

Modeling of Polish Intonation for Statistical-Parametric Speech Synthesis

Tomasz Kuczmarski

Prof. zw. Dr hab. Inż. Grażyna Demenko

Supervisor

Adam Mickiewicz University



Faculty of Modern Languages and Literature
Institute of Ethnolinguistics

May 18, 2022

Intonation

Definition

"The tonal structure of speech expressed by the melody produced by our larynx. It has a phonetic aspect, the fundamental frequency (F_0), and a grammatical (phonological) aspect." (Féry 2016)

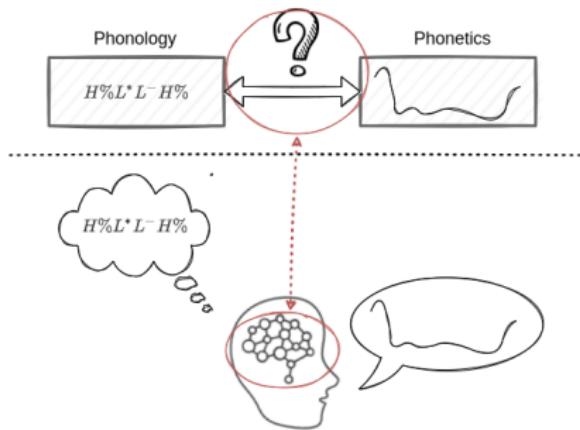
Intonation

Definition

"The tonal structure of speech expressed by the melody produced by our larynx. It has a phonetic aspect, the fundamental frequency (F_0), and a grammatical (phonological) aspect." (Féry 2016)

All definitions of intonation "are epistemological definitions, i.e., not a priori programmatic definitions, but a posteriori statements of a practice and methodology." (Rossi 2000)

Motivation



Motivation

- ✓ Dualistic gap between phonology and phonetics.
- ✓ Unification within a broader metatheory.
- ✓ Unknown nature of the mappings between mental categories and continuous contours of F_0 .
- ✓ How linguistic features of an utterance influence its F_0 contours.
- ✓ Need for a physicalist (neurobiological) model.
- ✓ Modern statistical-parametric speech synthesis provides a framework for experimentation and evaluation of such
- ✓ Remarkable properties of biologically-motivated Convolutional Neural Networks.

Objectives

- ① Build a robust biologically-inspired neural model of the probabilistic mapping between discrete low-level linguistic features of an utterance and its intonation contours (F_0 values).
- ② Build a state-of-the-art neural source-filter resynthesis framework for Polish read speech.
- ③ Deploy the intonation model within the resynthesis framework and measure the perceived naturalness of the output intonation contours, and to
- ④ operationalize the results of these measurements as an indicator of the model's robustness.
- ⑤ Develop a method to explain the relevance of the individual input linguistic features for the produced intonation contour.
- ⑥ Analyze how specific linguistic features contribute to the F_0 contours of an utterance.

Objectives

- ① Build a robust biologically-inspired neural model of the probabilistic mapping between discrete low-level linguistic features of an utterance and its intonation contours (F_0 values).
- ② Build a state-of-the-art neural source-filter resynthesis framework for Polish read speech.
- ③ Deploy the intonation model within the resynthesis framework and measure the perceived naturalness of the output intonation contours, and to
- ④ operationalize the results of these measurements as an indicator of the model's robustness.
- ⑤ Develop a method to explain the relevance of the individual input linguistic features for the produced intonation contour.
- ⑥ Analyze how specific linguistic features contribute to the F_0 contours of an utterance.

Hypotheses

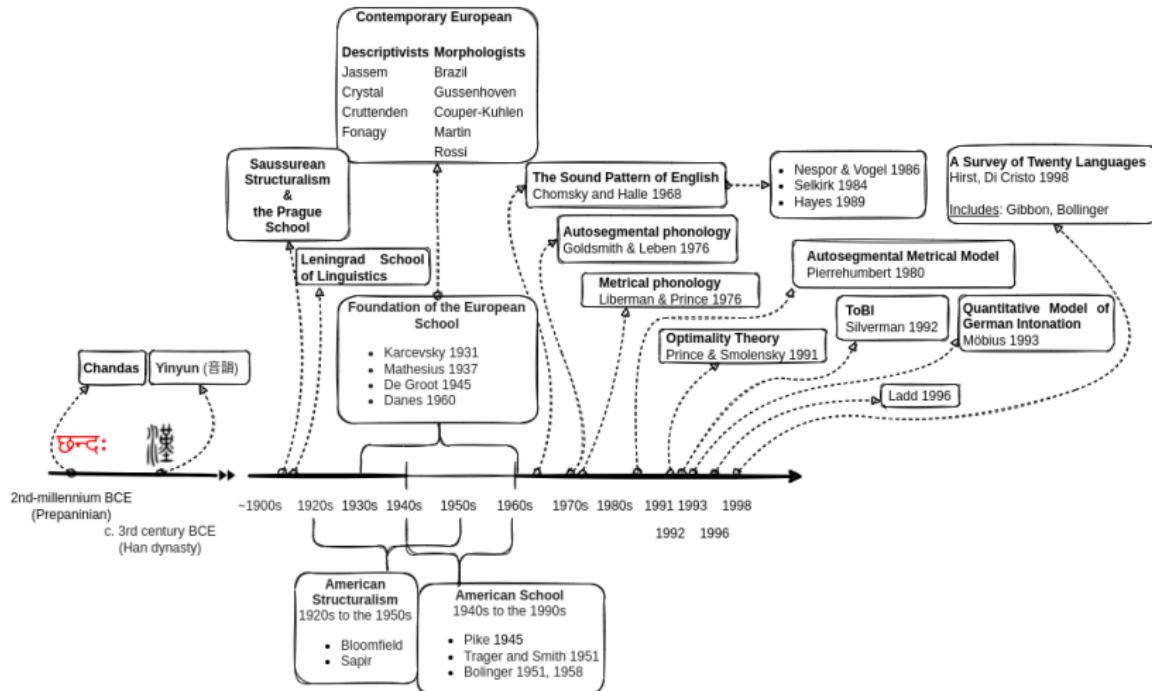
HYPOTHESIS 1: *The continuous F_0 contours of an utterance emerge from its discrete linguistic features through a series of successive probabilistic mappings into intermediate latent representations.*

HYPOTHESIS 2: *The biologically-inspired Deep Temporal Convolutional Network can be an effective model of these mappings and hence of Polish neutral read speech intonation in the context of statistical-parametric speech synthesis.*

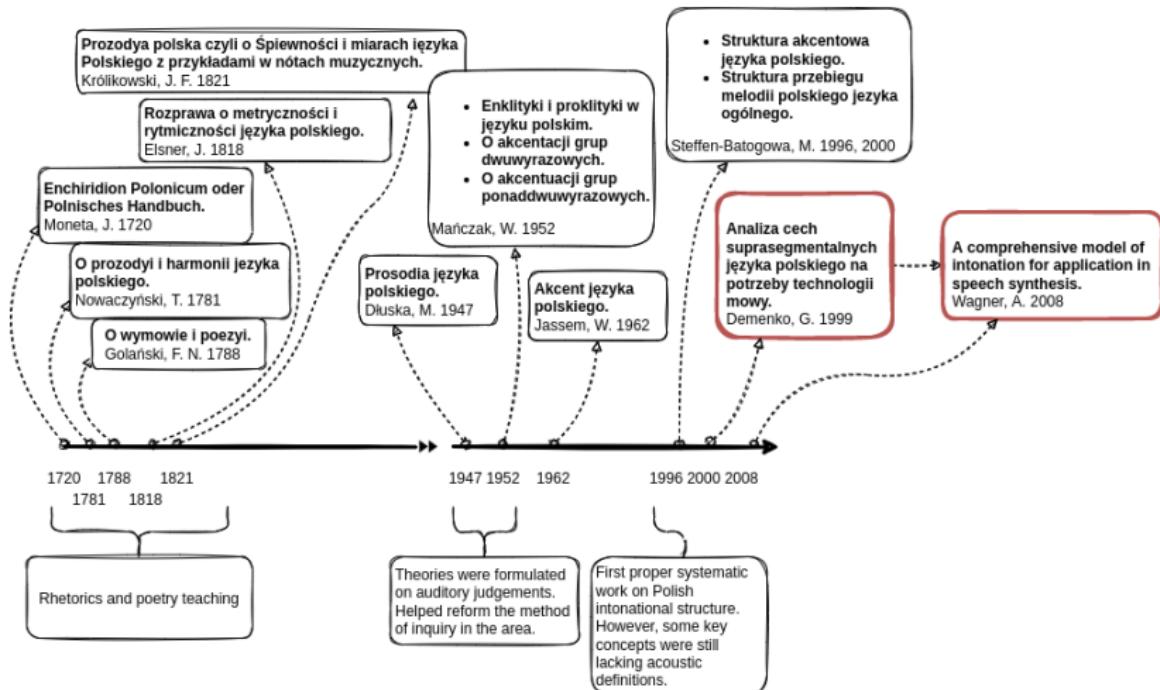
HYPOTHESIS 3: *The set of shallow linguistic features used in this thesis provides information which is sufficient for synthesis of natural sounding intonation in the context of statistical-parametric speech synthesis.*

HYPOTHESIS 4: *A Deep Temporal Convolutional Network can become an explanatory scientific model of mappings between linguistics features and the intonation of an utterance.*

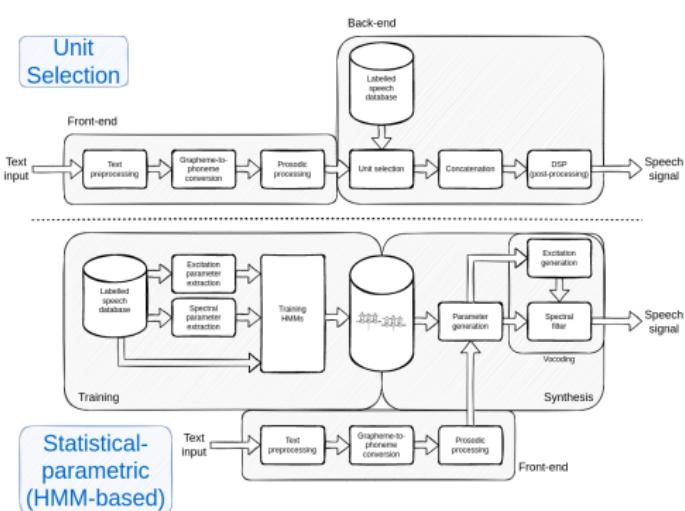
Background



Background



Speech synthesis



Speech synthesis - Wavenet

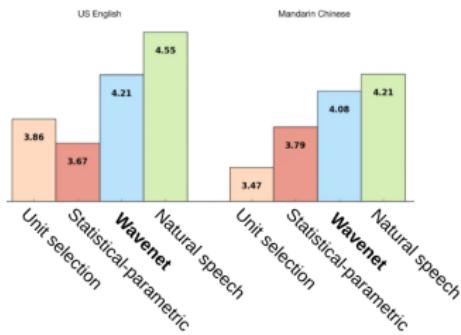
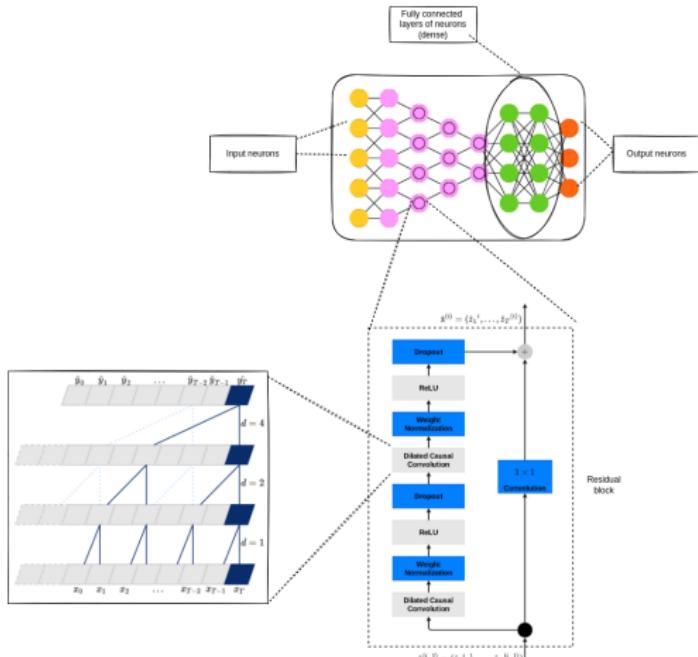


Figure: Google WaveNet evaluation results (Oord et al. 2016).

Wavenet belongs to a class of models known as **Convolutional Neural Networks (CNNs)** which are used mainly in the area of image recognition where they excel.



