

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 (F0), and a grammatical (phonological) aspect." (Féry 2016)

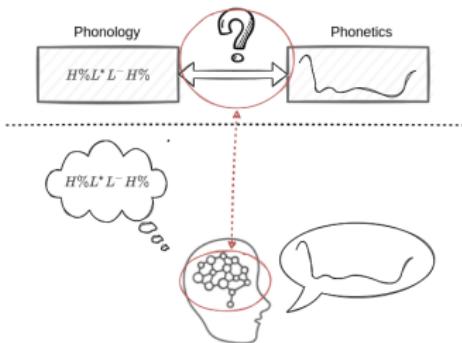
Intonation

Definition

"The tonal structure of speech expressed by the melody produced by our larynx. It has a phonetic aspect, the fundamental frequency (F0), 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

- ✓ Epistemological definition of intonation.
- ✓ Dualistic gap between phonology and phonetics.
- ✓ Unification within a broader metatheory.
- ✓ Unknown nature of the mappings between mental categories and continuous contours of $F0$.
- ✓ Explore how linguistic features of an utterance influence its $F0$ contours.
- ✓ Need for a physicalist (neurobiological) model.
- ✓ Modern statistical-parametric speech synthesis provides a framework for experimentation and evaluation of such models.

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 ($F0$ 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 $F0$ 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 ($F0$ 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 $F0$ contours of an utterance.

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.*

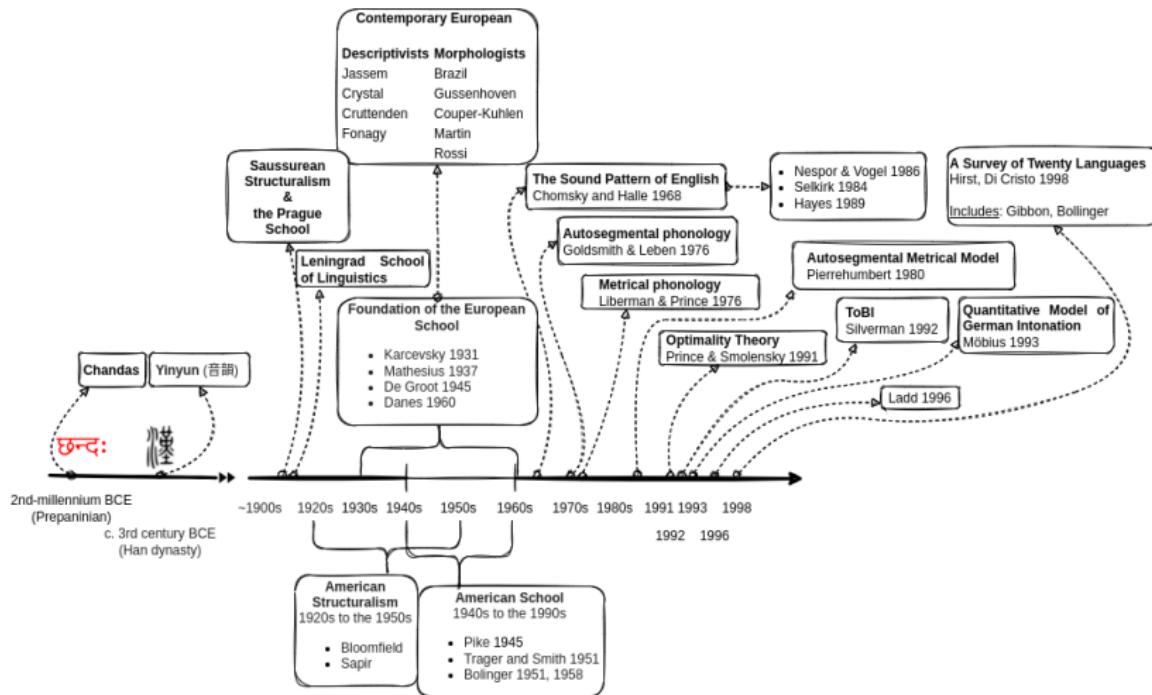
HYPOTHESIS 2: *The biologically-inspired Deep Temporal Convolutional Network can be an effective model of these mappings and hence of Polish neutral read speech intonation in the context of statistical-parametric speech synthesis.*

