

# 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)

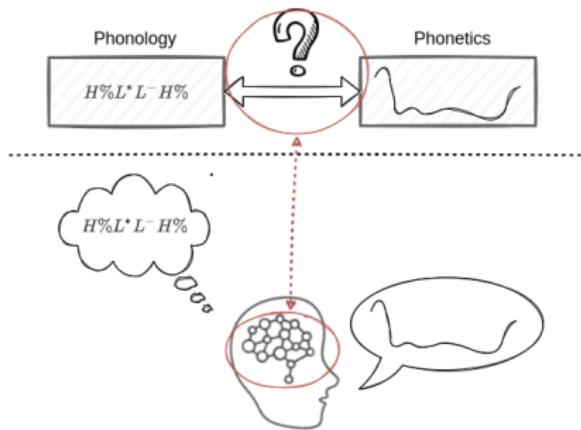
# Intonation

## Definition

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

- ✓ Dualistic gap between phonology and phonetics.
- ✓ Unification within a broader metatheory.
- ✓ Unknown nature of the mappings between mental categories and continuous contours of  $F0$ .
- ✓ 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
- ✓ 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 (*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.*

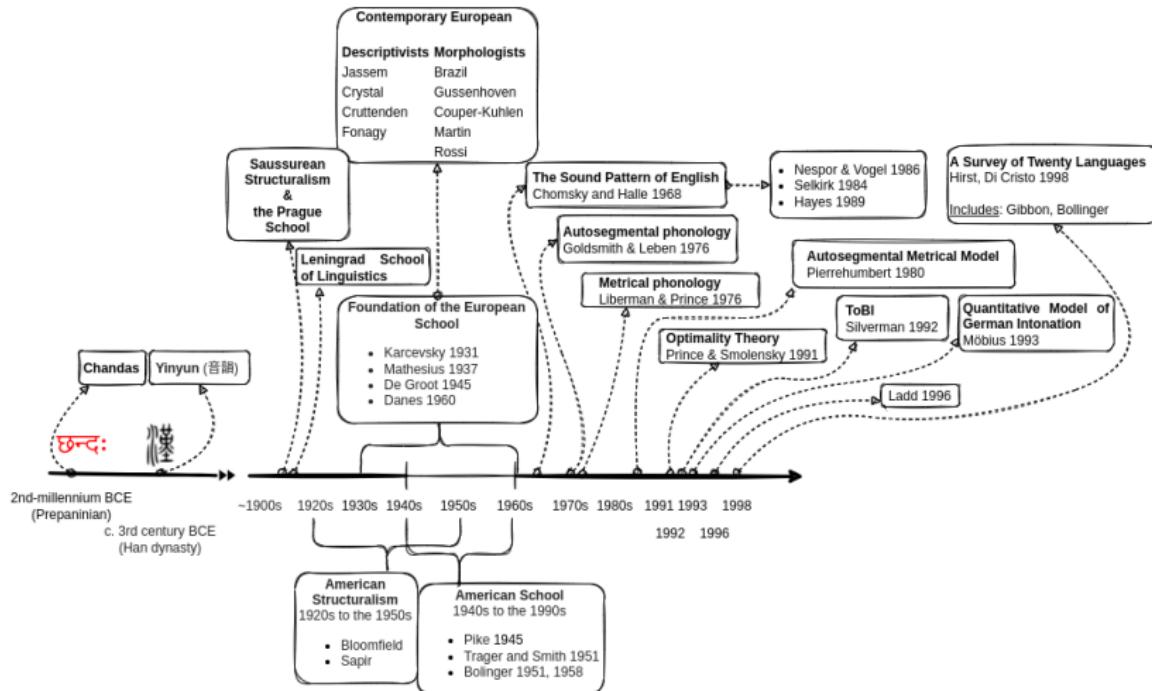
**HYPOTHESIS 2:** *The biologically-inspired Deep Temporal Convolutional Network can be an effective model of these mappings and hence of Polish neutral read speech intonation in the context of statistical-parametric speech synthesis.*

**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.*

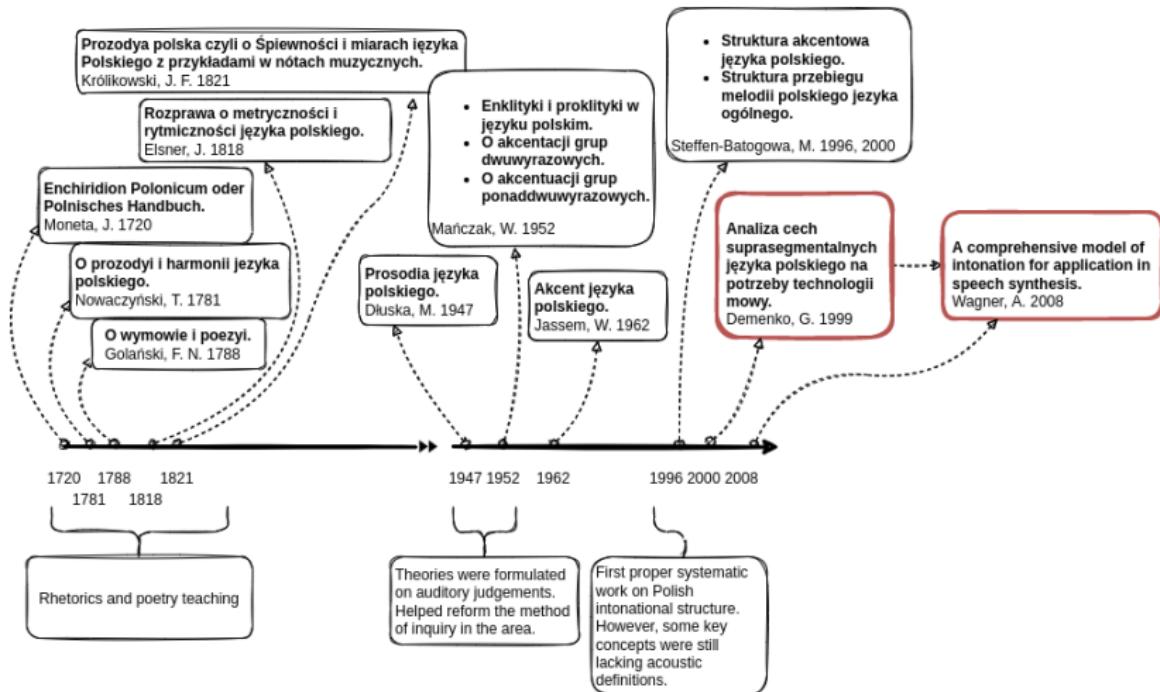
# Contributory Methodological Hypothesis