Visual processing in the human brain

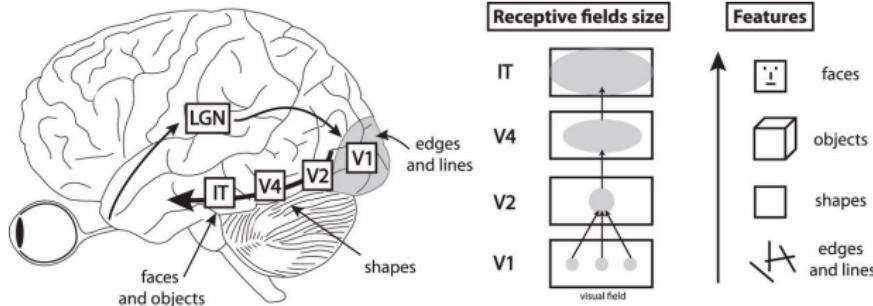


Figure: Hierarchical, feedforward visual processing in human brain. Adopted from (Manassi, Sayim, and Herzog 2013).

CNN feature maps

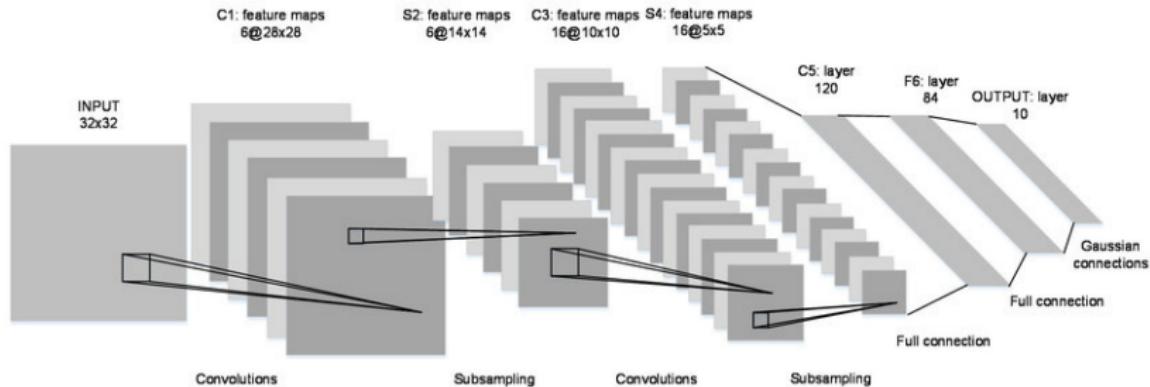
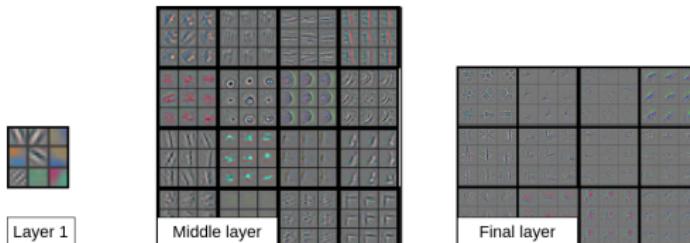


Figure: Image recognition Convolutional Neural Network (LeNet-5). Adopted from (LeCun et al. 1989).



Simple and complex cells in the visual cortex

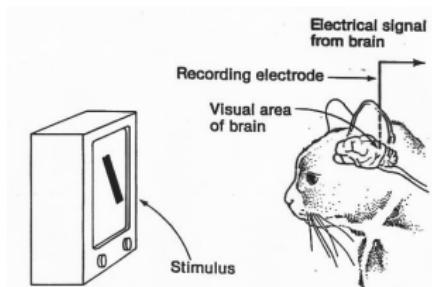


Figure: Famous Hubel and Wiesel cat experiment.
Adopted from (Hubel and Wiesel 1959).

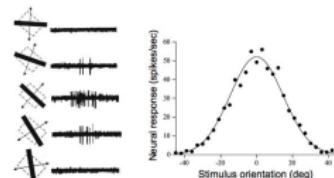


Figure: Neural response of simple cells. Adopted from (Hubel and Wiesel 1962).

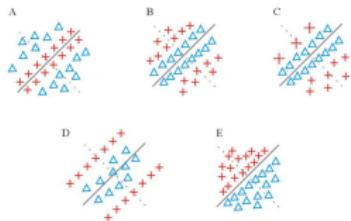


Figure: Simple receptive fields. Adopted from (Hubel and Wiesel 1962).

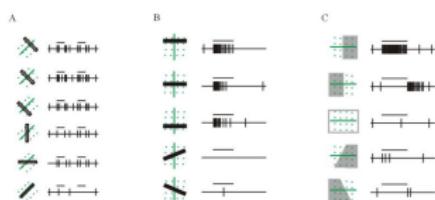
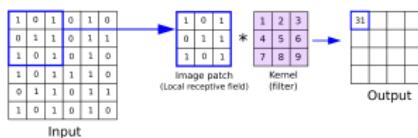


Figure: Three different types of complex receptive fields.
Adopted from (Hubel and Wiesel 1962).

Neurobiological foundations of CNNs



Original	Gaussian Blur	Sharpen	Edge Detection
$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$

Figure: Example of a 2-dimensional matrix convolution.

Figure: Examples of convoluting and image with different convolution kernels. Adopted from Wikipedia: (Kernel (image processing) 2021).

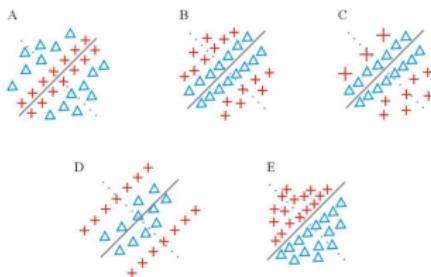


Figure: Simple receptive fields. Adopted from (Hubel and Wiesel 1962).

Simple cells perform edge and line detection which can be very effectively approximated with matrix **convolution**.

Neurobiological foundations of CNNs

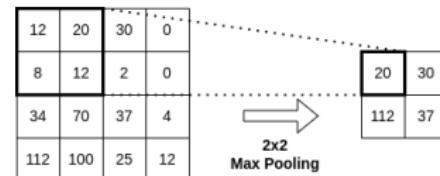


Figure: Example of a 2×2 max pooling matrix operation.

The computations that the **complex cells** perform are similar to the **max pooling** operation.

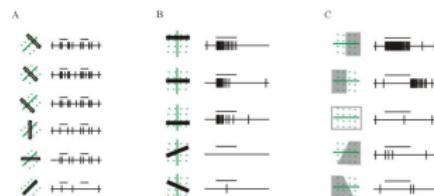


Figure: Three different types of complex receptive fields.
Adopted from (Hubel and Wiesel 1962).

Simple and complex cells in the auditory cortex

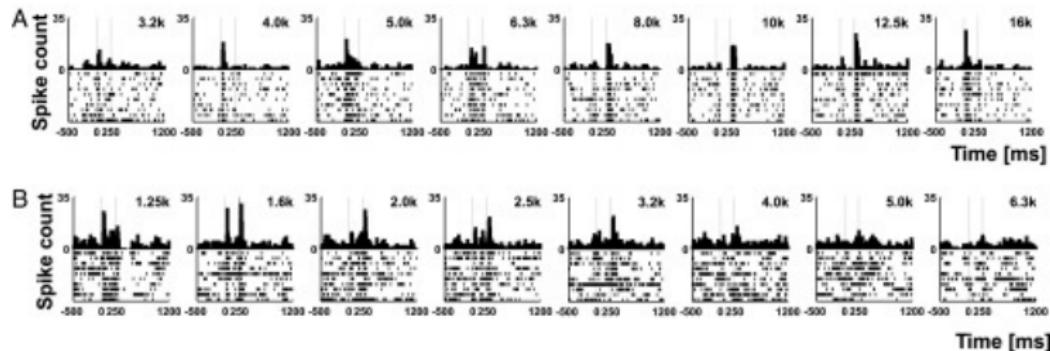


Figure: Responses of single neurons in primary auditory cortex (A1) of rhesus monkeys to band-passed noise (BPN) bursts centered at particular frequencies. Adopted (Tian, Kuśmirek, and Rauschecker 2013).

Auditory cortex also contains **simple (and complex) cells** with two-dimensional receptive fields that are similar to those in the visual cortex but which operate in the **spectrotemporal domain** instead.

Explanatory properties of CNNs

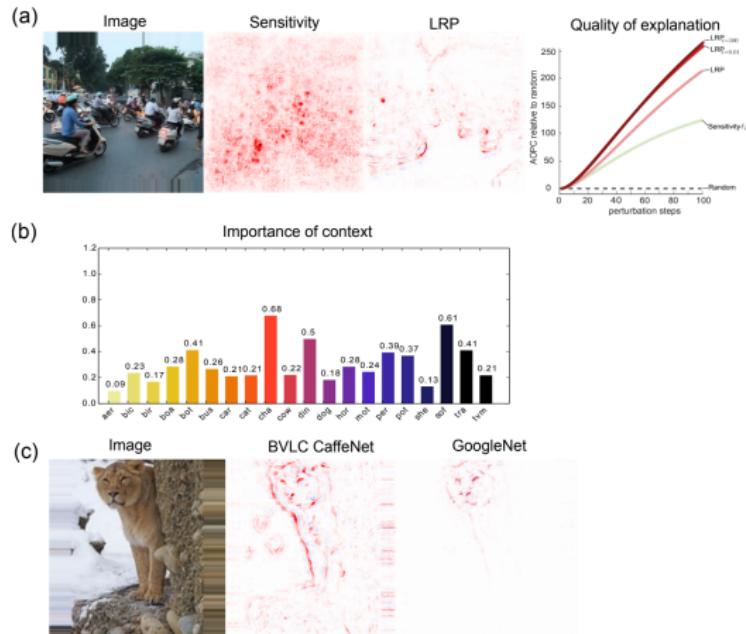


Figure: Results of sensitivity-based and relevance-based explainability methods. Based on (Samek et al. 2016)

