

VIBE PROJECT

Virtual Biomedical and STEM/STEAM Education

2021-1-HU01-KA220-HED-000032251



**Funded by
the European Union**



**PÉCSI TUDOMÁNYEGYETEM
UNIVERSITY OF PÉCS**

U.PORTO

**Politechnika
Śląska**

DEX
innovation centre



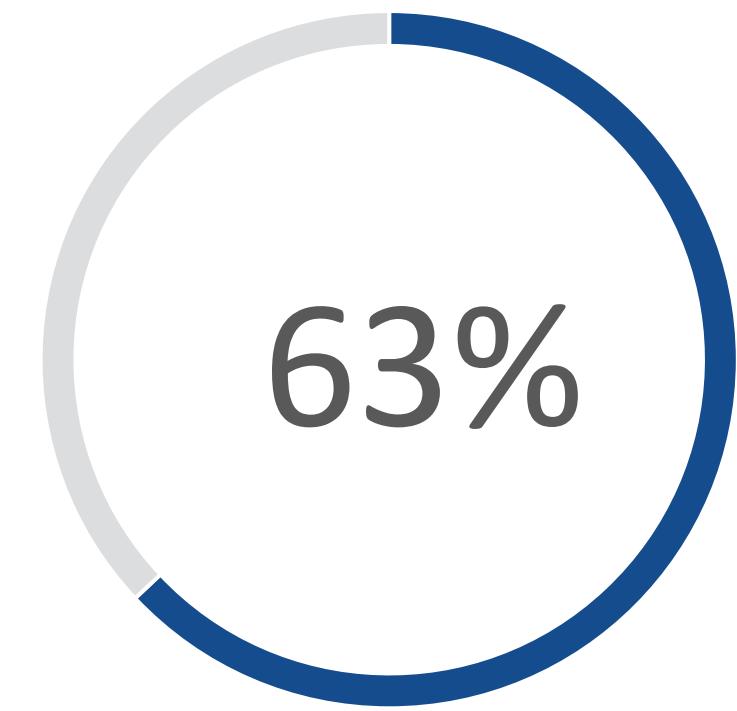
PATTERN RECOGNITION IN BIOMEDICAL ENGINEERING

ANALYSIS OF SPATIAL ACOUSTIC
MODELS OF SIBILANT SOUNDS IN
SIGMATISM DIAGNOSIS

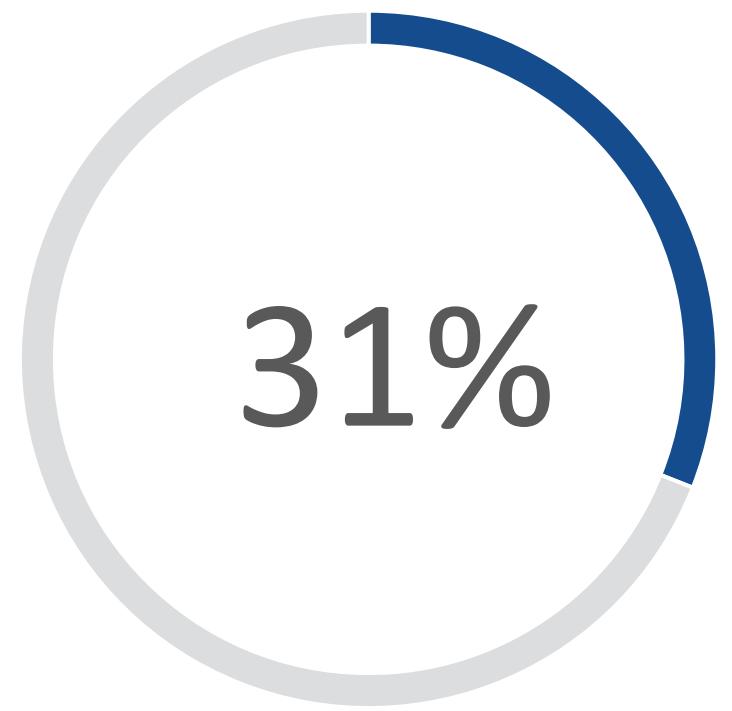


Speech disorders in Poland

3



Dyslalia among children
in Poland,



including sigmatism
(lisp).



