



Corpus of Regional
African American Language

Editors: Tyler Kendall & Charlie Farrington (University of Oregon)

CORAAL (Version 2021.07)

- ↳ CORAAL:ATL (Atlanta, GA 2017; Version 2020.05)
- ↳ CORAAL:DC
 - ↳ CORAAL:DCA (Washington, DC 1968; Version 2018.10.06)
 - ↳ CORAAL:DCB (Washington, DC 2016; Version 2018.10.06)
- ↳ CORAAL:LES (Lower East Side, NY 2009; Version 2021.07)
- ↳ CORAAL:PRV (Princeville, NC 2004; Version 2018.10.06)
- ↳ CORAAL:ROC (Rochester, NY 2016; Version 2020.05)
- ↳ CORAAL:VLD (Valdosta, GA, 2017; Version 2021.07)

CORAAL User Guide (as of July 2021)

This material is based upon work supported by the National Science Foundation under Grant No. BCS-1358724.

Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.



Table of Contents

Table of Contents	2
About the Corpus of Regional African American Language.....	4
Why AAL?.....	4
The Online Resources for African American Language (ORAAL) Project.....	5
CORAAL Supplements	5
Acknowledgments and Development Team	7
Contact Information	8
About Version 2021.07 (Change Log).....	9
About Version 2020.05 (Change Log).....	9
About Version 2018.10.06 (Change Log).....	9
Errata (v. 2018.10.06)	10
About Version 2018.04.06 (Change Log).....	10
Errata (v.2018.04.06)	10
About Version 2018.01.06 (Change Log).....	11
Obtaining and Using CORAAL.....	12
What comes with the corpus?	12
Obtaining the corpus	12
The CORAAL Explorer website and R functions	15
Terms of use.....	15
Citing the corpus	16
CORAAL Transcription.....	18
Transcription practices & conventions	18
A. Symbols and punctuation	19
B. Orthographic conventions	19
C. Disfluent speech	21
D. Other features	21
Redaction and participant anonymity	24
CORAAL transcription formats.....	26
CORAAL Component Details	27
CORAAL:ATL (Atlanta, GA 2017).....	27
CORAAL User Guide (as of July 2021)	2

CORAAL:DCA (Washington, DC 1968).....	28
CORAAL:DCB (Washington, DC 2016).....	29
CORAAL:LES (Lower East Side, NY 2008).....	30
CORAAL:PRV (Princeville, NC 2004).....	31
CORAAL:ROC (Rochester, NY 2016).....	32
CORAAL:VLD (Valdosta, GA 2017).....	33
CORAAL Metadata.....	34
Speaker and file labeling conventions.....	34
Metadata files and notes.....	34
Metadata notes (all CORAAL components).....	35
Metadata notes: ATL.....	38
Metadata notes: DCA.....	38
Metadata notes: DCB.....	39
Metadata notes: PRV.....	40
Metadata notes: LES.....	40
Metadata notes: ROC.....	41
Metadata notes: VLD.....	41
Projects Using CORAAL.....	41
Publications.....	41
References.....	43

About the Corpus of Regional African American Language

The Corpus of Regional African American Language (CORAAL) is the first public corpus of African American Language (AAL) data. CORAAL features recorded speech from regional varieties of AAL and includes the audio recordings along with time-aligned orthographic transcription.

CORAAL is a long-term corpus-building project conceived of in terms of several components. The core components of CORAAL focus on AAL in Washington DC, the nation's capital, a city with a long-standing African American majority, and the site of much early research on AAL (e.g. Fasold 1972; see Farrington & Schilling 2019). CORAAL:DC, first released in January 2018, is comprised of over 100 sociolinguistic interviews with AAL speakers in DC born between 1890 and 2005. CORAAL:DC consists of two sub-components, CORAAL:DCA and CORAAL:DCB. In addition to CORAAL:DC, CORAAL includes several smaller components to provide regional breadth. As of July 2021, there are five supplemental components: CORAAL:ATL, which includes 14 sociolinguistic interviews from speakers living in Atlanta, Georgia; CORAAL:LES, comprised of 10 sociolinguistic interviews of speakers from the Lower East Side of New York City; CORAAL:PRV, which includes 15 sociolinguistic interviews from the town of Princeville, a rural African American community in central North Carolina; CORAAL:ROC, which includes 14 sociolinguistic interviews from Rochester, a city in Western Upstate New York; and CORAAL:VLD, which includes 12 speakers from Valdosta, a small city in South Georgia. All CORAAL recordings have been anonymized and orthographically transcribed with time-alignment at the utterance level. Audio is available in high-quality uncompressed (.wav) format, and transcripts are available in three formats, Praat TextGrid (.TextGrid) files, ELAN (.eaf) files, and as plain text (.txt) files with tab-delimited fields. In March 2021, a phone-level aligned set of transcripts and the phonological model created with the Montreal Forced Aligner were released (see <http://lingtools.uoregon.edu/coraal/aligned/>).

The CORAAL team plans to release updates periodically, with the next release scheduled for later in 2021. This next release will primarily include updates to CORAAL's organization and metadata, which will affect many CORAAL components. A sixth supplemental component, CORAAL:DET, with recordings from the Detroit Dialect Study in 1966 (Wolfram 1969), is also planned for the near future. Additionally, a syntactically parsed version of a portion of the CORAAL data is also under development, which will include disfluency coding and part-of-speech tagging for approximately one million words. This is in progress and will be made available as soon as possible.

Why AAL?

AAL (often referred to as African American English, AAE, or African American Vernacular English, AAVE) has been a central object of study in North American linguistics and especially sociolinguistics for over 50 years (e.g., Labov, Cohen, Robins, & Lewis 1968; Wolfram 1969; Labov 1969, 1972; Fasold 1972; Bailey, Maynor, & Cukor-Avila 1991; Mufwene, Rickford, Bailey, & Baugh 1998; Rickford 1999; Poplack & Tagliamonte 2001; Green 2002; Wolfram & Thomas 2002; Yaeger-Dror & Thomas 2010; Rickford, Sweetland, Rickford, & Grano 2012; Lanehart 2015). Already by the 1990s, AAL was described as having

inspired more than five times as many sociolinguistic publications as any other ethnic or regional dialect (Schneider 1996:3).

From this extensive work, much is known about many structures of AAL varieties and a large body of research has investigated its origins (e.g., Kurath 1949; McDavid & McDavid 1951; Stewart 1968; Bailey et al. 1991; Poplack & Tagliamonte 2001; Wolfram & Thomas 2002) and current trajectories of change (e.g., Bailey & Maynor 1985, 1987; Labov 1987, 1998; Cukor-Avila 1995, 2001; Dayton 1996; Wolfram & Thomas 2002; Yaeger-Dror & Thomas 2010). Yet, there remain important questions about the origin of these varieties, their current and future development, and their relationship(s) to regional European American and other socioethnic varieties. There also continue to be a range of important social and educational applications of enhanced knowledge about the nature of AAL.

At the same time that AAL has been so extensively studied, it has remained massively underrepresented in terms of publicly available datasets and in terms of its use in general linguistic theory building (Green 2002; Kendall, Bresnan, & Van Herk 2011). Sociolinguists (and the field of linguistics more generally; see Berez-Kroeker, Holton, Kung, & Pulsifer 2017) have increasingly adopted models of data compilation in recent years that include data sharing and promoting data re-use, but thus far almost all AAL data has remained unavailable for wider, public sharing, due to ethical considerations or limitations from how the data were collected (e.g. participant consent; Warner 2014).

The availability of a public corpus of AAL is meant to enable new research and new uses. It provides access to primary data for a wider range of scholars, for example those who do not have access to field sites or to sociolinguistic data themselves (such as educational professionals and graduate students). It also seeks to support new “open science”-based approaches, where direct testing of competing theories or methodologies or reanalysis (see Rickford, Ball, Blake, Jackson, & Martin 1991; Kendall 2011) can be made on the same data.

The Online Resources for African American Language (ORAAL) Project

CORAAL is a publication of our larger, umbrella project, the Online Resources for African American Language (ORAAL) Project, housed at the Language Variation and Computation (LVC) Laboratory in the Department of Linguistics at the University of Oregon. The ORAAL website (<http://oraal.uoregon.edu/>; Kendall, McLarty, & Farrington 2020) seeks to act as a central web-based source for research and educational information about AAL. Please visit the ORAAL website for more general information about AAL, and to obtain CORAAL or find out about future updates.

CORAAL Supplements

In addition to the official components of CORAAL, beginning in May 2020, ORAAL also features a related set of materials, **CORAAL Supplements**. These CORAAL Supplements highlight newly available datasets that are available publicly but are not part of CORAAL’s official components. CORAAL Supplements represent recordings or selections from larger datasets important to the field of sociolinguistics. These CORAAL Supplements have been prepared and are curated by the CORAAL team. The first CORAAL Supplements represent important data in the study of prosody and intonation in AAL. This includes a collection of recordings from Bengt Loman’s (1967) *Conversations in a Negro American Dialect*. These

selections were made available by the Center for Applied Linguistics in 1967 with the release of *Conversations*. Coming soon are recordings from Elaine Tarone's (1972) dissertation fieldwork in Seattle, Washington, investigating the intonation patterns of African American and European American speakers from Seattle. More information for each supplement can be found on ORAAL at <https://oraal.uoregon.edu/coraal/supplements>.

Acknowledgments and Development Team

The Corpus of Regional African American Language (CORAAL) is part of the Online Resources for African American Language (ORAAL) Project at the University of Oregon, U.S.A. CORAAL and the larger ORAAL Project have been made possible by support from the U.S. National Science Foundation (Grant No. BCS-1358724 “Enhancing data and tools for research and education on African American English”), by the University of Oregon, and by the contributions of many people.

The ORAAL Project is centered in the Language Variation and Computation (LVC) Lab in the Department of Linguistics at the University of Oregon. The LVC Lab, directed by Tyler Kendall, is also home to other linguistic resources, including the NORM Vowel Normalization and Plotting Suite and Vowels.R Package (see <http://lingtools.uoregon.edu/norm/>; Thomas and Kendall 2007) and has developed the Sociolinguistic Archive and Analysis Project (SLAAP; Kendall 2007a, 2008), which houses sociolinguistic datasets through its website, <https://slaap.chass.ncsu.edu/>, hosted at North Carolina State University.

The main CORAAL development team over the years has consisted of Tyler Kendall, Charlie Farrington, Jason McLarty, Shelby Arnson, and Brooke Josler. Lucas Jensen, Emma Mullen, Chloe Tacata, Jaidan McLean, Deepika Viswanath, Savannah Ray, and Matthew Bauer have also contributed to the corpus transcription, annotation, and redaction. Fieldwork would not have been possible without the major contributions of Minnie Quarthey, Carlos Huff, Patrick Slay Brooks, Sharese King, Kara Becker, and Ryan Rowe. We also thank Ralph Fasold, Natalie Schilling, Charlotte Vaughn, Walt Wolfram, and Danica Cullinan for their many contributions to the project.

We are especially, and deeply, grateful to the very many individuals who participated in the sociolinguistic interviews for the corpus and the projects upon which it is based. Obviously, their generous contributions of their voices, their words, and their time, are the foundation of this project. Specifically, we thank:

Rafeeq Badde	Shanquette Dannah	Dwayne Lawson-
Russell Baker	Janet Davenport	Brown
Vernal Batts	Andrea Davis	Gary Lewis
Lamonte Belk	Angela Dorsey	Barbara Ligon
Devon Bennett	Bryanna Duncan	LaTanya Malloy
Darren Black	Renee Edelin	Keisha Matos-Joaquin
Jason Black	R’Mani Fitchett	Rhonice Miles
Nicholas Black	Devonte Gooding	Melvin Moore
Niya Black	Michelle Graham	Sharon Quarthey
Michelle Broadus	Theressa Green	Curtis Robinson
Sheila Brockington	Ayanna Holmes	Daejah Ross
Justin Brown	Domonique Inniss	Tavis Saunders
Robert Brown	Alonte Jenkins	Euchrist Smith
Robin Brown	Terri Johnson	Zymiah Speller
Christopher Burgman	Angelo Johnson	Gabrielle Spencer
Barbara Byrd	Jazmine Jones	Samuel Streater
Marco Coleman	LeVar Jones	Andrea Talbert
Yolanda Coppedge	Linda Jones	Mylz Taylor

Ted Thomas
Shirleen Thompson
Jamiris Tolbert
Devin Turner

Monique Van Buren
Khalil Vaughn
Teleia White
Brandon Williams

Tiffany Woodberry-
Black
Geo Wright

We also thank the very many other, anonymous participants who also contributed their speech and their stories to the corpus. (Participants were given the choice of whether they wished to be recognized by name or to remain anonymous – thus we only recognize by name those participants who asked not to remain anonymous; see [discussion of redaction in this document](#).) A number of students and research assistants at the University of Oregon have contributed to the project. In addition to the research assistants mentioned above, we thank students in T. Kendall's Spring 2016 Seminar on African American English. Finally, we thank colleagues who have acted as beta-testers and consultants as we have developed the corpus, including: Tricia Cukor-Avila, Jon Forrest, Jessi Grieser, Chris Hall, Nicole Holliday, Taylor Jones, Sharese King, John Rickford, and Tracey Weldon.