Dataset

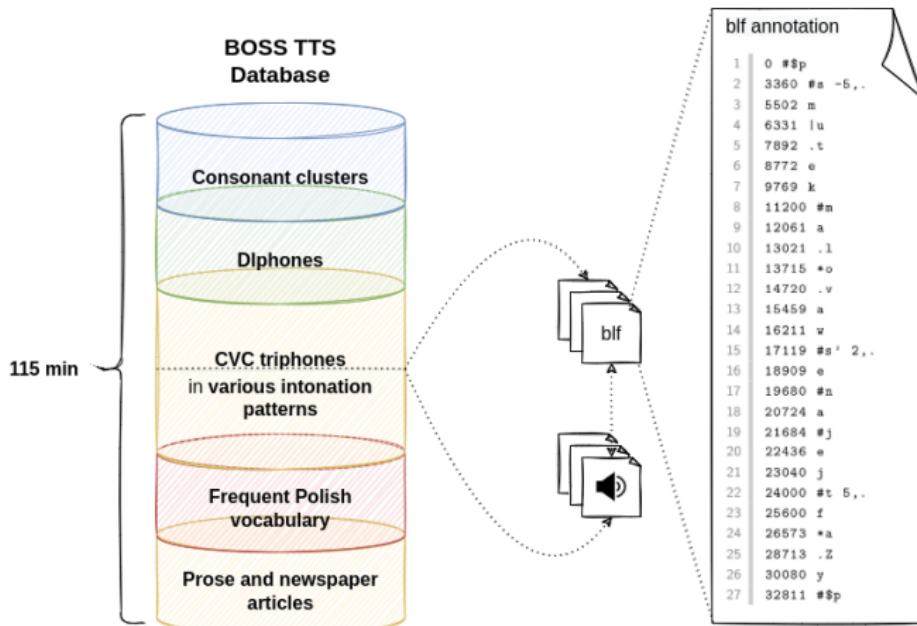
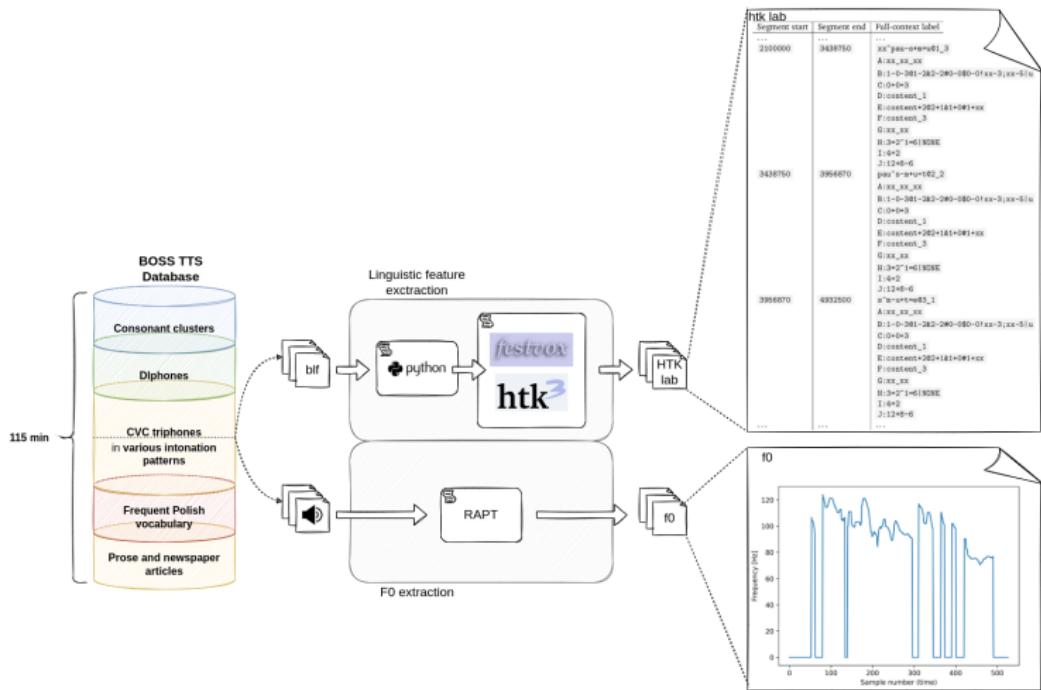
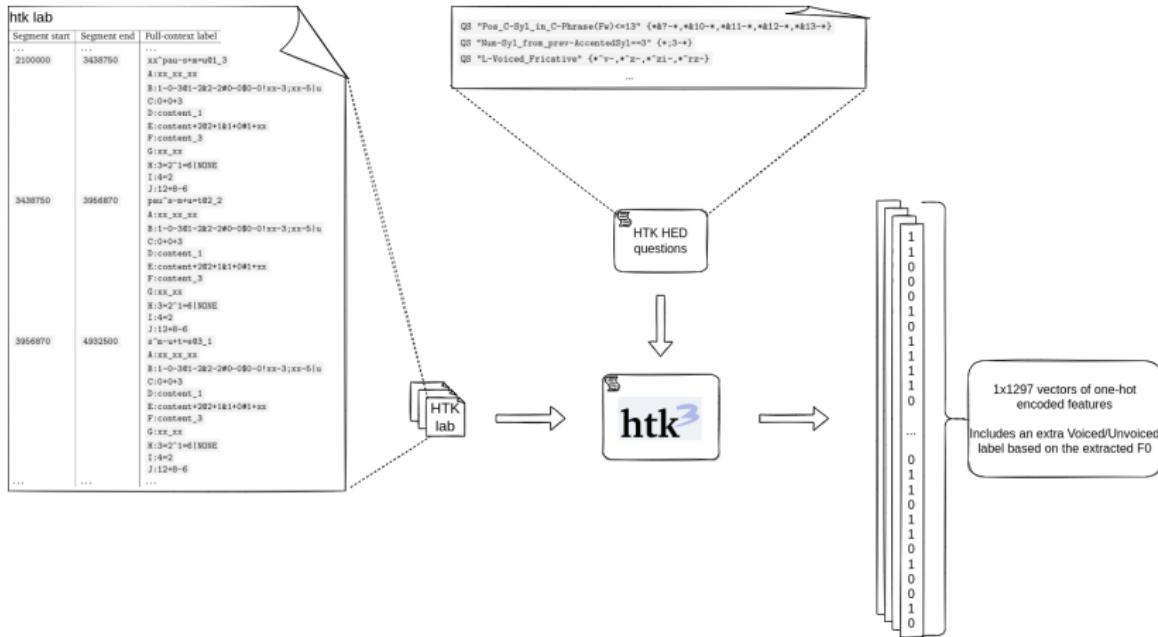


Figure: Speech corpus built originally for the purpose of the Polish BOSS unit selection synthesizer (Grażyna Demenko and Wagner 2007; Grażyna Demenko, Bachan, et al. 2008; G. Demenko, B. Möbius, and K. Klessa 2010; Grażyna Demenko, Katarzyna Klessa, et al. 2010).

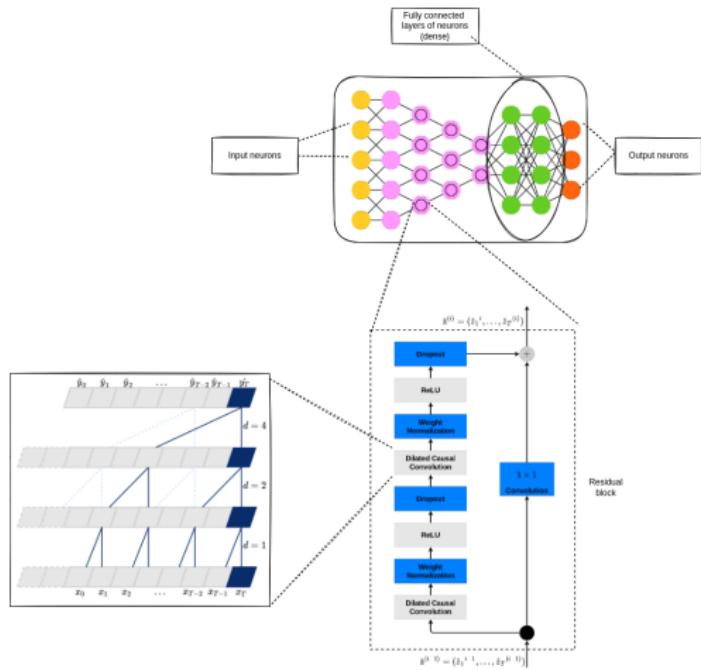
Data preprocessing



Feature extraction



Model implementation



TCN parameters

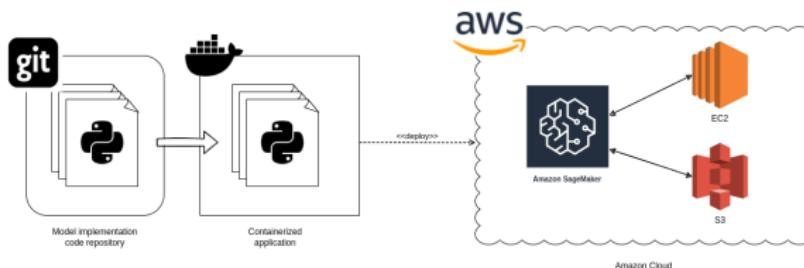
- 6 residual layers,
 - 64 convolution filters of length 2
 - 2^n dilations ($[1, 2, 4, 8, 16, 32]$)
 - 2 final fully connected dense layers
(ReLU activation; sizes 64 and 1)

Complete code repository

✓ https://github.com/mrsslacklines/intonation_synthesis

github.com/philipperemy/keras-tcn
(Remy 2020)

Model training infrastructure



- 64 CPUs; 488GB
- 8 NVIDIA Tesla V100 GPUs; 128GB

Memory requirements

A single batch of data, which is a $64 \times 1900 \times 1297$ vector of 8-byte boolean values, occupies 1.2617216 gigabytes of memory, and the model comprises of a total of 449,409 parameters (446,337 trainable and 3,072 non-trainable).



Training parameters:

- ADAM optimizer
- initial learning rate of 0.1,
- $\beta_1 = 0.9$,
- $\beta_2 = 0.999$,
- $\epsilon = 1e - 07$.
- **Loss metric:** Mean Squared Error (MSE)
- 200 epochs
- random 8:1:1 dataset split

F_0 inference and feature relevance analysis

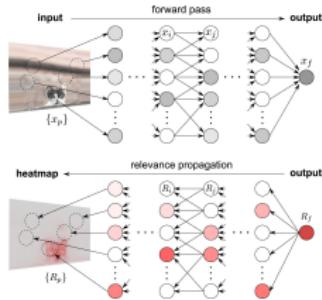
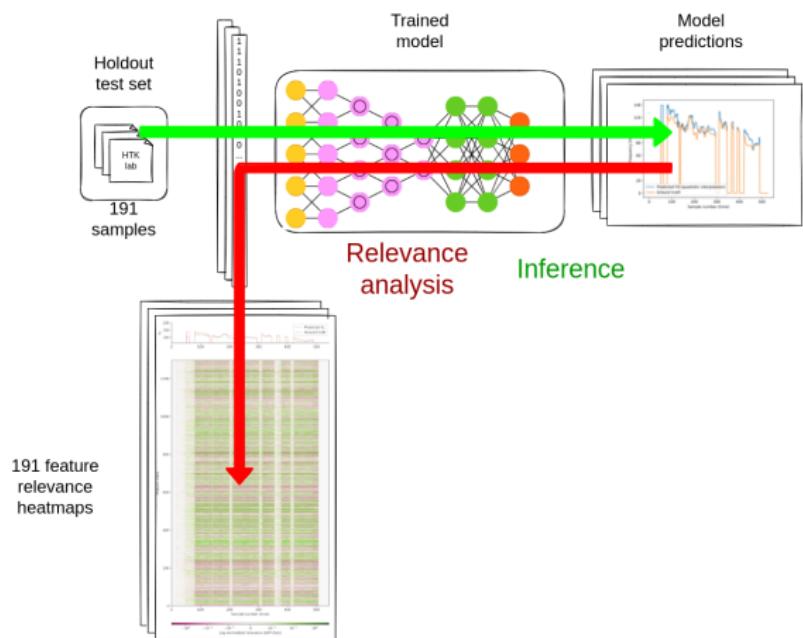


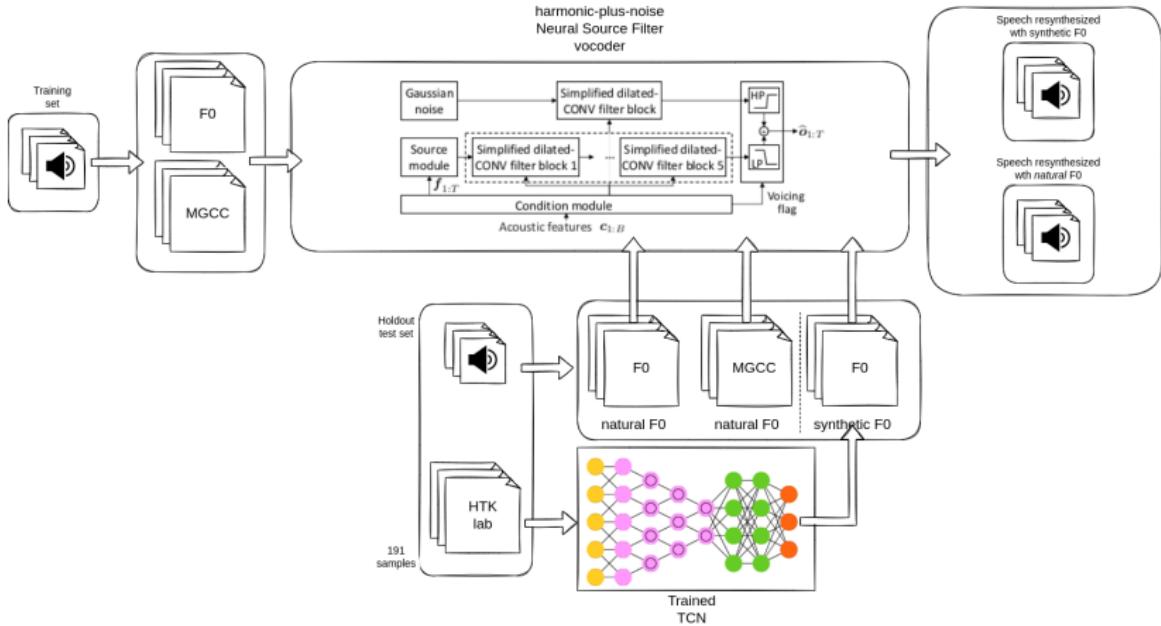
Figure: Computational flow of deep Taylor decomposition.
(Adopted from Montavon 2017).