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.*

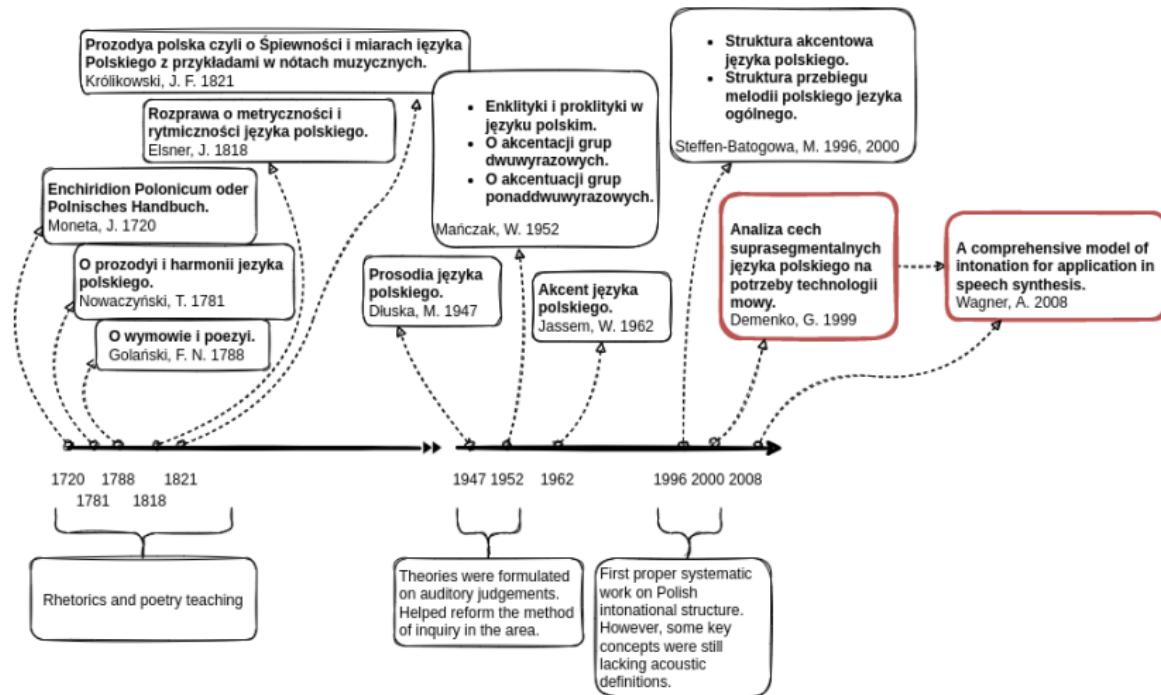
Contributory Methodological Hypothesis

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

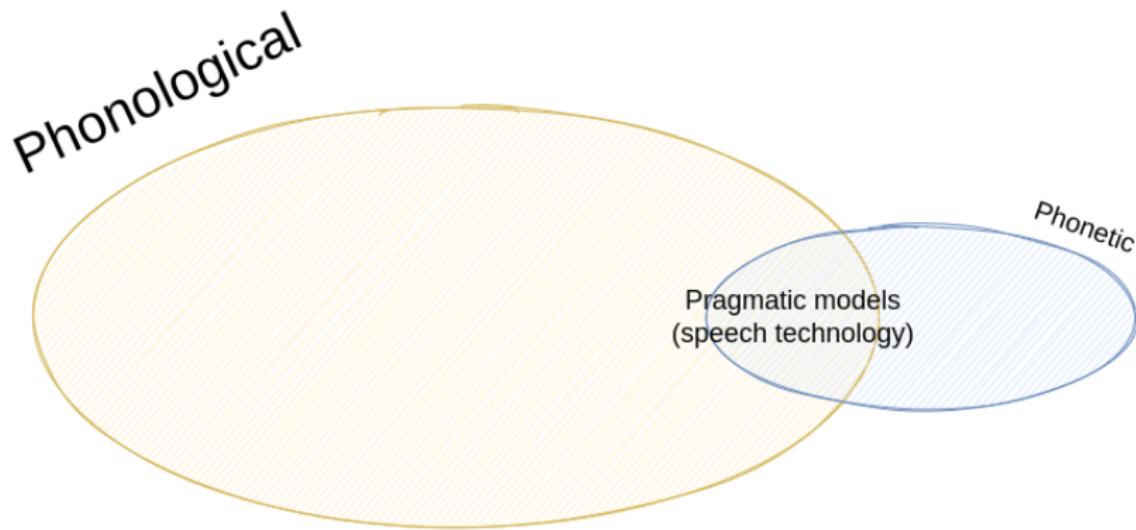
Background



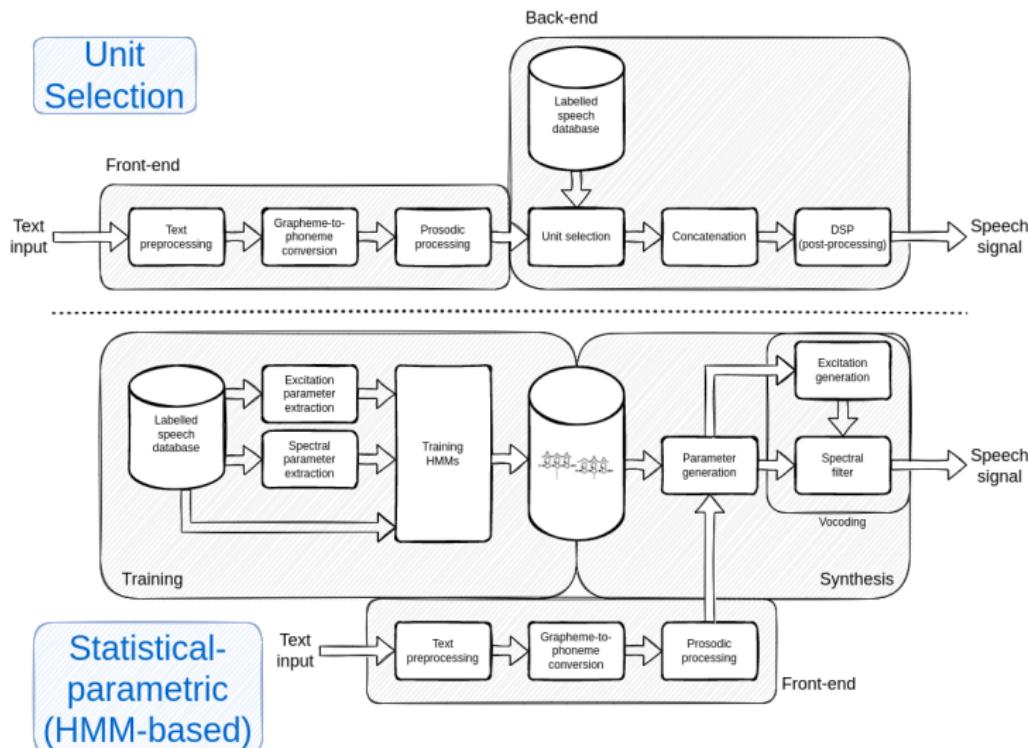
Background



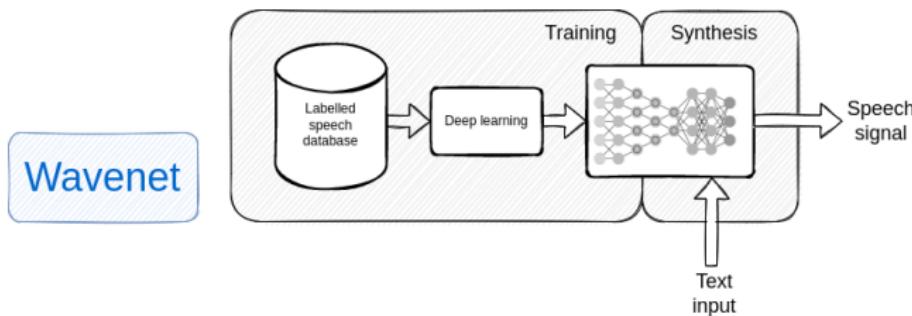
Intonation models



Speech synthesis



Speech synthesis - Wavenet



Wavenet

Wavenet belongs to a class of models known as Convolutional Neural Networks (CNNs).

Speech synthesis - Wavenet

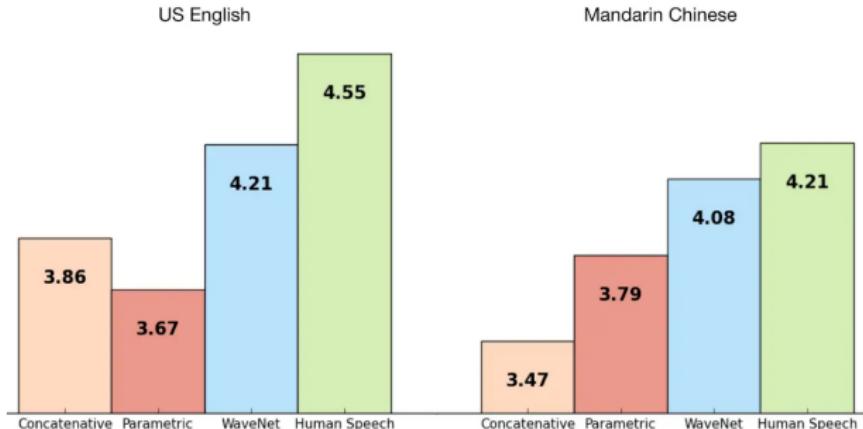


Figure: Google WaveNet evaluation results as compared with Google's best concatenative and parametric systems. (from van den Oord 2016).



Speech synthesis - Wavenet

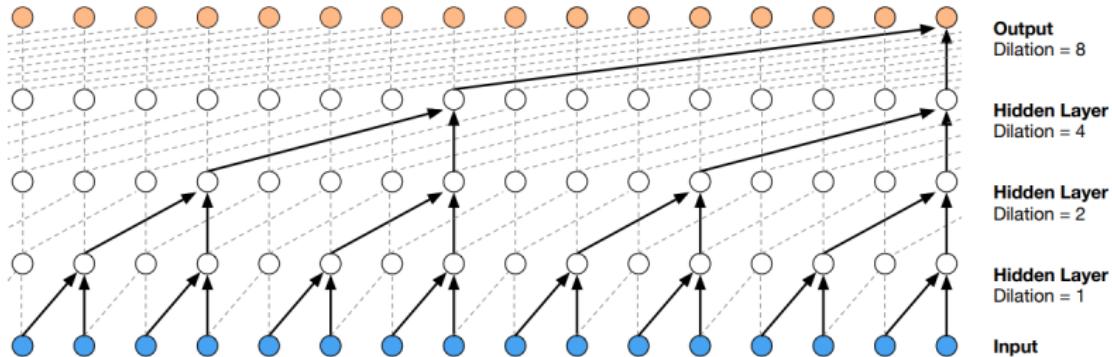


Figure: Dilated causal convolutions. (Adopted from the original WaveNet paper).