Contact Information

Please contact the CORAAL development team via:

Email: corpusofregionalAAL@gmail.com

Twitter: [@CorpusAAL](https://twitter.com/CorpusAAL)

You can also write the editors:

Dr. Tyler Kendall: tsk_at_uoregon_dot_edu

Dr. Charlie Farrington: crf_at_uoregon_dot_edu

Department of Linguistics
1290 University of Oregon
Eugene, OR 97403-1290 USA

About Version 2021.07 (Change Log)

Version 2021.07 is the fifth release of CORAAL. This release involves the addition of two subcomponents; all other components remain unchanged.

CORAAL:LES, v. 2021.07, the fourth subcomponent of CORAAL, and CORAAL:VLD, v. 2021.07, the fifth subcomponent of CORAAL, have been added. These components are expected to be relatively stable, with no additional speakers or audio files planned. For these two subcomponents, Age Group categories follow the groupings used for CORAAL:DC, with Age Group 2 representing ages 20 to 29, Age Group 3, ages 30 to 50, and Age Group 4, ages 51+. This is a difference from how ATL, PRV, and ROC were implemented, but allows for more comparability with the CORAAL:DC components. We anticipate updating the other components to match this practice in an upcoming release, to regularize the treatment of Age Group across all of CORAAL.

The components in the CORAAL User Guide are now presented in alphabetical order. The User Guide also now mentions the availability of phone-level aligned TextGrids, which were released by the CORAAL project in March 2021.

About Version 2020.05 (Change Log)

Version 2020.05 is the fourth release of CORAAL. Beginning with Version 2020.05, version numbers have been shortened to only represent the year and month of initial publication.

CORAAL:ATL, v. 2020.05, the third sub-component of CORAAL, has been added. This component is expected to be relatively stable, with no additional speakers or audio files planned, though we anticipate amending transcription files to improve accuracy if/when errors are discovered.

Six additional files have been added to the CORAAL:ROC component. These include one new speaker, ROC_se0_ag2_m_01, as well as additional recordings from ROC_se0_ag3_f_02.

CORAAL's web interface, the CORAAL explorer, has been updated to include ATL and the new ROC recordings.

There are no changes to CORAAL:DCA, CORAAL:DCB, and CORAAL:PRV. Hence, the individual version numbers for these components remain 2018.10.06. Errata is not included for this version because there were no changes to previously uploaded files.

About Version 2018.10.06 (Change Log)

Version 2018.10.06 is the third release of CORAAL.

Through a round of close consistency checking and editing, changes were made to every TextGrid file previously available (in v. 2018.04.06), including removing spaces that sometimes occurred at the beginning and ends of utterance as well as fixing typos (e.g. “<unintelligible>” changed to “<unintelligible>” for DCA_se2_ag1_m_03_1). This also means that all text files and ELAN files have been changed as well. Some inconsistencies in transcripts were fixed to more closely align to the transcription conventions (see that section). This affects all files in DCA, DCB, and PRV. Additionally, audio file processing has been re-run on ten files listed in the

errata below (redaction, amplitude normalization, and conversion). Due to the increased consistency across files, we recommend that all users replace prior versions of all CORAAL data with v. 2018.10.06 to have the most up-to-date transcriptions and audio.

CORAAL:DCB v. 2018.10.06 includes one additional speaker (five new files).

CORAAL:ROC, the second sub-component, has been added. We expect to add two speakers to this component in the next release.

Additionally, with this release comes the first version of a web-interface to CORAAL, the **CORAAL explorer website** (<http://lingtools.uoregon.edu/coraal/explorer/>) and R functions to work with CORAAL transcripts directly in R (http://lingtools.uoregon.edu/coraal/explorer/R/CORAAL_web.R).

Errata (v. 2018.10.06)

- Audio re-processing
 - DCA_se1_ag2_m_01_1, DCA_se1_ag3_f_01_1, DCA_se1_ag3_f_02_1, DCA_se1_ag4_m_01_1, DCA_se1_ag4_m_02_1, DCA_se3_ag2_f_04_1, DCA_se3_ag4_m_01_1, DCB_se2_ag3_m_03_1, DCB_se3_ag4_f_02_1, PRV_se0_ag2_m_02_1
- Numerous minor changes throughout all CORAAL transcripts (see change log above).

About Version 2018.04.06 (Change Log)

Version 2018.04.06 is the second release of CORAAL.

CORAAL:DCA data have not changed. However, updates were made to the metadata adding new columns to parallel changes to CORAAL:DCB (see just below) and to correct word counts listed for individual speakers.

CORAAL:DCB v. 2018.04.06 includes ten additional speakers who were not included in the opening release. New audio files are available in the download tar.gz parts 11-14. Additionally, one transcript file from CORAAL:DCB has been updated with corrections (see Errata below). Users will need to download those new audio tar.gz bundles as well as the new metadata file and new transcript tar.gz files to obtain the new data. (Audio files in parts 01-10 are unchanged so the content of those tar.gz parts are the same as for v. 2018.01.06.) CORAAL:DCB is not yet complete, but we anticipate completing the sample for the Summer 2018 release. The DCB metadata spreadsheet now includes a Primary Speaker column (Primary.Spkr) and columns indicating when files were added and last modified. The audio file is named with the primary speaker's code, but when there are other interviewees present, they are listed as non-primary speakers.

CORAAL:PRV, the first sub-component, has been added. This component is expected to be relatively stable, with no additional speakers or audio files planned, though we anticipate amending transcription files to improve accuracy if/when errors are discovered. For details about the sub-corpus, see CORAAL:PRV section below. The PRV metadata spreadsheet also includes a Primary Speaker column. The audio file is named with the primary speaker's code, but when there are other interviewees present, they are listed as non-primary speakers.

Errata (v.2018.04.06)

- DCB_se2_ag3_m_02_1 (.txt, .TextGrid, .eaf)

- Changed four instances of “this that and a third” to “this that and the third”

About Version 2018.01.06 (Change Log)

Version 2018.01.06 is the opening release of CORAAL.

CORAAL:DCA is expected to be relatively stable, with no additional speakers or audio files planned (although future updates will presumably amend transcription files to improve their accuracy).

CORAAL:DCB is not quite complete. CORAAL:DCB will have additional interviews released in a future update, with the goal of including at least two speakers per demographic cell (see CORAAL:DCB section below). We do not anticipate making substantial changes to the existing interviews/transcripts but we hope to complete the sample for the next update to the corpus.

Obtaining and Using CORAAL

What comes with the corpus?

CORAAL contains audio files along with corresponding orthographic transcription, time-aligned at the utterance level. Metadata files are available for each component that provide extensive information about the speakers (see metadata section below). Metadata files are plain text in a tab-delimited format; these can be read by any text editing software, but can also be loaded into spreadsheet software like MS Excel or software like R. Audio files are available in uncompressed .wav format (generally 44.1 kHz, 16 bit, mono). Transcripts are available in three formats – Praat TextGrids (.TextGrid), ELAN files (.eaf), and plain, tab-delimited text (.txt). All three formats contain identical information. As described in the transcription section, transcripts are created by the CORAAL development team directly as TextGrids in Praat. The Praat TextGrid files are automatically processed (by script) into text files and (by ELAN) into ELAN format.

As of March 2021, phone-level aligned Praat TextGrids and a language model are now available through the CORAAL download site (<http://lingtools.uoregon.edu/coraal/aligned/>). These were produced by the CORAAL project using the [Montreal Forced Aligner](#) (McAuliffe et al. 2017). Initial alignment and creation of the CORAAL MFA language model was completed in November 2018 using CORAAL version 2018.10.06. In June 2019, CORAAL was re-aligned using this pre-trained model. Note that this version of CORAAL includes all speakers who were available in June 2019. This includes all speakers from DCA, DCB, and PRV. For ROC, two speakers who were added in version 2020.05 are not included in the alignment. The LES and VLD speakers are not yet available at the time of this publication.

Obtaining the corpus

All components of CORAAL are available for download from its home at <http://oraal.uoregon.edu/coraal>. (See the following subsections for other ways to access the corpus.) The corpus is organized by sub-component and, due to their large sizes, each sub-component is broken down into several parts, which are then compressed using standard “tar” and “gnu zip” compression. Audio files are contained in numerous separate tar.gz files. Transcripts are organized into tar.gz files by file type. Metadata files are available as uncompressed plain text files as these have smaller file sizes. You can download the corpus by saving the individual parts over your web-browser. Any standard decompression software should readily be able to decompress the downloaded files. We suggest you move all of the files from their individual parts’ folders to a single folder, as the individual part folders, especially for the audio parts, are not meaningful other than as ways to organize the files for downloading.

Please note that CORAAL is quite large (≈ 27.5 GB for v. 2021.07) and can take a long time to download. ATL, in its compressed format for download, is 1.7 GB; DCA, compressed, is 7.5 GB; DCB, compressed, is 8.1 GB; LES, compressed, is 1.8 GB; PRV, compressed, is 3.3 GB; ROC, compressed, is 2.9 GB; and VLD, compressed is 2.2 GB.

Quicker way: Automate the download

A text file containing a list of all files for all of the components of the current version of CORAAL is available here:

<https://tinyurl.com/coraalfiles>

(This is a shortcut to http://lingtools.uoregon.edu/coraal/coraal_download_list.txt)

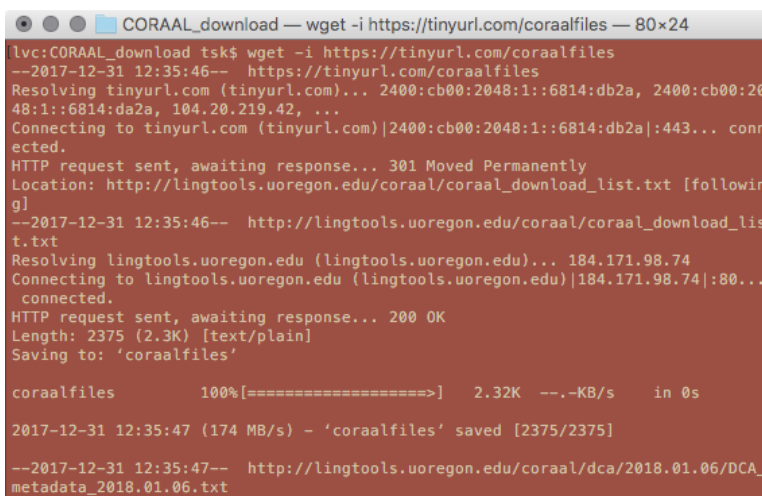
You can use this file to automate the downloading of the entire corpus, such as through one of the techniques described below. (It will probably still take a while but at least you can avoid manually downloading every file!) Here are some suggested ways you can automate this process:

Use wget software

An easy way to download all of the files is by using the popular open-source “wget” application. Download “wget” if you don’t already have it (see below). Then from your system’s shell (e.g. Mac Terminal application), navigate to the location you wish to save all of the files and run the following:

```
wget -i https://tinyurl.com/coraalfiles
```

This will download all of the files for the corpus. Here is a screenshot of a Mac Terminal after executing the wget command (the first line shows the wget command, everything below that is wget doing its magic):



```
CORAAL_download — wget -i https://tinyurl.com/coraalfiles — 80x24
[vc:CORAAL_download tsk$ wget -i https://tinyurl.com/coraalfiles
--2017-12-31 12:35:46-- https://tinyurl.com/coraalfiles
Resolving tinyurl.com (tinyurl.com)... 2400:cb00:2048:1::6814:db2a, 2400:cb00:20
48:1::6814:da2a, 104.20.219.42, ...
Connecting to tinyurl.com (tinyurl.com)|2400:cb00:2048:1::6814:db2a|:443... conn
ected.
HTTP request sent, awaiting response... 301 Moved Permanently
Location: http://lingtools.uoregon.edu/coraal/coraal_download_list.txt [followin
g]
--2017-12-31 12:35:46-- http://lingtools.uoregon.edu/coraal/coraal_download_lis
t.txt
Resolving lingtools.uoregon.edu (lingtools.uoregon.edu)... 184.171.98.74
Connecting to lingtools.uoregon.edu (lingtools.uoregon.edu)|184.171.98.74|:80...
connected.
HTTP request sent, awaiting response... 200 OK
Length: 2375 (2.3K) [text/plain]
Saving to: 'coraalfiles'

coraalfiles      100%[=====] 2.32K  --.-KB/s   in 0s

2017-12-31 12:35:47 (174 MB/s) - 'coraalfiles' saved [2375/2375]

--2017-12-31 12:35:47-- http://lingtools.uoregon.edu/coraal/dca/2018.01.06/DCA_
metadata_2018.01.06.txt
```

You may find easy-to-install versions of “wget” for your operating system by searching the web for something like “download wget for mac”, but as of this writing the following links are helpful starters. For Mac and Linux, using a package manager (like “homebrew” on Mac, see <https://brew.sh> and note that wget is the example on the homebrew website!) is an easy way to install the software.