INNvestigate Neural Networks!
(Alber et al. 2019)

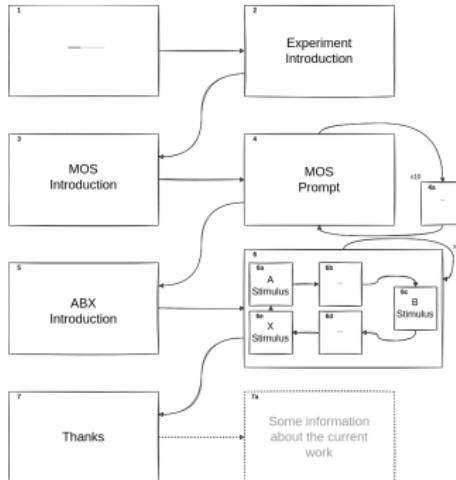
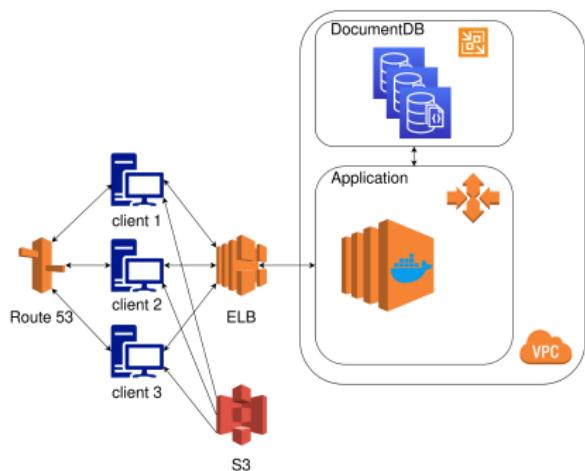
[github.com/albermax/
innvestigate](https://github.com/albermax/innvestigate)

Resynthesis



Neural Source Filter Vocoder from (Wang, Takaki, and Yamagishi 2019).

Perceptual evaluation experiment



Available at:
fonetyka.cudaniewidły.org/experiment
Code at:
github.com/mrslacklines/listening_experiments

F_0 inference results

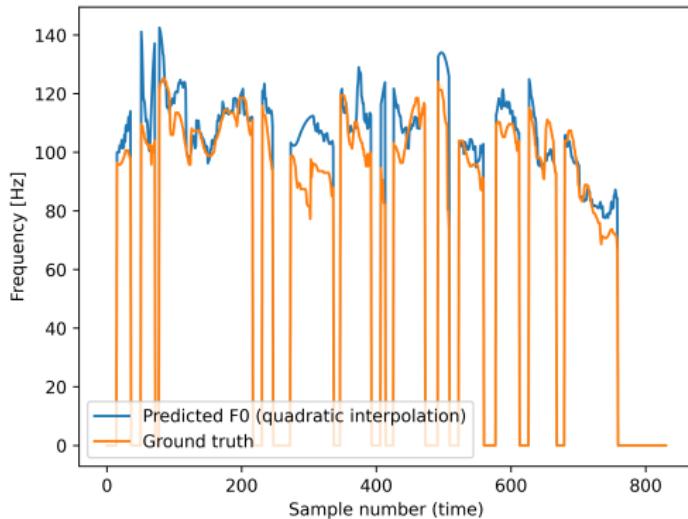


Figure: Result of F_0 prediction for "Lokatorzy znaleźli się w podbramkowej sytuacji i musieli się wyprowadzić" (*The tenants found themselves in a difficult situation and had to move out*).

F_0 inference results

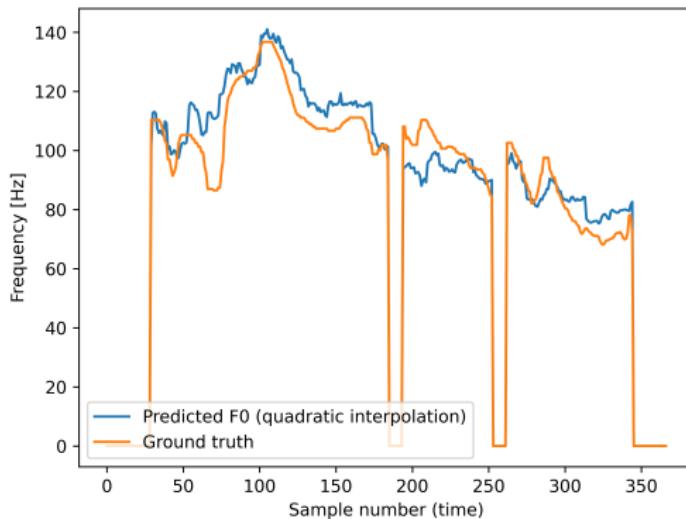


Figure: Result of F_0 prediction for "Powodzenie nie jest gwarantowane" (*Success is not guaranteed*).

F_0 inference results

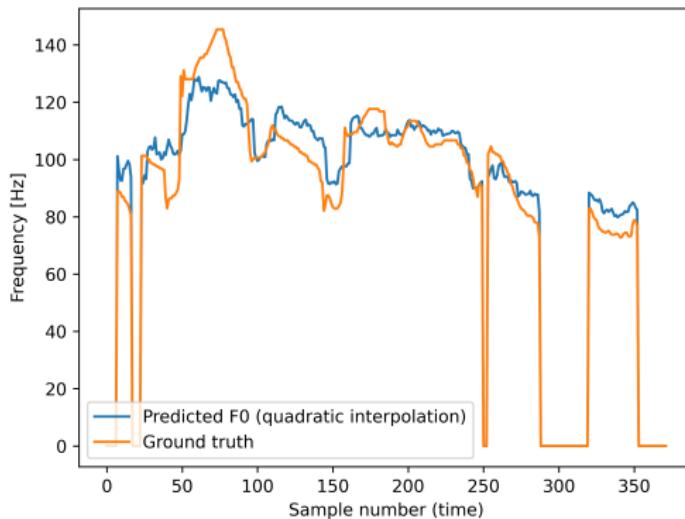


Figure: Result of F_0 prediction for "Gaduła była bardzo nieznośna" (*Gabby was very annoying.*).

F_0 inference results

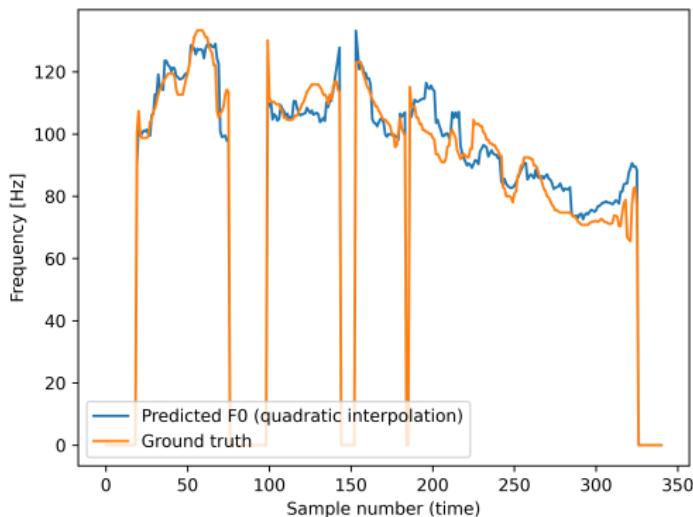


Figure: Result of F_0 prediction for "Może przyniosą też gorzałę" (*Maybe they will bring booze too.*).

F_0 inference results

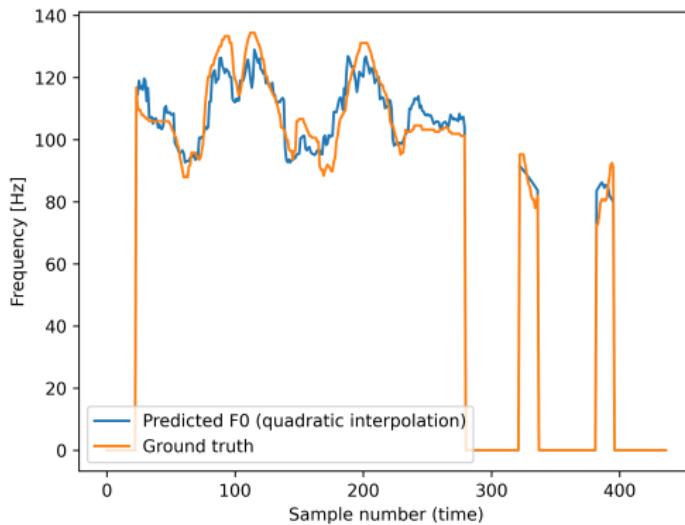


Figure: Result of F_0 prediction for "To jest ważna godzina dla nas wszystkich" (*This is an important hour for all of us*).

F_0 inference results

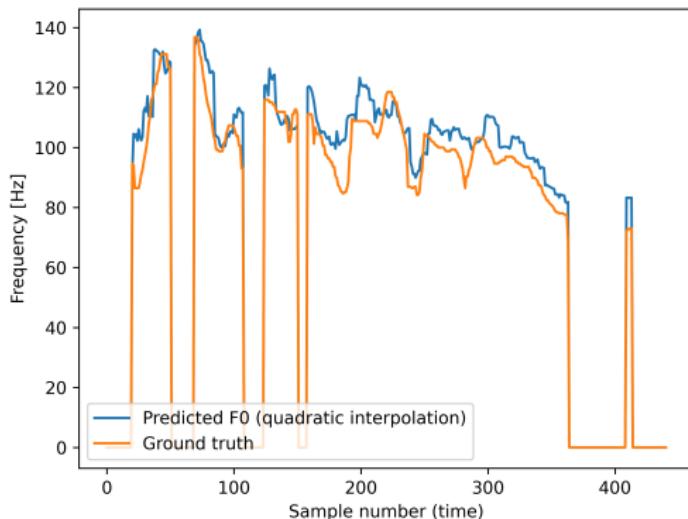


Figure: Result of F_0 prediction for "Myślę, że chleb razowy będzie najlepszy" (*I think that a wholemeal bread will be the best*).

F_0 inference results

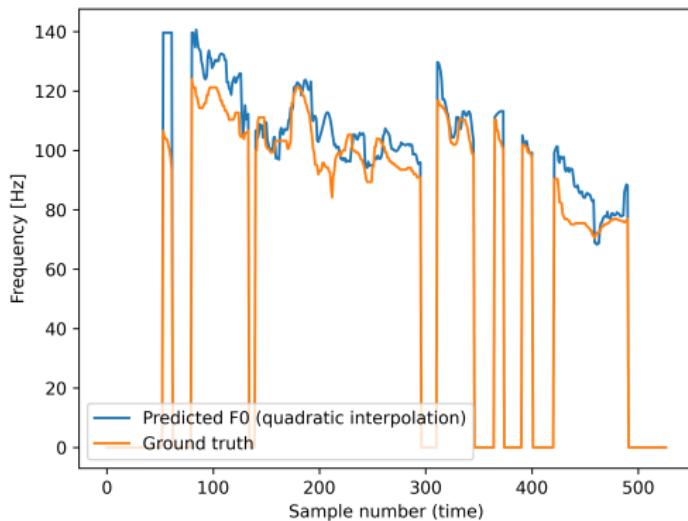


Figure: Result of F_0 prediction for "Słyszałam odgłos zbliżającego się pociągu" (*I heard the sound of an approaching train*).

F_0 inference results

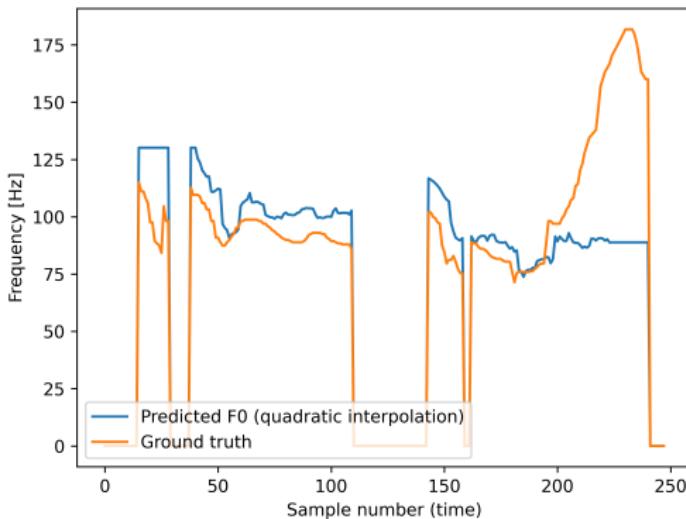


Figure: Result of F_0 prediction for "Czy to był łatwy dobór?" (*Was it an easy choice?*).

F_0 inference results

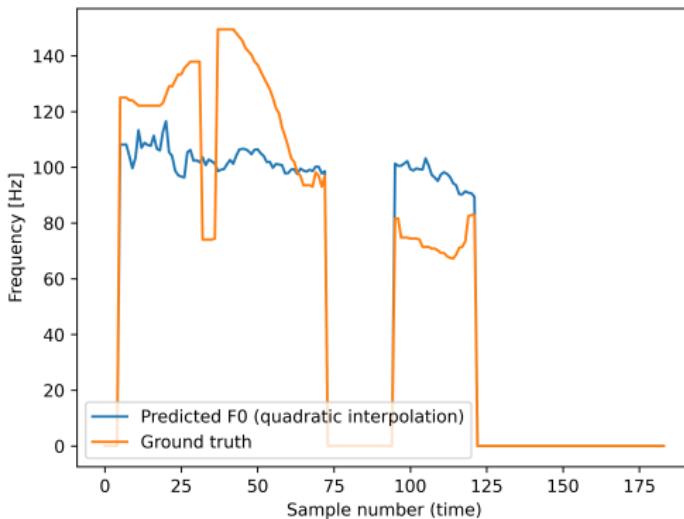


Figure: Result of F_0 prediction for "To Majka" (*This is Majka*).

