

# 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

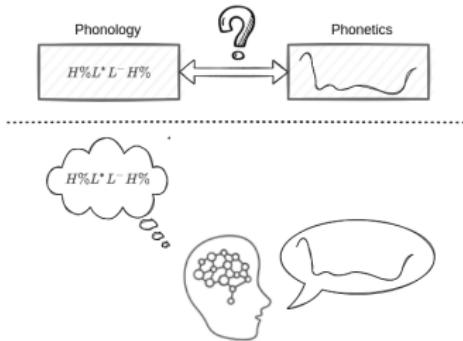
*"The tonal structure of speech expressed by the melody produced by our larynx. It has a phonetic aspect, the fundamental frequency (F0), and a grammatical (phonological) aspect." (Féry 2016)*

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)

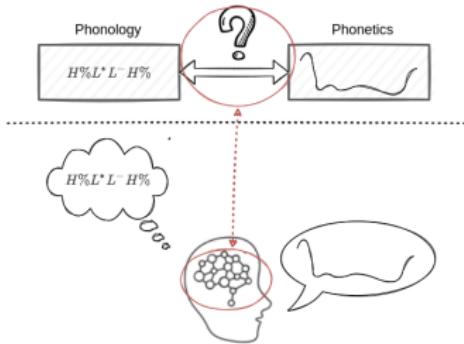
# Motivation

## 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.
- ✓ Need for a physicalist (neurobiological) model.
- ✓ Modern statistical-parametric speech synthesis provides a framework for experimentation and evaluation of such models.



# 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  $F0$ .
- ✓ 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.*

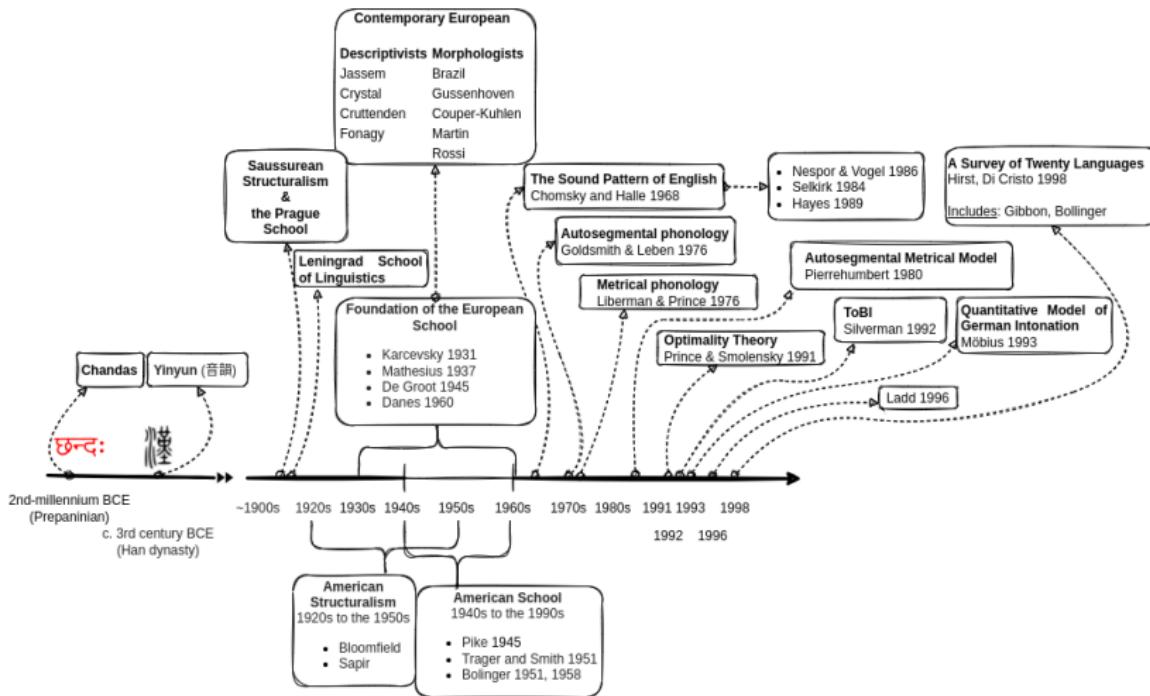
**HYPOTHESIS 2:** *The biologically-inspired Deep Temporal Convolutional Network can be an effective model of these mappings and hence of Polish neutral read speech intonation in the context of statistical-parametric speech synthesis.*