Sigmatism

A lisp (sigmatism) is a speech impairment in which a person misarticulates sibilants (in Polish):

/s,z,ts,dz/

/ʃ,ʒ,ʈʃ,ɖʒ/

/tʂ,ʐ,ʈʂ,ɖʐ/

4



www.Thinkstockphotos.co
m

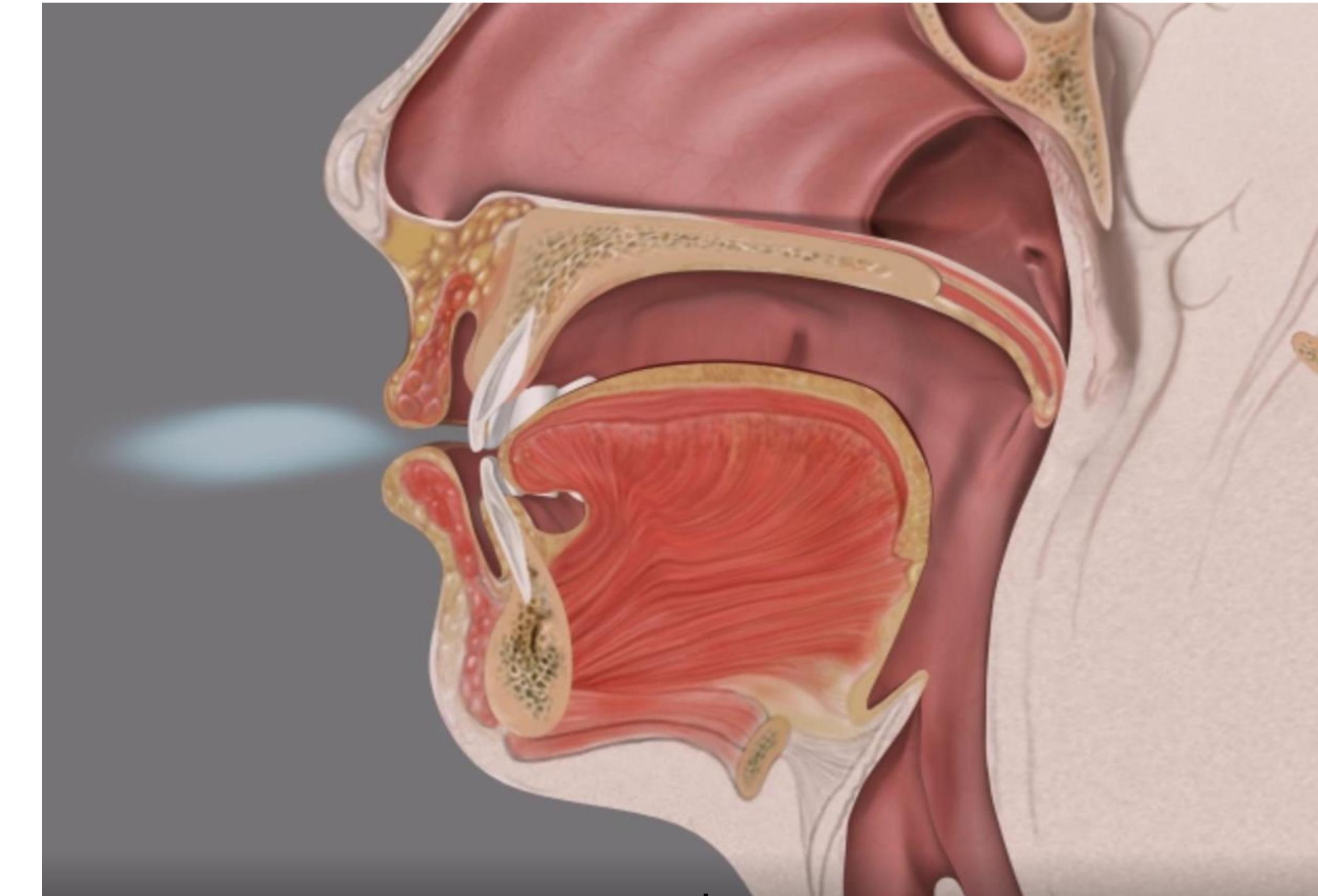
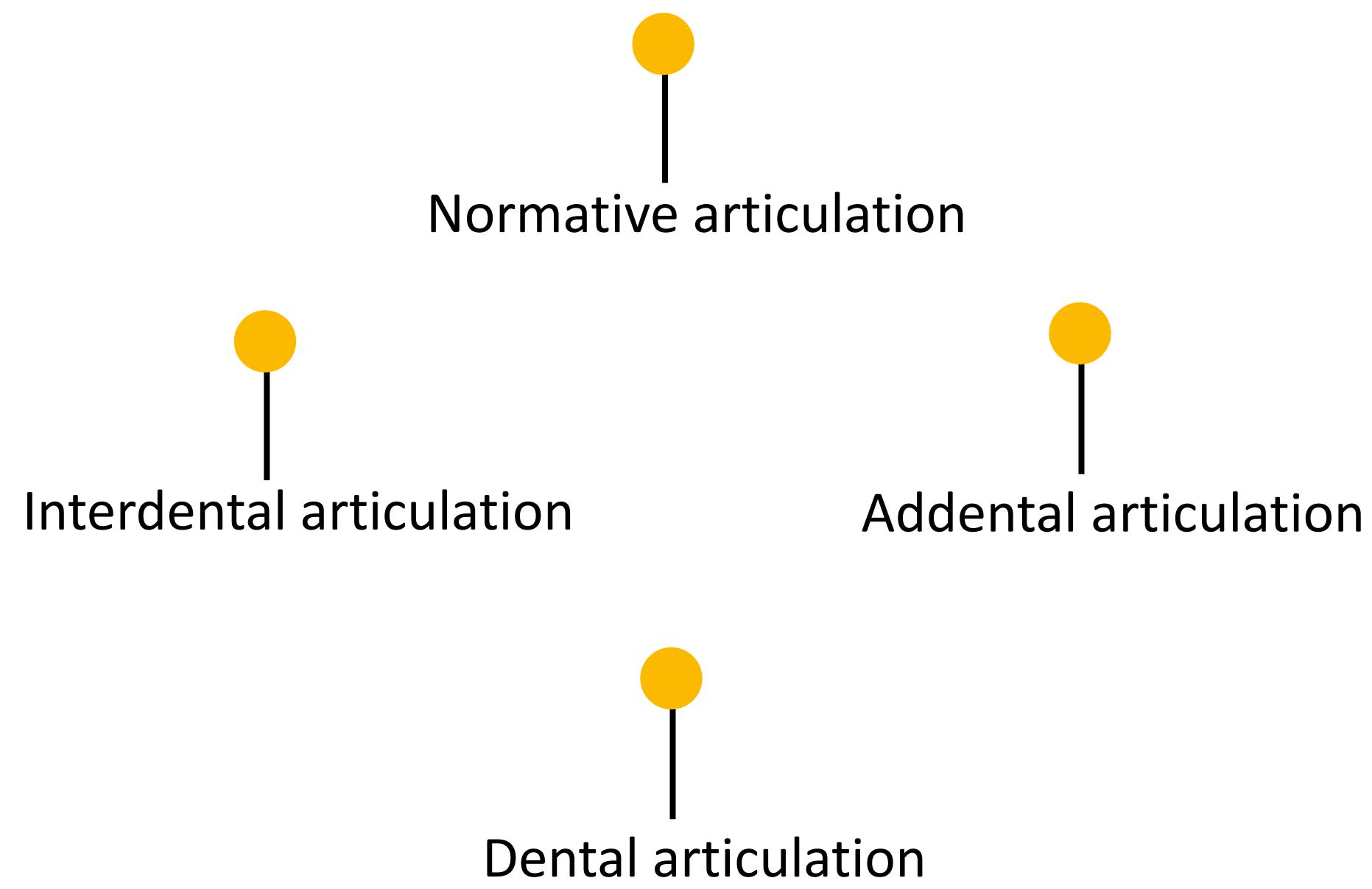
Causes