**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.*

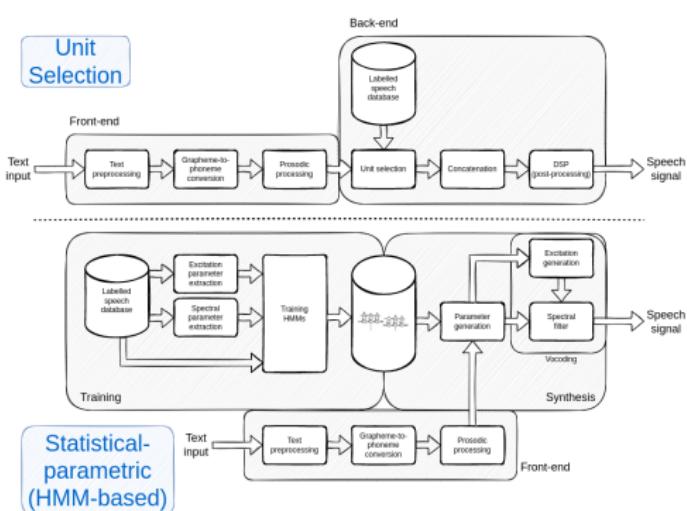
# Background



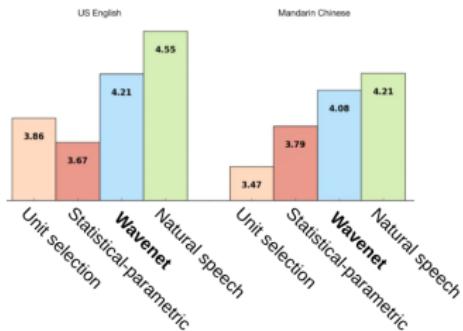
# Background



# Speech synthesis

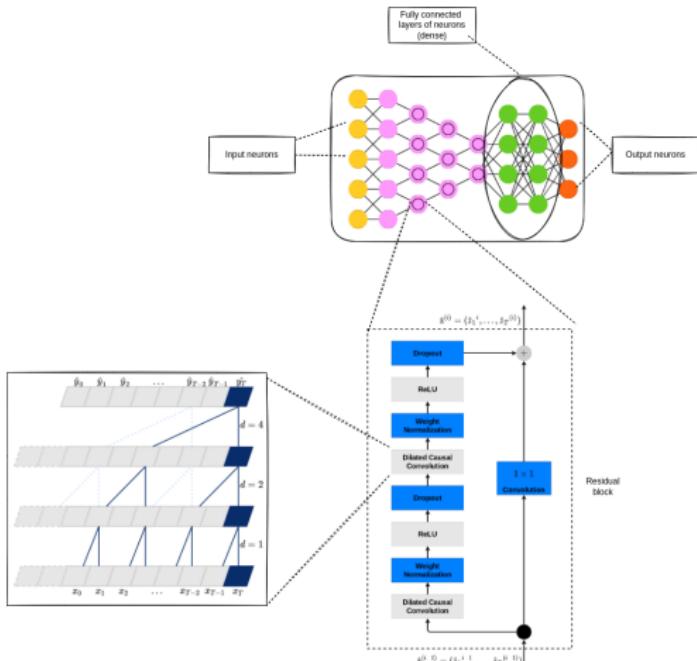


# Speech synthesis - Wavenet



**Figure:** Google WaveNet evaluation results.  
(from van den Oord 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

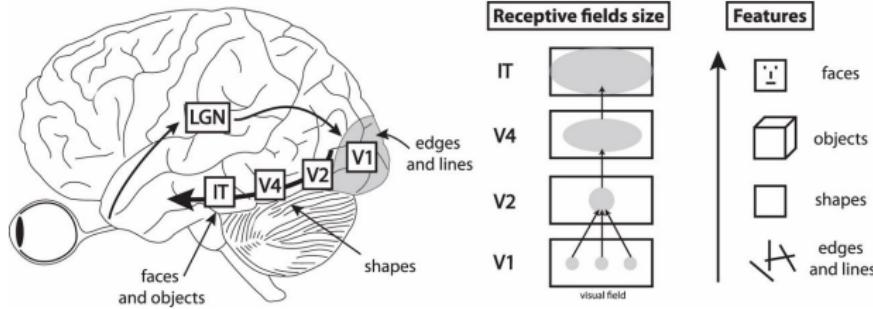


Figure: Hierarchical, feedforward visual processing in human brain. (Adopted from Manassi et al. 2013)

# CNN feature maps

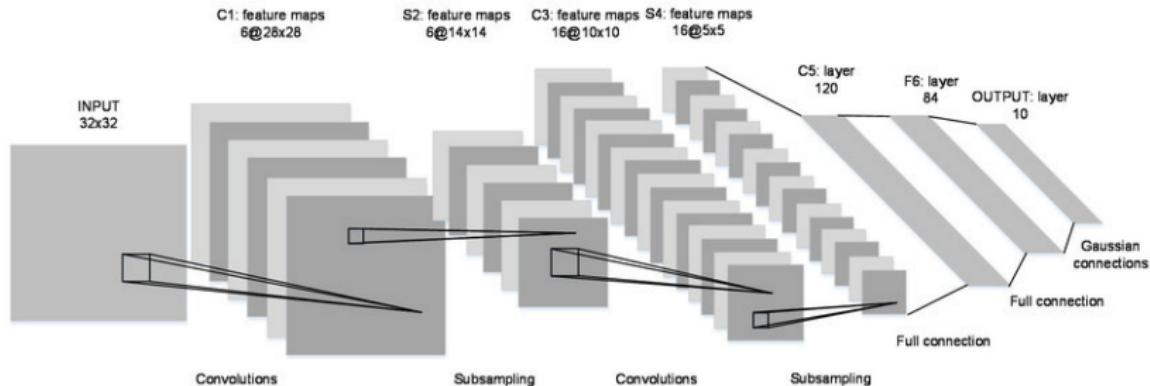
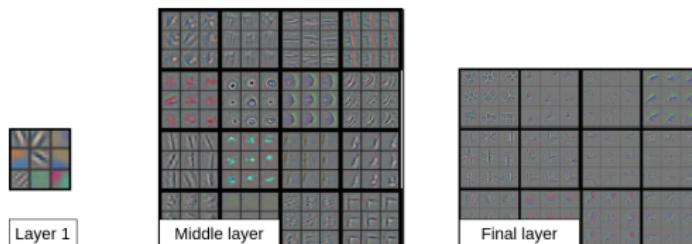
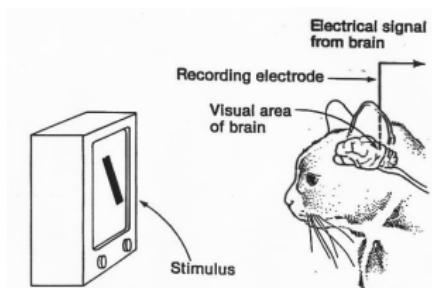


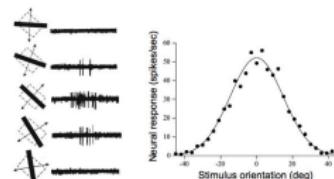
Figure: Image recognition Convolutional Neural Network (LeNet-5). (Adopted from LeCun et al. 1989)



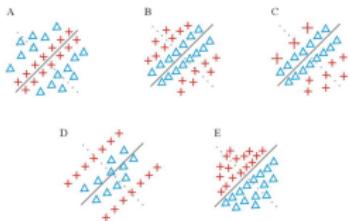
# Simple and complex cells in the visual cortex



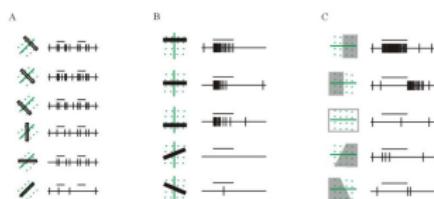
**Figure:** Famous Hubel and Wiesel cat experiment.  
(adopted from Hubel and Wiesel 1959).



**Figure:** Neural response of simple cells. (adopted from Hubel and Wiesel 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 of CNNs

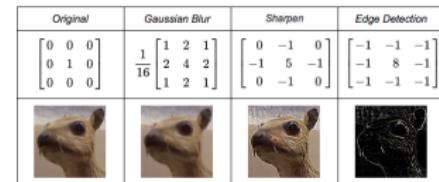
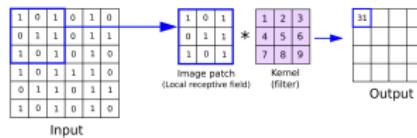


Figure: Example of a 2-dimensional matrix convolution.

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

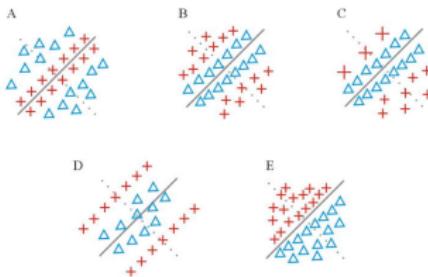


Figure: Simple receptive fields. (adopted from Hubel and Wiesel 1962).

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

# Neurobiological foundations of CNNs

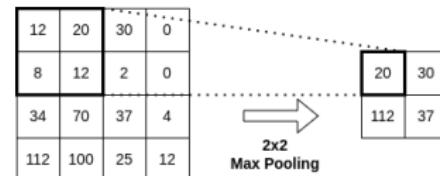


Figure: Example of a  $2 \times 2$  max pooling matrix operation.

