

Singing Labeling Data Tool Design Document

Group members:

Nandith Narayan nnarayan2018@my.fit.edu

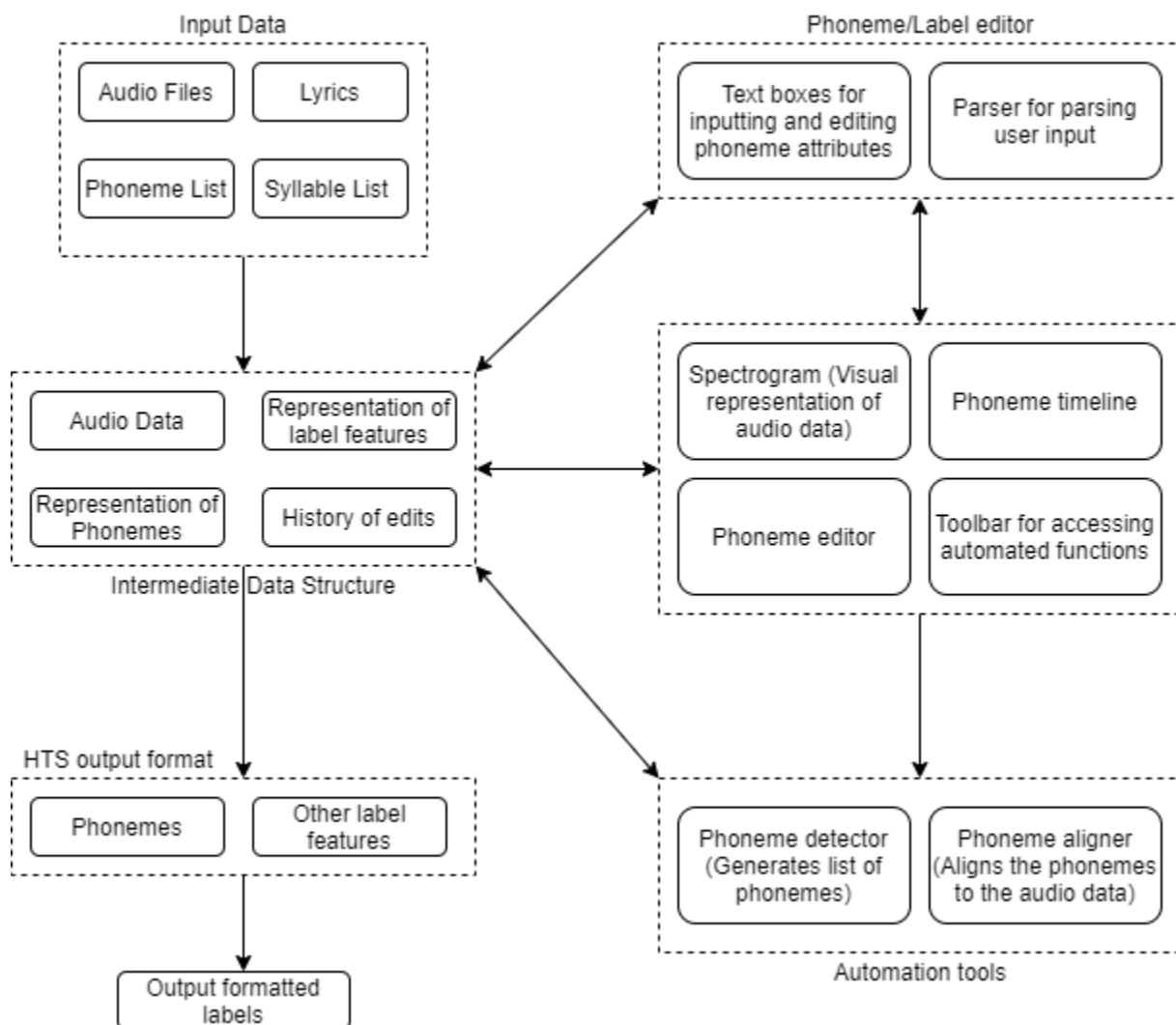
Avinash Persaud apersaud2018@my.fit.edu

Carlos Cepeda ccepeda2018@my.fit.edu

Advisor:

Dr. William Shoaff

System Architecture Diagram:

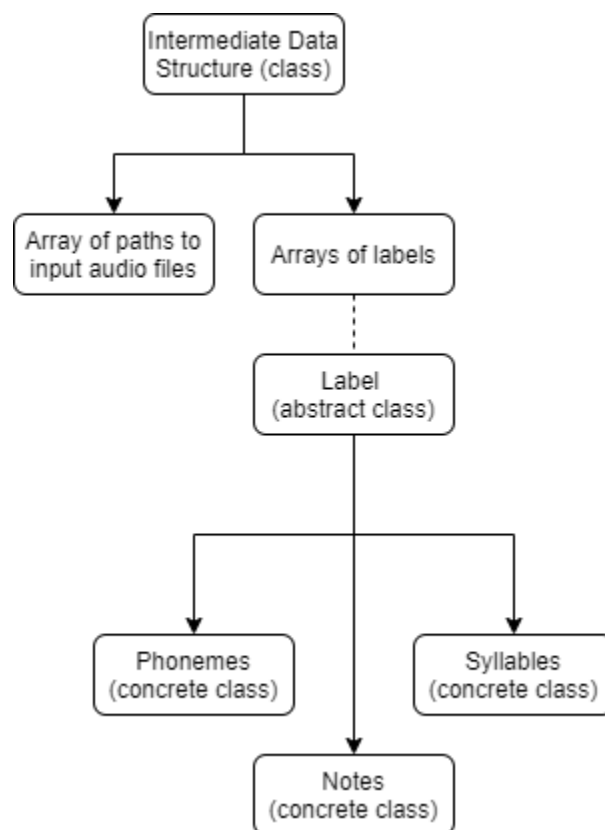


Input Data:

The user can import data into the project at any time. This data can be audio data, a list of phonemes, a list of notes, a list of syllables, and the lyrics of the song(s). The user will be presented with a Graphical User Interface to select which files to import into the project when they press the import data button.

Intermediate Data Structure:

The intermediate data structure will store all the data required by the tool and the data inputted by the user. This will be used to let the user save a project to load it later and resume labeling. It will be a class which contains arrays of instances of label classes. These label classes will all inherit from an abstract label class. So, every different type of label will have its own class. The audio files will be represented by their path.



Upon clicking the save button, the contents of the intermediate data structure will be serialized and written to files within a project folder. Upon loading a project, all the data from the saved files will be read and parsed into the intermediate data structure. The data will be stored in a text format and the audio files will be stored as a path to the file rather than the file itself.

Phoneme/Label Editor

This editor will contain text boxes for the user to input phoneme attributes. The Graphical User Interface for this editor will contain buttons to add/remove attributes as well as edit them. All the text from every text box will be concatenated and parsed as a whole. Parsing will use Antlr v4 as a parser generator.

phonemes[0].start * 10000000

phonemes[0].end * 10000000

phonemes[0].type

@

phonemes[-2]

^

phonemes[-1]

-

phonemes[0]

+

phonemes[1]

Insert dynamic entry

On failure insert:

xx

Name:

HTS Song Full Label

Repeat for each:

Phoneme

V

Cancel

Save

Example grammar:

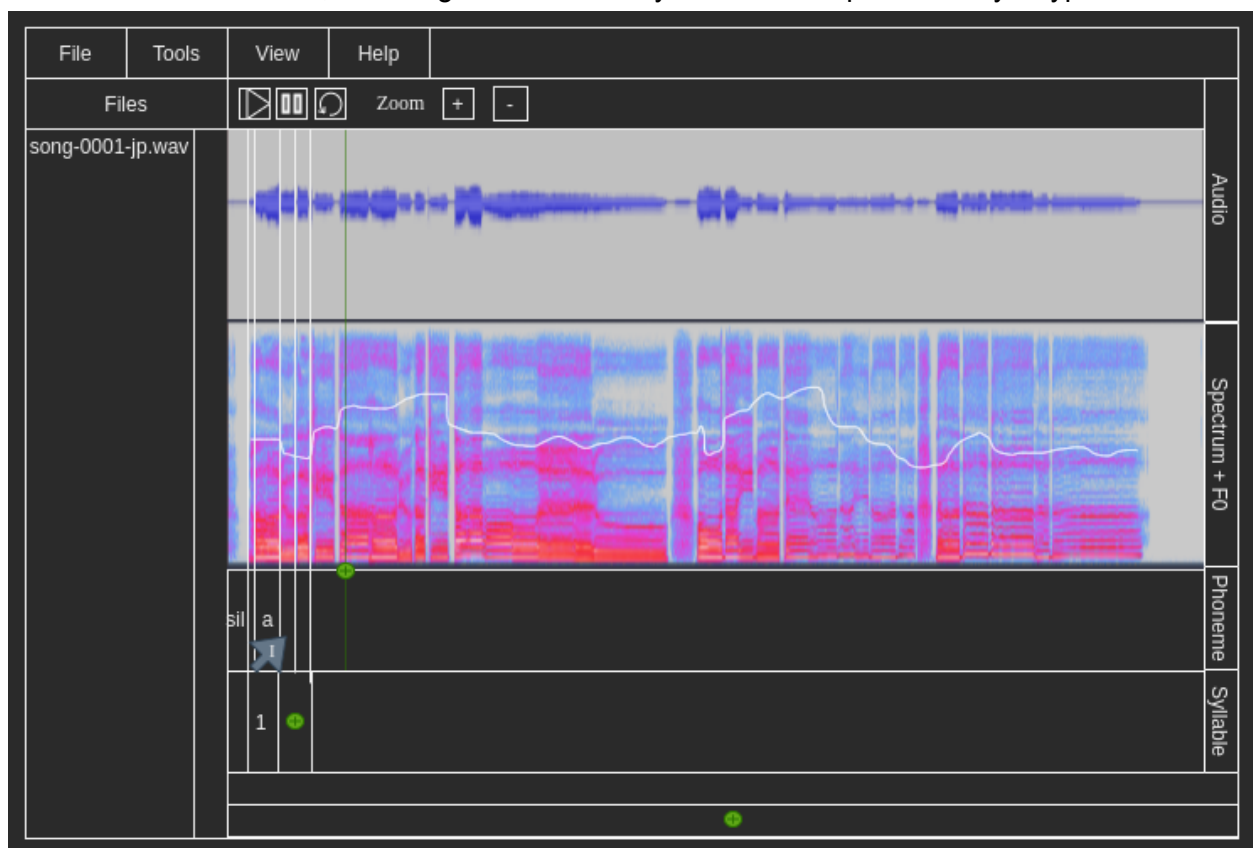
S -> statement*
 statement -> assign | operation
 assign -> label (dot attribute)? = value
 label -> type openSqrBracket index closeSqrBracket
 attribute -> identifier
 operation -> label (dot attribute)? binaryOp value
 binaryOp -> '*' | '-' | '+' | '/' | '^' | '%'
 value -> integer
 type -> identifier
 index -> integer
 dot -> '.'
 openSqrBracket -> '['
 closeSqrBracket -> ']
 identifier -> [a-zA-Z_]*
 integer -> '-'? [0-9]*

Automation tools

The user will be able to click on the toolbar and select from a list of tools for automation. These tools consist of automatic phoneme detection, and automatic phoneme alignment. Automatic phoneme detection will generate a list of phonemes. Automatic phoneme alignment will align those phonemes to the audio data.

Graphical User Interface

The graphical user interface will contain multiple scenes. These scenes are listed as follows, the phoneme editor, the primary labeling scene, the import data scene, and the export scene. The primary labeling scene will display a graphical representation of the raw audio data as well as a spectrogram of the frequencies present in the audio data. It will also contain a timeline for adding, removing, editing, and duplicating phonemes. Other feature layers can be added and edited in a similar matter, including user defined layers based on premade layer types.



The import data scene will open up a window that will let the user select which files they desire to import into the project. The export data scene will open up a window that will let the user select where they want the output to be written to, as well as the format to output the data as. The user will be able to use shortcuts like ctrl+c, ctrl+v, and ctrl+d to copy, paste, duplicate phonemes on the timeline. The user will be able to use the scroll wheel to change the scaling of the timeline and audio data.