Mac OSX: <https://www.fossmint.com/install-and-use-wget-on-mac/> (or from source: <http://osxdaily.com/2012/05/22/install-wget-mac-os-x/>)

Windows: <https://eternallybored.org/misc/wget/> (or, longer version: <https://builtvisible.com/download-your-website-with-wget/>)
Linux: e.g. <https://www.tecmint.com/10-wget-command-examples-in-linux/>

Use curl software

Many operating systems, like Mac OSX, have a program called “curl” already installed. “curl” can also be used to download the corpus quickly, although it cannot read a list of files as readily as “wget”. A fast way using “curl” involves also using a second program “xargs”. You can test to see if your system has “curl” and “xargs” by going to your system’s shell (e.g. Terminal on Mac) and executing:

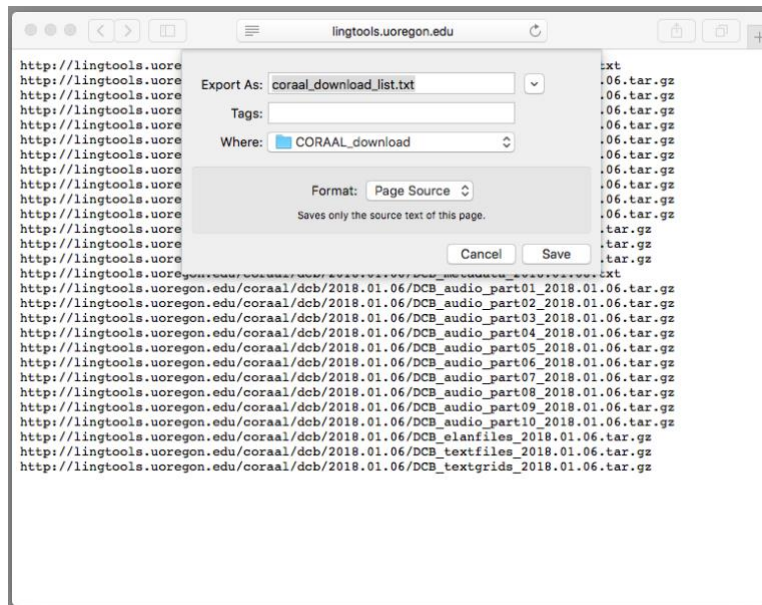
```
which curl  
which xargs
```

If you have these programs, your system will respond by telling you where they are. If you don’t, your system will simply show you another prompt.

Here, you will need to manually download the download_list.txt files before you can proceed. Go to <https://tinyurl.com/coraalfiles> and save this file to your computer as a text file (e.g. coraal_download_list.txt). Then, from your system’s shell (e.g. Terminal on Mac), navigate to the location you wish to save all of the files and execute:

```
xargs -n 1 curl -O < /path/to/coraal_download_list.txt
```

This will download each of the files in the list. Here are screenshots of saving the coraal_download_list.txt file (for CORAAL version 2018.01.06) from a web browser:



And executing the xargs/curl commands in Mac Terminal (in this screenshot, the first two files have transferred and curl is 6% of the way through downloading the third file):

```

CORAAL_download — curl + xargs -n 1 curl -O — 80x24
lvc:CORAAL_download tsk$ which curl
/usr/bin/curl
lvc:CORAAL_download tsk$ which xargs
/usr/bin/xargs
lvc:CORAAL_download tsk$ ls
coraal_download_list.txt
lvc:CORAAL_download tsk$ xargs -n 1 curl -O < coraal_download_list.txt
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
   Dload  Upload  Total   Spent    Left     Speed
100 40106  100 40106    0     0  1551k      0  --:--:--  --:--:--  --:--:-- 1566k
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
   Dload  Upload  Total   Spent    Left     Speed
100 688M  100 688M    0     0  79.4M      0  0:00:08  0:00:08  --:--:-- 77.4M
% Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
   Dload  Upload  Total   Spent    Left     Speed
 6 596M   6 37.0M    0     0  74.2M      0  0:00:08  --:--:--  0:00:08 74.1M

```

The CORAAL Explorer website and R functions

Beginning with CORAAL v. 2018.10.06, we have created a web-interface for “exploring” CORAAL. The CORAAL Explorer website is available at <http://lingtools.uoregon.edu/coraal/explorer/>. The online interface currently has two primary features, a browse page (<http://lingtools.uoregon.edu/coraal/explorer/browse.php>) and a search page (<http://lingtools.uoregon.edu/coraal/explorer/search.php>). The browse feature provides text and audio for each CORAAL file (and also allows users to download individual files). The search feature provides a web-based front end to a set of R functions for working with the corpus. The pages are (hopefully) relatively straightforward to use, but users can contact CORAAL developers with questions or problems. We anticipate that the web pages will continue to be developed and new versions will be released on a rolling basis (i.e. not at the same rate as the periodic CORAAL data updates).

An R script, which supports the direct downloading and creating of R data structures for CORAAL is available at http://lingtools.uoregon.edu/coraal/explorer/R/CORAAL_web.R. You can save that file to disk to use or you can simply load it into R over the internet. Execute the following to make CORAAL and some helper functions available in your R session:

```
source("http://lingtools.uoregon.edu/coraal/explorer/R/CORAAL_web.R")
```

This will show a short help message including instructions for how to download CORAAL directly into your R session (using the `coraal.webget.data()` function). Note that you need the Rcurl library installed in R to use the web-based download feature in R.

Terms of use

The Corpus of Regional African American Language (CORAAL), its data, and websites are available for free, public use for research purposes. CORAAL is available under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-NC-SA 4.0) License. This means you are free to use and reuse the corpus for non-commercial purposes, but that you must cite the original corpus and any derivative versions of CORAAL you develop and



wish to share with others must be distributed using the same license. A summary of the license is available on the Creative Commons website at <https://creativecommons.org/licenses/by-nc-sa/4.0/> and the full license is available at <https://creativecommons.org/licenses/by-nc-sa/4.0/legalcode>.

CORAAL: DCB, PRV, ATL, and VLD are available under additional licenses, such as those for commercial purposes. For more information about these other license options, please contact Tyler Kendall (tsk_at_uoregon_dot_edu).

Citing the corpus

If you use the corpus, we ask that you cite the corpus. Below are suggested citations for CORAAL and its available subcomponents. More generally, we urge you to learn about the *Austin Principles of Data Citation*, which provide guidelines for citation and attribution of linguistic data. The Austin Principles are described at <http://site.uit.no/linguisticsdatacitation/>.

Recommended Citation and Version Number for the main CORAAL project:

- Kendall, Tyler and Charlie Farrington. 2021. *The Corpus of Regional African American Language*. Version 2021.07. Eugene, OR: The Online Resources for African American Language Project. [<https://oraal.uoregon.edu/coraal>]

Recommended Citations and Version Number for CORAAL:ATL (2017):

- Farrington, Charlie, Tyler Kendall, Patrick Slay Brooks, Lucas Jenson, Chloe Tacata, and Jaidan McLean. 2020. *The Corpus of Regional African American Language: ATL (Atlanta, GA 2017)*. Version 2020.05. Eugene, OR: The Online Resources for African American Language Project.

Recommended Citations and Version Number for CORAAL:DCA (1968):

- Kendall, Tyler, Ralph Fasold, Charlie Farrington, Jason McLarty, Shelby Arnson, and Brooke Josler. 2018. *The Corpus of Regional African American Language: DCA (Washington DC 1968)*. Version 2018.10.06. Eugene, OR: The Online Resources for African American Language Project.
- Fasold, Ralph. 1972. *Tense marking in Black English: A linguistic and social analysis*. Arlington, VA: Center for Applied Linguistics. [<https://eric.ed.gov/?id=ED129065>]

Recommended Citation and Version Number for CORAAL:DCB (2016):

- Kendall, Tyler, Minnie Quartey, Charlie Farrington, Jason McLarty, Shelby Arnson, and Brooke Josler. 2018. *The Corpus of Regional African American Language: DCB (Washington DC 2016)*. Version 2018.10.06. Eugene, OR: The Online Resources for African American Language Project.

Recommended Citations and Version Number for CORAAL:LES (2008):

- Becker, Kara, Charlie Farrington, Tyler Kendall, Chloe Tacata, and Jaidan McLean. 2021. *The Corpus of Regional African American Language: LES (Lower East Side, NY 2008)*. Version 2021.07. Eugene, OR: The Online Resources for African American Language Project.

- Becker, Kara. 2010. *Regional dialect features on the Lower East Side of New York City: Sociophonetics, ethnicity, and identity*. Ph.D. dissertation. New York, New York: New York University. [<https://search.proquest.com/openview/5684221c59338341afe0941c1d462532/>]

Recommended Citations and Version Number for CORAAL:PRV (2004):

- Rowe, Ryan, Walt Wolfram, Tyler Kendall, Charlie Farrington, and Brooke Josler. 2018. *The Corpus of Regional African American Language: PRV (Princeville, NC 2004)*. Version 2018.10.06. Eugene, OR: The Online Resources for African American Language.
- Rowe, Ryan. 2005. *The development of African American English in the oldest Black town in America: Plural -s absence in Princeville, North Carolina*. MA Thesis. Raleigh: North Carolina State University. [<http://repository.lib.ncsu.edu/handle/1840.16/711>]

Recommended Citations and Version Number for CORAAL:ROC (2016):

- King, Sharese, Charlie Farrington, Tyler Kendall, Emma Mullen, Shelby Arnson, and Lucas Jenson. 2020. *The Corpus of Regional African American Language: ROC (Rochester, NY 2016)*. Version 2020.05. Eugene, OR: The Online Resources for African American Language Project.
- King, Sharese. 2018. *Exploring social and linguistic diversity across African Americans from Rochester, New York*. Ph.D. dissertation. Palo Alto, CA: Stanford University. [<https://searchworks.stanford.edu/view/12739840>]

Recommended Citation and Version Number for CORAAL:VLD (2017):

- Quartey, Minnie, Charlie Farrington, Tyler Kendall, Lucas Jenson, Chloe Tacata, and Jaidan McLean. 2020. *The Corpus of Regional African American Language: VLD (Valdosta, GA 2017)*. Version 2021.07. Eugene, OR: The Online Resources for African American Language Project.

Recommended Citation and Version Number for CORAAL: MFA-Aligned:

- Farrington, Charlie and Tyler Kendall. 2019. *The Corpus of Regional African American Language: MFA-Aligned*. Version 2019.06. Eugene, OR: The Online Resources for African American Language Project. [<http://lingtools.uoregon.edu/coraal/aligned/>]

CORAAL Transcription

Transcription practices & conventions

This section outlines the transcription conventions used in the creation of CORAAL. Our practices, and thus parts of this section, have been adapted from those developed for the Sociolinguistic Archive and Analysis Project (SLAAP), described in the SLAAP User Guide, Version 0.96, June 2009, available here: <https://slaap.chass.ncsu.edu/userguide/>.

Transcripts are created entirely in the Praat software (<http://praat.org>; Boersma and Weenink 2014) using TextGrid annotation objects. (Other transcription formats, like ELAN, are derived from these original TextGrid files.) In a TextGrid transcript, each speaker is represented in an interval tier. Occasionally an additional interval tier is used to transcribe speech by an interloper. In some cases, e.g. for some of the recordings in CORAAL:DCA the corpus development team had earlier transcript documents to work with, but even in these cases the entire transcript was re-keyed following the conventions outlined here. Earlier transcripts were referred to for spelling or for words that could not be determined by the transcription team from the audio.

Following SLAAP conventions, transcripts align text to speech at a per-utterance level, where utterances are defined as uninterrupted speech sounds by the same individual, with utterances delimited at pauses. For CORAAL, the criteria for how transcribers delimit utterances is set at a pause length of 60-70 milliseconds (see Kendall 2009, 2013).

For each audio file, three rounds of transcription and editing were completed. The first round of transcription was completed by one of two undergraduate research assistants, who were the primary transcribers between 2015 and 2017. The second round of transcription – a thorough editing of the TextGrid – was done by a graduate student in linguistics, who listened to the entire audio file while reviewing the transcript, making corrections where necessary. The third round created a Redaction (RD) tier which time-stamped portions of the transcript/audio that needed to be redacted, while also cleaning up any remaining inconsistencies. Despite our attempts to have maximally clean and accurate transcripts, we expect some degree of disagreement on transcript accuracy. Transcripts are always the (in process) product of individual analysis (Edwards 2001; Kendall 2008).

A primary goal of a time-aligned transcript is to act as a proxy to the original recording, to allow for easy searching and browsing of the recording. It is not to make an exact, textually accurate representation of the speech. Along these lines, orthographic transcription conventions use simple orthography and standard-like spelling. As a general rule, morphosyntactic variants (e.g. *was* for *were*) are transcribed, but phonological variants (such as *velar nasal fronting* (*in* for *ing*), *r*-lessness, and dialectal vowel qualities) are not.