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 (van den Oord 2016), i.e.:

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

Speech synthesis - Wavenet

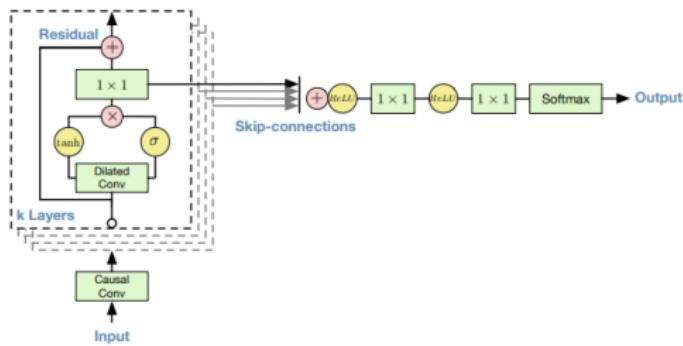


Figure: Residual and skip connections from a stack of k gated convolutional layers (Adopted from the original WaveNet paper).

Gated convolutional layers:

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

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.

Neurobiological foundations

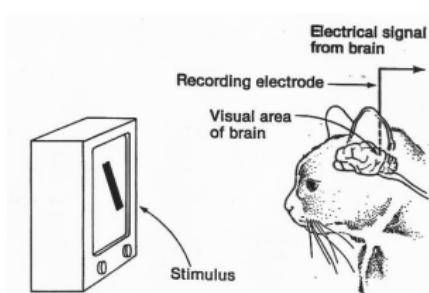


Figure: Famous Hubel and Wiesel cat experiment. (adopted from Hubel and Wiesel 1959).

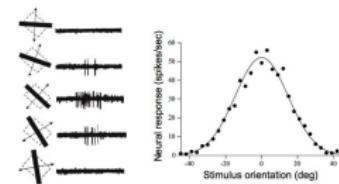


Figure: Neural response of simple cells. (adopted from Hubel and Wiesel 1968).

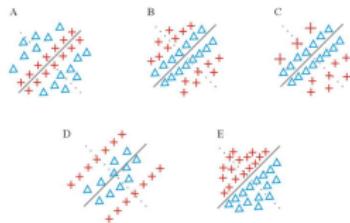


Figure: Simple receptive fields. (adopted from Hubel and Wiesel 1962).