F_0 inference results

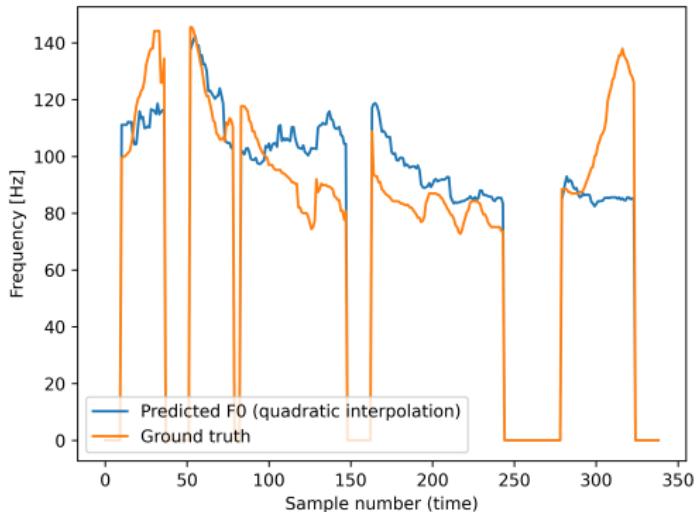


Figure: Result of F_0 prediction for "Na czym polega kandyzacja?" (*How does candying work?*).

Objective evaluation results

MAE	0.015382
MSE	0.002495
NRMSE	0.009534
RMSE	0.048128

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \tilde{y}_i)^2$$

Subjective evaluation results - MOS

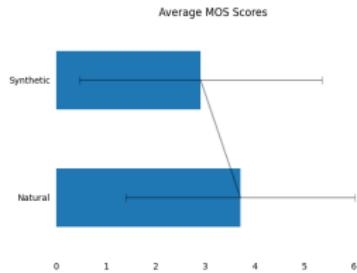


Figure: Mean Opinion Score-based evaluation results.

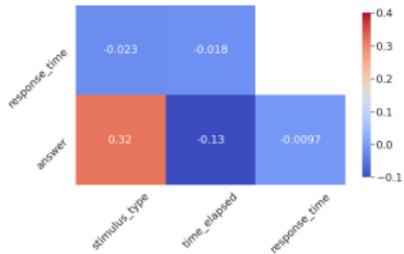


Figure: Mean Opinion Score-based evaluation parameters correlation matrix.

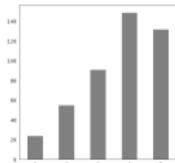


Figure: Mean Opinion Score-based evaluation total numbers of specific scores for natural stimuli.

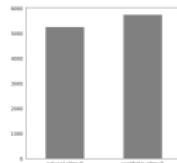


Figure: Mean Opinion Score-based evaluation mean response times for synthetic and natural stimuli.

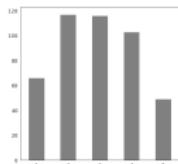


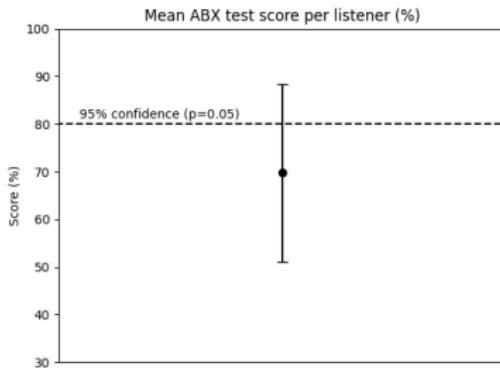
Figure: Mean Opinion Score-based evaluation total numbers of specific scores for synthetic stimuli.

Subjective evaluation results – ABX

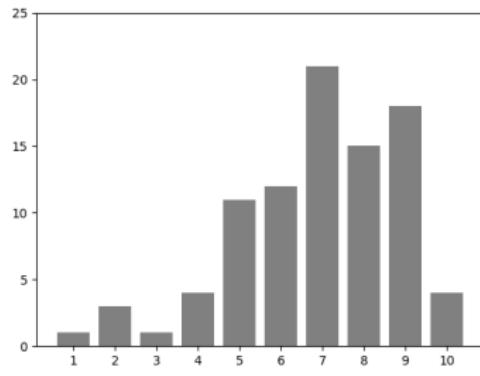
χ^2 test

NULL HYPOTHESIS (H_0): *There are no perceptually significant differences between resynthesized recordings with synthetic and natural F_0 .*

HYPOTHESIS (H_a): *There are perceptually significant differences between resynthesized recordings with synthetic and natural F_0 .*



Number of individuals with n-correct scores out of 10



Subjective evaluation results – ABX

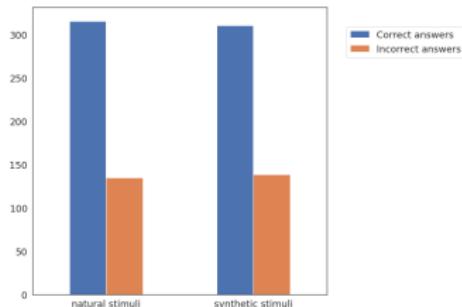


Figure: ABX experiment number of correct and incorrect answers in case of natural and synthetic stimuli.

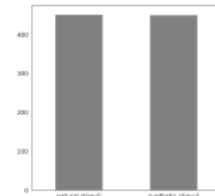


Figure: ABX experiment mean response times for synthetic and natural stimuli.

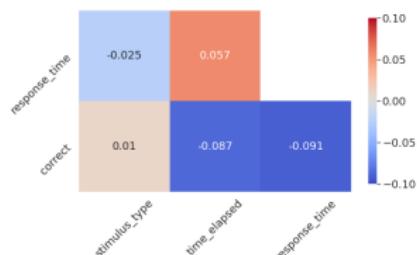


Figure: ABX experiment parameters correlation matrix.

Feature relevance analysis results

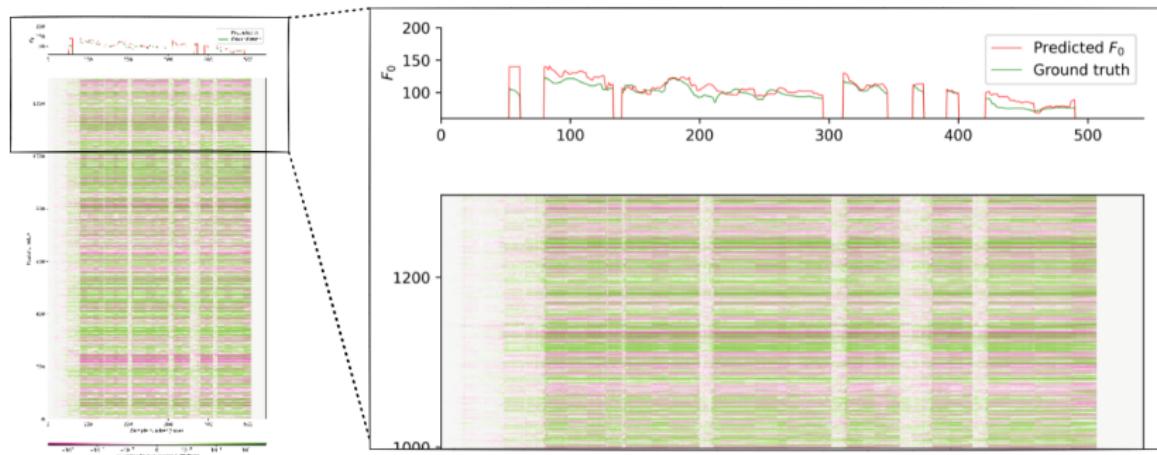
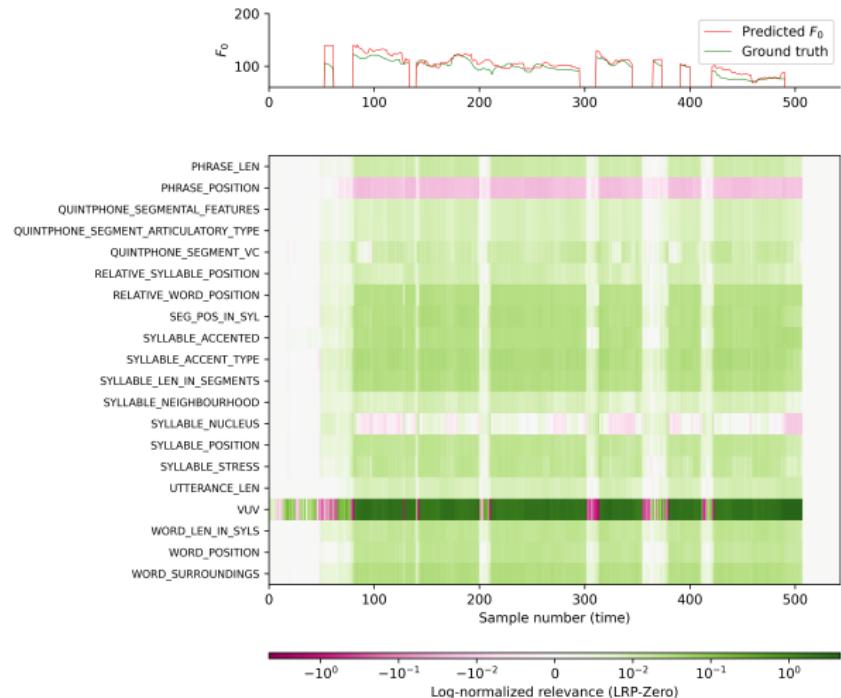


Figure: Fundamental frequency predictions for "Słyszałam odgłos zbliżającego się pociągu" (*I heard the sound of an approaching train*) aligned with feature relevance heatmap.

Feature relevance analysis results

Figure: Fundamental frequency predictions for "Słyszałam odgłos zbliżającego się pociągu" (*I heard the sound of an approaching train*) aligned with relevance heatmap for general feature groups.



Feature relevance analysis results

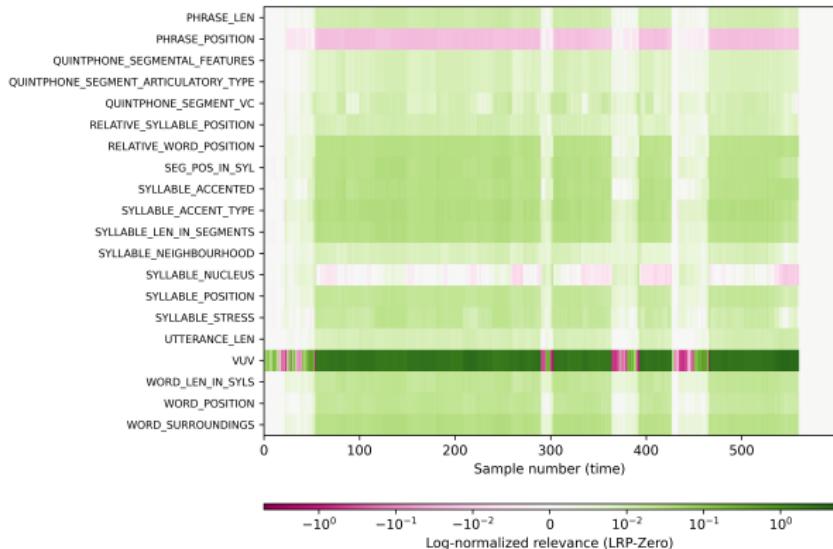
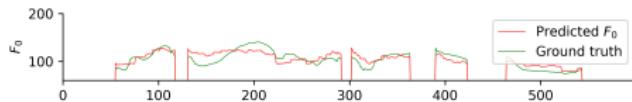
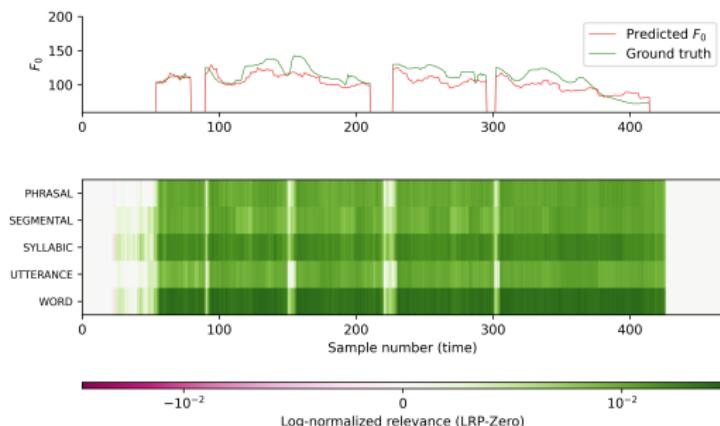


Figure: Fundamental frequency predictions for "Waluta jeny - w języku greckim - jest cenna" (*The currency yen, in Greek, is valuable*) aligned with relevance heatmap for general feature groups.