At the same time, the transcript text attempts to accurately account for all the “noises” of speech, such as laughter, filled pauses (like “uh” or “um”), and restarts (e.g. “I- I- I di- didn’t mean to”) as well as misspoken words (e.g. “/brack/ in the seventies”). Standard-like capitalization and punctuation is used, with the hyphen, -, used to indicated lexical and morphosyntactic restarts, as well as incomplete intonation. Silent pauses (of course) are not described or coded, as they are represented in all cases by empty intervals. Since a number of speech sounds (e.g. “Mm-hm”) do not have codified spellings while some extremely common productions do have agreed upon non-standard written forms (e.g. “I’m’a”), transcription conventions attempt to standardize possible spellings.

The following subsections outline the specific transcription conventions. Part A provides information on symbols and punctuation used in the corpus. Part B provides orthographic conventions and examples for commonly encountered non-standard words and constructions. Part C gives conventions for disfluent speech, and Part D gives examples of other features encountered in CORAAL (including miscellaneous topics that came up for our transcription team).

A. Symbols and punctuation

Special symbols

- The basic symbols used in transcription follow the conventions used by SLAAP.

[...] contains overlapping speech, e.g.,

```
Speaker A: So [I went-]
Speaker B: [You did] what?
```

/ ... / represents several categories, including

- Inaudible or unintelligible speech
- Redacted speech
- Misspoken words

< ... > contains non-linguistic sounds, such as <cough>

(...) contains line-level notes, such as (laughing)

- Standard punctuation is used for ease of transcription and readability. Punctuation includes periods (.), commas (,), question marks (?), exclamation points (!), all of which are used as normal in English (prosody). Apostrophes (') are used as they are generally used in English. Dashes (-) are limited to disfluent speech and compounds.
 - I considered myself- Oh! Horseshoes! I played quite a bit of horseshoes at that age. (DCA_se3_ag3_m_01)
- In initial versions of the corpus, quoted speech is not represented orthographically (e.g. by use of quotation marks). This may change in future versions, but determining what speech is direct vs. indirectly reported speech vs. other “quoted speech” is not trivial and highly interpretive.

B. Orthographic conventions

General Notes

- Spelling and capitalization follow standard English writing practices.
- All numbers are written out as complete words.
 - I've been in Washington for twenty-two years. (DCA_se3_ag3_m_04)
 - When aught years (e.g. 2008) are referred to as '08, o eight is transcribed.

- Transcribers write compounds as two words, or a hyphenated word.
- Abbreviations are avoided except for personal titles (e.g. Mr., Mrs., Dr.). Junior (e.g. Thomas Junior) is not abbreviated.
- Acronyms are written without punctuation (e.g. TV, DC). Letters that are pronounced are separated by dashes. ASAP is always transcribed in caps. An apostrophe is included for plural acronyms (e.g. TV's, A's).
 - And I think his name is spelled B-O-O-N-K. (DCB_se3_ag1_f_01)
- Question marks (?) are used in transcripts as normal in English, based on a combination of prosodic and syntactic cues.

Reduced forms

- In CORAAL, several reduced constructions have orthographic representations. Transcribers used this list when fuller (more conventionally standard forms) are *reduced*. Otherwise full forms are transcribed.
- In the case of *have* reduction, conventional contractions (e.g. must've, would've, etc.) are also transcribed when necessary.
- Some commonly reduced constructions (e.g. used to, kind of, sort of, out of) are never orthographically represented as reduced. In the following table, reduced forms are listed as transcribed in CORAAL. Rows without “reduced form” entries should always be transcribed in their full form, regardless of pronunciation.

Category	CORAAL Representation (Full Form)	CORAAL Representation (Reduced Form)	Notes
<i>have</i> reduction	must have	musta	
	would have	woulda	for wouldn'ta, transcribe wouldn't have
	should have	shoulda	
	could have	coulda	
	might have	mighta	
<i>to</i> reduction	going to	gonna ~ I'm'a	gonna includes [gon]/[gõ] variants
	have to	hafta	
	used to		Always transcribed as used to
	trying to	tryna	
	supposed to	sposta	sposta includes [po] variants
	fixing to	finna	
	got to	gotta	
	want to	wanna	
	ought to	oughta	
syllable reduction	because	cause	
	until	til	
	about		
	remember		
	around		

Other reduction	talking about		e.g. quotative talmbout (Vaughn-Cooke 1976; Jones 2016) transcribed as talking about
	them	`em	
	let me	lemme	
	what do you/what are you	whatchu/ whatcha	The difference between whatchu and whatcha is the final vowel (whatchu is more common in CORAAL)
	got you	gotcha	

C. Disfluent speech

Restarts

- Speaker restarts are indicated with a single dash.
 - That was- that was some good times. (DCB_se2_ag2_m_03)
- When a filled pause precedes a restart, the restart should be indicated before the filled pause.
 - And so, she went to- um, she was born in Georgia. (ROC_se0_ag2_f_04)

Filled Pauses

- Only *uh*, *um*, and clear cases of *ah* are used. More than one filled pause in a row is not treated as a restart. A comma should precede and follow each filled pause.
 - It'd be at the, um, other campus, in Largo. (DCB_se1_ag1_f_01)
 - So, when I moved to East Capitol, uh, mm, you know, I had to wait for a bunch of stuff when I first got there, had to set up my room. (DCB_se3_ag1_f_01)
 - Ah, I think the former Grammar teacher that I was just telling you about, um, was very good. (DCA_se3_ag3_m_01)

Mispronounced Words

- Mispronounced words are transcribed as they are pronounced, in slashes.
 - Christopher /Folumbus/ I mean Columbus. (DCA_se2_ag1_f_03)

D. Other features

Discourse markers

- Utterance final *so*. is transcribed with a following period.
 - And there's sort of like a slump there, you know, so. (DCA_se2_ag2_f_03)
- Lip Smacking/Teeth Sucking, *if* relevant to the present discourse, are transcribed as <ts>.
 - Like, try to be funny, so. <ts> I gave him my number. But a long story short, we started talking. (DCB_se1_ag1_f_03)

Overlap

- Speaker overlap is noted by the use of square brackets, for all parties to the overlap. The overlap markers are always placed at word boundaries.
 - DCB_int_01: Or is he [retired? No?]
 - DCB_se1_ag2_f_02: [No, he was] in the reserve but he just went back to [active.] Yeah.

- o DCB_int_01: [Oh.]

Unintelligible/Inaudible Speech:

- Slashes are used to enclose sections of unsure transcription. Transcribers often place “best guesses” within the slashes, or write /unintelligible/ for unintelligible talk or /inaudible/ for inaudible talk. For unintelligible talk of less than three syllables, transcribers can also use question marks, /??./, within the slashes to indicate each syllable of unintelligible speech.
 - o DCA_se3_ag3_m_01: And here again now, I can appreciate that /unintelligible/, [but at] that time, I thought he was being very unreasonable. I was only at the circus.
 - o DCA_int_01: [Yeah.]
 - o DCA_se3_ag3_m_01: [<laugh>]
 - o DCA_int_01: [/inaudible/.]

Non-linguistic/meta-linguistic noises:

- Noises like laughter, hand clapping, and throat clears are often indicated by short descriptions enclosed within angle brackets. These are only used to describe actual noises, not features like voice quality.
 - o DCA_int_04: What about the kids that go there? Can you tell me anything about them?
 - o DCA_se3_ag1_f_06: Well we’re all different. <laugh>
 - o DCA_int_04: Mm-hm.
- Other examples used in CORAAL
 - o <cough>, <clears throat>, <laugh>, <yawns>, <snap>, <sound effect>, <grumbles>, <inhale>, <exhale>, <microphone feedback>, <clap>, <ts>, <pp>, <imitates music>,
 - <ts> represents teeth suck. Transcribers were asked to transcribe <ts> if it appeared to be relevant to the present discourse, and related to communication. <ts> covers a range of sounds.
 - <pp> represents a bilabial trill interjection.
 - o *speech in slashes as well as angle brackets, if overlapping, can also be within square brackets
 - [I mean /unintelligible/] you know, you- you gonna have your trickery, <clears throat> in the government as far as that goes. (DCB_se1_ag3_m_02)

Line-level notes:

- Notes can be included by the use of parentheses in a transcribed utterance. Features like voice quality are noted this way. These do not have to be within square brackets.
 - o (whispered) (while chewing) (laughing) (atypical pronunciation) (singing) (rapping) (coughing) (breathy)
 - (atypical pronunciation) (e.g. in DCB_se1_ag3_m_02) is used when the speaker is describing an accent specific pronunciation.
 - Baltimore, they be like, what up dog (atypical pronunciation) (DCB_se1_ag3_m_02)

Redaction

- Slashes and a redaction code are used to obscure real names, addresses, places of work, and schools.

- Redaction codes are as follows, with the # being the number of syllables obscured:
- RD-WORK, RD-SCHOOL, RD-ADDRESS, RD-NAME
 - e.g. Winston High School → /RD-SCHOOL-4/
- In CORAAL, all redacted codes have been replaced with tones, which were generated based on the mean pitch and amplitude of the speech being redacted.

Morphophonological vs. morphosyntactic differences

- When it’s not clear that the process is phonological or morphosyntactic (e.g. a result of consonant cluster reduction vs. an unmarked verb), transcribers were instructed to err on the side of using standard orthographical conventions (i.e. transcribed the *d* in “The boy named Bill”).
- But, for clearly morphological/syntactic features, transcribers were instructed to transcribe as close to the audio as possible.
- Examples of common AAL morphosyntactic features included in transcripts are:

<i>Category</i>	<i>Example</i>	<i>Notes</i>
Third person singular –s absence	So she run over to her bag (DCA_se2_ag3_m_01)	
Possessive –s absence	We was over my uh, father mother house. (DCA_sel_ag1_f_01)	
Possessive <i>they</i>	The young girls today wanna be friends with they kids. (DCB_sel_ag3_f_02)	Not r-less <i>their</i>

Common interjections/filler words

- Since several speech sounds (e.g., “mm-hm”) don’t have codified spellings. The table below gives orthographic guidelines for common interjections and filler words.

<i>Transcription</i>	<i>Notes</i>
Uh-huh	Positive/neutral
Uh-uh	Negative
Nuh-uh	Negative (nasalized uh-uh)
Mm-hm	Positive/neutral
Mm	Positive/neutral
Mm-mm	Negative
Okay	
Mkay	
Yep/Yup	Both are transcribed
Nah/Naw	Both are transcribed
Oh	[oo]
Ooh	[u]
Ayo	[ejo]
Hoo	[hu]

CORAAL lexical conventions and dialect specific items

- The following are several common, dialect specific lexical items that appear in the corpus that (1) don't have a standardized spelling or (2) might be unfamiliar to some CORAAL users. In some cases, these are features of African American Language.

<i>Transcription</i>	<i>Notes</i>
aight	Reduction of <i>alright</i>
aks	Metathesis of <i>ask</i>
and them	Associative plural marker <i>nem/dem</i>
ay!	Common exclamation ([eɪ])
`bacca	Reduced <i>tobacco</i> in Princeville, Valdosta
bih	Reduced <i>bitch</i>
bougie	Perceived as upscale
brazy	Means crazy
bruh	Variant of <i>bro</i>
cuz	Reduced <i>cousin</i> . Not reduced <i>because</i>
`em	Reduced <i>them</i>
fella	Variant of <i>fellow</i>
go go	Popular style of music in DC
hisself	Regularized <i>hisself</i>
jai	Means really/very in DC. [dʒaɪ]
effed up	Not <i>F-ed up</i>
mama	Not spelled <i>momma</i>
mumbo sauce	Popular condiment found in DC
murk	Means murder
ratchet	Negative evaluative term
shorty	Female companion
turnt	Variant of turned meaning excited, e.g., turnt up
wilding	To act wild ([waɪlɪŋ])
wont	Past tense <i>wont</i> in Princeville
youngin	Transcribed as –in