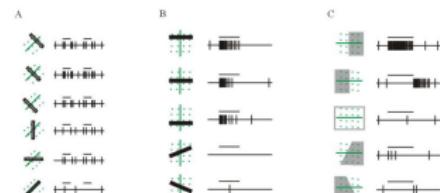


Figure: Three different types of complex receptive fields. (adopted from Hubel and Wiesel 1962).

Neurobiological foundations

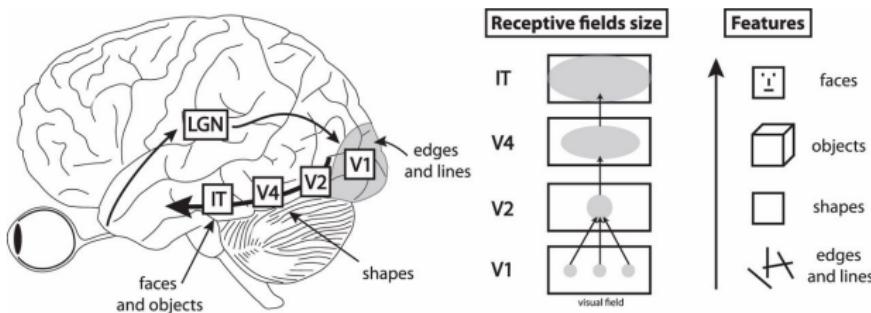


Figure: Hierarchical, feedforward visual processing in human brain. (Adopted from Manassi et al. 2013)

Neurobiological foundations

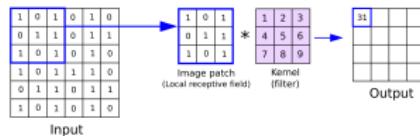


Figure: Example of a 2-dimensional matrix convolution.

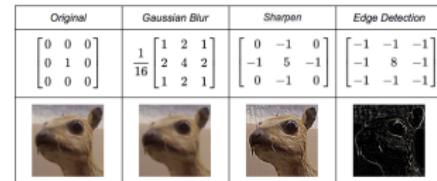


Figure: Examples of convoluting and image with different convolution kernels. (Adopted from the Wikipedia).

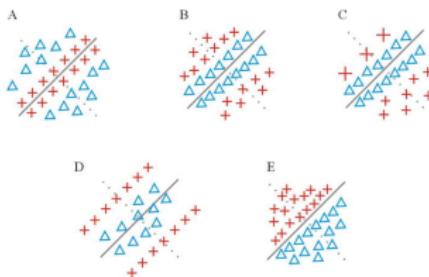


Figure: Simple receptive fields. (adopted from Hubel and Wiesel 1962).

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

Neurobiological foundations

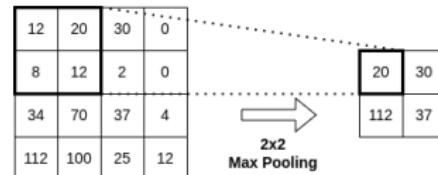


Figure: Example of a 2x2 max pooling matrix operation.

The function of the complex cells can be well approximated by the max pooling operation.

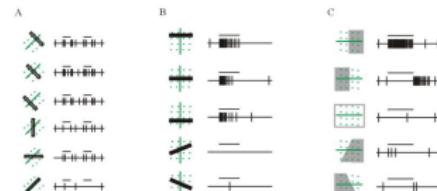


Figure: Three different types of complex receptive fields.
(adopted from Hubel and Wiesel 1962).

Neurobiological foundations

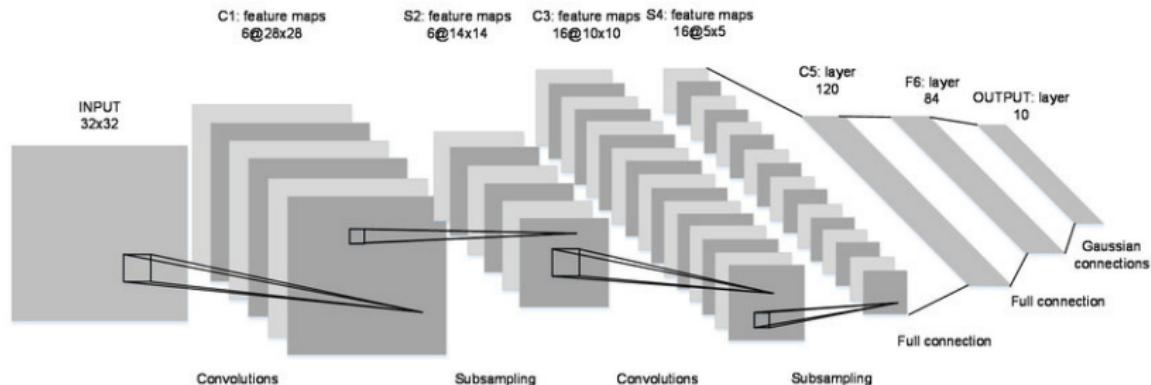
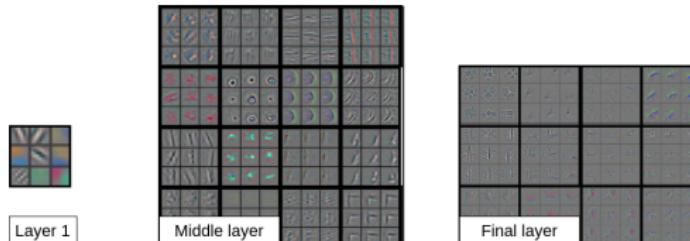


Figure: Image recognition Convolutional Neural Network (LeNet-5). (Adopted from LeCun et al. 1989)



