately 9,000 hours of recorded American and British speech, including the British National Corpus<sup>22</sup> and holdings of Penn Linguistic Data Consortium, together with transcriptions. The goal of the project was to assess the challenges of working with large-scale digital audio collections of spoken word. To this end, the researchers decided that a first important step in addressing these challenges was to develop a technology for forced alignment between the audio and the transcripts that would yield a phonemic transcription with detailed timing information, including the start and end of every vowel, consonant, and word. They used this data to consider sex differences in conversational speaking rates, the differences in phrasal speaking rates across genres, dialects, and languages (Coleman et al. 2011). The Harvesting Speech Datasets for Linguistic Research on the Web project comprised hundreds of 25-second snippets from podcasts, news broadcasts, and public and educational lectures gathered across the web. The researchers collected the datasets “to evaluate hypothesized correlations between acoustic form and grammatical and contextual features, and to identify the particular acoustic features (such as pitch, duration, intensity, or vowel quality) that are significant in marking prosodic distinctions.”<sup>9</sup>

This same literature indicates that the accuracy and efficacy of these methods remains inconclusive and that scholars are in the early days of epistemological questions surrounding accuracy and efficacy, especially with historical speech recordings (Joudrey, Taylor, and Wisser. 2018; Marin-ganti 2017; Mascaro 2011; Svenonius and McGarry 1993; Xu et al. 2020). No matter how scientists manipulate damping ratios, gain, frequencies, spectra, and pitch energy, transcription accuracy is influenced by recording quality, accents, and the presence of background noise (among other signals). Data gathering, cleaning, structuring, and feature selection are other roadblocks that Williford and Henry (2012) identify in the SALAMI, Mining a Year of

Speech, and Harvesting Speech Datasets for Linguistic Research on the Web projects. In addition, recent work in the broader field of acoustics, speech, and signal processing relies on features generated or learned as a byproduct of training large-scale deep networks for some general tasks like acoustic scene classification or source identification (Baevski et al. 2020; Cramer et al. 2019), but inductive bias—assumptions about the data that are encoded in the model to learn the target function and to generalize beyond training data—persist, obfuscated in black-box methods that can often resist critical intervention. Indeed, in the textbook *Speech and Audio Signal Processing: Processing and Perception of Speech and Music* (2011), Ben Gold, Nelson Morgan, and Dan Ellis describe the history of speech transmission in terms of a balance between meaning and matter in which infrastructural resources are also a significant variable: “Speech,” they write, “conveys emphasis, emotion, personality, etc., and we still don’t know how much bandwidth is needed to transmit these kinds of information” (21). In addition, the analysis of such large datasets requires computing resources to which very few scholars have access (Clement et al. 2016).<sup>10</sup>

This sample of distant listening research marks promising inquiries that could potentially alter scholarship with audio, but in their points of interference across scholarly disciplines, these projects also expose how much scholars need to learn together to do this work effectively. Scholars must not only reflect on scientific practices from humanistic perspectives or on the nature of objective versus subjective reasoning. We should also understand these perspectives diffractively as constructive and destructive interferences: superimposed patterns of thinking at the same point in time that, like interfering sound waves, amplify or silence ideas as they dissolve or smooth into one another. Scientists and humanities scholars approaching the distant listening problem space from different epistemological paradigms are concerned with intellectual bandwidths too. What is possible to know? What questions are productive or interesting to ask? Is this work worth the effort and resources needed?

### **Entelechic and Agential Meaning Making Is Premised by a Co-constructed Field of Possible Meanings**

Paradigm shifts across instantiated protocols and practices can be daunting, especially when large-scale concerns impacting institutions of scholarly

practice (universities, funding agencies, scholarly societies, and publishing) are always evolving. Christa Williford and Charles Henry (2012) evaluated projects from the first two cohorts of Digging into Data projects for the Council on Library and Information Resources and determined that distant reading and distant listening required changes in how stakeholders (researchers, administrators, scholarly societies, academic publishers, research libraries, and funding agencies) understood research practices. Their recommendations include expanding concepts of research and research data, embracing interdisciplinarity and collaboration, increasing training, adopting new models for sharing credit among collaborators, and sharing resources among institutions, reenvisioning scholarly publication, and generally making greater, sustained institutional investments in human infrastructure and cyberinfrastructure (Williford and Henry 2012, 5–6). These recommendations are premised on large financial investments and the idea that distant listening is a methodology that will produce results of interest. This is one path forward.

From my experience, I sense another path. What I would suggest for humanities scholars interested in knowing more about the potentials and pitfalls of distant listening is to learn more about the process of distant listening and knowledge production in the following key areas.