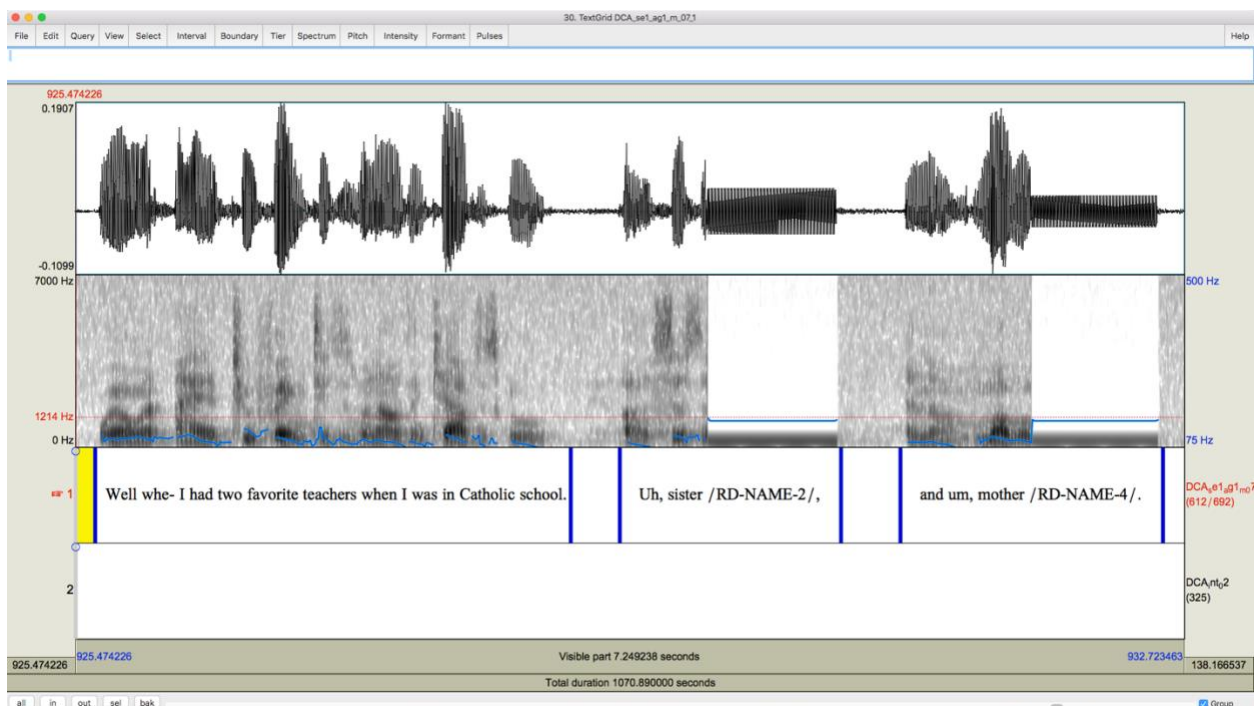
Redaction and participant anonymity

A guiding principle in the development of CORAAL is to protect the anonymity of its participants. Participants who were interviewed specifically for the corpus project (e.g. for CORAAL: ATL, DCB, and VLD) were given the choice in the consent process of whether they wished to be recognized by name, with the default being that they will not be named. A large number of participants (the majority) did ask to be recognized by name and we acknowledge them by name in the acknowledgments section above. (We are incredibly grateful to all of our participants, named and unnamed.) For participants not interviewed by the project team (e.g. CORAAL: DCA, LES, PRV, and ROC), we do not disclose any names, unless there is an explicit (i.e. documented) permission given by the participant. (*If you are a participant in the*

corpus and did not get recognized by name, but wish to, please contact the project team and we will be happy to send you a new consent form, where you can give us this permission.)

Our redaction process involved several steps. During the first round of transcription, transcribers were asked to mark different categories of sensitive information, such as names, street addresses, places of work, and other kinds of personally identifying information. Additionally, transcribers marked the numbers of syllables of the item(s) to be redacted. Redaction codes are described in Part D of the Transcription practices and conventions section, above. The third round of transcription involved the creation of a redaction tier in Praat, where boundaries were placed directly around the portion of the interview to be redacted. The amount of material redacted varies widely by interview. Some interviews have only one or two redacted utterances while others have a great many.

Once completed, redaction ‘bleeps’, which were generated based on the mean pitch and amplitude of the speech being redacted, replaced the sensitive information. An example of two redacted utterances in an interview from CORAAL:DCA is shown below in a screenshot from Praat.



Note that in a few cases, we have excised portions of interviews. This is occasionally by the request of the participant and it is occasionally done as a decision of the project team based on the content of the interview. There are a few cases where we have excised content even though the participant gave us permission to include it. For example, in DCB_se1_ag4_f_01, a passage from 537.8 to 671.8 is excised because of a graphic personal story. This is marked in the Notes section of the metadata spreadsheet. Other excised content includes reading passages and word lists. *Although they are not included in the published corpus, many of the excised portions may still be available for bona fide research use upon request.*

CORAAL transcription formats

Transcripts are available in three formats:

- Praat TextGrids (.TextGrid): Each speaker is on a separate interval tier. Occasionally transcripts have a “misc” (miscellaneous) tier with additional information or transcription of substantial interlopers.
- Plain text (.txt): These files are automatically generated from Praat TextGrids, using a script in Praat. They are formatted as tab-delimited text so they can be opened as text, as Excel spreadsheets, or in R, etc.; they contain the same timestamp information as TextGrids.
- ELAN files (.eaf): ELAN versions of the transcripts were automatically generated from the Praat TextGrids, using ELAN’s “Import Multiple Files As...” feature. (Note that ELAN files are unmodified output from ELAN’s batch conversion tool; the files are not linked to their corresponding audio files and may need to be processed in other ways to maximize their usability in ELAN.)

While they are not published as a part of the official corpus (at this time), the CORAAL development team has also generated phone- and word-level alignments for the transcripts from CORAAL v. 2018.10.06 using the Montreal Forced Aligner (version 1.1) (McAuliffe et al. 2017). The aligned TextGrids and trained language model are available at <http://lingtools.uoregon.edu/coraal/aligned/>.

In addition to releasing the corpus in the formats available here, we are also working on a syntactically parsed version for a large portion of the corpus (similar to Tortora et al.’s ongoing work on [AAPCAppE](#) and [CUNY-CoNYCE](#)). This work is in progress, and we will release those annotated versions of the data once they are completed.

The R functions in http://lingtools.uoregon.edu/coraal/explorer/R/CORAAL_web.R can create a few different versions of the CORAAL transcripts as R (data.frame) objects. These include both utterance-level transcripts (as are provided in the main data files) and turn-level transcripts, where speaker utterances are collapsed on a turn-by-turn basis.

CORAAL Component Details

CORAAL:ATL (Atlanta, GA 2017)

Authors: Charlie Farrington, Tyler Kendall, Patrick Slay Brooks, Lucas Jenson, Chloe Tacata, and Jaidan McLean
Release Date: May 2020 (v. 2020.05)
Interviewer: Patrick Slay Brooks

About CORAAL:ATL

CORAAL:ATL consists of 13 primary speakers across 14 audio files, collected in 2017 and 2018 by Patrick Slay Brooks, a music producer in Atlanta (www.slayinrecords.com/), specifically for CORAAL. Speakers represent a modern friendship network in Atlanta, GA. Atlanta has been described as a “black mecca” in the South (Hobson 2010), especially in the context of the so-called reverse Great Migration, the movement of African Americans from Northern and Western cities back to the (urban) South. Brooks has a friendship group that highlights a diversity of experiences in Atlanta. Speakers range from being born and raised in Atlanta, to growing up in places like New York City, Washington DC, and Los Angeles, CA. As with all sub-components, see metadata for speaker details.

Speakers were interviewed by Brooks for CORAAL to fill a 2 x 2 demographic matrix. In file naming, like with CORAAL:PRV and CORAAL:ROC, the socioeconomic group is listed as “0” (e.g., ATL_se0_ag1_m_01_1) to denote no focus on socioeconomic groups (not to indicate a group lower than 1). We have attempted to capture and include in the metadata broad information about speakers’ demographic backgrounds, such as length of residence and other places lived, but leave questions of interpretation up to end users.

CORAAL:ATL data

The 14 audio files are 44.1 kHz, 16 bit, mono in WAV format, totaling 8.6 hours and 93.5K words. Interviews were recording on a Zoom H5 recorder, with either a lapel microphone or an internal microphone, between 2017 and 2018. Interviews are sociolinguistic styled interviews and conversations on topics such as life in Atlanta, and the interviewee’s neighborhood, schooling, and work history.

Speaker numbers are listed in each cell.

	Socio-Economic Group 0	
	<i>Female</i>	<i>Male</i>
Age Group 1 (under 29)	3	5
Age Group 2 (30 to 50)	2	3

CORAAL:DCA (Washington, DC 1968)

Authors: Tyler Kendall, Ralph Fasold, Charlie Farrington, Jason McLarty, Shelby Arnson, and Brooke Josler

Release Date: January 6, 2018 (v. 2018.01.06)

Update Dates: April 6, 2018 (v. 2018.04.06) [see errata]
October 6, 2018 (v. 2018.10.06) [see errata]

Interviewers: Ralph Fasold, Walt Wolfram, Carolyn Cunningham, Virginia Lundstrom, Veronica Johnson, Roger Shuy, James Goines, and Gail Marble

About CORAAL:DCA

CORAAL:DCA consists of 68 speakers across 74 recordings, originally collected as part of Ralph Fasold's foundational study of African American Language in Washington, DC (Fasold 1972). The speakers were recorded between March 1968 and August 1969, with dates of birth ranging from 1891 to 1958. The 68 speakers selected for CORAAL are not the exact same set of speakers analyzed by Fasold (1972). We have selected speakers from Fasold's interviews to best represent four age groups and three social class groups, although a balanced demographic matrix is not possible given the emphasis of the original project on young speakers. The youngest age group has additional speakers for two reasons: there are lots of these speakers in Fasold's data and their interviews tend to be shorter, so extra speakers were included to increase the amount of total data available for the demographic group. The social class groups are not completely analogous to Fasold's groups, which are based on the Index of Status Characteristics, but are meant to capture broad social strata.

CORAAL:DCA data

The 74 audio files are 44.1 kHz, 16 bit, mono in WAV format, totaling 34.0 hours and 333.5K words. The data reflects a sociolinguistic interview style in the Labovian tradition, with interview topics including games, school, and favorite movies, among others.

Speaker numbers are listed in each cell.

	Socio-Economic Group 1 (≈LWC)		Socio-Economic Group 2 (≈UWC)		Socio-Economic Group 3 (≈MC)	
	<i>Female</i>	<i>Male</i>	<i>Female</i>	<i>Male</i>	<i>Female</i>	<i>Male</i>
Age Group 1 (under 19)	5	8	7	6	6	6
Age Group 2 (20 to 29)	1	1	0	3	5	3
Age Group 3 (30 to 50)	2	1	0	3	1	4
Age Group 4 (51 and over)	0	2	1	1	0	2

CORAAL:DCB (Washington, DC 2016)

Authors: Tyler Kendall, Minnie Quartey, Charlie Farrington, Jason McLarty, Shelby Arnson, and Brooke Josler
Release Date: January 6, 2018 (v. 2018.01.06)
Update Dates: April 6, 2018 (v. 2018.04.06) [several speakers added]
October 6, 2018 (v. 2018.10.06) [one speaker added; see errata]
Interviewers: Minnie Quartey and Carlos Huff

About CORAAL:DCB

CORAAL:DCB currently consists of 48 primary speakers across 63 audio files, collected specifically for CORAAL. The speakers were recorded between July 2015 and December 2017. Speakers were collected through a friend of a friend network to fill a 4 x 3 demographic matrix, as was done for DCA. The socioeconomic groups here are meant to capture broad social strata; the qualitative labels are simple descriptors to help orient users around the ordering. These are not meant to represent theoretically motivated socioeconomic assessments of individuals. They are also not intended to be perfectly analogous to Fasold’s classifications. There are theoretical and practical issues comparing socioeconomic indices in the DC community 50 years apart. We have tried to capture and include in the metadata broad information about speakers’ demographic backgrounds, but leave questions of interpretation up to end-users.

CORAAL:DCB data

The 63 audio files are 44.1 kHz, 16 bit, mono in WAV format, totaling 46.0 hours and 515K words. Most interviews were recorded at a higher audio quality (48 kHz, 32 bit, stereo) but down-sampled for distribution. Interviews are sociolinguistic styled interviews on topics such as life in Washington DC, and the interviewee’s neighborhood, schooling, and work history. (In most cases, interviewers also collected two reading passages, a word list, and metalinguistic commentary although these are not transcribed, redacted, or otherwise included in CORAAL; these may be made available at a later date.)

Speaker numbers are listed in each cell. We seek to eventually include, at minimum, two speakers per demographic cell. In a few cases, more than two speakers are available for a cell, due to the availability of additional persons in these demographic groups.

	Socio-Economic Group 1 (≈WC)		Socio-Economic Group 2 (≈LMC)		Socio-Economic Group 3 (≈UMC)	
	<i>Female</i>	<i>Male</i>	<i>Female</i>	<i>Male</i>	<i>Female</i>	<i>Male</i>
Age Group 1 (under 19)	3	3	1	1	1	1
Age Group 2 (20 to 29)	3	3	2	1	1	0
Age Group 3 (30 to 50)	3	3	2	3	2	2
Age Group 4 (51 and over)	1	2	5	1	2	2

CORAAL:LES (Lower East Side, NY 2008)

Authors: Kara Becker, Charlie Farrington, Tyler Kendall, Chloe Tacata, and Jaidan McLean

Release Date: July 2021 (v. 2021.07)

Interviewers: Kara Becker

About CORAAL:LES

CORAAL:LES consists of 10 primary speakers across 15 audio files, collected by Kara Becker as part of her dissertation research on the Lower East Side (LES) of Manhattan while at New York University. The LES was the site of William Labov's (1966) early sociolinguistic research on English in New York City. Demographically, the neighborhood changed dramatically between the 1960s and the time of Becker's (2010) research. As of 2010, most residents (~66%) were non-white with a median income of almost half of the median income for the borough of Manhattan (Becker 2014).