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

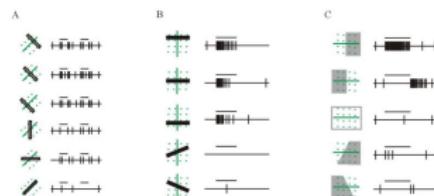
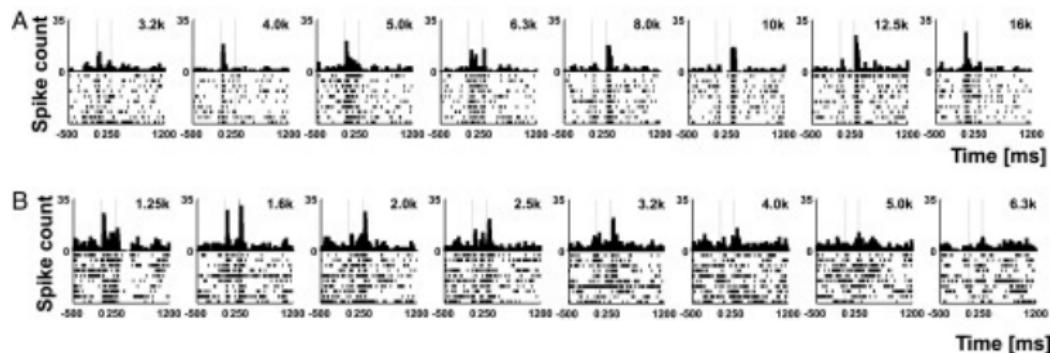


Figure: Three different types of complex receptive fields.  
(adopted from Hubel and Wiesel 1962).

# 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 from Tian et al. 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

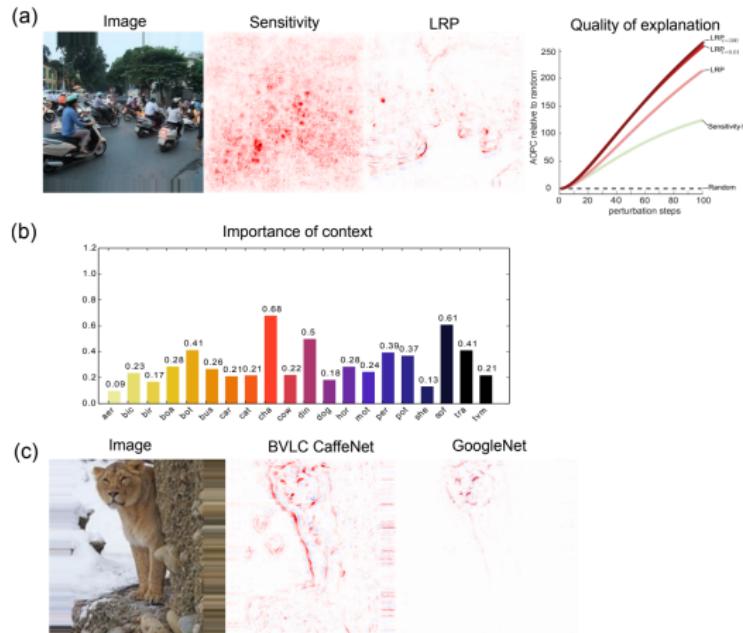
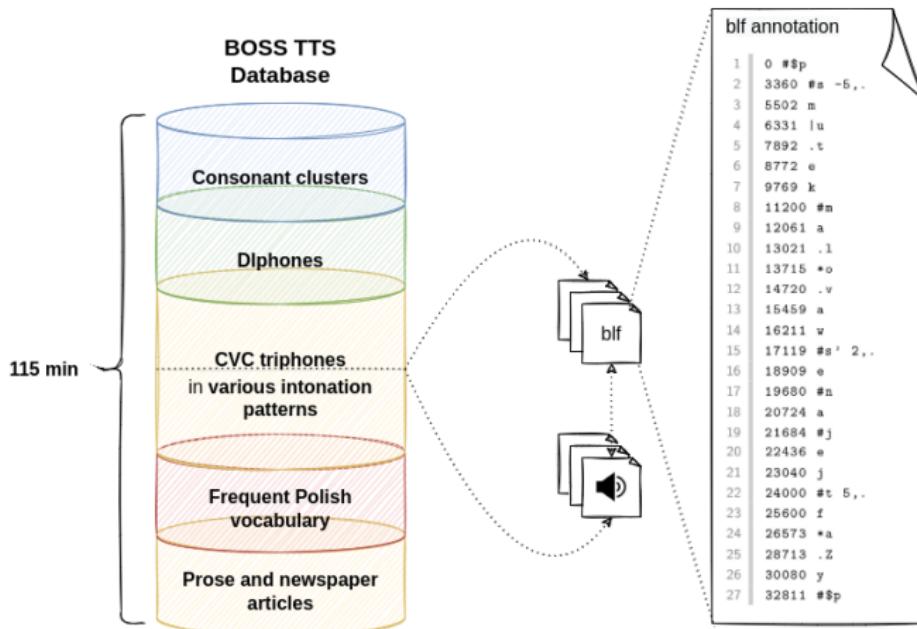


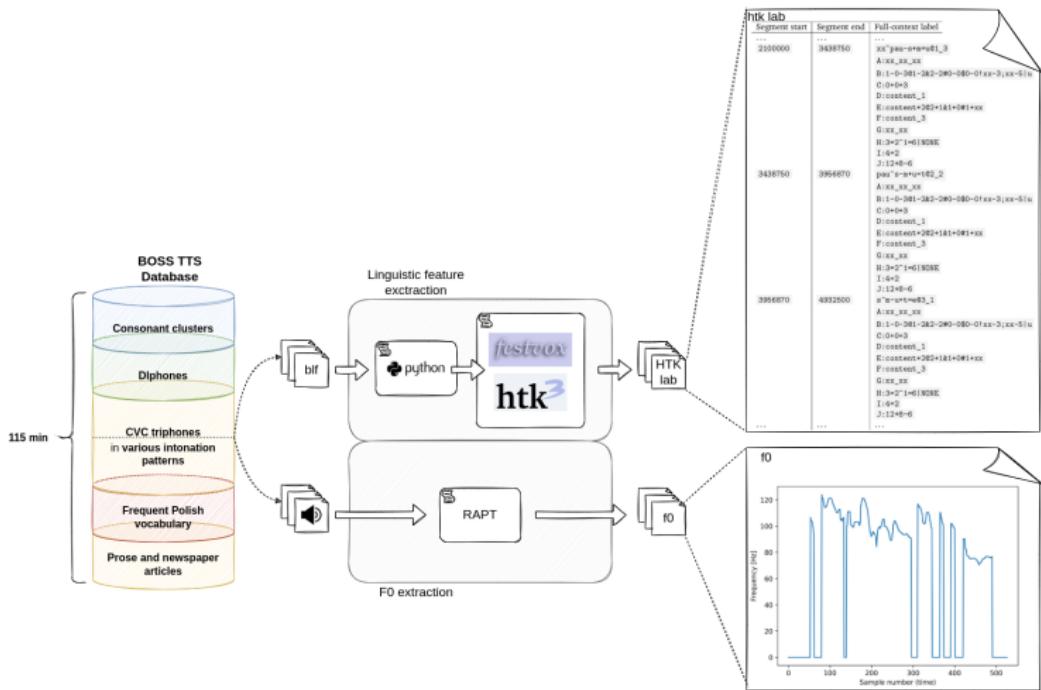
Figure: Results of sensitivity-based and relevance-based explainability methods. (Based on Samek 2016)