Neural auditory processing

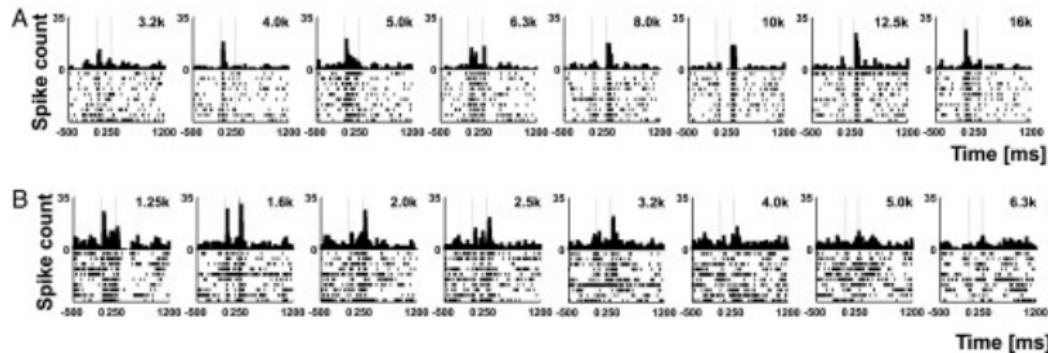


Figure: Responses of single neurons in primary auditory cortex (A1) of rhesus monkeys to band-passed noise (BPN) bursts centered at particular frequencies.. (Adopted from Tian et al. 2013)

Auditory cortex also contains simple cells with two dimensional receptive fields that are similar to those in the visual cortex but operate in the spectrotemporal domain.

Explanatory properties of CNNs

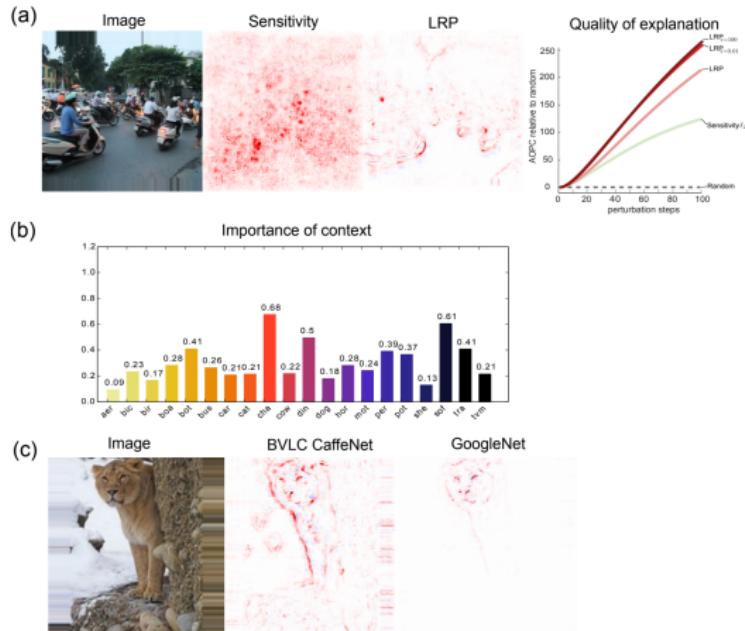


Figure: Results of sensitivity-based and relevance-based explainability methods. (Based on Samek 2016)

Deep Taylor Decomposition

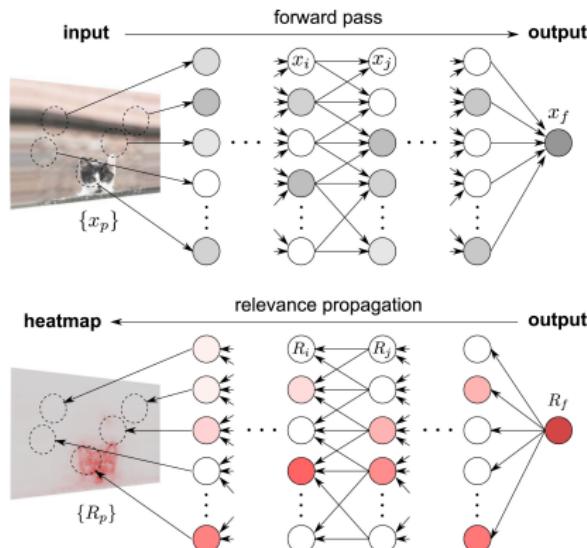


Figure: Computational flow of deep Taylor decomposition. (Adopted from Montavon 2017).

The relevance in this framework can be defined as:

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

Dataset

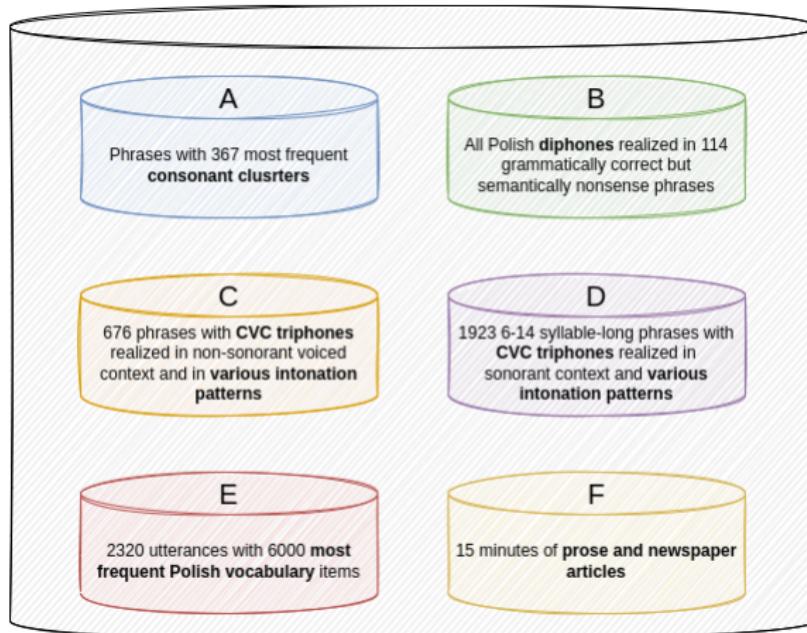


Figure: Speech corpus built originally for the purpose of the Polish BOSS unit selection synthesizer (Demenko, Bachan, Möbius 2008; Demenko, Klessa, Szymański, Breuer, Hess 2010).