HTS output format

The HTS output format represents the labels as regular expressions. The specification for this format is as follows:

$p_1 \sim p_2 - p_3 + p_4 = p_5 @ p_6 - p_7$
 /A: a₁ a₂ a₃ /B: b₁ b₂ b₃ @ b₄ b₅ & b₆ b₇ # b₈ b₉ \$ b₁₀ b₁₁ ! b₁₂ b₁₃ ; b₁₄ b₁₅ | b₁₆ /C: c₁ + c₂ + c₃
 /D: d₁ d₂ /E: e₁ + e₂ @ e₃ + e₄ & e₅ + e₆ # e₇ + e₈ /F: f₁ f₂
 /G: g₁ g₂ /H: h₁ = h₂ ^ h₃ = h₄ | h₅ /I: i₁ = i₂
 /J: j₁ + j₂ - j₃

p ₁	the phoneme identity before the previous phoneme
p ₂	the previous phoneme identity
p ₃	the current phoneme identity
p ₄	the next phoneme identity
p ₅	the phoneme after the next phoneme identity
p ₆	position of the current phoneme identity in the current syllable (forward)
p ₇	position of the current phoneme identity in the current syllable (backward)
a ₁	whether the previous syllable stressed or not (0: not stressed, 1: stressed)
a ₂	whether the previous syllable accented or not (0: not accented, 1: accented)
a ₃	the number of phonemes in the previous syllable
b ₁	whether the current syllable stressed or not (0: not stressed, 1: stressed)
b ₂	whether the current syllable accented or not (0: not accented, 1: accented)
b ₃	the number of phonemes in the current syllable
b ₄	position of the current syllable in the current word (forward)
b ₅	position of the current syllable in the current word (backward)
b ₆	position of the current syllable in the current phrase (forward)
b ₇	position of the current syllable in the current phrase (backward)
b ₈	the number of stressed syllables before the current syllable in the current phrase
b ₉	the number of stressed syllables after the current syllable in the current phrase
b ₁₀	the number of accented syllables before the current syllable in the current phrase
b ₁₁	the number of accented syllables after the current syllable in the current phrase
b ₁₂	the distance per syllable from the previous stressed syllable to the current syllable
b ₁₃	the distance per syllable from the current syllable to the next stressed syllable
b ₁₄	the distance per syllable from the previous accented syllable to the current syllable
b ₁₅	the distance per syllable from the current syllable to the next accented syllable
b ₁₆	name of the vowel of the current syllable
c ₁	whether the next syllable stressed or not (0: not stressed, 1: stressed)
c ₂	whether the next syllable accented or not (0: not accented, 1: accented)
c ₃	the number of phonemes in the next syllable
d ₁	gpos (guess part-of-speech) of the previous word
d ₂	the number of syllables in the previous word
e ₁	gpos (guess part-of-speech) of the current word
e ₂	the number of syllables in the current word
e ₃	position of the current word in the current phrase (forward)
e ₄	position of the current word in the current phrase (backward)
e ₅	the number of content words before the current word in the current phrase
e ₆	the number of content words after the current word in the current phrase
e ₇	the distance per word from the previous content word to the current word
e ₈	the distance per word from the current word to the next content word
f ₁	gpos (guess part-of-speech) of the next word
f ₂	the number of syllables in the next word
g ₁	the number of syllables in the previous phrase
g ₂	the number of words in the previous phrase
h ₁	the number of syllables in the current phrase
h ₂	the number of words in the current phrase
h ₃	position of the current phrase in this utterance (forward)
h ₄	position of the current phrase in this utterance (backward)
h ₅	TOBI endtone of the current phrase
i ₁	the number of syllables in the next phrase
i ₂	the number of words in the next phrase
j ₁	the number of syllables in this utterance
j ₂	the number of words in this utterance
j ₃	the number of phrases in this utterance

Time Level User Interface

The main project level User Interface will display the spectrogram and waveform of the audio data as well as feature layers in a timeline format. By scrolling with the mouse wheel, the time scale can be changed. The spectrogram and waveform display will be generated dynamically. Feature layers will be overlaid on top of the spectrogram. The time dependent parts of the UI will be drawn using a `time_scale` value that will be modified when the user scrolls the mouse wheel or manually sets it in a dropdown menu.