**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.*

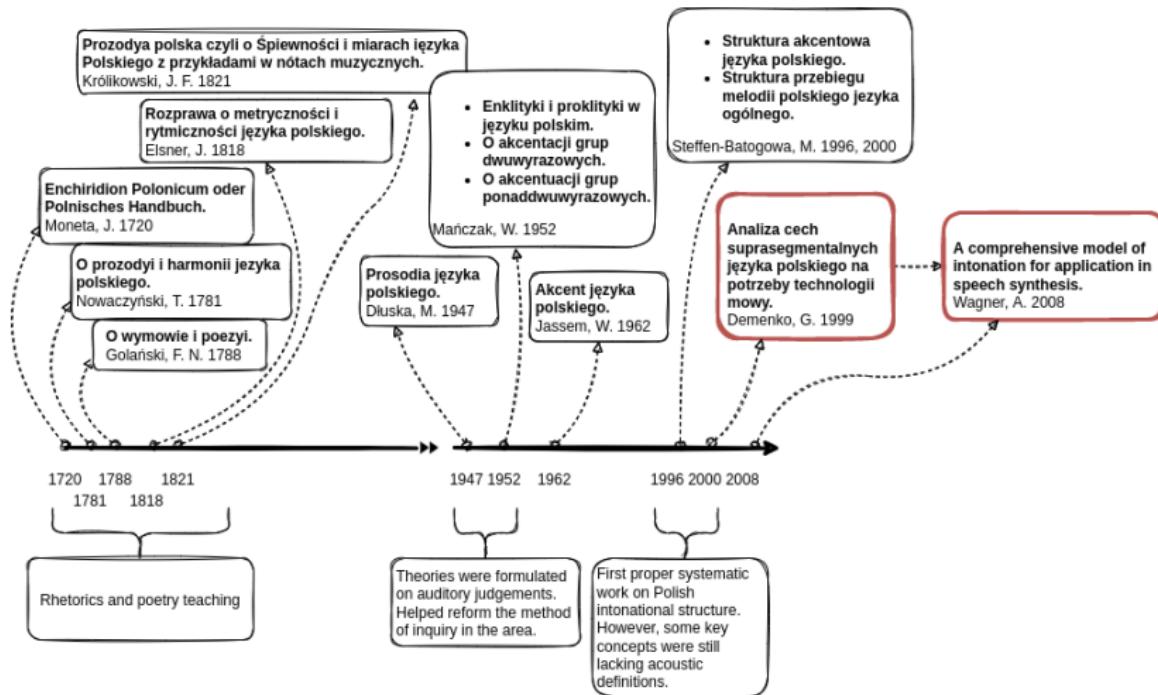
## Contributory Methodological Hypothesis