## Dataset

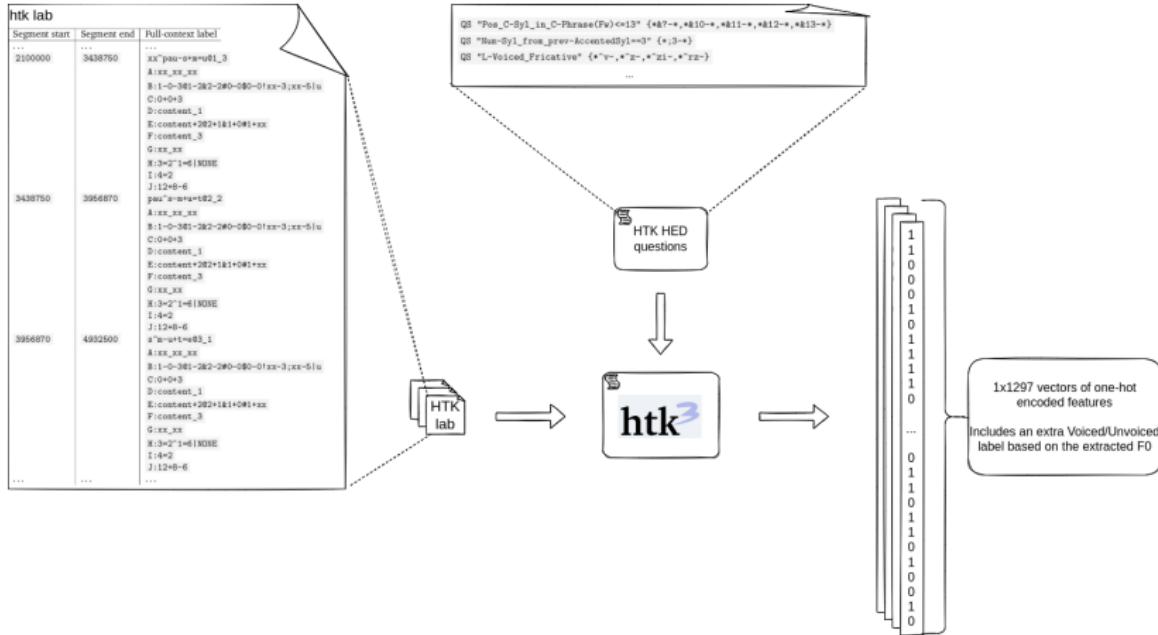


**Figure:** Speech corpus built originally for the purpose of the Polish BOSS unit selection synthesizer (Demenko, Bachan, Möbius 2008; Demenko, Klessa, Szymański, Breuer, Hess 2010).

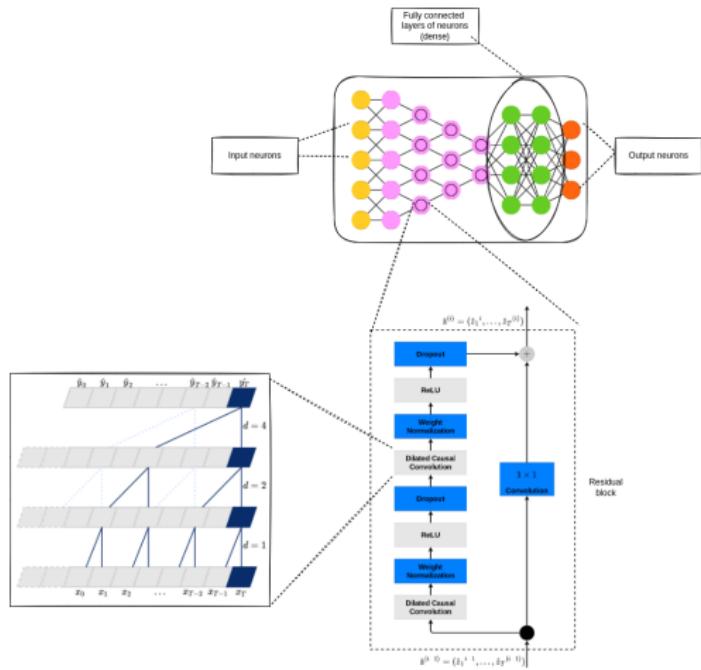
## Data preprocessing



## Feature extraction



## Model implementation



## TCN parameters

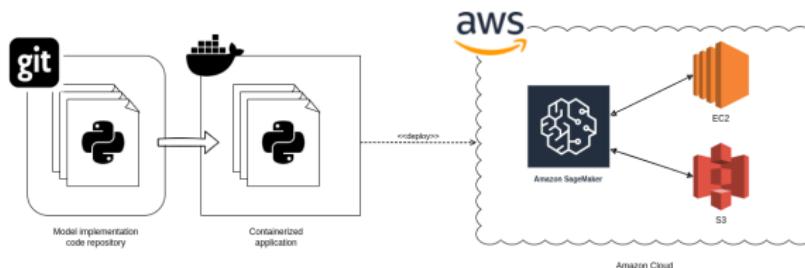
- 6 residual layers,
  - 64 convolution filters of length 2
  - $2^n$  dilations ( $[1, 2, 4, 8, 16, 32]$ )
  - 2 final fully connected dense layers  
(ReLU activation; sizes 64 and 1)

## Complete code repository

✓ [https://github.com/mrsslacklines/intonation synthesis](https://github.com/mrsslacklines/intonation-synthesis)

[github.com/philipperemy/keras-tcn](https://github.com/philipperemy/keras-tcn)

# Model training infrastructure



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

## Memory requirements

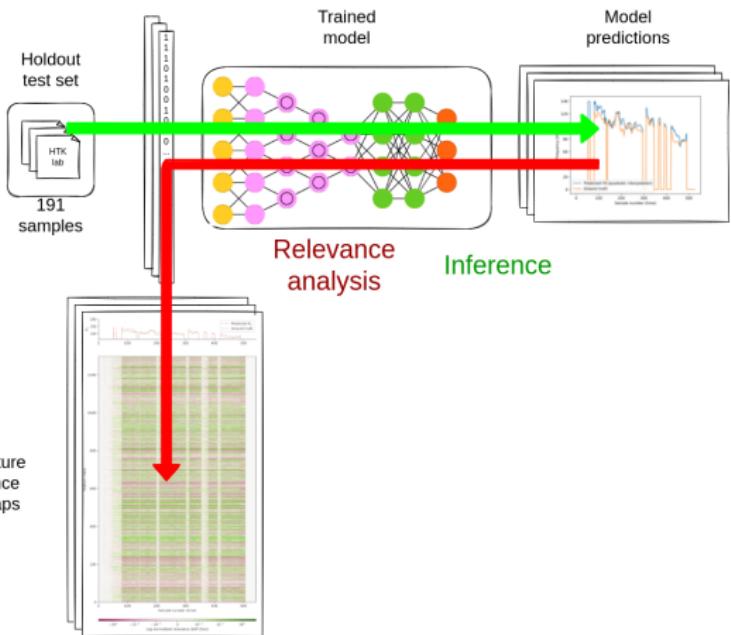
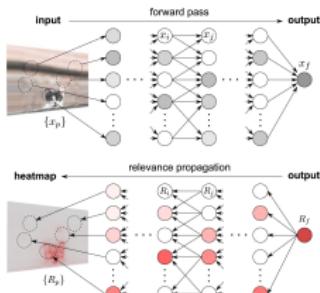
A single batch of data, which is a  $64 \times 1900 \times 1297$  vector of 8-byte boolean values, occupies 1.2617216 gigabytes of memory, and the model comprises of a total of 449,409 parameters (446,337 trainable and 3,072 non-trainable).



## 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

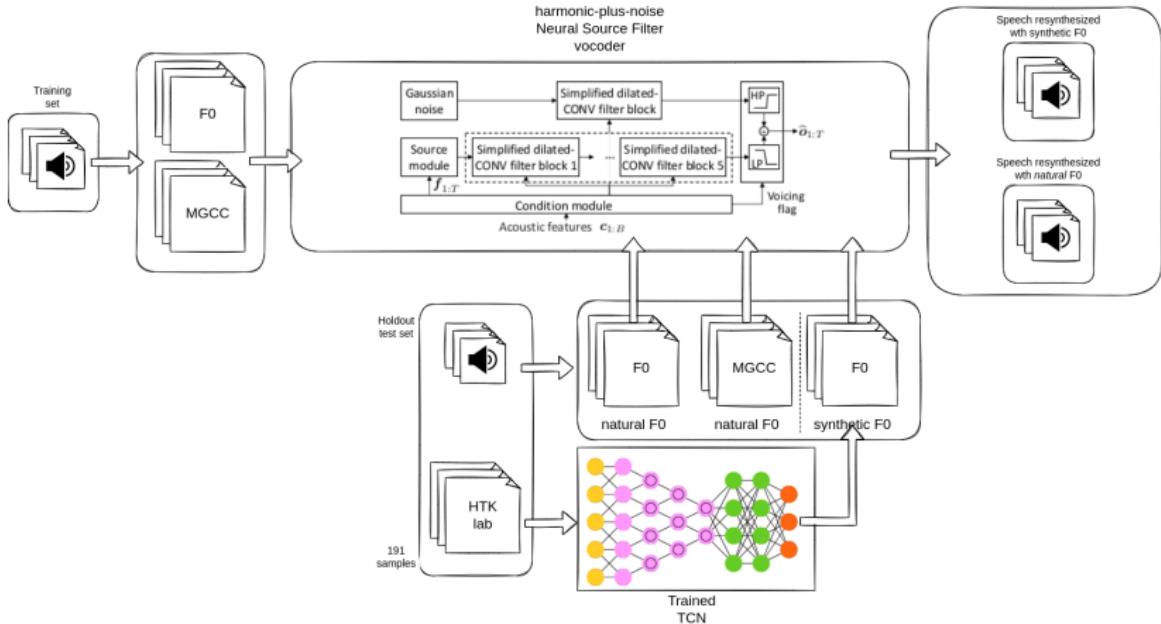
# F0 inference and feature relevance analysis