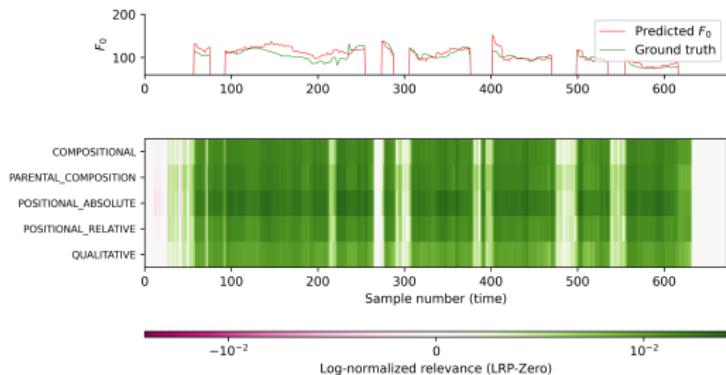
Feature relevance analysis results

Figure: Fundamental frequency predictions for "Musisz dojrzeć do tego, by to zrozumieć" (*You have to grow up to understand it*) aligned with relevance heatmap for features grouped by the level of utterance.



Feature relevance analysis results

Figure: Fundamental frequency predictions for "Przepłynęłam na grzbiecie siedemnaście długości basenu" (*I swam seventeen pool lengths on my back.*) aligned with relevance heatmap for features grouped by the type of relation.



Feature relevance analysis results

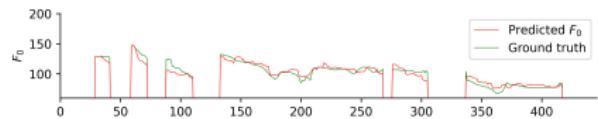
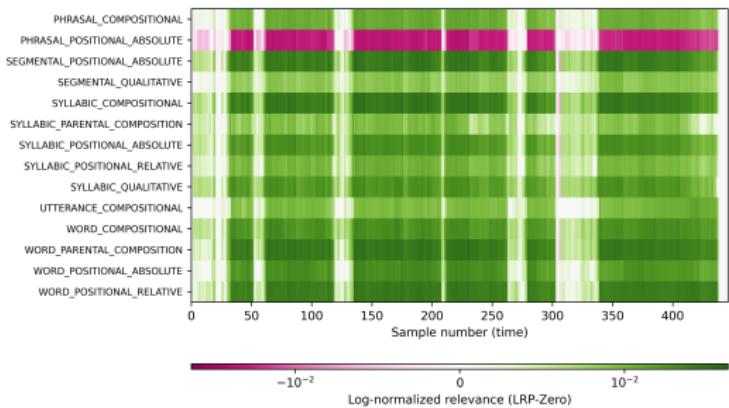
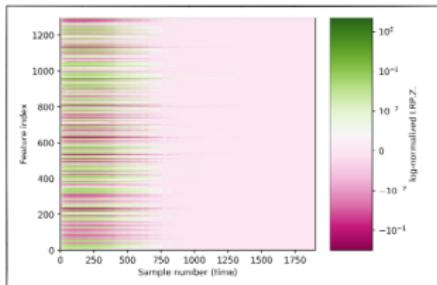


Figure: Fundamental frequency predictions for "Ciocia pracuje w urzędzie państwowym" (*Aunt works in a government office*) aligned with relevance heatmap for features grouped by the type of relation and the level of utterance.

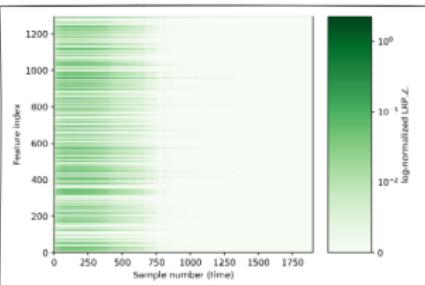
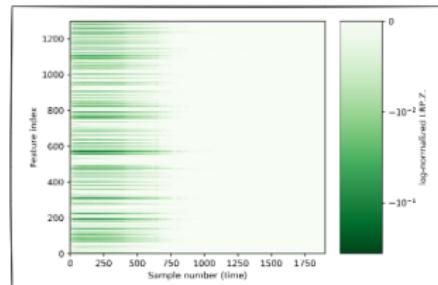
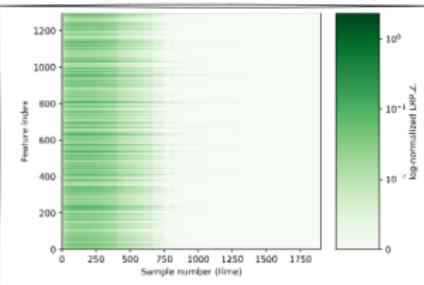


Feature relevance analysis results

Mean



Mean (from absolute sum)

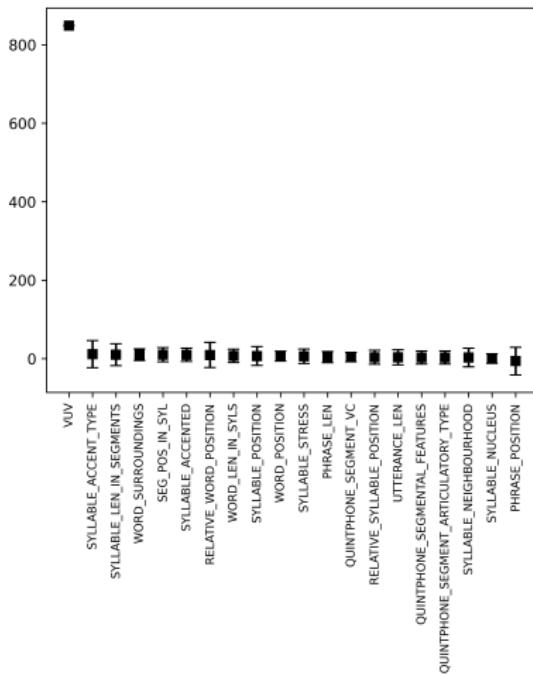


Mean (negative-only values)

Mean (positive-only values)

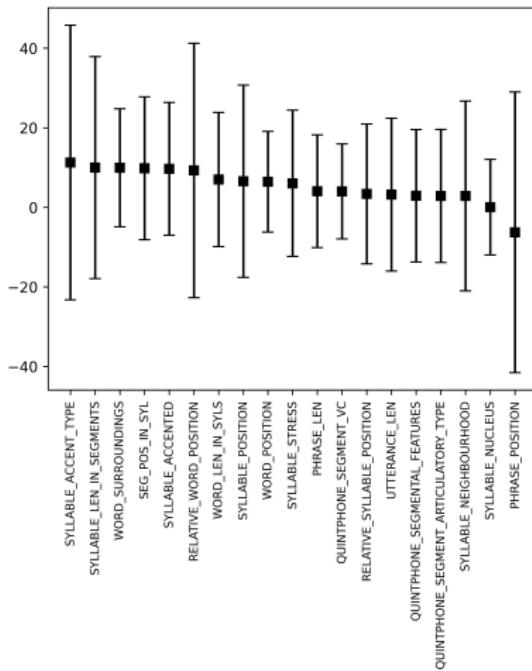
Feature relevance analysis results

Figure: Regular mean relevance per feature group. The y-axis scaling due to the high relevance of V/UV renders comparatively flat plots for other features.



Feature relevance analysis results

Figure: Regular mean relevance per feature group with the most relevant feature (V/UV) excluded.



Feature relevance analysis results

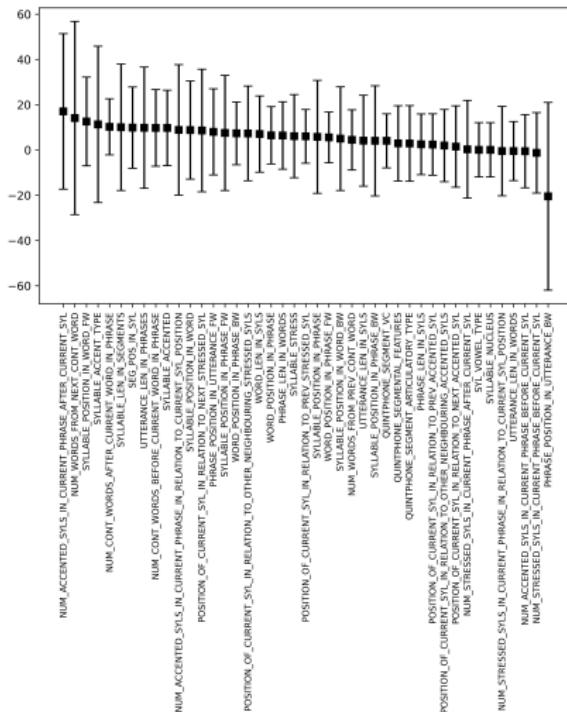


Figure: Regular mean relevance per feature group with the most relevant feature (V/UV) excluded. Medium granularity of feature groups.

Feature relevance analysis results

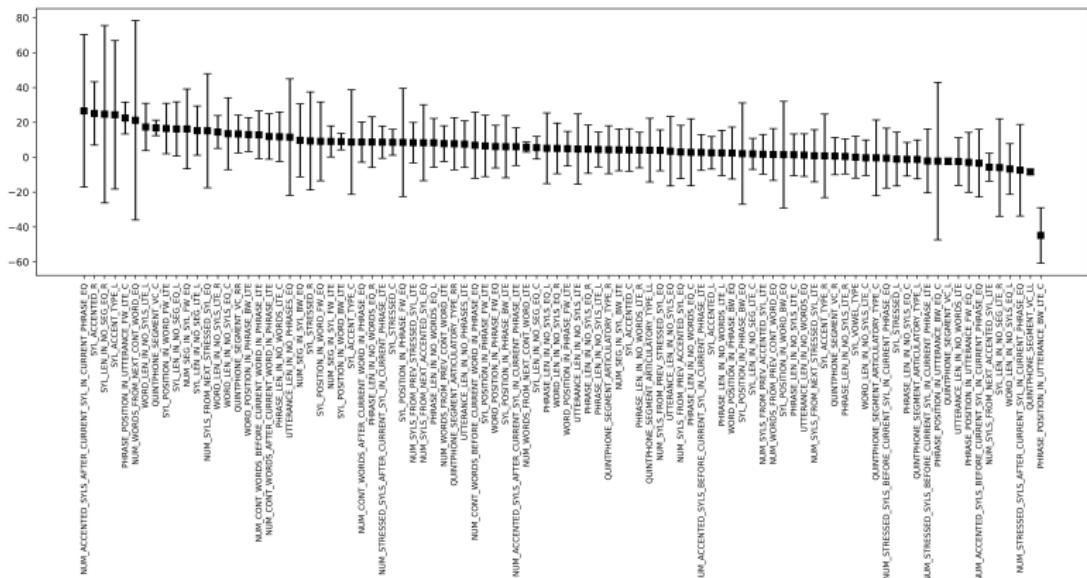
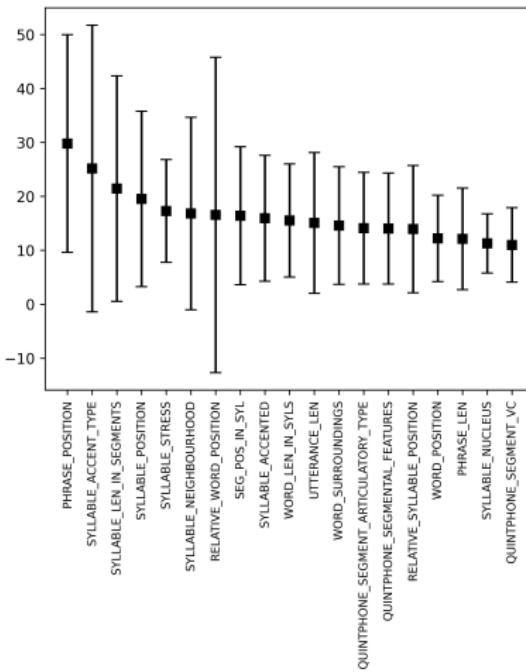


Figure: Regular mean relevance per feature group with the most relevant feature (V/UV) excluded. High granularity of feature groups.