**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.*

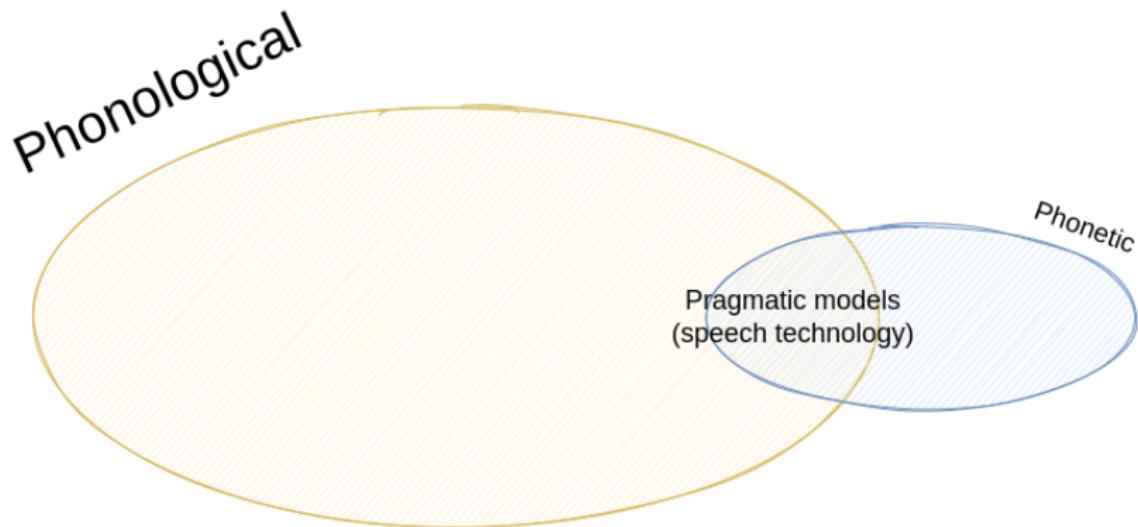
## Background



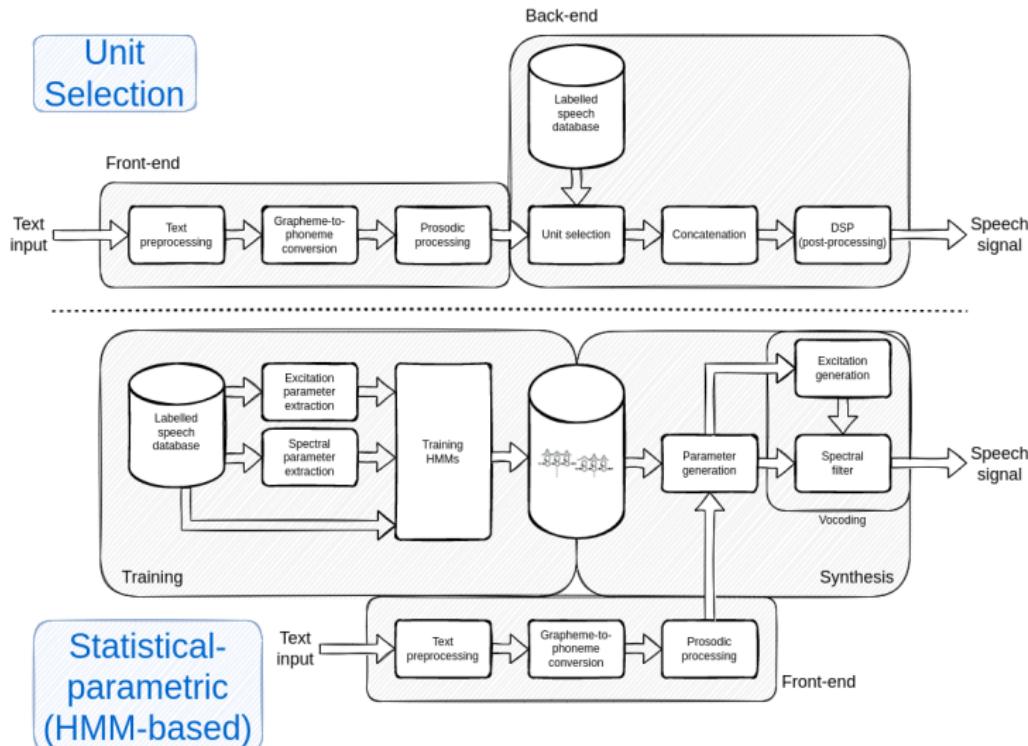
# Background



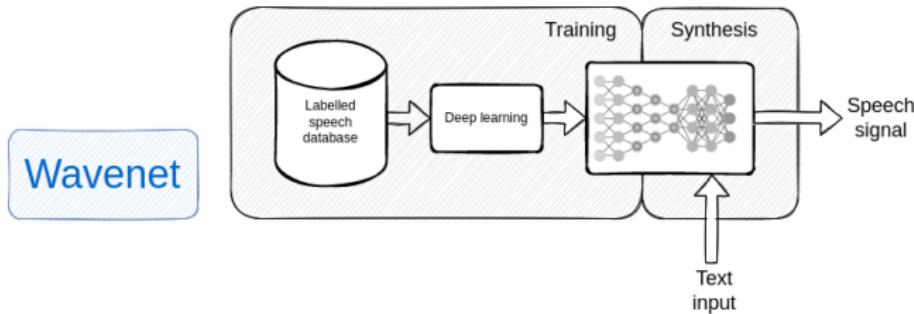
# Intonation models



# Speech synthesis

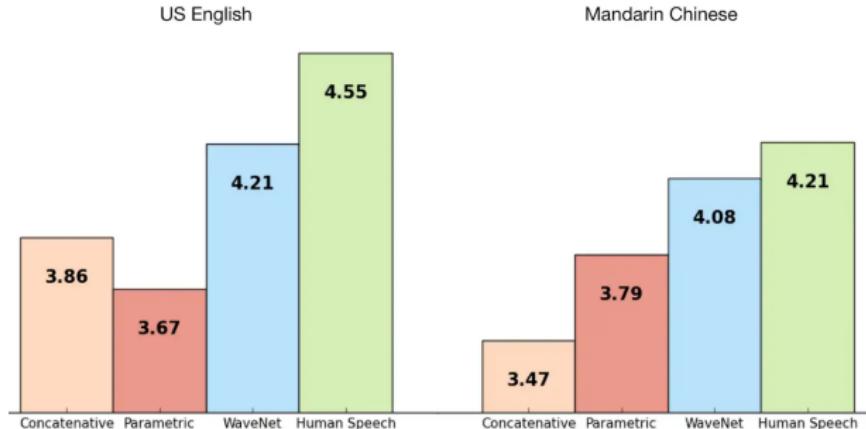


## Speech synthesis - Wavenet



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

## Speech synthesis - Wavenet



The idea behind WaveNet is based on a PixelCNN Image Generator also developed by van den Oord (2016) at Google's DeepMind.

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



# Neurobiological foundations

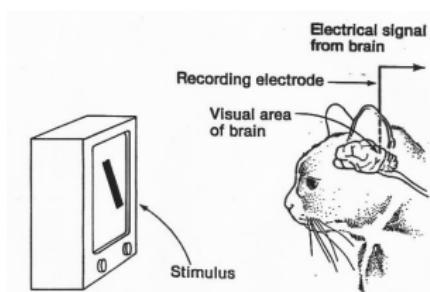


Figure: Famous Hubel and Wiesel cat experiment. (adopted from Hubel and Wiesel 1959).