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

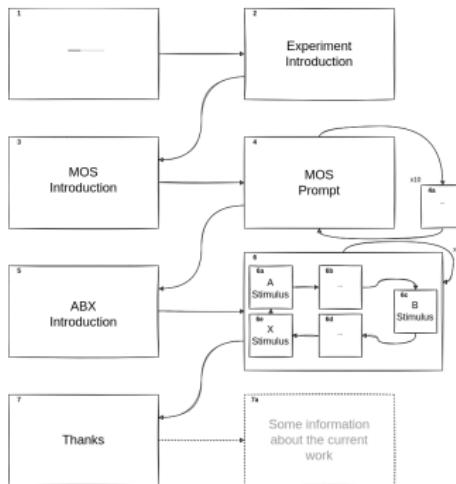
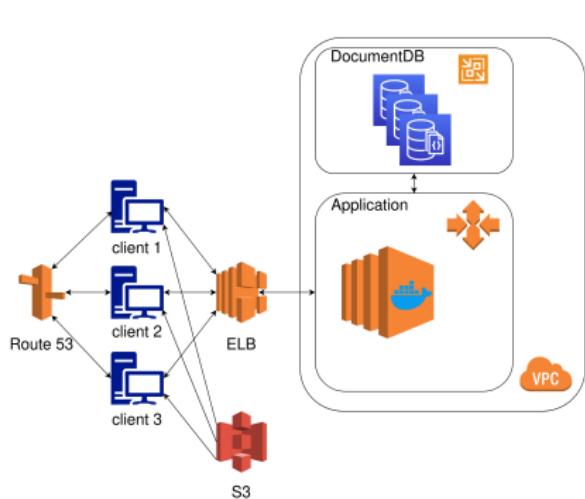
INNvestigate Neural Networks!  
(Alber et al. 2019)  
[github.com/albermax/  
innvestigate](https://github.com/albermax/innvestigate)

# Resynthesis



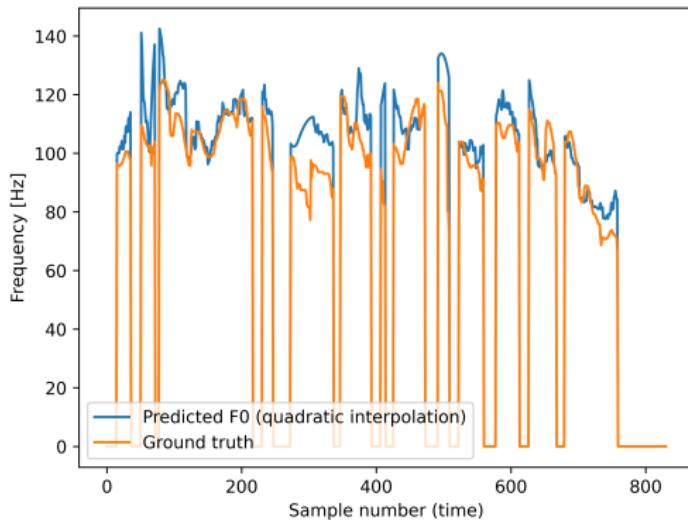
Neural Source Filter Vocoder from (X. Wang et al. 2019)

# Perceptual evaluation experiment



**Available at:**  
[fonetyka.cudaniewidły.org/experiment](http://fonetyka.cudaniewidły.org/experiment)  
**Code at:**  
[/github.com/mrslacklines/listening\\_experiments](https://github.com/mrslacklines/listening_experiments)

## F<sub>0</sub> 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*).

## F0 inference results

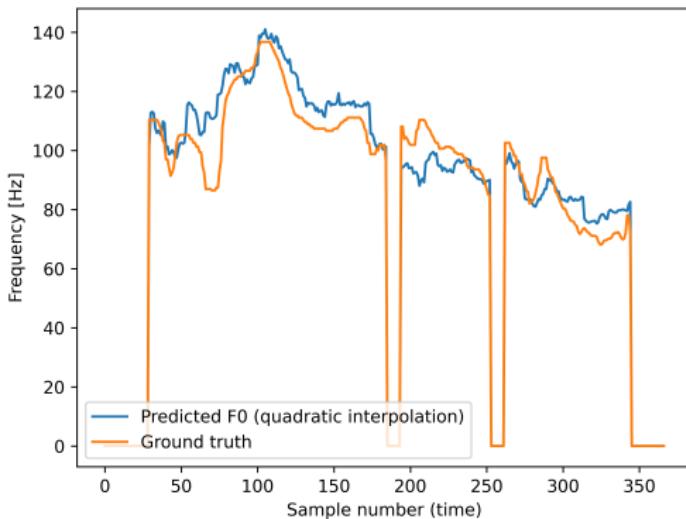


Figure: Result of  $F_0$  prediction for "Powodzenie nie jest gwarantowane" (*Success is not guaranteed*).

## F0 inference results

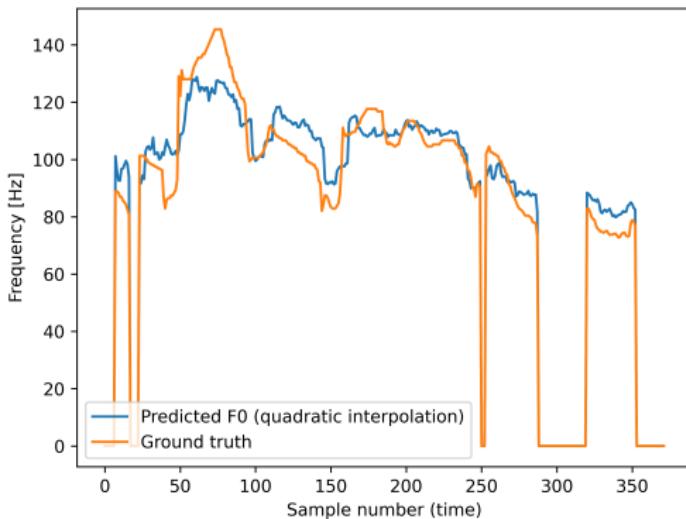


Figure: Result of  $F_0$  prediction for "Gaduła była bardzo nieznośna" (Gabby was very annoying.).

## F0 inference results

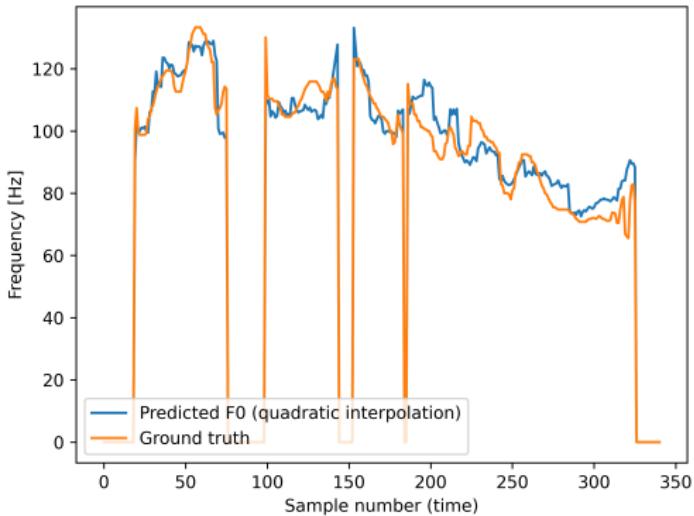
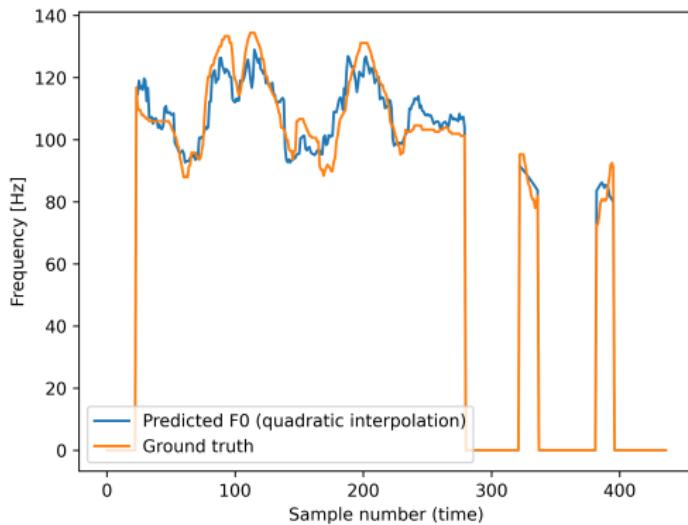