- abnormalities in the structure and functioning of the articulation organs, anatomical changes,
- low motor efficiency of the articulation organs,
- incorrect structure and functioning of the hearing organ.



Types of articulations*

5



www.komlogo.pl/fonem

Normative realization of the sound /s/

*due to the location of the apex in the oral cavity

The aim of the research conducted

MAIN AIM

DEVELOPMENT OF SPATIAL ACOUSTIC MODELS OF SIBILANTS AND THEIR
USE IN THE CLASSIFICATION OF NORMATIVE AND NON-NORMATIVE
SIBILANT ARTICULATIONS

6

SECONDARY AIM

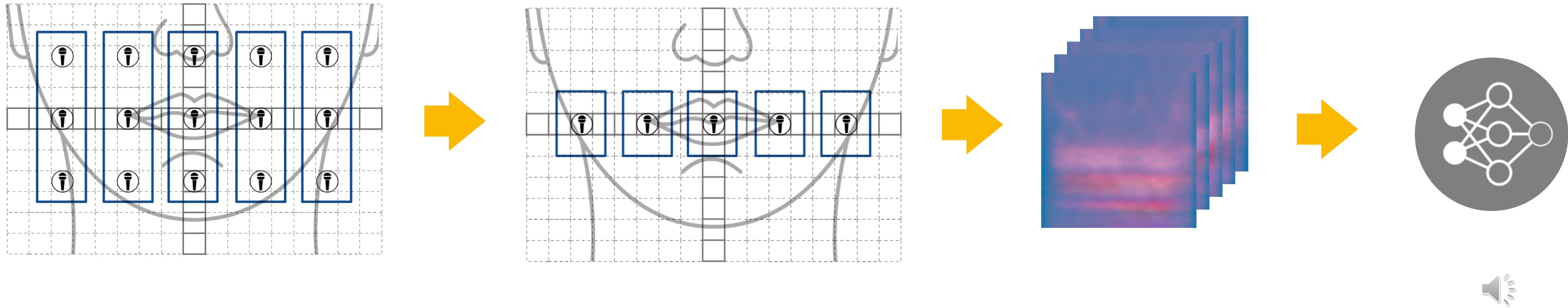
DESIGN AND CONSTRUCTION OF A MEASUREMENT STATION ENABLING
SPATIAL ACQUISITION OF THE SPEECH SIGNAL



Research thesis

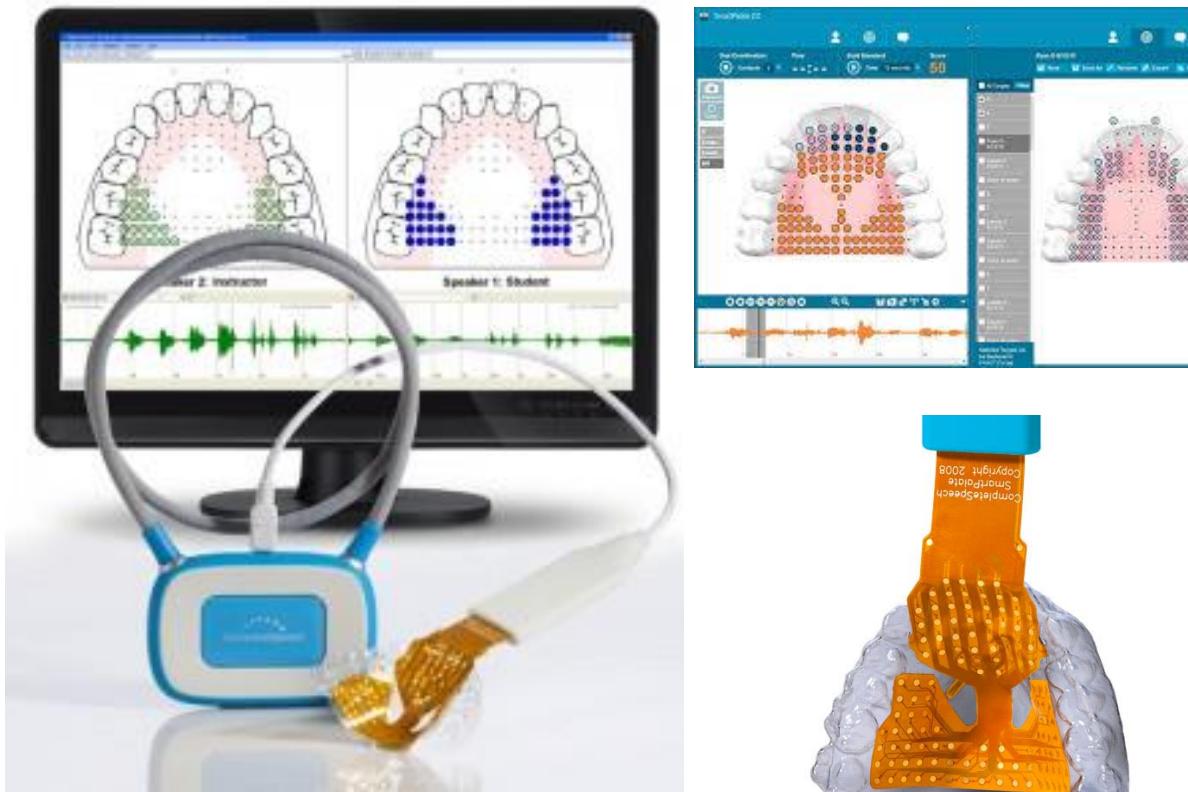
THE USE OF SPATIAL SPEECH SIGNAL PROCESSING METHODOLOGY ENABLES THE RECOGNITION OF SIGMATISM AND DETERMINATION OF ITS TYPE

