

Nkululeko: A Python package to predict speaker characteristics with a high-level interface.

Felix Burkhardt ^{1,2*} and Bagus Tris Atmaja ^{3*}

¹ audEERING GmbH, Germany ² TU Berlin, Germany ³ National Institute of Advanced Industrial Science and Technology (AIST), Japan * These authors contributed equally.

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

Summary

Nkululeko (Burkhardt et al., 2022) is open-source software written in Python and hosted on GitHub. It is predominantly a framework for audio-based machine learning explorations without the need to write Python code, and is strongly based on machine learning packages like sklearn (Pedregosa et al., 2011) and pytorch (Chaudhary et al., 2020). The main features are: training and evaluation of labelled speech databases with state-of-the-art machine learning approach and acoustic feature extractors, a live demonstration interface, and the possibility to store databases with predicted labels. Based on this, the framework can also be used to check on bias in databases by exploring correlations of target labels, like, e.g. depression or diagnosis, with predicted, or additionally given, labels like age, gender, signal distortion ratio or mean opinion score.

Design Choices

The program is intended for novice people interested in speaker characteristics detection (e.g., emotion, age, and gender) without proficient in (Python) programming language. Its main target is for education and research with the main features as follows:

- Finding good combinations of variables, e.g., acoustic features, models (classifier or regressor), feature standardization, augmentation, etc., for speaker characteristics detection (e.g., emotion);
- Characteristics of the database, such as distribution of gender, age, emotion, duration, data size, and so on with their visualization;
- Inference of speaker characteristics from a given audio file or streaming audio (can be said also as “weak” labeling for semi-supervised learning).

Hence, one should be able to use Nkululeko after installing and preparing/downloading their data in the correct format.

```
$ nkululeko.MODULE_NAME --config CONFIG_FILE.ini
```

How does it work?

nkululeko is a command line tool written in Python, best used in conjunction with the Visual Studio code editor (but can be run stand-alone). To use it, a text editor is needed to edit the experiment configuration. You would then run nkululeko like this:

```
$ nkululeko.explore --config conf.ini
```

and inspect the results afterward; they are represented as images, texts, and even a fully automatically compiled PDF report written in latex.

nkululeko's data import format is based on a simple CSV formalism, or alternatively, for a more detailed representation including data schemata, audformat.¹ Basically, to be used by nkululeko, the data format should include the audio file path and a task-specific label. Optionally, speaker ID and gender labels help with speech data. An example of a database labelled with emotion is

```
file, speaker, gender, emotion
x/sample.wav, s1, female, happy
...
```

As the main goal of nkululeko is to avoid the need to learn programming, experiments are specified by means of a configuration file. The functionality is encapsulated by software *modules* (interfaces) that are to be called on the command line. We list the most important ones here:

- **nkululeko**: do machine learning experiments combining features and learners
- **demo**: demo the current best model on the command line
- **explore**: perform data exploration (used mainly in this paper)
- **augment**: augment the current training data. This could also be used to reduce bias in the data, for example, by adding noise to audio samples that belong to a specific category.
- **** aug_train**: augment the training data and train the model with the augmented data.
- **predict**: predict features like signal distortion ratio, mean opinion score, arousal/valence, age/gender (for databases that miss this information), with deep neural nets models, e.g. as a basis for the *explore* module.
- **segment**: segment a database based on VAD (voice activity detection)
- **ensemble**: ensemble several models to improve performance

The configuration (INI) file consists of a set of key-value pairs that are organised into several sections. Almost all keys have default values, so they do not have to be specified.

Here is a sample listing of an INI file with database section:

```
[EXP]
name = explore-androids
[DATA]
databases = ['androids']
androids = /data/androids/androids.csv
target = depression
labels = ['depressed', 'control']
samples_per_speaker = 20
min_length = 2
[PREDICT]
sample_selection = all
targets = ['pesq', 'sdr', 'stoi', 'mos']
[EXPL]
value_counts = [['gender'], ['age'], ['est_sdr'], ['est_pesq'], ['est_mos']]
[REPORT]
latex = androids-report
```

As can be seen, some of the values simply contain Python data structures like arrays or dictionaries. Within this example, an experiment is specified with the name *explore-androids*, and a result folder with this name will be created, containing all figures and textual results, including an automatically generated Latex and PDF report on the findings.