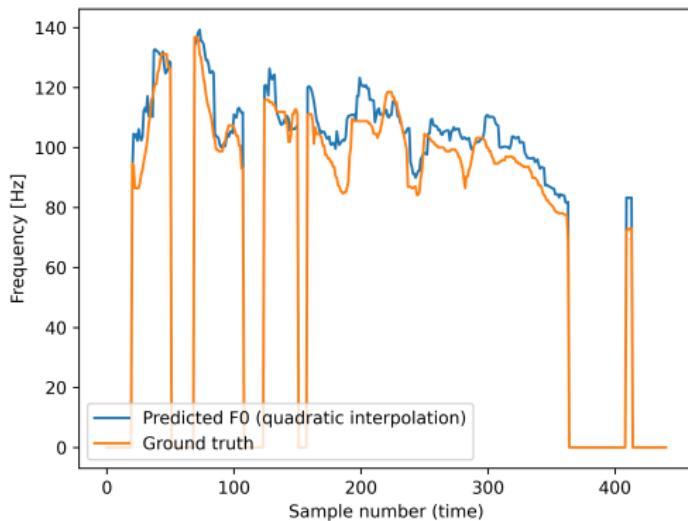
Figure: Result of  $F_0$  prediction for "Może przyniosą też gorzałę" (*Maybe they will bring booze too.*).

## F0 inference results



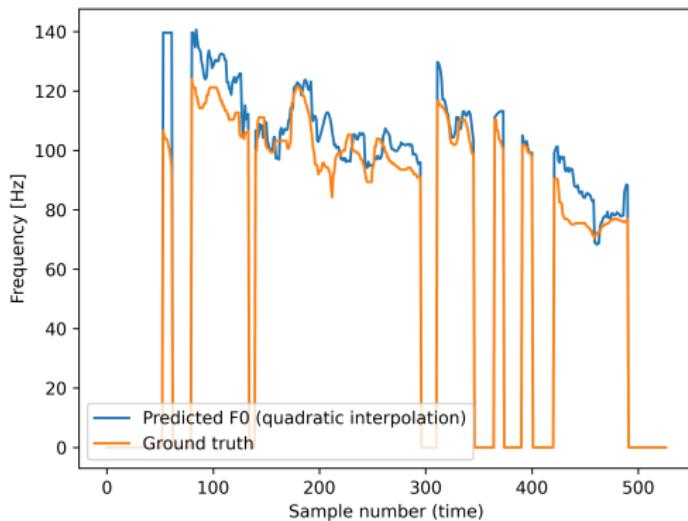
**Figure:** Result of  $F_0$  prediction for "To jest ważna godzina dla nas wszystkich" (*This is an important hour for all of us*).

## F0 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<sub>0</sub> inference results



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

## F0 inference results

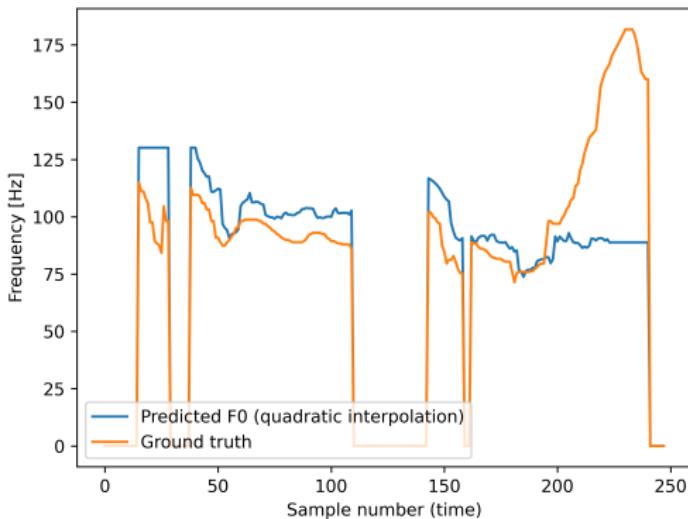


Figure: Result of  $F_0$  prediction for "Czy to był łatwy dobór?" (*Was it an easy choice?*).

## F<sub>0</sub> inference results

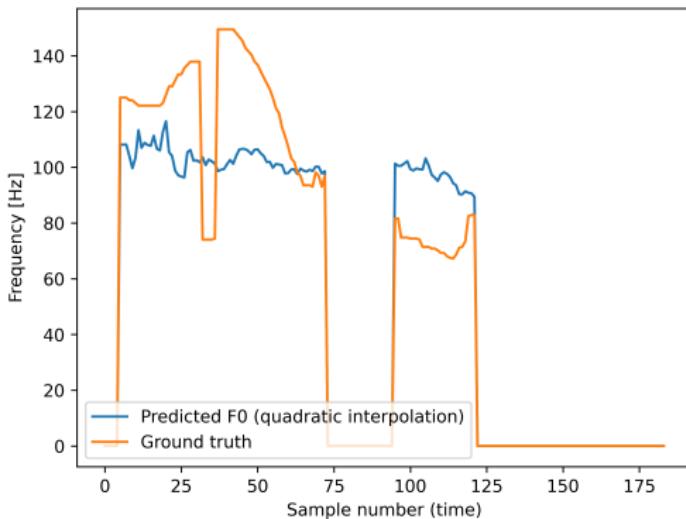


Figure: Result of  $F_0$  prediction for "To Majka" (*This is Majka*).

## F0 inference results

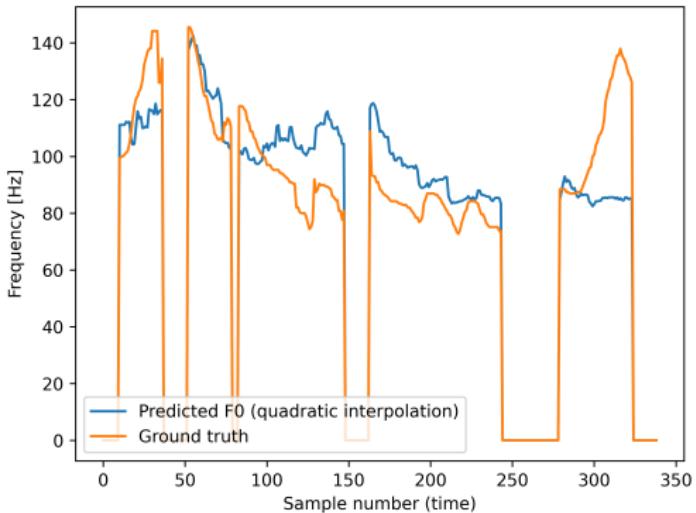


Figure: Result of  $F_0$  prediction for "Na czym polega kandyzacja?" (*How does candying work?*).

## Subjective evaluation results - MOS

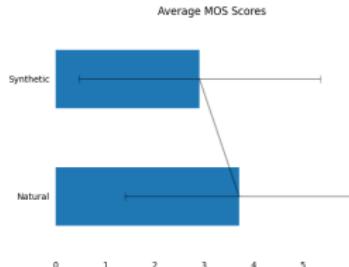


Figure: Mean Opinion Score-based evaluation results.

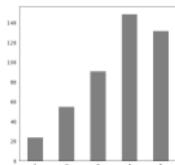


Figure: Mean Opinion Score-based evaluation total numbers of specific scores for natural stimuli.

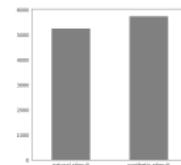


Figure: Mean Opinion Score-based evaluation mean response times for synthetic and natural stimuli.

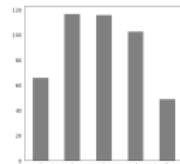


Figure: Mean Opinion Score-based evaluation total numbers of specific scores for synthetic stimuli.

