and /d/ and sonorants like /ɹ/. In fact, the model tends to over-estimate its confidence in these segments and compress all of the phones that it is not as confident about into a tiny margin in between the ones that it is more confident in. This is an undesirable behaviour, but we encountered this same tendency, to minimize certain phones in models used on this task, for regular speech. Beyond this problem, the model does place certain boundaries in nearly the same position as the hand-alignment, which is encouraging, as placing boundaries with high accuracy is necessary for the model to be useful in automating the alignment task.
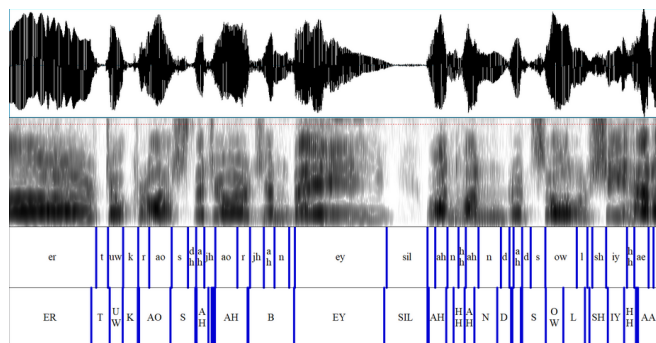


**Figure 1:** Sample alignment of a section of a song. From top to bottom: Waveform, spectrogram, target transcription, model transcription.

## 4  Discussion

One important question that arises from our results is why, despite achieving good frame-wise accuracy rates in the training set, the alignments produced by the model are so bad in certain places. Manual inspection of the alignment suggests that the model seems to almost ignore certain phones entirely, which is most likely due to it not having as much confidence in being able to identify those phones correctly when compared to sonorants and bust releases. If its confidence is spread between many classes, it will have a hard time deciding which one it should select, and may ignore those lower probability phones in favour of items on which the model reaches a much higher level of confidence. For this reason it may be very useful for us to look at the way that we translate the labels outputted by the model into a Praat TextGrid, as there might be a way to do this that takes into consideration the fact that some classes of phone will have overall lower levels of probability than others due to them appearing more similar to other options. Minimum and maximum duration constraints, for example, on the alignment algorithm should help ameliorate this behavior.

Future work should focus on determining why the model is making classification mistakes. A confusion matrix or inverse layer maximization may help indicate what kinds of mistakes the network is making.

## 5  Conclusion

In this project, we achieved our main initial goals of creating a model that is able to classify phones when presented with singing as well as produce a time-aligned Praat TextGrid that can be compared to the original audio track. Some of our other goals, such as using various audio pre-processing methods to compare their effectiveness, we did not achieve. Despite the difficulty in creating from scratch a dataset that would be sufficiently large to adequately train a model for this challenging task, we were able to create a model that is able to classify phones with some degree of accuracy. Excitingly, this suggests to us that this task, although more difficult than the same task performed on regular speech, can be handled in a similar way and with similar levels of success.

Ultimately, we hope to improve the model and the automatic phone alignment further, until it can be packaged as a bespoke application for use in research. We also plan to explore the many questions about how our model can classify sung phonemes, and what differences and similarities it might have with humans, as we continue to test and improve upon it.

## References

[1] Jiahong Yuan and Mark Liberman. Speaker identification on the SCOTUS corpus. *Journal of the Acoustical Society of America*, 123(5):3878, 2008.

[2] Michael McAuliffe, Michaela Socolof, Sarah Mihuc, Michael Wagner, and Morgan Sonderegger. Montreal Forced Aligner: Trainable text-speech alignment using Kaldi. In *Interspeech 2017*, pages 498–502, 2017.

[3] Rong Gong and Xavier Serra. Singing voice phoneme segmentation by hierarchically inferring syllable and phoneme onset positions. In *Proc. Interspeech 2018*, pages 716–720, 2018.

[4] "Moses and Frances Asch Collection". Sound Studies Institute: University of Alberta, May 2018.

[5] O.J.Abbott; Artist; Edith Fawke; Lining Notes. Irish and British Folk Songs from the Ottawa Valley, 1961.

[6] Karen James; Artist; Samuel Gesser; Lining Notes. Karen James, 1961.

[7] Paul Boersma and David Weenink. Praat, a system for doing phonetics by computer, 2018.

[8] John S Garofolo, Lori F Lamel, William M Fisher, Jonathan G Fiscus, and David S Pallett. DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM. NIST speech disc 1-1.1. *NASA STI/Recon technical report n*, 93, 1993.

[9] Matthew C. Kelley and Benjamin V. Tucker. A comparison of input types to a deep neural network-based forced aligner. In *Proc. Interspeech 2018*, pages 1205–1209, 2018.