¹<https://audeering.github.io/audformat/>

64 The *DATA* section sets the location of the database and specifies filters on the sample, in this
 65 case limiting the data to 20 samples per speaker at most and at least 2 seconds long. In this
 66 section, the split sets (training, development, and test) are also specified. There is a special
 67 feature named *balance splits* that lets the user specify criteria that should be used to stratify
 68 the splits, for example, based on signal distortion ratio.

69 With the *predict* module, specific features like, for example, signal distortion ratio or mean
 70 opinion score are to be predicted by deep learning models. The results are then used by a
 71 following call to the *explore* module to check whether these features, as well as some ground
 72 truth features (*age* and *gender*), correlate with the target variable (*depressed* in the given
 73 example) in any way.

74 The *nkululeko* configuration can specify further sections:

- 75 ■ **FEATS** to specify acoustic features (e.g. *opensmile* (Eyben et al., 2010) or deep learning
 76 embeddings; e.g. *wav2vec 2.0* (Baevski et al., 2020)) that should be used to represent
 77 the audio files.
- 78 ■ **MODEL** to specify statistical models for regression or classification of audio data.

79 Statement of need

80 Open-source tools are believed to be one of the reasons for accelerated science and technology.
 81 They are more secure, easy to customise and transparent. There are several open-source
 82 tools that exist for acoustic, sound, and audio analysis, such as *librosa* (McFee et al., 2015),
 83 *TorchAudio* (Yang et al., 2021), *pyAudioAnalysis* (Giannakopoulos, 2015), *ESPNET* (Watanabe
 84 et al., 2018), and *SpeechBrain* (Ravanelli et al., 2021). However, none of them are specialised
 85 in speech analysis with high-level interfaces for novices in the speech processing area.

86 One exception is *Spotlight* (Suwelack, 2023), an open-source tool that visualises metadata
 87 distributions in audio data. An existing interface between *nkululeko* and *Spotlight* can be
 88 used to combine the visualisations of *Spotlight* with the functionalities of *Nkululeko*.

89 *Nkululeko* follows these principles:

- 90 ■ **Minimum programming skills:** the only programming skills required are to prepare the
 91 data in the correct (CSV) format and to run the command line tool. For *AUDFORMAT*,
 92 no preparation is needed.
- 93 ■ **Standardised data format and label:** the data format is based on CSV and *AUFORMAT*,
 94 which is a widely used format for data exchange. The standard headers are like 'file',
 95 'speaker', 'emotion', 'age', and 'language' but also can be customised.
- 96 ■ **Replicability:** the experiments are specified in a configuration file, which can be shared
 97 with others including the splitting of training, development, and test partition. All results
 98 are stored in a folder with the same name as the experiment.
- 99 ■ **High-level interface:** the user specifies the experiment in an INI file, which is a simple
 100 text file that can be edited with any text editor. The user does not need to write Python
 101 code for experiments.
- 102 ■ **Transparency:** as CLI, *nkululeko* *always output debug*, in which info, warning, and error
 103 will be displayed in terminal (and should be easily understood). The results are stored in
 104 the experiment folder for further investigations and are represented as images, texts, and
 105 even a fully automatically compiled PDF report written in latex.

Usage in Existing Research

Nkululeko has been used in several research projects since its first appearance in 2022 (Burkhardt et al., 2022). The following list gives an overview of the research papers that have used Nkululeko:

- (?): this paper reported a database development of synthesized speech for basic emotions and its evaluation using Nkululeko toolkit.
-
-

Acknowledgements

We acknowledge support from these various projects:

- European SHIFT (*MetamorphoSiS of cultural Heritage Into augmented hypermedia assets For enhanced accessibiliTy and inclusion*) project (Grant Agreement number: 101060660);
- European EASIER (*Intelligent Automatic Sign Language Translation*) project (Grant Agreement number: 101016982);
- Project JPNP20006 commissioned by the New Energy and Industrial Technology Development Organization (NEDO), Japan;
- Project 24K02967 from the Japan Society for the Promotion of Science (JSPS).