Speakers were selected for CORAAL to fill a 2 x 3 demographic matrix. For LES, less comprehensive samples are targeted than for CORAAL:DC. We do not focus on socioeconomic strata, but focus on providing a distribution across gender and age groups. In file naming, the socioeconomic group is listed as "0" (e.g., LES_se0_ag2_m_01) to denote no focus on socioeconomic groups (not to indicate a group lower than 1). We have attempted to capture and include in the metadata broad information about speakers' demographic backgrounds, but leave questions of interpretation up to end users. Following CORAAL's current file conventions (see the [Change Log](#)), there is intentionally no Age Group 1.

CORAAL:LES data

The 15 audio files are 44.1 kHz, 16 bit, mono in WAV format, totaling 8.4 hours and 102.2K words. Interviews were recorded on a Zoom H4 digital recorder with a lavalier microphone between August 2008 and September 2009. Interviews come from an oral history and ethnographic fieldwork project (Becker 2010), with a range of topics including growing up and living on the LES, defining community, gentrification, and discussing the future of the neighborhood.

Speaker numbers are listed in each cell. These age groupings differ somewhat from Becker (2010) due to our focus on age and not year of birth. This CORAAL sub-component includes a selection of the African American interviews collected on the Lower East Side. Additionally, Becker (2010) interviewed Chinese, Jewish, Puerto Rican and white residents of the Lower East Side.

	Socio-Economic Group 0	
	<i>Female</i>	<i>Male</i>
Age Group 2 (20 to 29)	2	1
Age Group 3 (30 to 50)	1	2
Age Group 4 (51 and over)	2	2

CORAAL:PRV (Princeville, NC 2004)

Authors: Ryan Rowe, Walt Wolfram, Tyler Kendall, Charlie Farrington, and Brooke Josler

Release Date: April 6, 2018 (v. 2018.04.06)

Update Date: October 6, 2018 (v. 2018.10.06) [see errata]

Interviewers: Ryan Rowe, Jeanine Carpenter, Drew Grimes, Kristy D'Andrea

About CORAAL:PRV

CORAAL:PRV consists of 16 primary speakers across 32 audio files, collected by Ryan Rowe, Walt Wolfram, and colleagues for the North Carolina Language and Life Project. Princeville, NC is the oldest town incorporated by African Americans in the U.S. Many community members can trace their families back to the original founders of the town. The speakers here were recorded between August 2003 and June 2004. As of the 2000 census, African Americans composed 97% of the population (Rowe 2005, see also Kendall 2007b and Kendall & Wolfram 2009).

Speakers were selected for CORAAL from the larger dataset to fill a 2 x 3 demographic matrix. For PRV, as well as additional upcoming corpus sub-components, less comprehensive samples are targeted than for CORAAL:DC. We do not focus on socioeconomic strata, but focus on providing a distribution across gender and age groups. In file naming, the socioeconomic group is listed as “0” (e.g., PRV_se0_ag1_m_01) to denote no focus on socioeconomic groups (not to indicate a group lower than 1). This is meant simply to maintain the file naming structure of CORAAL:DC. We have attempted to capture and include in the metadata broad information about speakers’ demographic backgrounds, but leave questions of interpretation up to end users.

CORAAL:PRV data

The 32 audio files are 44.1 kHz, 16 bit, mono in WAV format, totaling 13.9 hours and 156.1K words. Interviews were recorded on cassette tape and transferred to digital formats in 2007 and 2008 at NC State University. Interviews are sociolinguistic styled interviews on topics such as life in Princeville, schooling, and Hurricane Floyd, which left much of the town underwater in 1999.

Speaker numbers are listed in each cell. This CORAAL sub-component includes only a selection of the total interviews collected in Princeville. There are currently no plans to transcribe more speakers. In a few cases, more than two speakers are available for a cell, due to the availability of additional persons in these demographic groups.

	Socio-Economic Group 0	
	<i>Female</i>	<i>Male</i>
Age Group 1 (under 29)	2	2
Age Group 2 (30 to 50)	3	2
Age Group 3 (51 and over)	4	3

CORAAL:ROC (Rochester, NY 2016)

Authors: Sharese King, Charlie Farrington, Tyler Kendall, Emma Mullen, Lucas Jenson, and Shelby Arnson

Release Date: October 6, 2018 (v. 2018.10.06)

Update Date: May 2020 (v. 2020.05) [one speaker added]

Interviewer: Sharese King

About CORAAL:ROC

CORAAL:ROC consists of 15 primary speakers across 19 audio files, collected in 2016 and 2017 by Sharese King as part of her dissertation research in Rochester, New York. Rochester is city on Lake Ontario, in Monroe County in western New York state. Since the early twentieth century, Rochester has been home to a large African American population (see King 2018).

Speakers were provided by King for CORAAL from a larger dataset to fill a 2 x 3 demographic matrix. For ROC, less comprehensive samples are targeted than for CORAAL:DC. We do not focus on socioeconomic strata, but focus on providing a distribution across gender and age groups. In file naming, like with CORAAL:PRV, the socioeconomic group is listed as “0” (e.g., ROC_se0_ag1_m_01_1) to denote no focus on socioeconomic groups (not to indicate a group lower than 1). We have attempted to capture and include in the metadata broad information about speakers’ demographic backgrounds, but leave questions of interpretation up to end users.

CORAAL:ROC data

The 19 audio files are 44.1 kHz, 16 bit, mono in WAV format, totaling 13.2 hours and 138.9K words. Interviews were recorded on a Zoom H2N recorder, with an AudioTechnica AT831b microphone, and a MiniJack to XLR audio cable between 2016 and 2017, often at the homes of the interviewees. Interviews are sociolinguistic styled interviews on topics such as life in Rochester, schooling, as well as metalinguistic questions about perception of Rochester accents.

Speaker numbers are listed in each cell. In two cases, more than two speakers are available for a cell, due to the availability of additional persons in these demographic groups.

	Socio-Economic Group 0	
	<i>Female</i>	<i>Male</i>
Age Group 1 (under 29)	3	3
Age Group 2 (30 to 50)	4	1
Age Group 3 (51 and over)	2	2

CORAAL:VLD (Valdosta, GA 2017)

Authors: Minnie Quartey, Charlie Farrington, Tyler Kendall, Chloe Tacata, and Jaidan McLean

Release Date: July 2021 (v. 2021.07)

Interviewer: Minnie Quartey

About CORAAL:VLD

CORAAL:VLD consists of 12 primary speakers across 14 audio files, collected specifically for CORAAL between collected in 2017 and 2019 by Minnie Quartey. Quartey, also the primary fieldworker for CORAAL:DCB, recorded these recordings in her hometown of Valdosta, GA with friends and family. Valdosta, the county seat in Lowndes County in South Georgia, has a population of 56,000, approximately 53% African American, as of the 2019 US Census estimate.

Speakers were interviewed by Quartey for CORAAL to fill a 3 x 2 demographic matrix. In file naming, the socioeconomic group is listed as “0” (e.g., VLD_se0_ag2_f_01_1) to denote no focus on socioeconomic groups (not to indicate a group lower than 1). We have attempted to capture and include in the metadata broad information about speakers’ demographic backgrounds, such as length of residence and other places lived, but leave questions of interpretation up to end users. Following CORAAL’s current file conventions (see the [Change Log](#)), there is intentionally no Age Group 1.

CORAAL:VLD data

The 14 audio files are 44.1 kHz, 16 bit, mono in WAV format, totaling 11.5 hours and 112K words. Interviews were recording on a Marantz PMD661 MKII, with a Shure SM93 lapel microphone. Interviews are sociolinguistic styled interviews on topics such as life in Valdosta, personal histories, and high school sports.

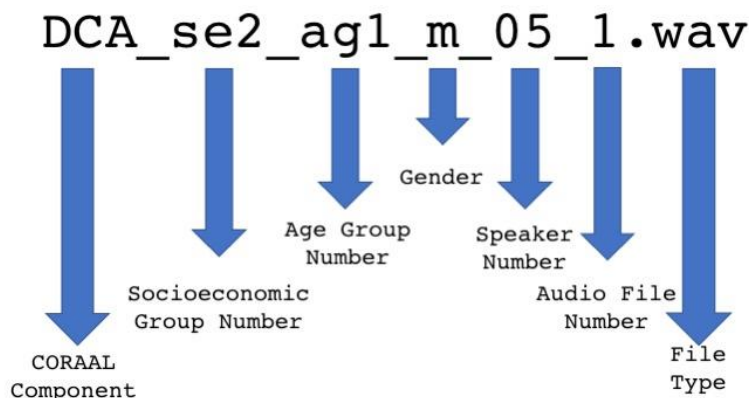
Speaker numbers are listed in each cell.

	Socio-Economic Group 0	
	<i>Female</i>	<i>Male</i>
Age Group 2 (20 to 29)	2	1
Age Group 3 (30 to 50)	2	3
Age Group 4 (51 and over)	2	2

CORAAL Metadata

Speaker and file labeling conventions

Speaker and file names in the corpus are labeled systematically.



For example, file `DCA_se2_ag1_m_05_1.wav` is an audio (WAV) file for DCA, the Washington, DC 1968 component of CORAAL. The file's primary speaker is in socioeconomic group 2 (`se2`), age group 1 (`ag1`; this is the youngest age group in DCA, see metadata information below), male (`m`) number 5 (i.e. the fifth speaker in the cell of the demographic matrix). The final 1 indicates the audio file number. As noted in the descriptions of the smaller sub-components, `se0` is used to notate that a speaker is uncategorized for socioeconomic group (not that the speaker is in group 0). For gender, three codes are used, `f`, for female, `m`, for male, and `n`, for non-binary (as of Version 2021.07, there are no non-binary speakers in CORAAL).

Most speakers are contained in just one audio file but occasionally, such as for `DCA_se2_ag1_m_05`, there are multiple files (this speaker has two files: `DCA_se2_ag1_m_05_1.wav` & `DCA_se2_ag1_m_05_2.wav`). As discussed above, for each WAV file there is a corresponding Praat TextGrid, ELAN file, and plain text file with identical file labels. (Again, all three of these transcript versions contain identical information and are derived from the Praat Textgrid.)

Metadata files and notes

Each component of CORAAL has its own metadata file that contains a range of information about the recordings and their speakers. These files are tab-delimited text files that can be readily opened in a spreadsheet program, like MS Excel, or in R. The metadata files are `.txt` files downloadable with the rest of each corpus component's files. For example, metadata for DCA are in the file labeled `DCA_metadata_2018.10.06.txt` (for version 2018.10.06).

All files in CORAAL have been anonymized and have been trimmed only to contain conversation portions of the sociolinguistic interviews. All of the files are stored (in original formats) on the Sociolinguistic Archive and Analysis Project (SLAAP; <https://slaap.chass.ncsu.edu/>), and SLAAP often contains more files from a sociolinguistic fieldwork project than just those included in CORAAL. The SLAAP codes for the files (e.g.

SLAAP.Spkr) are provided in the metadata file when appropriate. For DCA, SLAAP codes are reflective of the original codes used in Fasold (1972), included here for backward compatibility. For DCB, SLAAP codes were given to recordings as they were completed. For PRV, SLAAP codes are reflective of the codes used when the audio tapes were digitized and uploaded to SLAAP in 2007 and 2008.

Several categories in the metadata spreadsheets apply to all CORAAL components, while others apply only to specific components. Metadata notes are provided below, organized by those that apply to all components first, with information specific to sub-components in the sections following. For DCA, most speaker files obtained from Ralph Fasold came with an Informant Data Sheet (IDS), where much of the metadata information comes from. The IDS collected basic demographic data, such as sex, age, address, birthplace, parents' birthplace, as well as a social class section (see below for discussion). For ATL, DCB, and VLD, interviewers were asked to complete a similar Interview Report Form (IRF) for each speaker, which collected the same kinds of demographic information as Fasold's IDS. In addition to general demographic information, the IRF contains additional interview notes (e.g., interruptions, background noise, etc.) as well as topics covered over the course of the interview. For LES, PRV, and ROC, some speaker information was gathered from the metadata made available by the primary interviewer, while other information was obtained from the content of the interviews themselves.

Metadata notes (all CORAAL components)

- In the following table, each column provided in the tab-delimited text file is outlined. The majority of the categories apply to all components, but some categories require further information about how the categories were coded in each corpus component.