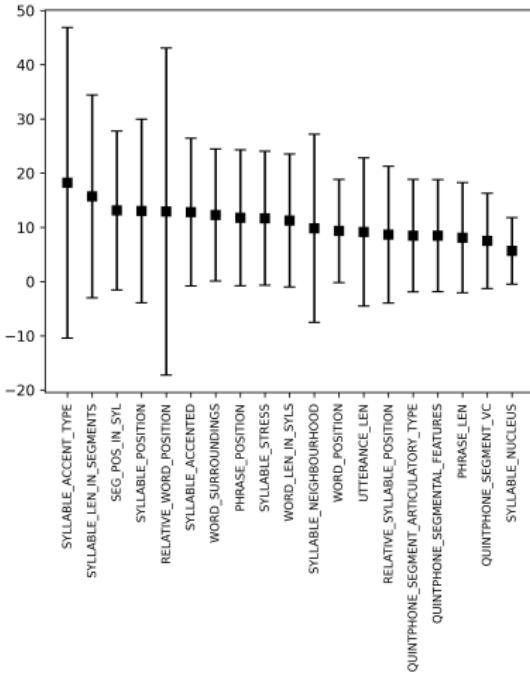
Feature relevance analysis results

Figure: Absolute mean relevance per feature group with the most relevant feature (V/UV) excluded.



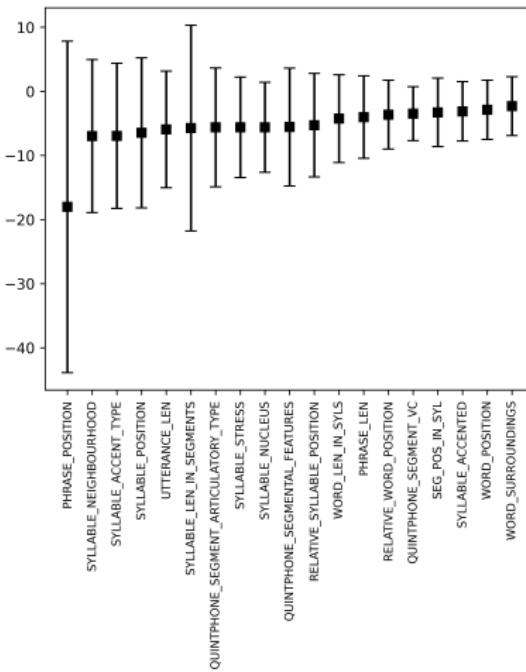
Feature relevance analysis results

Figure: Positive relevance-only mean per feature group with the most relevant feature (V/UV) excluded.



Feature relevance analysis results

Figure: Negative relevance-only mean per feature group with the most relevant feature (V/UV) excluded.



Main Hypotheses

HYPOTHESIS 1: *The continuous F_0 contours of an utterance emerge from its discrete linguistic features through a series of successive probabilistic mappings into intermediate latent representations.*

HYPOTHESIS 2: *The biologically-inspired Deep Temporal Convolutional Network can be an effective model of these mappings and hence of Polish neutral read speech intonation in the context of statistical-parametric speech synthesis.*

HYPOTHESIS 3: *The set of shallow linguistic features used in this thesis provides information which is sufficient for synthesis of natural sounding intonation in the context of statistical-parametric speech synthesis.*

HYPOTHESIS 4: *A Deep Temporal Convolutional Network can become an explanatory scientific model of mappings between linguistics features and the intonation of an utterance.*

- ✓ Generalized neural model of the mappings based on auditory cortical processing,
- ✓ Linguistically agnostic,
- ✓ Mapping model spontaneously converges towards a set of most relevant linguistic features as might be anticipated based on generally accepted linguistic theory (eg. syllabic features),
- ✓ Functional performance and naturalness of the produced F_0 contours as demonstrated through the subjective evaluation procedure,
- ✓ Substantial size of results produced based on the relevance analysis of model output.

Results

The complete results in the form of CSV files, charts, listings, etc. are available in the code repository at:
github.com/mrslacklines/intonation_synthesis

Contribution and possible applications

- ✓ Successful application of LRP explainability algorithm to a temporal signal modelling modelling (speech),
- ✓ Evaluation of a biologically-motivated model of the mappings between abstract linguistic categories and the concrete contours of F_0 ,
- ✓ Evidence for the effectiveness of Deep Neural Networks as explanatory scientific models and contribution to their popularization in linguistic research,
- ✓ Application in computational neuroscience research,
- ✓ Modest step towards the unification of theories of intonation,
- ✓ Methodology for determining the most relevant features that could help optimize the implementation of practical DNN-based models and the preparation of training data sets,
- ✓ Evidence that could support or undermine some of the current assumptions in phonological grammars of intonation,
- ✓ Complete open-source implementation available in public repositories including quickstart instructions.

Code repository

https://github.com/mrslacklines/intonation_synthesis

Improvements and future plans

- ✓ Size, coverage and quality of the dataset,
- ✓ Better data work, exploratory analysis (e.g. feature correlations)
- ✓ Better objective metric,
- ✓ Comparison of multiple model instances,
- ✓ Hyperparameter optimization,
- ✓ Account for the influence of the biomechanics of the articulatory tract on the final contour of the intonation,
- ✓ Measure relevance given the accuracy of the predictor,
- ✓ Redesign the model based on a single coherent neural theory and provide a better justification of the model architecture,
- ✓ Run experiments on features from other levels of language (such as semantics, syntax, pragmatics) or for expressive speech,
- ✓ Test other variants of the LRP algorithm,
- ✓ Verify the relevance of features with an ablation study,
- ✓ Design better subjective evaluation experiment,
- ✓ Analyze the intermediate latent features extracted by the deep convolutional filters.



Thank you!

References I

-  Alber, Maximilian et al. (2019). "iNNvestigate neural networks!" In: *Journal of Machine Learning Research* 20.93, pp. 1–8. URL: <http://jmlr.org/papers/v20/18-540.html>.
-  Demenko, G., B. Möbius, and K. Klessa (2010). "Implementation of Polish speech synthesis for the BOSS system". In: *Bulletin of the Polish Academy of Sciences. Technical Sciences* 58.3, pp. 371–376.
-  Demenko, Grażyna (1999). *Analiza cech suprasegmentalnych języka polskiego na potrzeby technologii mowy*. Wydawnictwo Naukowe UAM.
-  Demenko, Grażyna, Jolanta Bachan, et al. (2008). "Development and evaluation of Polish speech corpus for unit selection speech synthesis systems". In: *Ninth Annual Conference of the International Speech Communication Association*.
-  Demenko, Grażyna, Katarzyna Klessa, et al. (2010). "Polish unit selection speech synthesis with BOSS: extensions and speech corpora". In: *International journal of speech technology* 13.2, pp. 85–99.

References II

-  Demenko, Grażyna and Agnieszka Wagner (2007). "Prosody annotation for unit selection TTS synthesis". In: *Archives of acoustics* 32.1, pp. 25–40.
-  Féry, Caroline (2016). *Intonation and Prosodic Structure*. Key Topics in Phonology. Cambridge University Press. DOI: [10.1017/9781139022064](https://doi.org/10.1017/9781139022064).
-  Gibbon, Dafydd (1976). *Perspectives of intonation analysis*. Herbert Lang Bern.
-  Hubel, D. H. and T. N. Wiesel (1959). "Receptive fields of single neurones in the cat's striate cortex". In: *The Journal of Physiology* 148.3, pp. 574–591. DOI: <https://doi.org/10.1113/jphysiol.1959.sp006308>. eprint: <https://physoc.onlinelibrary.wiley.com/doi/pdf/10.1113/jphysiol.1959.sp006308>. URL: <https://physoc.onlinelibrary.wiley.com/doi/abs/10.1113/jphysiol.1959.sp006308>.

References III

-  – (1962). “Receptive fields, binocular interaction and functional architecture in the cat's visual cortex”. In: *The Journal of Physiology* 160.1, pp. 106–154. DOI: <https://doi.org/10.1113/jphysiol.1962.sp006837>. eprint: <https://physoc.onlinelibrary.wiley.com/doi/pdf/10.1113/jphysiol.1962.sp006837>. URL: <https://physoc.onlinelibrary.wiley.com/doi/abs/10.1113/jphysiol.1962.sp006837>.
-  Kernel (image processing) (2021). *Kernel (image processing)* — Wikipedia, The Free Encyclopedia. [Online; accessed 23.01.2021]. URL: [https://en.wikipedia.org/wiki/Kernel_\(image_processing\)](https://en.wikipedia.org/wiki/Kernel_(image_processing)).
-  LeCun, Y. et al. (1989). “Backpropagation Applied to Handwritten Zip Code Recognition”. In: *Neural Computation* 1.4, pp. 541–551. DOI: [10.1162/neco.1989.1.4.541](https://doi.org/10.1162/neco.1989.1.4.541). eprint: <https://doi.org/10.1162/neco.1989.1.4.541>. URL: <https://doi.org/10.1162/neco.1989.1.4.541>.
-  Manassi, M., B. Sayim, and M. H. Herzog (2013). “When crowding of crowding leads to uncrowding”. In: *Journal of vision* 13.13, pp. 10–10.

References IV

-  Möbius, Bernd (1995). "Components of a quantitative model of German intonation". In: *Proceedings of ICPHS*. Vol. 95, pp. 108–115.
-  Möbius, Bernd, Matthias Pätzold, and Wolfgang Hess (1993). "Analysis and synthesis of German F0 contours by means of Fujisaki's model". In: *Speech Communication* 13.1-2, pp. 53–61.
-  Möbius, Bernd and Jan PH Van Santen (2000). "A quantitative model of F0 generation and alignment". In: *Intonation*. Springer, pp. 269–288.
-  Montavon, Grégoire et al. (2017). "Explaining nonlinear classification decisions with deep Taylor decomposition". In: *Pattern Recognition* 65, pp. 211–222.
-  Oord, Aaron van den et al. (2016). *WaveNet: A Generative Model for Raw Audio*. arXiv: 1609.03499 [cs.SD].
-  Remy, Philippe (2020). *Temporal Convolutional Networks for Keras*.
<https://github.com/philipperemy/keras-tcn>.

References V

-  Botinis, Antonis, ed. (2000). *Intonation: Past, Present, Future*. Dordrecht: Springer Netherlands, pp. 13–52. ISBN: 978-94-011-4317-2. DOI: 10.1007/978-94-011-4317-2_2. URL: https://doi.org/10.1007/978-94-011-4317-2_2.
-  Samek, Wojciech et al. (2016). “Interpreting the predictions of complex ML models by layer-wise relevance propagation”. In: *arXiv preprint arXiv:1611.08191*.
-  Tian, Biao, Paweł Kuśmierk, and Josef P. Rauschecker (2013). “Analogues of simple and complex cells in rhesus monkey auditory cortex”. In: *Proceedings of the National Academy of Sciences* 110.19, pp. 7892–7897.
-  Wang, Xin, Shinji Takaki, and Junichi Yamagishi (2019). “Neural source-filter waveform models for statistical parametric speech synthesis”. In: *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 28, pp. 402–415.

Extra slides

Speech synthesis - Wavenet

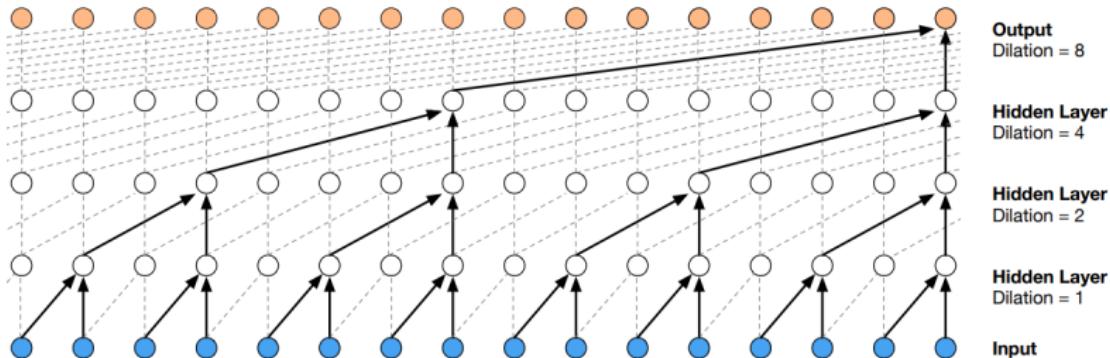


Figure: Dilated causal convolutions. Adopted from (Oord et al. 2016).