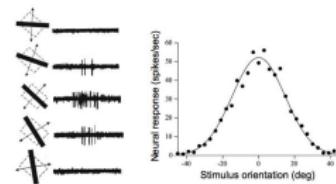


Figure: Neural response of simple cells. (adopted from Hubel and Wiesel 1968).

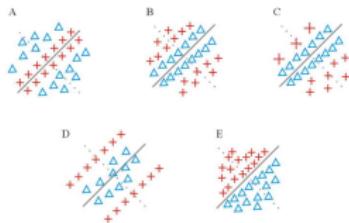


Figure: Simple receptive fields. (adopted from Hubel and Wiesel 1962).

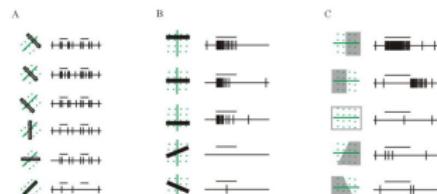


Figure: Three different types of complex receptive fields. (adopted from Hubel and Wiesel 1962).

## Neurobiological foundations

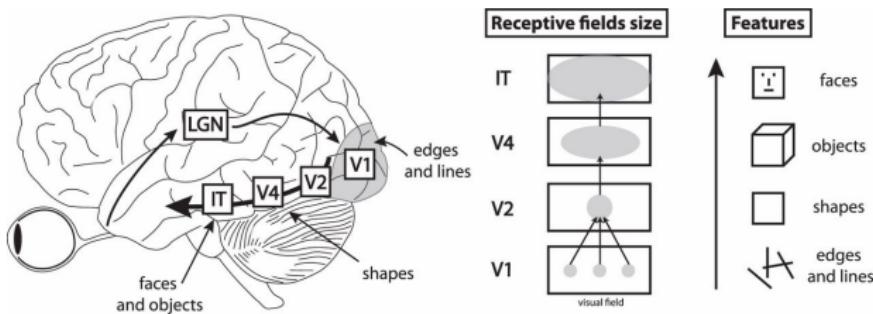


Figure: Hierarchical, feedforward visual processing in human brain. (Adopted from Manassi et al. 2013)

# Neurobiological foundations

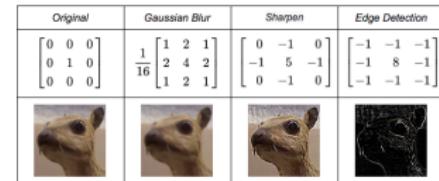
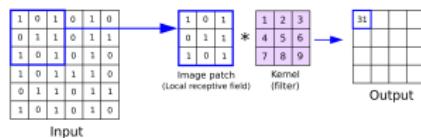


Figure: Example of a 2-dimensional matrix convolution.