### **Audio Data Curation**

What is included in audio data? Knowledge about the technologies of audio, of computational methods, and of humanist inquiry to do new kinds of research in this area is essential. As audio-based “virtual assistant technologies” become more prevalent in our phones, computers, cars, homes, and places of work, leisure, and culture, humanists must have a better understanding of audio and machine learning to address how the data these technologies collect do and do not represent the complexities of people and culture. Scholars and academics must do research and design curriculum topics of interest to the humanities to attract other humanities scholars to this work and to address these and other areas of study and practice in the humanities.<sup>11</sup>

### **Analysis Tools**

Interfaces and tools that facilitate audio analysis can be highly technical in nature, asking users to toggle damping ratios, gain, frequencies, spectra,



energy, and pitch energy for a diverse range of sounds from music to speech to bird calls. Some of the most powerful tools for audio analysis are proprietary and owned by large companies such as IBM or Amazon, but free tools are powerful too, such as Audacity, Praat, and SonicVisualizer. These tools offer what could be called “learning interfaces,” in which practitioners are able to play with sonic parameters on a single audio file to learn how different features change sound and how it is visualized. In order to take advantage of the methods these tools deploy, however, scholars require technical knowledge of how tools visualize sound amplitude and frequencies as waves and as spectrograms and how these measurements map on to what resonates in the humanities.<sup>12</sup>

### **Accuracy Thresholds**

When is good enough good enough? Audio is always mediated in some format in real time through some mechanism. Historical recordings can be difficult to hear due to poor recording quality or noise in the sonic environment. In addition, underresourced and overworked information professionals without time to listen to every audio file in their collections produce descriptive metadata, which is used for indexing as well as many of the machine-learning projects listed above. When is “good enough” metadata good enough for discovery? While libraries and archives must adhere to professional guidelines that dictate precision and accuracy (often for very good reasons around privacy and copyright),<sup>13</sup> scholars, students, and the public can play with developing transcripts and annotations for audio recordings by creating their own projects, editions, and playlists. Tools and methodologies for engaging the public in describing audio visual collections can both ameliorate the backlog of undescribed audio collections and complicate what “fidelity” to an original sound recording means.

### **Scalability**

Scalability is also a factor in accuracy. The more data that informs a machine-learning model, the more accurate the model. Institutions with more data, storage, and processing power are best suited to conduct local and large-scale audiovisual (AV) analyses. What scale is enough and how much does it cost? For HiPSTAS projects, which included processing approximately 6,000 hours of poetry performances from PennSound, among other collections, I had access to the supercomputing cluster Stampede at the Texas Advanced

Computing Center at the University of Texas (Clement and McLaughlin 2016). The size of the cluster used to be described in Texas parlance on the website as “Each Stampede node is like a beefy desktop computer” (Texas Advanced Computing Center). At the time, a desktop computer typically had two or four processing cores; Stampede had 522,080.<sup>14</sup> In 2017, Google released AudioSet, which is likely the largest dataset of audio samples available for generating models (Gemmeke et al. 2017).<sup>15</sup> Currently, AudioSet includes 2,084,320 human-labeled 10-second sound clips drawn from YouTube videos and 632 audio event classes (e.g., “music,” “speech,” “vehicle,” “creak,” and “gargling”). Google also has the advanced processing power to run analyses across these clips, having recently announced at its March 2023 annual Google I/O developer conference an AI supercomputer with 26,000 graphics processing units (GPUs). For reference, most motherboards on a typical computer allow up to four GPUs, and the world’s fastest public supercomputer has 37,000 GPUs. The cost of this work is not limited to the financial cost of the machinery. Studies show there is also an environmental price to running large data centers that needs further exploration (Hogan 2015).

### **Sustainability**

What are local, national, and global scale issues for sustaining literacy, usability, and accuracy? How does this work fit back into the preservation and discovery infrastructures already in place in archives, libraries, classrooms? Standards such as the newest International Image Interoperability Framework (IIIF) guidelines for audiovisual materials can open the closed circle of authority around institutionally driven descriptive practices.<sup>16</sup> IIIF uses the WC3 Web Annotation Data Model standard (W3C 2017) for browsers to facilitate sharing digital image and AV data across technology systems. Third-party software gives users new kinds of access to images and AV, allowing for viewing, zooming, comparing, manipulating, and working with annotations.<sup>17</sup> With IIIF, users can reference audiovisual artifacts linked from Libraries, Archives, and Museums (LAMs) into software that allows them to annotate AV in new ways without impacting the institution’s presentation of the item. In these cases, researchers can discover, compare, refer, sample, illustrate, and represent their interpretations of these cultural heritage objects, which in turn encourages their broader use.

Done thoughtfully, distant listening research around data curation, analysis tools, accuracy, scalability, and sustainability can expose knowledge production as a messy sociotechnical process of interferences that are amplified in moments when meaning making and world building resonate because of or despite personal histories, institutional sociopolitics, and the materials and technologies (modalities, media, and devices) that produce, reproduce, store, and play archival records.