The causality is expressed through the joint probability of the modeled waveform $\vec{x} = \{x_1, \dots, x_T\}$ being factorized as a product of conditional probabilities of all previous timesteps (Oord et al. 2016), i.e.:

$$p(\vec{x}) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1}) \quad (1)$$

Speech synthesis - Wavenet

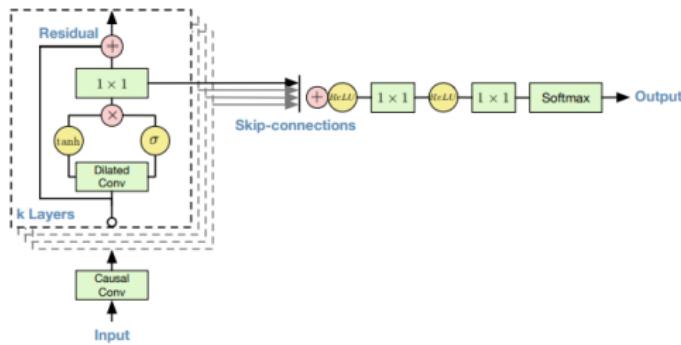


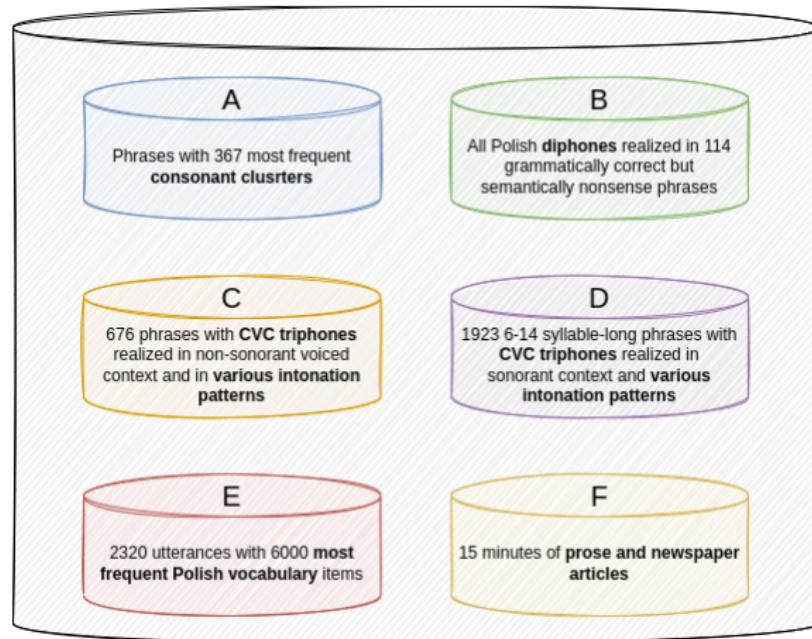
Figure: Residual and skip connections from a stack of k gated convolutional layers Adopted from (Oord et al. 2016).

Gated convolutional layers:

$$\vec{z} = \tanh \left(W_{f,k} * \vec{x} \right) \odot \sigma \left(W_{g,k} * \vec{x} \right), \quad (2)$$

where $*$ denotes a convolution operator, \odot denotes an element-wise multiplication operator, $\sigma(\cdot)$ is a sigmoid function, k is the layer index, f and g denote filter and gate, respectively, and W is a learnable convolution filter.

Dataset



Stress and accent type labels

%	rising accent realized by F_0 rise on post accented syllable/syllables or F_0 interval between accented and post accented vowels
,	rising accent realized by F_0 change (rise on accented syllable)
"	falling accent realized by F_0 fall on post accented syllable/syllables or F_0 interval between accented and post accented vowels
&	falling accent realized by F_0 change (fall on accented syllable)
	rising-falling accents with rise-fall shape of F_0 movement on accented vowel
*	level accent realized by F_0 interval between preaccented and accented vowels; near zero slope of fundamental frequency
<	level accent realized only by differences in duration between preaccented, accented and postaccented vowels

Figure: Stress and accent labels used in the original Polish BOSS speech corpus.

Stress and accent type labels

-5, .	Intonation on the first word in a sentence with falling accent F (or level accent L). In most cases it is used for declarative sentences or wh-questions. Mark on the first phoneme of the first word in the sentence
-5, ?	Intonation on the first word in a sentence with rising accent R. It can be used in different complex sentences. Mark on the first phoneme of the first word in the sentence
5, .	Intonation on the last word in sentence with falling accent F (or level accent L). In most cases it is used for declarative sentences or wh-questions. Mark on the first phoneme of the last word in the sentence
5, ?	Intonation on the last word in a sentence with rising accent R. In most cases it is used for yes-no questions. Mark on the first phoneme of the last word in the sentence
5, !	Intonation on the last word in a sentence with falling accent F. In most cases it is used for exclamatory sentences. Mark on the first phoneme of the last word in the sentence.
2, ?	Intonation on the last word in the phrase with rising accent R. In most cases it is used for continuation phrases. Mark on the first phoneme of the last word in the phrase.
2, ..	Intonation on the last word in the phrase with falling accent F (or level accent L). In most cases it is used in declarative phrases in complex sentences. Mark on the first phoneme of the last word in the sentence.

Figure: Prosodic phrase boundary labels used in the original Polish BOSS speech corpus.

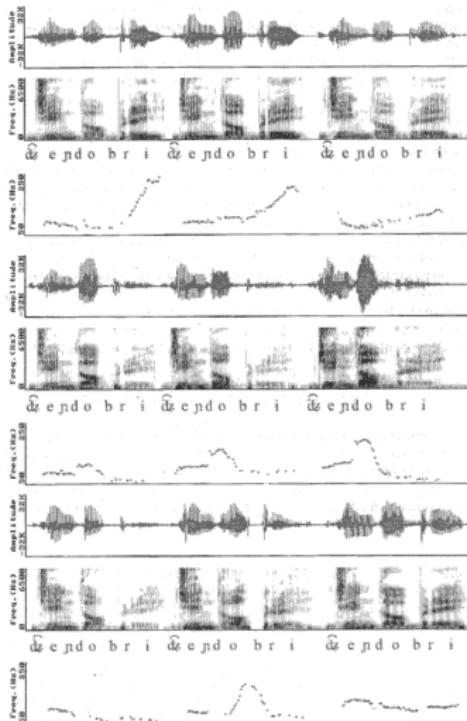


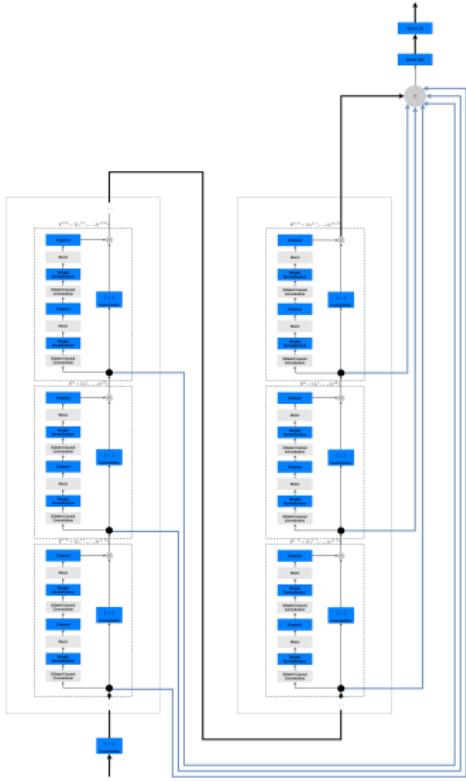
Figure: Acoustic realizations of the 9 different accents. Adopted from (Grażyna Demenko 1999).

Feature set

Question type	Segments	
Vowel	{i, ɛ, e, ɔ, u, ɯ, [schwa]}	Number of preceding/succeeding segments in the previous/current/next syllable is equal to/less than or equal to 0-7
Consonant	{g, ɣ, p, b, t, d, k, g, k̚, g̚, f, v, s, ʂ, z, ʐ, m̚, m̚, r, x, c, dz, cz, drz, cl, tʂl, n, n̚, ʂl, ng, l, r, v, w, j, ʐl}	Previous/current/next syllable is stressed
Stop	{g, ɣ, p, b, t, d, k, g̚, ʂ}	Previous/current/next syllable has accent X (where X is one of the ToBI accents described above)
Nasal	{v, w, ʐl, n, m, n̚, ʂl}	Number of preceding/succeeding segments in the next syllable is equal to/less than or equal to 0-7
Frikative	{t, v, ʂ, zl, xl, z, ʐz, rz, x̚}	Forward/backward position of the current syllable in current word is equal to/less than or equal to 0-7
Front	{e, ɛ, i, y, t̚, v, p, b, m, v, w̚}	Forward/backward position of the current syllable in current phrase is equal to/less than or equal to 0-20
General	{[schwa], n, ʂ, t, d, n, ʂl, x, x̚, n, r, t, d, dz, rz, rz, cz, drz, c, dzl, cl, ʂl}	Number of stressed syllables before/after the current syllable in current phrase is equal to/less than or equal to 0-12
Back	{o, ɔ, x, ʐ, ʂl, gl, ʂg, x, ʂd}	Number of accented syllables before/after the current syllable in current phrase is equal to/less than or equal to 0-12
Front Vowel	{e, ɛ, i, y}	Number of accented syllables before/after the current syllable in current phrase is equal to/less than or equal to 0-7
Central Vowel	{a, [schwa]}	Number of syllables from previous/next stressed syllable is equal to/less than or equal to 0-5
Back Vowel	{o, ɔ}	Number of syllables from previous/next accented syllable is equal to/less than or equal to 0-16
High Vowel	{i, ɛ}	Current syllable nucleus is a non-vowel, vowel, front vowel, central vowel, back vowel, high vowel, medium vowel, low vowel, rounded vowel, unrounded vowel, [i], [e], [a], [o], [u], [y], [schwa]
Medium Vowel	{e, ɔ}	Number of syllables in the previous/current/next word is equal to/less than or equal to 0-7
Low Vowel	{a}	Forward/backward position of the current word in the current phrase is equal to/less than or equal to 0-13
Rounded Vowel	{o, ɔ, ʐl}	Number of content words before/after the current word in the current phrase is equal to/less than or equal to 0-9
Unrounded Vowel	{a, ɛ, i, ʂl}	Number of words from previous/next content word is equal to/less than or equal to 0-5
Xlowl (e.g. Altvowel)	{i, ɛ, e, ɔ, u, ɯ, [schwa]}	Number of syllables in the previous/current/next phrase is equal to/less than or equal to 0-20
Unvoiced Consonant	{g, ɣ, p, t, k, k̚, f, v, s, ʂ, x, ʐ, m̚, m̚, r, x̚, dz, dzl, n, m, ʂl, l, r, v, w̚, j, ʐl}	Number of words in the previous/current/next phrase is equal to/less than or equal to 0-15
Voiced Consonant	{b, d, g, ɣl, v, z, ʐz, rz, xz, drz, dzl, n, m, ʂl, l, r, v, w̚, j, ʐl}	Forward/backward position of the current phrase in the utterance is equal to/less than or equal to 0-4
Front Consonant	{t, v, f, p, b, m, v, w̚}	Number of syllables in the utterance is equal to/less than or equal to 0-28
General Consonant	{t, d, s, ʂl, x, ʐz, n, r, l, t̚, d̚, rz, cz, drz, c, dzl, cl, ʂl}	Number of words in the utterance is equal to/less than or equal to 0-13
Back Consonant	{g, ɣ, k, g̚, k̚, gl, ʂg, x̚}	Number of phrases in the utterance is equal to/less than or equal to 0-4
Fortis Consonant	{g, ɣ, cz, t, k, p, s, ʂz, n, r, cl, ʂl}	
Lensis Consonant	{d̚, v, ʂ, b, ʐz, x, d, dzl, dz, gl, ʂl}	
Neither F or L	{n, ʂ, x, ʂl, ng, l, r, v, w̚, j, ʐl}	
Voiced Stop	{b, d, ɣl}	
Unvoiced Stop	{p, t, k, ʂl}	
Front Stop	{b, p}	
General Stop	{t̚, t̚}	
Back Stop	{g, ʂ, ʂl}	
Voiced Frikative	{v, z, ʐz}	
Unvoiced Frikative	{t̚, s, ʂl, m̚, x̚}	
Front Frikative	{t̚, v}	

Figure: Segmental features for a quintphone-wide context.

Figure: Non-segmental features.



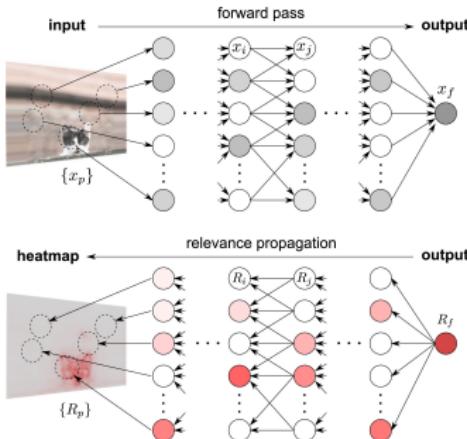


Figure: Computational flow of deep Taylor decomposition. Adopted from (Montavon et al. 2017).

The relevance in this framework can be defined as:

$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k \quad (3)$$