## Subjective evaluation results - MOS

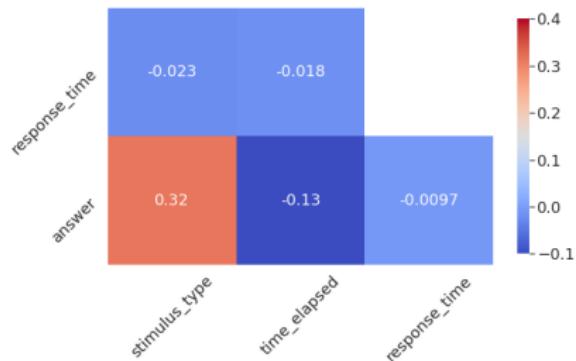
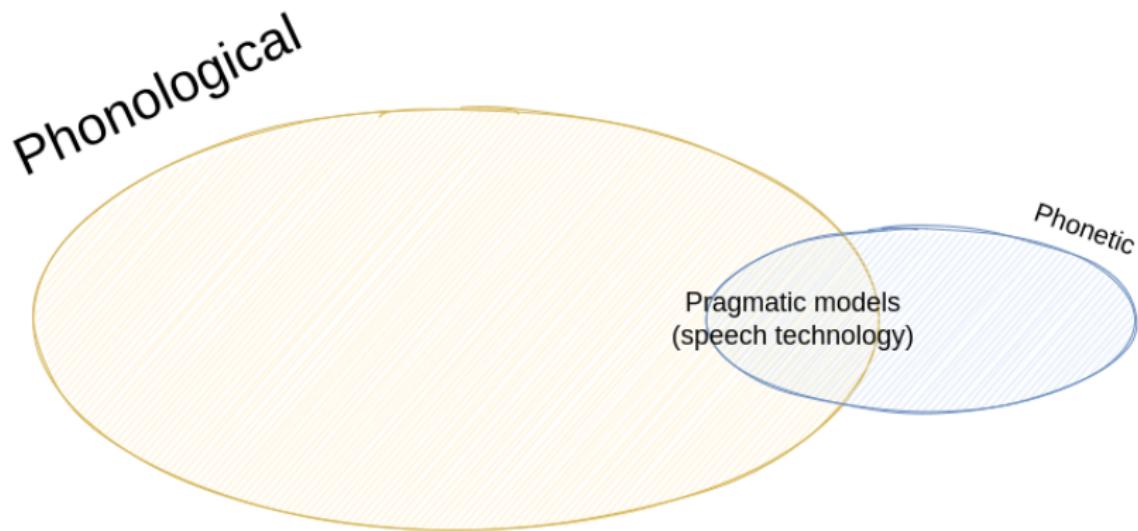


Figure: Mean Opinion Score-based evaluation parameters correlation matrix.

Thank you.

# Intonation models



## Speech synthesis - Wavenet

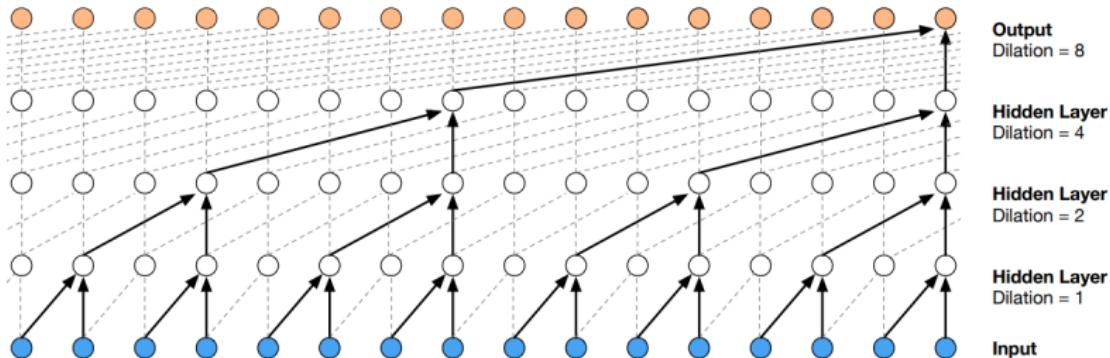
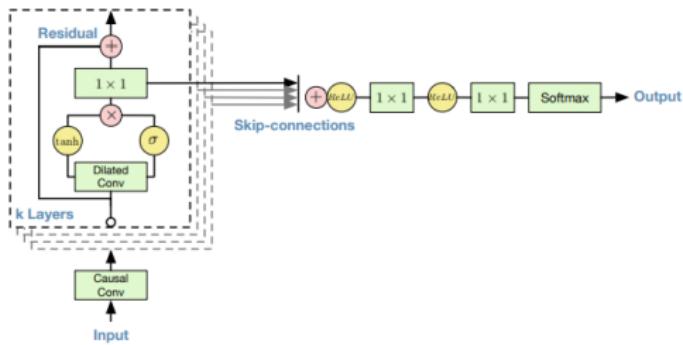


Figure: Dilated causal convolutions. (Adopted from the original WaveNet paper).

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



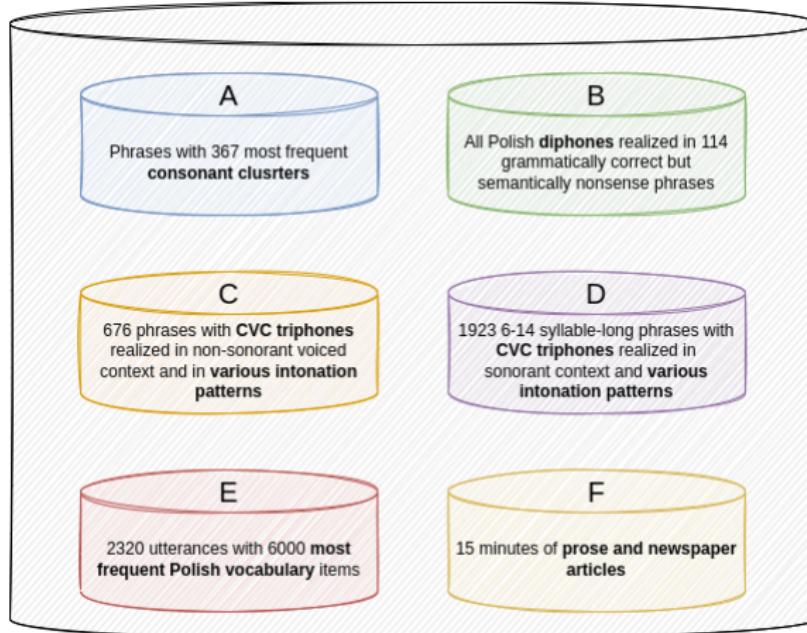
**Figure:** Residual and skip connections from a stack of  $k$  gated convolutional layers (Adopted from the original WaveNet paper).

Gated convolutional layers:

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

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



**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).

### 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.

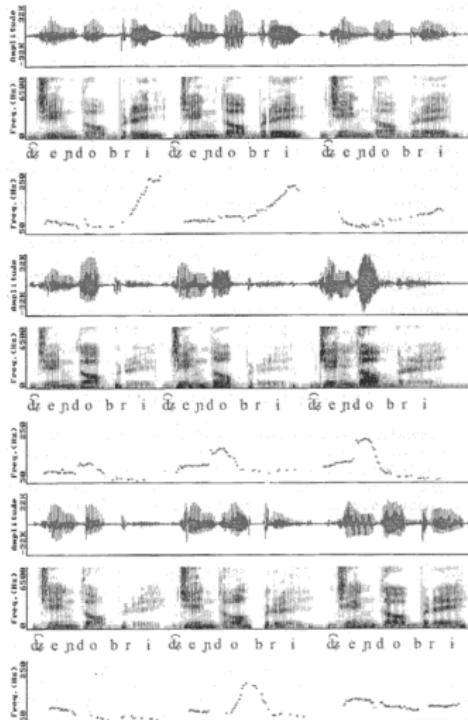
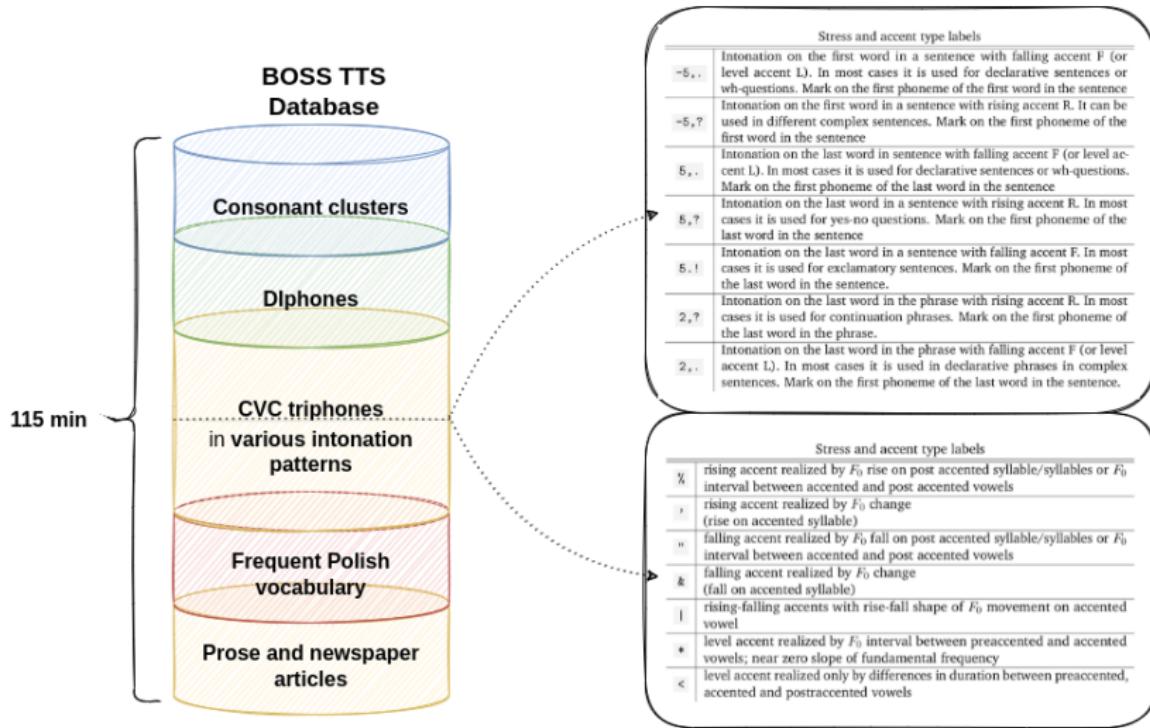


Figure: Acoustic realizations of the 9 different accents. Adopted from (Demenko 1999).

# Dataset

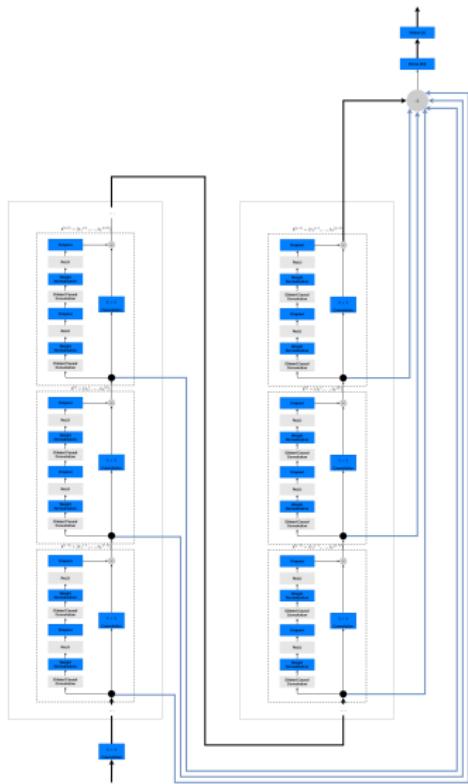


# Feature set

| Question type        | Segments  |  |
|----------------------|---|--|
| Vowel                | {i, y, e, a, o, 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, g̊, p, b, t, d, k, l, g, kl, g̊l,<br>f, v, s, z, x, zl, m, n, rz, x,<br>c, dz, cz, drz, cl, dzl, n, n,<br>x̄l, ng, l, r, v, w̄, j, jj̄} | Previous/current/next syllable is stressed   |
| Stop                 | {g, g̊, p, b, t, d, k, l, g̊l}  | Previous/current/next syllable has accent X<br>(where X is one of the ToBI accents described above)  |
| Nasal                | {vn, jj̄, n, m, n̄, nḡ}  | Number of preceding/succeeding segments in the next syllable is equal to/less than or equal to 0-7   |
| Frikative            | {t̄, v, s, xl, zl, nz, rz, x̄l}   | Forward/backward position of the current syllable in current word is equal to/less than or equal to 0-7  |
| Front                | {e, ɔ, y, t̄, v, p, b, m, v, w̄}  | Forward/backward position of the current syllable in current phrase is equal to/less than or equal to 0-20   |
| General              | {[schwa], n, t̄, v, d, n, xl, s, x, xl,<br>n, r, t̄, t, d, zz, rz, cz, drz,<br>c, dn, cl, dz̄}  | Number of stressed syllables before/after the current syllable in current phrase is equal to/less than or equal to 0-12  |
| Back                 | {o, u, x, g, xl, gl, ng, x, gx̄}  | 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, ɔ, 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, u}  | Number of syllables from previous/next accented syllable is equal to/less than or equal to 0-16  |
| High Vowel           | {i, y, u}   | 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, u}  | 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, e, i, y}  | Number of words from previous/next content word is equal to/less than or equal to 0-5  |
| Xlowl (e.g. Alvherr) | {i, y, e, a, o, u, [schwa]}   | Number of syllables in the previous/current/next phrase is equal to/less than or equal to 0-20   |
| Unvoiced Consonant   | {g̊, g̊l, p, t, k, kl, f, v, s, rz,<br>x, c, cz, cl̄}   | Number of words in the previous/current/next phrase is equal to/less than or equal to 0-15   |
| Voiceless Consonant  | {b, d, g, gl̄, v, z, xl, rz, dz̄,<br>drz, dzl̄, n, m, nl̄, ng, l, r, v, w̄, j, jj̄}   | Forward/backward position of the current phrase in the utterance is equal to/less than or equal to 0-4   |
| Front Consonant      | {t̄, v, f, p, b, m, v, w̄}  | Number of words in the utterance is equal to/less than or equal to 0-28  |
| General Consonant    | {t̄, d, s, xl, zl, n, r, l,<br>t, d, zz,<br>rz, cz, drz, c, dz, cl, dz̄}  | Number of words in the utterance is equal to/less than or equal to 0-13  |
| Back Consonant       | {g, g̊, k, g̊l, nḡ, x̄l}   | Number of phrases in the utterance is equal to/less than or equal to 0-4   |
| Fortis Consonant     | {g̊, cz, f, k, p, s, nz, t, cl,<br>c, xl̄}  |  |
| Lensis Consonant     | {drz, v, g, b, rz, x, d, dn̄l̄,<br>dz, gl̄, zl̄}  |  |
| Neither F or L       | {n, x, xl̄, nḡ, l, r, v, w̄,<br>j, jj̄}  |  |
| Voiced Stop          | {b, d, g̊}  |  |
| Unvoiced Stop        | {p, t, k, g̊l}  |  |
| Front Stop           | {b, p}  |  |
| General Stop         | {d, t̄}   |  |
| Back Stop            | {g̊, x, gl̄}  |  |
| Voiced Frikative     | {v, z, xl̄, rz̄}  |  |
| Unvoiced Frikative   | {t̄, s, xl̄, nz̄, x̄l̄}   |  |
| Front Frikative      | {t̄, v}   |  |

Figure: Segmental features for a quintphone-wide context.

Figure: Non-segmental features.



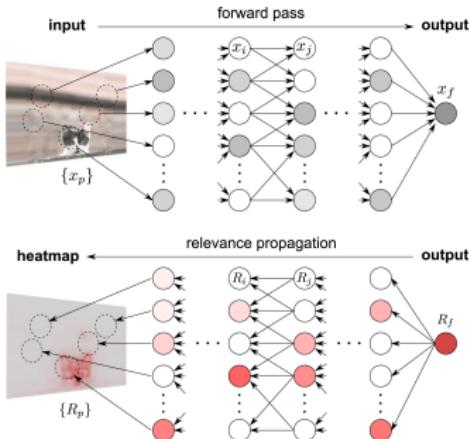


Figure: Computational flow of deep Taylor decomposition. (Adopted from Montavon 2017).

The relevance in this framework can be defined as:

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