Dataset

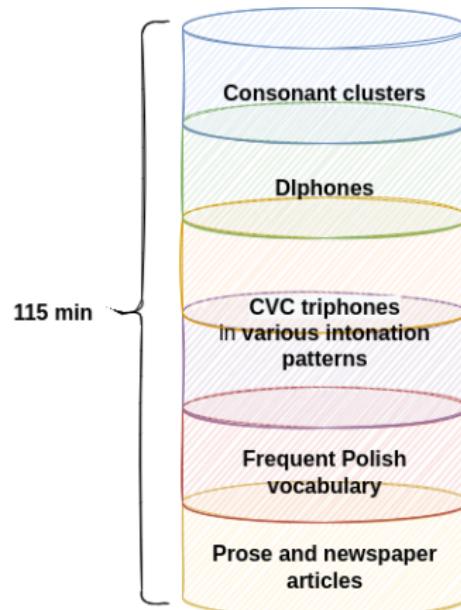
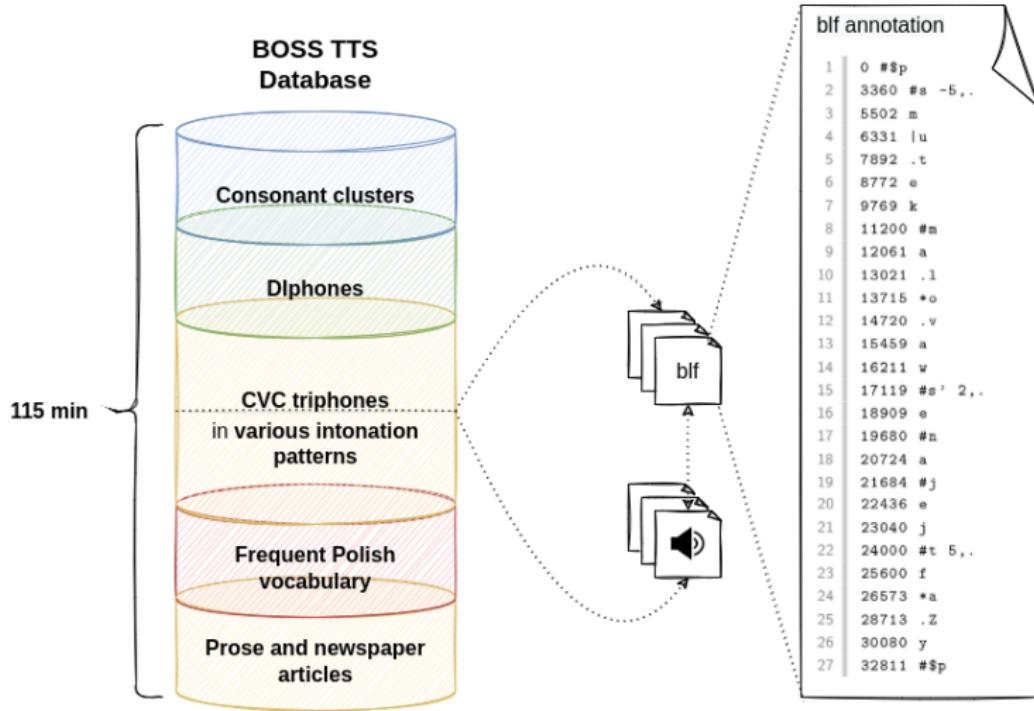
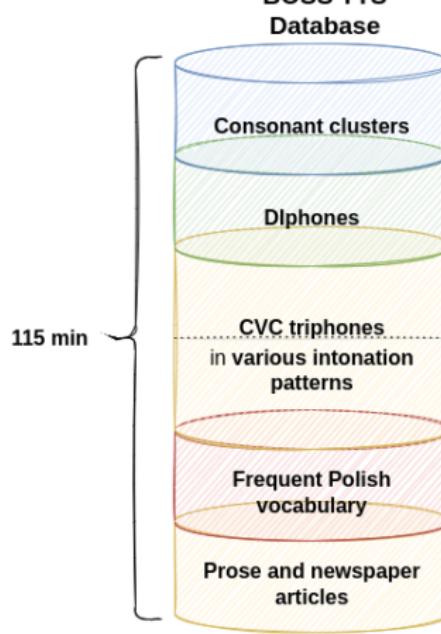


Figure: Speech corpus built originally for the purpose of the Polish BOSS unit selection synthesizer (Demenko, Bachan, Möbius 2008; Demenko, Klessa, Szymański, Breuer, Hess 2010).