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

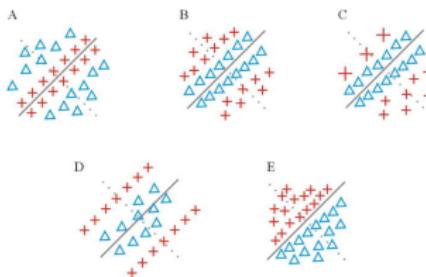


Figure: Simple receptive fields. (adopted from Hubel and Wiesel 1962).

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

# Neurobiological foundations

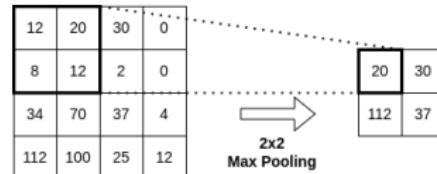


Figure: Example of a 2x2 max pooling matrix operation.

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

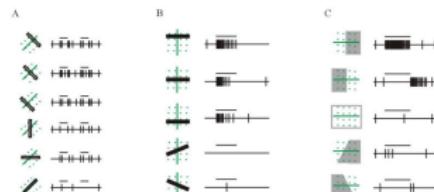


Figure: Three different types of complex receptive fields.  
(adopted from Hubel and Wiesel 1962).

## Neurobiological foundations

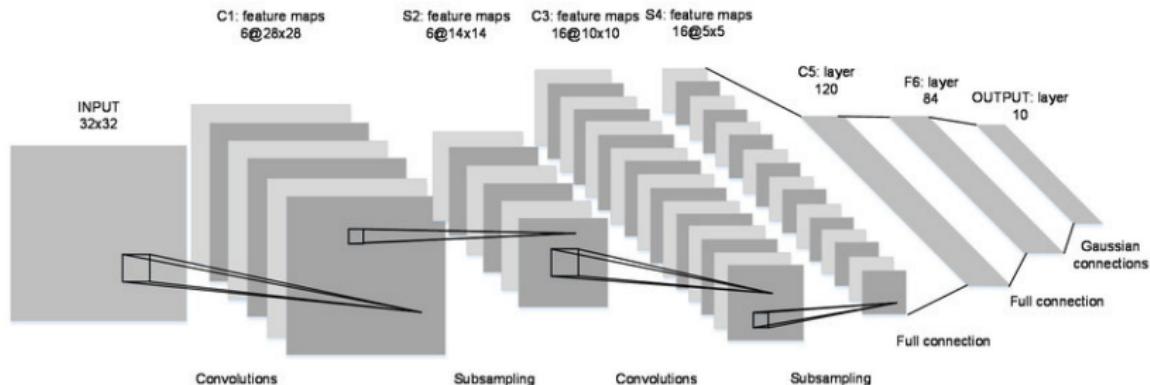
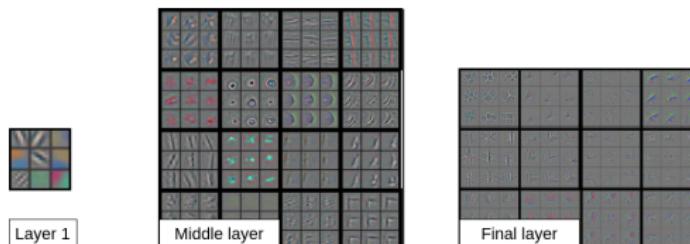


Figure: Image recognition Convolutional Neural Network (LeNet-5). (Adopted from LeCun et al. 1989)