<i>Field</i>	<i>Description and Examples</i>
CORAAL.Sub	CORAAL component three-character code (DCA; DCB; PRV; ROC; ATL)
Version.Created	Version number of the recording's initial release (e.g. v. 2018.01.06, v. 2018.04.06, v. 2018.10.06, v. 2020.05).
Version.Modified	When modifications are made to the transcripts, the most up-to-date version number is listed here. Any modifications will also be listed in the Change Log found earlier in this guide.
CORAAL.Spkr	Speaker code for CORAAL (e.g. DCB_se1_ag1_f_01; PRV_se0_ag3_m_01).
CORAAL.File	File name used for audio file and transcription files (e.g. DCB_se1_ag1_f_01_1; PRV_se0_ag3_m_01_1).
Audio.Folder	Name of folder that houses the specific audio file (e.g. DCA_audio_part01_2018.10.06)
Tarball	Name of compressed folder with extension (tar.gz) (e.g. DCA_audio_part01_2018.10.06.tar.gz)
Primary.Spkr	In DCB & PRV there are several interviews where there are other interviewees present in addition to a primary interviewee. In the metadata, there is a row for each interviewee present in the interview. <i>Yes</i> means the speaker was the <u>primary</u> interviewee, <i>no</i> means the speaker was <u>secondary</u> .
SLAAP.Collection	Alternate codes were used in SLAAP. DCA==wds, DCB=wdc, PRV==prv, VLD==sga.
SLAAP.Spkr	SLAAP speaker code (e.g. wdc003).

SLAAP.Interview	SLAAP interview file code (e.g. wdc0030d).
Gender	Speaker's gender: Male, Female, Non-binary.
Age	Speaker's actual age in years (currently ranging from 12 to 83).
Age.Group	<p>CORAAL:DCA, DCB, LES, VLD AG1 = -19 AG2 = 20 to 29 AG3 = 30 to 50 AG4 = 51+</p> <p>CORAAL:PRV and CORAAL:ROC (three age groups): AG1 = -29 AG2 = 30 to 50 AG 3 = 51+</p> <p>CORAAL:ATL (two age groups): AG1 = -29 AG2 = 30 to 50</p> <p><i>Notes:</i> LES and VLD follow the Age Groups from CORAAL:DC but only contain speakers from AG2, AG3, and AG4.</p> <p>We anticipate updating the Age Group designation for ATL, PRV, and ROC to match DCA, DCB, LES, and VLD to regularize the treatment of Age Group across all CORAAL components.</p>
Year.of.Birth	Speaker's actual year of birth (currently ranging from 1891 to 2005).
Year.of.Interview	<p>ATL interviews took place in either 2017 or 2018. DCA interviews took place in either 1968 or 1969. DCB interviews took place between 2015 and 2017. LES interviews took place between 2008 and 2009. PRV interviews took place in either 2003 or 2004. ROC interviews took place in either 2016 or 2017. VLD interviews took place between 2017 and 2019.</p>
CORAAL.SEC.Group	<p>See DCA and DCB Metadata notes for specific information on how CORAAL.SEC was coded.</p> <p>For sub-components without coding for socioeconomic group (ATL, PRV, ROC, and VLD), SEC.Group is always listed as 0 (e.g., PRV_se0_ag2_f_02).</p>
Education	<p>For DCA, this is what was provided on the IDS, variability in what was reported is a result of different fieldworkers.</p> <p>For ATL, DCB and VLD, this is what was provided on the IRF.</p> <p>For LES, PRV, and ROC this was primarily determined through the information found in the sociolinguistic interview.</p>

Edu.Group	Categories are collapsed from the Education field into groupings: College (completed college) Elementary School (completed elementary school) Graduate School (completed graduate school) High School (completed high school, or equivalent) Some College (attended college but did not graduate) Some High School (attended high school but did not graduate) Student_college (current college student) Student_hs (current high school student) Student_ms (current middle school student).
Occupation	For DCA, this is provided on the IDS (Student; Postal worker; Engineering). For ATL, DCB, VLD, this is provided on the IRF. For LES and ROC, this was determined through information provided by the primary interviewer. For PRV, this was primarily determined through the information provided in the sociolinguistic interview, but occasionally available in the metadata available on SLAAP.
Region.in.City	CORAAL:DC only. Quadrant in DC where the participant resided at the time of the interview (NW; NE; SE; SW). "DC" is used if participant resided in multiple quadrants for an unknown amount of time.
LOR	Length of Residence in current location. CORAAL:DC as listed on the IDS/IRF (e.g. 20 years). For ATL, LES, ROC, and VLD this is provided where available. Not included in PRV metadata because information is not consistent across speakers. Most speakers have lived in the Edgecombe County area their whole life.
LOR.Percent	LOR/Age (e.g. $20 (LOR)/25(Age) = 80\% LOR.Percent$).
Other.Places.Lived	LOR is not available for Other Places Lived. If the LOR.Percent is 100 and there is something listed in this category, assume this is under a year of residence.
Relationship.To.Others.In.Corpus	Occasionally, relationships are made clear in the interview, or is written on the IDS/IRF. The format lists the CORAAL.Code of the speaker, with the relationship in parentheses: e.g. DCB_se1_ag3_f_02 (mother) indicates that DCB_se1_ag3_f_02 is this speaker's mother.
Guardian.1.Birthplace	Mother birthplace.
Guardian.1.Birthplace.State	Mother birthplace by state.
Guardian.1.Education	Mother education level; only available for DCA.
Guardian.1.Occupation	Mother occupation.
Guardian.2.Birthplace	Father birthplace.
Guardian.2.Birthplace.State	Father birthplace by state.
Guardian.2.Education	Father education level; only available for DCA.
Guardian.2.Occupation	Father occupation.
Guardian.Notes	Occasionally, Guardian.1 and Guardian.2 do not correspond to mother and father. This column describes that and any other relevant information.
Interviewer.Code	CORAAL code for interviewers, (e.g., DCA_int_01). In DCB, speaker DCB_se1_ag3_m_01 is also an interviewer. His CORAAL.Code is also used as his Interviewer.Code.

	In PRV, there is one primary interviewer (PRV_int_01), and three secondary interviewers. ATL, LES, ROC, and VLD only have one primary interviewer each (ATL_int_01, LES_int_01, ROC_int_01 and VLD_int_01).
Interviewer.Initials	For ATL, DCA, LES, PRV, ROC and VLD initials are included (e.g., RF; WW; RR; SK; PB, etc.). For DCB, SLAAP codes are included here (e.g., wdc000; wdc001).
Interviewer.Gender	Female, Male.
Interviewer.Ethnicity	Ethnicity of interviewer.
Interviewer.Age	Age is approximated into five year ranges.
Interviewer.Relationship	In DCB, information provided by interviewer. Collapsed into three categories: Acquaintance; Close relationship; No previous relationship . In DCA, LES, PRV, and ROC, no previous relationship is assumed, but this is not definitive. In ATL and VLD all interviews were collected because of a previous relationship between the interviewer and interviewee.
Recording.Equipment	Equipment used during the recording (if known).
Microphone	Microphone(s) used during the recording (if known).
Stereo.Mono	Number of channels used during the recording (if known).
Bit.Rate	Bit rate used during the recording (if known).
Sampling.Rate	Sampling rate used during the recording (if known).
Source.Device	Source device used during the recording (if known).
Dig.Capture.System	System used in digitization (if used).
Dig.Capture.Device	Device used in digitization (if used).
Dig.Capture.Software	Software used in digitization (if used).
Dig.Sampling.Rate	44.1kHz.
Dig.Bit.Rate	16 bit.
Dig.Channels	Mono.
CORAAL.Length.of.Transcript	Length of transcripts in seconds (listed for Primary.Spkr only).
CORAAL.Word.Count	Number of words in each transcript (listed for Primary.Spkr only).
Is.Misc.Tier	Lists whether a miscellaneous tier was included in the CORAAL release.
Notes	Describes the activity and time (in seconds) where a misc tier is relevant, as well as any general notes regarding recording.

Metadata notes: ATL

Recordings were made on a Zoom H5 digital recorder primarily using a lapel mic (Shure SM93 microphone) or the built-in microphone in a variety of settings in the Atlanta area. Atlanta files were recorded in mono with a sampling rate of 44.1kHz and a bit rate of 16.

Metadata notes: DCA

Recordings were made on a reel-to-reel recorder of unknown make and model in a variety of settings (at the Center for Applied Linguistics, the participant's home, a local church, etc.). Digitization occurred at North Carolina State University in 2013, using an AEC A-6300 Reel-to-Reel player, captured on Audacity with a Sound Devices USBPre 2 preamp, in an effort supervised by Michael J. Fox. Most files were digitized with sampling rates primarily at 48kHz, but a few were digitized at 44.1kHz. For consistency across CORAAL, all audio files were converted to mono with a sampling rate of 44.1kHz and a bit rate of 16.

<i>Field</i>	<i>Description and Examples</i>
In.Fasold.1972	This column reports whether the speaker is examined in Fasold (1972). The 68 speakers selected for CORAAL are <i>not</i> the exact same set of speakers analyzed by Fasold (1972), which only used a selection of the entire dataset of 90 speakers. We have selected speakers from Fasold’s interviews to best represent four age groups and three social class groups.
Fasold.ISC.Group	The categories are based on Warner et al.’s (1960) Index of Status Characteristics (ISC). The actual number value is provided in <i>ISC.Fasold</i> . More information about this rating system can be found in Fasold (1972:17-21). An example ISC is shown below.
Fasold.SEC.Group	This category is Charlie Farrington’s attempt to collapse the detailed ISC groupings into workable class categories, paying attention to number of speakers per group. The categories are: <ul style="list-style-type: none"> ▪ Lower Working Class (CORAAL.SEC.Group 1) ▪ Upper Working Class (CORAAL.SEC.Group 2) ▪ Lower Middle Class (CORAAL.SEC.Group 3) ▪ Upper Middle Class (CORAAL.SEC.Group 3) ▪ Upper Class (CORAAL.SEC.Group 3).
Occ.Score.Fasold	ISC Occupation Score. See Fasold.ISC.Group above, and Fasold (1972: 17-21) for more information.
House.Score.Fasold	ISC House Type Score. See Fasold.ISC.Group above, and Fasold (1972: 17-21) for more information.
Dwelling.Area.Score.Fasold	ISC Dwelling Area Score. See Fasold.ISC.Group above, and Fasold (1972: 17-21) for more information.
ISC.Fasold	ISC Score calculated by Fasold. See rating example below.

- *Fasold ISC Example*
 - The categories are based on Warner et al.’s (1960) Index of Status Characteristics (ISC). The actual number value is provided in *ISC.Fasold*. More information about this rating system can be found in Fasold (1972:17-21).

Sample ISC (Speaker # DCA_se1_ag1_m_06)

	<i>Rating</i>	<i>Weight</i>	<i>Total</i>
<i>Occupation</i>	7	5	35
<i>House Type</i>	6	4	24
<i>Dwelling Area</i>	5	3	15
I.S.C. 74 (Lower Working)			

Metadata notes: DCB

Almost all interviews for DCB were recorded on a Marantz PMD661 MK II digital recorder, using SHURE SM93 and Microflex MX 100 microphones. These recordings were made in stereo, with a bit rate of 24 (PCM-24), and a sampling rate of 48kHz. The last few interviews were recorded using a SHURE SM93 microphone alone, in mono, with a bit rate of 24 (PCM-24), and a sampling rate of 44.1kHz. There is one interview, recorded in early 2015, using an Olympus LS11 digital recorder and an Audio Technica AT8010 Omnidirectional Condenser microphone. Stereo files were converted to mono and down sampled to 44.1kHz/16 bit for the CORAAL release (this was done using the sox software).

<i>Field</i>	<i>Description and Examples</i>
CORAAL.SEC.Group	<p>These SEC groups are estimations based on the fieldworker’s knowledge of the individuals and the African American community in Washington D.C. We’ve done our best to be careful in making group assignments, but these group assignments are meant for balancing the corpus and as a heuristic. They are not meant to be taken as the outcome of a sociological/socio-economic analysis of the speakers.</p> <ul style="list-style-type: none"> ▪ Group 1 – roughly correlates to Working Class ▪ Group 2 – roughly correlates to Lower Middle Class ▪ Group 3 – roughly correlates to Upper Middle Class.

Metadata notes: PRV

Recordings were made on cassette recorders (most likely Marantz PMD222 devices, with a Sony ECM-44B lavalier microphone) in a variety of settings in the Princeville community (the participant’s home; fire house, etc.). Digitization occurred at North Carolina State University in 2006/2007 by graduate student employees using a Tascam CC-222 MKIII to convert the tape to compact disc format, which was then captured on a Mac OS using Audacity or, in a few cases, Praat. Princeville files were digitized with sampling rate of 44.1kHz and a bit rate of 16. All PRV audio files are mono with a sampling rate of 44.1kHz and a bit rate of 16.

Metadata notes: LES

Recordings were made on a Zoom H4 digital recorder using a lavalier microphone primarily in offices or the homes of participants. Lower East Side files were recorded in mono with a sampling rate of 44.1kHz and a bit rate of 16.

We are able to include metadata for CORAAL:LES some additional information available from Becker (2010).