We thank audeERING GmbH for partial funding.

References

- Baevski, A., Zhou, Y., Mohamed, A., & Auli, M. (2020). wav2vec 2.0: A framework for self-supervised learning of speech representations. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, & H. Lin (Eds.), *Advances in neural information processing systems* (Vol. 33, pp. 12449–12460). Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2020/file/92d1e1eb1cd6f9fba3227870bb6d7f07-Paper.pdf>
- Burkhardt, F., Wagner, J., Wierstorf, H., Eyben, F., & Schuller, B. (2022). Nkululeko: A tool for rapid speaker characteristics detection. *2022 Language Resources and Evaluation Conference, LREC 2022*, 1925–1932. ISBN: 9791095546726
- Chaudhary, A., Chouhan, K. S., Gajrani, J., & Sharma, B. (2020). *Deep learning with PyTorch*. <https://doi.org/10.4018/978-1-7998-3095-5.ch003>
- Eyben, F., Wöllmer, M., & Schuller, B. (2010). openSMILE – the munich versatile and fast open-source audio feature extractor. *MM'10 - Proceedings of the ACM Multimedia 2010 International Conference*, 1459–1462. <https://doi.org/10.1145/1873951.1874246>
- Giannakopoulos, T. (2015). pyAudioAnalysis: An open-source python library for audio signal analysis. *PLoS One*, 10(12), 1–17. <https://doi.org/10.1371/journal.pone.0144610>
- McFee, B., Raffel, C., Liang, D., Ellis, D., McVicar, M., Battenberg, E., & Nieto, O. (2015). librosa: Audio and Music Signal Analysis in Python. *Proc. 14th Python Sci. Conf., Scipy*, 18–24. <https://doi.org/10.25080/majora-7b98e3ed-003>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D.,

- 146 Brucher, M., Perrot, M., & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python.
147 *Journal of Machine Learning Research*, 12, 2825–2830.
- 148 Ravanelli, M., Parcollet, T., Plantinga, P., Rouhe, A., Cornell, S., Lugosch, L., Subakan, C.,
149 Dawalatabad, N., Heba, A., Zhong, J., Chou, J.-C., Yeh, S.-L., Fu, S.-W., Liao, C.-F.,
150 Rastorgueva, E., Grondin, F., Aris, W., Na, H., Gao, Y., ... Bengio, Y. (2021). *SpeechBrain:*
151 *A general-purpose speech toolkit*. <https://arxiv.org/abs/2106.04624>
- 152 Suwelack, S. (2023). Spotlight. In *GitHub repository*. [https://github.com/Renumics/](https://github.com/Renumics/spotlight/)
153 [spotlight/](https://github.com/Renumics/spotlight/); GitHub.
- 154 Watanabe, S., Hori, T., Karita, S., Hayashi, T., Nishitoba, J., Unno, Y., Soplin, N. E. Y.,
155 Heymann, J., Wiesner, M., Chen, N., Renduchintala, A., & Ochiai, T. (2018). ESPNet:
156 End-to-end speech processing toolkit. *Proc. Annu. Conf. Int. Speech Commun. As-*
157 *soc. INTERSPEECH, 2018-Sept*(September), 2207–2211. [https://doi.org/10.21437/](https://doi.org/10.21437/Interspeech.2018-1456)
158 [Interspeech.2018-1456](https://doi.org/10.21437/Interspeech.2018-1456)
- 159 Yang, S., Chi, P.-H., Chuang, Y.-S., Lai, C.-I. J., Lakhota, K., Lin, Y. Y., Liu, A. T., Shi,
160 J., Chang, X., Lin, G.-T., Huang, T.-H., Tseng, W.-C., Lee, K., Liu, D.-R., Huang,
161 Z., Dong, S., Li, S.-W., Watanabe, S., Mohamed, A., & Lee, H. (2021). SUPERB:
162 Speech Processing Universal PERformance Benchmark. *Interspeech 2021*, 1194–1198.
163 <https://doi.org/10.21437/Interspeech.2021-1775>

DRAFT