## Neural auditory processing

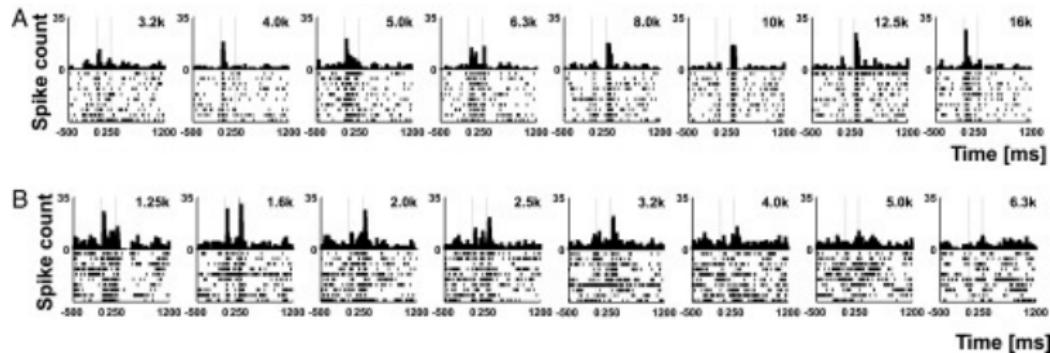


Figure: Responses of single neurons in primary auditory cortex (A1) of rhesus monkeys to band-passed noise (BPN) bursts centered at particular frequencies.. (Adopted from Tian et al. 2013)

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

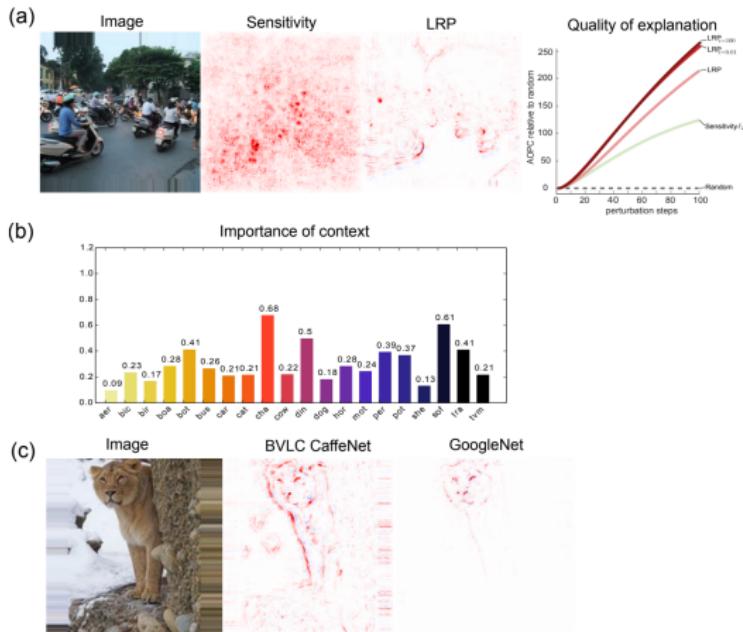


Figure: Results of sensitivity-based and relevance-based explainability methods. (Based on Samek 2016)

# Deep Taylor Decomposition

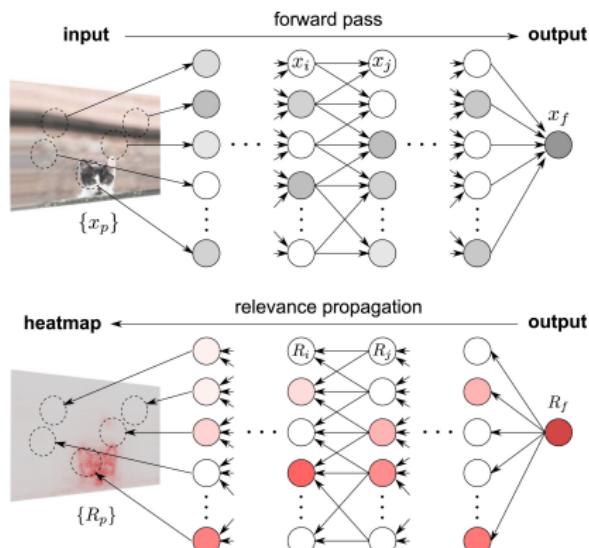


Figure: Computational flow of deep Taylor decomposition. (Adopted from Montavon 2017).

The relevance in this framework can be defined as:

$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k \quad (1)$$

Thank you.