Dataset



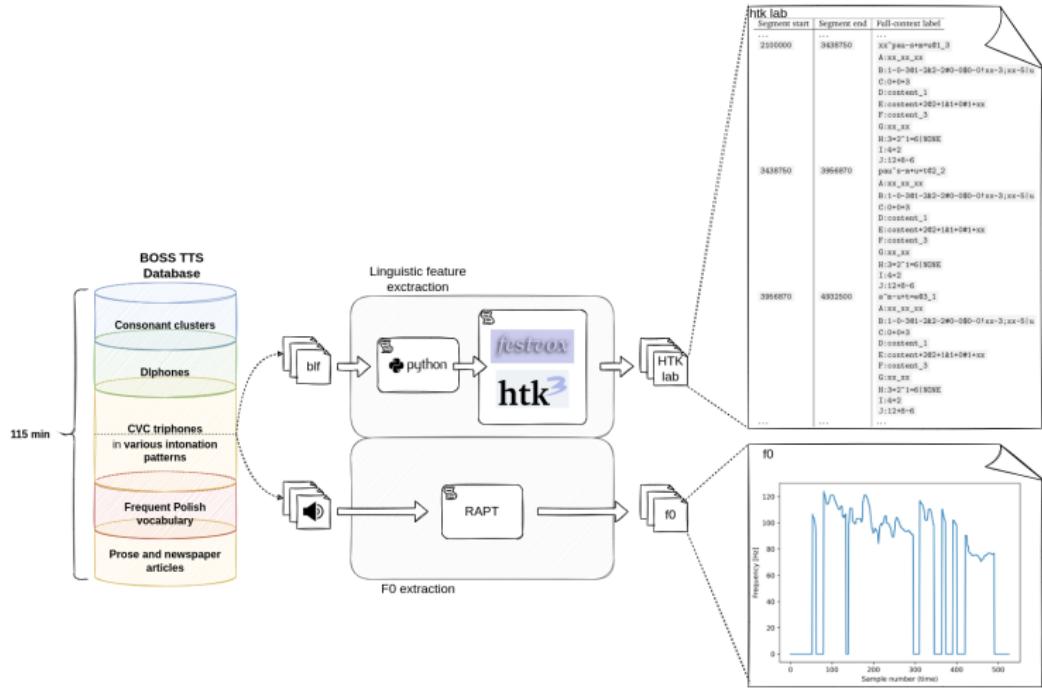
Dataset



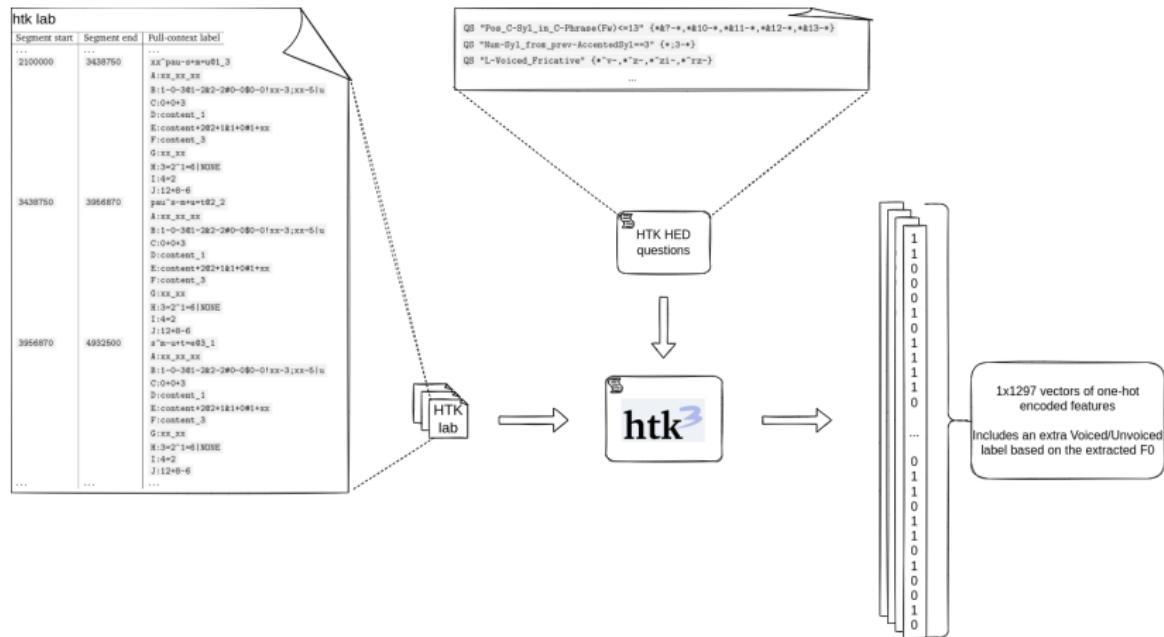
Stress and accent type labels	
-5..	Intonation on the first word in a sentence with falling accent F (or level accent L). In most cases it is used for declarative sentences or wh-questions. Mark on the first phoneme of the first word in the sentence
-5,7	Intonation on the first word in a sentence with 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 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 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.