These days I prefer close listening as a research methodology. Since 2018, I have been leading a new HiPSTAS project, called AVAnnotate.<sup>18</sup> Since monks included commentary on medieval manuscripts, annotations have been an essential humanities method for adding context and meaning to cultural objects for use in research, teaching, and publication (Clement and Fischer 2021). The goal of AVAnnotate is to facilitate sharing annotations on AV archives through a sustainable workflow that leverages IIIF standard for AV materials and simplifies the production of standards-based, user-generated, online projects that provide sustainable and much-needed commentary and context around underused and culturally sensitive AV collections.<sup>19</sup> These projects—which resemble AV-centered “editions” or “exhibits”—are a series of web pages, hosted on GitHub, that feature an audio or video recording linked from a library or archive that can be played in the context of user-generated, time-stamped annotations, alongside introductory material and an index of concepts and terms, all of which provide content for searching, browsing, and organizing recordings. Existing projects include curricula for recorded interviews with jailed student protestors during the Civil Rights movement, a bilingual edition of Radio Venceremos programs (the rebel radio station that broadcast during the Salvadoran Civil War, 1981–1992), a documentary of decades of events at the Furious Flower Poetry Center (the nation’s first academic center for Black poetry), as well as a set of oral histories from the Syilx Okanagan Peoples.<sup>20</sup> I have also used AVAnnotate to create a compendium of recordings with annotations for this book, including those on Hurston’s recordings at the Library of Congress, Ellison’s recordings at Harvard, Sexton’s recordings at the Harry Ransom Center and the Schlesinger Library, and Anzaldúa’s tarot readings at the Nettie Lee Benson Latin American Collection.<sup>21</sup>

The AVAnnotate project helps me put into action some of the goals with which I began *Dissonant Records*. Reflecting standards for simplicity and

sustainability that encourage easy-to-use and lightweight technical infrastructures, the AVAnnotate workflow is based on a minimalist computing approach to development: it does not require a heavy investment in a deep software stack. We define success in AVAnnotate in terms of increasing development, elevating awareness, and promoting sustainability around issues of access to AV in libraries and archives; scholarly, pedagogical, and public use of AV collections; standards for AV access and engagement; and the social and political contexts surrounding AV access and engagement. Our primary goal for a new kind of AV ecosystem for public knowledge is to open paths for responsible and sustainable collaborations between LAMs and the public by providing a free and easy way to collaboratively create and share annotations that help describe collections and make them more accessible for interpretation.

I believe that supporting close listening practices accelerates the continued use of AV and encourages information professionals, scholars, students, and the public to produce and discuss knowledge about the remarkable and varied AV artifacts archives hold. Hurston's tale about the laws of sound remains relevant: if a collection of sound forms can become so definite and concrete as a tool, the biggest questions are not how, but to what end? A war-ax or a food tool?

This is a section of [doi:10.7551/mitpress/14976.001.0001](https://doi.org/10.7551/mitpress/14976.001.0001)

# Dissonant Records

## Close Listening to Literary Archives

By: Tanya E Clement

### Citation:

*Dissonant Records: Close Listening to Literary Archives*

By: Tanya E Clement

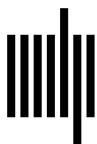
DOI: 10.7551/mitpress/14976.001.0001

ISBN (electronic): 9780262379229

Publisher: The MIT Press

Published: 2024

The open access edition of this book was made possible by generous funding and support from MIT Press Direct to Open



The MIT Press

© 2024 Massachusetts Institute of Technology

This work is subject to a Creative Commons CC-BY-NC-ND license.

This license applies only to the work in full and not to any components included with permission. Subject to such license, all rights are reserved. No part of this book may be used to train artificial intelligence systems without permission in writing from the MIT Press.



The MIT Press would like to thank the anonymous peer reviewers who provided comments on drafts of this book. The generous work of academic experts is essential for establishing the authority and quality of our publications. We acknowledge with gratitude the contributions of these otherwise uncredited readers.

This book was set in Stone Serif and Stone Sans by Westchester Publishing Services.

Library of Congress Cataloging-in-Publication Data

Names: Clement, Tanya Elizabeth, author.

Title: Dissonant records : close listening to literary archives /  
Tanya E. Clement.

Description: Cambridge, Massachusetts : The MIT Press, 2024. | Series:  
Media origins | Includes bibliographical references and index.

Identifiers: LCCN 2023052594 (print) | LCCN 2023052595 (ebook) |  
ISBN 9780262548724 (paperback) | ISBN 9780262379236 (epub) |  
ISBN 9780262379229 (pdf)

Subjects: LCSH: Hurston, Zora Neale—Archives. | Ellison, Ralph—Archives. |  
Sexton, Anne, 1928–1974—Archives. | Anzaldúa, Gloria—Archives. |  
Literature—Archival resources. | Sound archives—Social aspects. |  
Literature—Research. | Listening—Social aspects. | Intermediality.

Classification: LCC CD973.2 .C58 2024 (print) | LCC CD973.2 (ebook) |  
DDC 929.1072—dc23/eng/20240402

LC record available at <https://lcn.loc.gov/2023052594>

LC ebook record available at <https://lcn.loc.gov/2023052595>