<i>Field</i>	<i>Description and Examples</i>
Becker.Class.Score	<p>As described by Becker (2010: 73-80), a social class score was created for each individual based on three measures: education, occupation, and housing type. Each measure included a six-point scale describing a given score in detail. The African American speaker group ranges from 4 to 18, with a mean of 10. For analysis, Becker, divided the entire sample (N=64) into five evenly distributed class ranks.</p>
Becker.Generation	<p>As described in Becker (2010: 80-82), a generation value is given for each speaker to identify how long their family has lived in the New York City area. This value is independent of Length of Residence (LOR) and speaker age.</p> <ul style="list-style-type: none"> ▪ 1.5 – born elsewhere, but moved before age 5, both parents from elsewhere ▪ 2 – born in NYC, both parents from elsewhere ▪ 2.5 – born in NYC with one local parent ▪ 3 – born in NYC, with either of the following holding as well: both parents were local, or 1 hyperlocal parent (a hyperlocal parent is someone who was born in NYC and whose parents were also born in NYC) and 1 parent from elsewhere

	<ul style="list-style-type: none"> ▪ 3.5 – born in NYC with 1 or more grandparents born in NYC and other parent local ▪ 4 – born in NYC with all grandparents born here” <p>The African American speaker group has a mean generation of 2.5.</p>
--	--

For the two components in CORAAL v. 2021.07 (CORAAL:LES and CORAAL:VLD), Age.Group is organized slightly differently than it was for ATL, PRV, and ROC. Beginning with CORAAL:LES and CORAAL:VLD, the Age.Group age ranges match those in CORAAL:DC, with speakers distributed across Age Groups 2-4. Age Group 1 is not used because the components do not include speakers under the age of 20.

Metadata notes: ROC

Recordings were made on a Zoom H2N digital recorder, with an AudioTechnica AT831b microphone, connected with a MiniJack to XLR audio cable in a variety of settings in the Rochester area. Rochester files were recorded with a sampling rate of 44.1kHz and a bit rate of 16.

Metadata notes: VLD

Recordings were made on a Marantz PMD 661 MKII digital recorder primarily using a Shure SM93 lapel microphone in a variety of settings in the Valdosta area. Valdosta files were mostly recorded in mono with a sampling rate of 44.1kHz and a bit rate of 16, though some raw files were recorded in stereo with a sampling rate of 48kHz, but these were down-sampled to match the rest of CORAAL.

For the two components in CORAAL v. 2021.07 (CORAAL:LES and CORAAL:VLD), Age.Group is organized slightly differently than it was for ATL/PRV/ROC. Beginning with CORAAL:LES and CORAAL:VLD, the Age.Group age ranges match those in CORAAL:DC, with speakers distributed across Age Groups 2-4. Age Group 1 is not used because the component does not include speakers under the age of 20.

Projects Using CORAAL

We may not be able to keep an accurate list of projects and publications that have used CORAAL indefinitely, but for now we attempt to provide a list of current and published studies that have used CORAAL here. You can also view CORAAL citations on [Google Scholar](#). *Please let us know if you are using CORAAL and want to be added to this list!*

Publications

Aranson, Shelby & Charlie Farrington. 2017. Twentieth century sound change in Washington DC African American English. *Penn Working Papers in Linguistics* 23(2): Article 2.
<https://repository.upenn.edu/pwpl/vol23/iss2/2/>

- Cukor-Avila, Patricia & Ashley Balcazar. 2019. Exploring grammatical variation in the Corpus of Regional African American Language. *American Speech* 94(1): 36-53. DOI: <https://doi.org/10.1215/00031283-7321989>
- Farrington, Charlie. 2018. Incomplete neutralization in African American English: The case of final consonant voicing. *Language Variation and Change* 30(3): 361-383. DOI: <https://doi.org/10.1017/S0954394518000145>
- Farrington, Charlie. 2019. Language variation and the Great Migration: Regionality and African American Language. Doctoral dissertation, University of Oregon.
- Farrington, Charlie & Natalie Schilling. 2019. Contextualizing the Corpus of Regional African American Language, D.C.: AAL in the Nation's Capital. *American Speech* 94(1): 21-35. DOI: <https://doi.org/10.1215/00031283-7308060>
- Forrest, Jon & Walt Wolfram. 2019. The status of (ING) in African American Language: A quantitative analysis of social factors and internal constraints. *American Speech* 94(1): 72-90. DOI: <https://doi.org/10.1215/00031283-7308049>
- Grieser, Jessica A. 2019. Investigating topic-based style shifting in the classic sociolinguistic interview. *American Speech* 94(1): 54-71. DOI: <https://doi.org/10.1215/00031283-7322011>
- Holliday, Nicole R. 2019. Variation in question intonation in the Corpus of Regional African American Language. *American Speech* 94(1): 110-130. DOI: <https://doi.org/10.1215/00031283-7308038>
- Kendall, Tyler, & Valerie Fridland. 2021. *Sociophonetics*. Cambridge University Press.
- Kendall, Tyler, Charlotte Vaughn, Charlie Farrington, Kaylynn Gunter, Jaidan McLean, Chloe Tacata, & Shelby Arnson. 2021. Considering performance in the automated and manual coding of sociolinguistic variables: Lessons from variable (ING). *Frontiers in Artificial Intelligence*, 4, 43. DOI: <https://doi.org/10.3389/frai.2021.648543>
- Koenecke, Allison, Andrew Nam, Emily Lake, Joe Nudell, Minnie Quartey, Zion Mengesha, Connor Toups, John R. Rickford, Dan Jurafsky, & Sharad Goel. 2020. "Racial Disparities in Automated Speech Recognition." *Proceedings of the National Academy of Sciences* 117 (14): 7684–89. DOI: <https://doi.org/10.1073/pnas.1915768117>
- McLarty, Jason, Taylor Jones, & Christopher Hall. 2019. Corpus-based sociophonetic approaches to post-vocalic R-lessness in African American Language. *American Speech* 94(1): 91-109. DOI: <https://doi.org/10.1215/00031283-7362239>
- Quartey, Minnie & Natalie Schilling. 2019. Shaping 'connected' vs. 'disconnected' identities in narrative discourse in DC African American Language. *American Speech* 94(1): 131-147. DOI: <https://doi.org/10.1215/00031283-7322000>
- Tan, Samson, Shafiq Joty, Lav R. Varshney, and Min-Yen Kan. 2020. Mind Your Inflections! Improving NLP for Non-Standard English with Base-Inflection Encoding. *ArXiv:2004.14870 [Cs]*, April. <http://arxiv.org/abs/2004.14870>
- Tanner, James, Morgan Sonderegger, Jane Stuart-Smith, & Josef Fruehwald. 2019. Toward "English" Phonetics: Variability in the Pre-consonantal Voicing Effect Across English Dialects and Speakers. *Frontiers in Artificial Intelligence* 3, 38. DOI: <https://doi.org/10.3389/frai.2020.00038>

References

- Bailey, Guy & Natalie Maynor. 1985. The present tense of be in Southern Black folk speech. *American Speech* 60:195-213.
- Bailey, Guy & Natalie Maynor. 1987. Decreolization? *Language in Society* 16:449-74.
- Bailey, Guy, Natalie Maynor, & Patricia Cukor-Avila. eds. 1991. *The Emergence of Black English: Text and Commentary*. Amsterdam and Philadelphia: John Benjamins.
- Becker, Kara. 2010. *Regional dialect features on the Lower East Side of New York City: Sociophonetics, ethnicity, and identity*. Ph.D. dissertation. New York, New York: New York University.
- Berez-Kroeker, Andrea, Gary Holton, Susan Smyth Kung, & Peter Puslifier. Reproducible research in linguistics: toward a data-driven science of language. Paper presented at the 2017 Annual Meeting of the Linguistics Society of America. Austin, TX.
- Boersma, Paul & David Weenink. 2014. Praat: Doing phonetics by computer. [Software]
- Fasold, Ralph W. 1972. *Tense Marking in Black English*. Washington, DC: Center for Applied Linguistics.
- Green, Lisa. 2002. *African American English: A Linguistic Introduction*. Cambridge, U.K.: Cambridge University Press.
- Jones, Taylor. 2016. AAE talmbout: An overlooked verb of quotation. *Penn Working Papers in Linguistics* 22: Article 11.
- Kendall, Tyler. 2007a. Enhancing sociolinguistic data collections: The North Carolina Sociolinguistic Archive and Analysis Project. *Penn Working Papers in Linguistics* 13(2):15-26.
- Kendall, Tyler. 2007b. “The people what makes the town”: The semiotics of home and town spaces in Princeville, NC. *The North Carolina Folklore Journal* 54(1): 33-53.
- Kendall, Tyler. 2008. On the history and future of sociolinguistic data. *Linguistic and Language Compass* 2:332-351.
- Kendall, Tyler. 2011. Corpora and from a sociolinguistic perspective (Corpora sob uma perspectiva sociolinguística). In Stefan Th. Gries (ed.), *Corpus studies: Future directions*, special issue of *Revista Brasileira de Linguística Aplicada* 11(2):361-389.
- Kendall, Tyler, Jason McLarty & Charlie Farrington. 2020. *ORAAL: Online Resources for African American Language*. Eugene, OR: Online Resources for African American Language Project. <http://oraal.uoregon.edu/>
- Kendall, Tyler & Walt Wolfram. 2009. Local and external language standards in African American English. *Journal of English Linguistics*, 37(4): 305-330.
- Kendall, Tyler, Joan Bresnan, & Gerard Van Herk 2011. The dative alternation in African American English: Researching syntactic variation and change in a conglomerated corpus. *Corpus Linguistics and Linguistic Theory*, 7(2):229-244.
- Kurath, Hans. 1949. *A Word Geography of the Eastern United States*. Ann Arbor: University of Michigan Press.
- Labov, William, Paul Cohen, Clarence Robins & John Lewis. 1968. *A Study of the Non-Standard English of Negro and Puerto Rican Speakers in New York City*. Washington, D.C.: United States Office of Education Final Report, Research Project 3288.
- Labov, William. 1966. *The Social Stratification of English in New York City*. Washington DC: Center for Applied Linguistics.

- Labov, William 1969. *The Logic of Nonstandard English*. In James Alatis (ed.), Monograph Series on Languages and Linguistics, No. 22 [20th Annual Round Table: Linguistics and the Teaching of Standard English to Speakers of Other Languages or Dialects], Washington, DC: Georgetown University Press. 1-43.
- Labov, William. 1972. *Language in the Inner City: The Black English Vernacular*. Philadelphia, PA: University of Pennsylvania Press.
- Labov, William. 1987. Contribution to: Are black and white vernaculars diverging? Papers from the NWAV XIV panel discussion. *American Speech* 68(2):5-12.
- Lanehart, Sonja (ed.). 2015. *The Oxford Handbook of African American Language*. Oxford, UK: Oxford University Press.
- McAuliffe, Michael, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger (2017). Montreal Forced Aligner. Version 1.0.0. Computer program. <http://montrealcorpus.tools.github.io/Montreal-Forced-Aligner/>.
- McDavid, R. I. & V. G. McDavid. 1951. The relationship of the speech of negroes to the speech of whites. *American Speech* 26:3-17.
- Mufwene, Salikoko, John Rickford, Guy Bailey, & John Baugh (eds.). 1998. *African-American English: Structure, History and Use*. London: Routledge.
- Poplack, Shana & Sali Tagliamonte. 2001. *African American English in the Diaspora*. Malden/Oxford: Blackwell.
- Rickford, John R. 1999. *African American English: Features, Evolution, and Educational Implications*. Malden/Oxford: Blackwell.
- Rickford, John R., Arnetha Ball, Renee Blake, Raina Jackson, & Nomi Martin. 1991. Rappin on the copula coffin: Theoretical and methodological issues in the analysis of copula variation in African-American Vernacular English. *Language Variation and Change* 3: 103-132.
- Rickford, John R., Julie Sweetland, Angela Rickford, & Thomas Grano. 2012. *African American, Creole, and Other Vernacular Englishes in Education: A Bibliographic Resource*. New York: NCTE-Routledge Research Series.
- Rowe, Ryan. 2005. The development of African American English in the Oldest Black Town in America: Plural -s Absence in Princeville, NC. Master's Thesis. Raleigh: North Carolina State University.
- Schneider, Edgar W. (ed.). 1996. *Focus on the USA*. Philadelphia/Amsterdam: John Benjamins.
- Stewart, William A. 1968. Continuity and change in American Negro dialects. *The Florida FL Reporter* 6:3-4,14-16, 18.
- Thomas, Erik R. & Tyler Kendall. 2007. NORM: The Vowel Normalization and Plotting Suite. Online resource. <http://lingtools.uoregon.edu/norm/>
- Tortora, Christina, Beatrice Santorini, Michael Montgomery, & Frances Blanchette. 2012. A hands-on introduction to the Audio-Aligned and Parsed Corpus of Appalachian English (AAPCAppE). Workshop at Southeastern Conference on Linguistics 79. Lexington, KY.
- Warner, Natasha. 2004. Sharing of data as it relates to human subjects issues and data management plans. *Language and Linguistics Compass* 8: 512-518.
- Wolfram, Walter A. 1969. *A Sociolinguistic Description of Detroit Negro Speech*. Washington, DC: Center for Applied Linguistics.
- Wolfram, Walt & Erik R. Thomas. 2002. *The Development of African American English*. Malden/Oxford: Blackwell.

Yaeger-Dror, Malcah & Erik R. Thomas (eds). 2010. *African American Speakers and their Participation in Local Sound Changes: A Comparative Study*. Durham, NC: Duke University Press.