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

Data preprocessing



Feature extraction



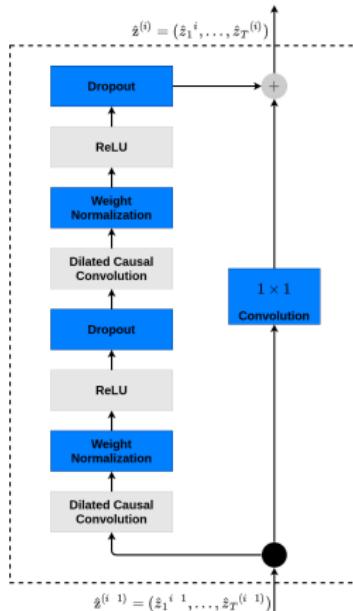
Feature set

Question type	Segments	
Vowel	(e, y, ε, ə, u, schwa)	Number of preceding/succeeding segments in the previous/current/next syllable is equal to/less than or equal to 0-7
Consonant	(g, p, b, v, d, k, g, kl, gl, f, v, s, z, sɪ, z, zl, m, rr, r, s, z, dz, cz, drz, cɪ, dɪ, n, n, nɪ, ng, l, r, w, wɪ, j, jɪ)	Previous/current/next syllable is stressed
Stop	(gɒ, p, b, t, d, k, g)	Previous/current/next syllable is accent
Nasal	(m, ŋ, n, w, n̩, m̩)	Previous/current/next syllable has accent X (where X is one of the ToBI accents described above)
Fricative	(f, v, z, s, sɪ, z, zl, m, rr, x)	Number of preceding/succeeding segments in the next syllable is equal to/less than or equal to 0-7
Front	(e, ɜ, y, ʌ, v, ɒ, b, m, w, v̩)	Forward/backward position of the current syllable in current word is equal to/less than or equal to 0-7
General	(əχnɒ, ə, ɜ, d, n, sɪ, z, zl, z, r, ɜ, t, d, rr, z, cz, drz, c, dz, cl, drɪl)	Forward/backward position of the current syllable in current phrase is equal to/less than or equal to 0-20
Back	(ə, u, ʌ, g, kl, gl, ng, z, ɒŋ)	Number of stressed syllables before/after the current syllable in current phrase is equal to/less than or equal to 0-12
Front Vowel	(e, ɜ, y)	Number of accented syllables before/after the current syllable in current phrase is equal to/less than or equal to 0-12
General Vowel	(ə, schwa)	Number of accented syllables before/after the current syllable in current phrase is equal to/less than or equal to 0-7
Back Vowel	(ə, u)	Number of syllables from previous/next stressed syllable is equal to/less than or equal to 0-5
High Vowel	(i, y, ɜ)	Number of syllables from previous/next accented syllable is equal to/less than or equal to 0-16
Medium Vowel	(ə, ɜ)	Current syllable nucleus is a non-vowel, vowel, front vowel, central vowel, back vowel, high vowel, medium vowel, low vowel, rounded vowel, unrounded vowel, [l], [ɛ], [a], [o], [u], [y], [schwa]
Low Vowel	(ə)	Number of syllables in the previous/current/next word is equal to/less than or equal to 0-7
Rounded Vowel	(ə, ɜ)	Forward/backward position of the current word in the current phrase is equal to/less than or equal to 0-13
Unrounded Vowel	(ə, e, ɪ, ɜ)	Number of content words before/after the current word in the current phrase is equal to/less than or equal to 0-9
XVowel (e.g. Alvowel)	(i, y, ə, e, ɔ, u, schwa)	Number of words from previous/next content word is equal to/less than or equal to 0-5
Unvoiced Consonant	(p, ph, t, k, kl, f, v, s, z, zl, x, c, ts, cɪ, dɪ)	Number of syllables in the previous/current/next phrase is equal to/less than or equal to 0-13
Voiced Consonant	(b, d, g, gɪ, v, z, zl, rr, dz, drz, dzl, n, n̩, m, m̩, ng, l, r, v, w, wɪ, j, jɪ)	Number of words in the utterance is equal to/less than or equal to 0-13
Front Consonant	(f, v, ɜ, p, b, m, w, v̩)	Number of words in the utterance is equal to/less than or equal to 0-13
General Consonant	(t, d, s, ə, sɪ, z, zl, n, r, ɜ, ʌ, t, d, sr,	Number of words in the utterance is equal to/less than or equal to 0-13
Back Consonant	(r, rr, c, drz, c, dz, cɪ, dɪ)	Number of words in the utterance is equal to/less than or equal to 0-13
Plosis Consonant	(p, g, k, kl, gɪ, z, x)	Number of words in the utterance is equal to/less than or equal to 0-13
Lens Consonant	(tr, v, g, b, v̩, z, ə, dsl, dz, gl, zɪ)	Number of words in the utterance is equal to/less than or equal to 0-13
Neigther F or L	(m, ə, sɪ, ng, l, r, w, wɪ, j, jɪ)	Number of syllables in the previous/current/next phrase is equal to/less than or equal to 0-20
Voiced Stop	(b, d, g)	Number of words in the previous/current/next phrase is equal to/less than or equal to 0-15
Unvoiced Stop	(p, t, k, ɒŋ)	Forward/backward position of the current phrase in the utterance is equal to/less than or equal to 0-4
Front Stop	(t, p)	Number of syllables in the utterance is equal to/less than or equal to 0-28
Central Stop	(t, t̩)	Number of words in the utterance is equal to/less than or equal to 0-13
Back Stop	(k, ɒŋ)	Number of syllables in the utterance is equal to/less than or equal to 0-13
Voiced Fricative	(v, z, ɜ, rr)	Number of words in the utterance is equal to/less than or equal to 0-13
Unvoiced Fricative	(f, s, ə, m, x)	Number of syllables in the utterance is equal to/less than or equal to 0-13
Front Fricative	(f, v)	Number of words in the utterance is equal to/less than or equal to 0-13

Figure: Segmental features for a quintphone-wide context.

Figure: Non-segmental features.

Model implementation



Complete code repository

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

Figure: Segmental features for a quintphone-wide context.

Thank you.

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.

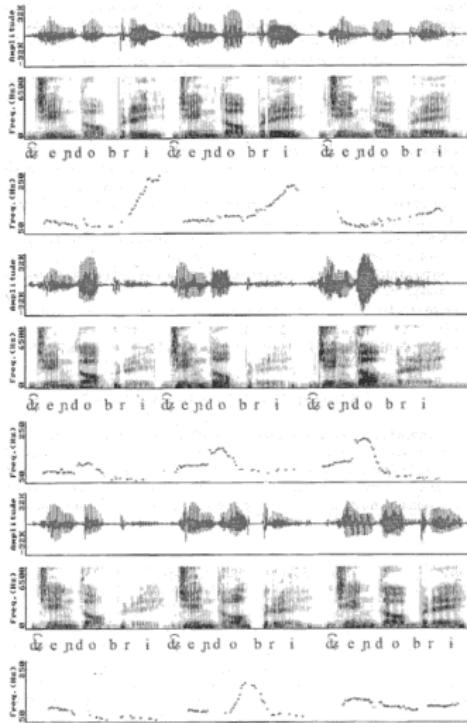


Figure: Acoustic realizations of the 9 different accents. Adopted from (Demenko 1999).