7



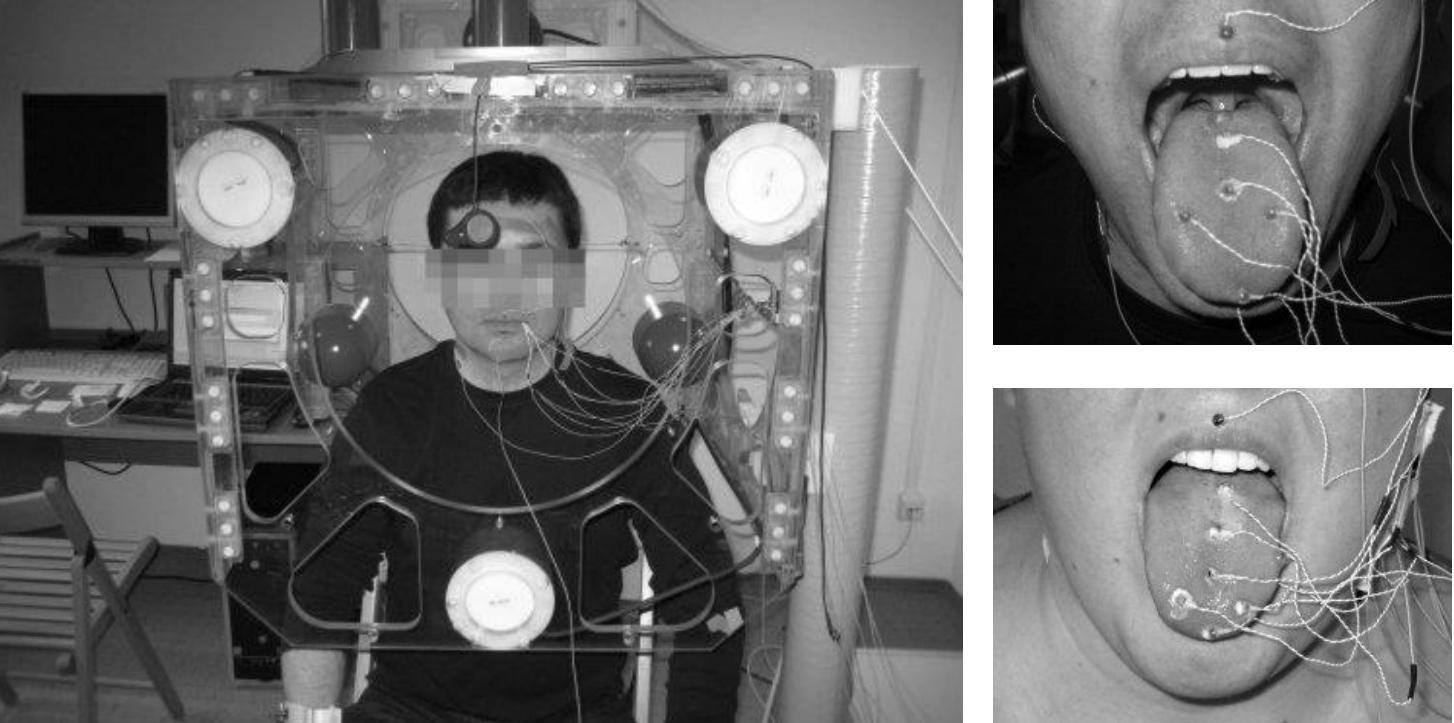
Measurement methods

Electropalatography

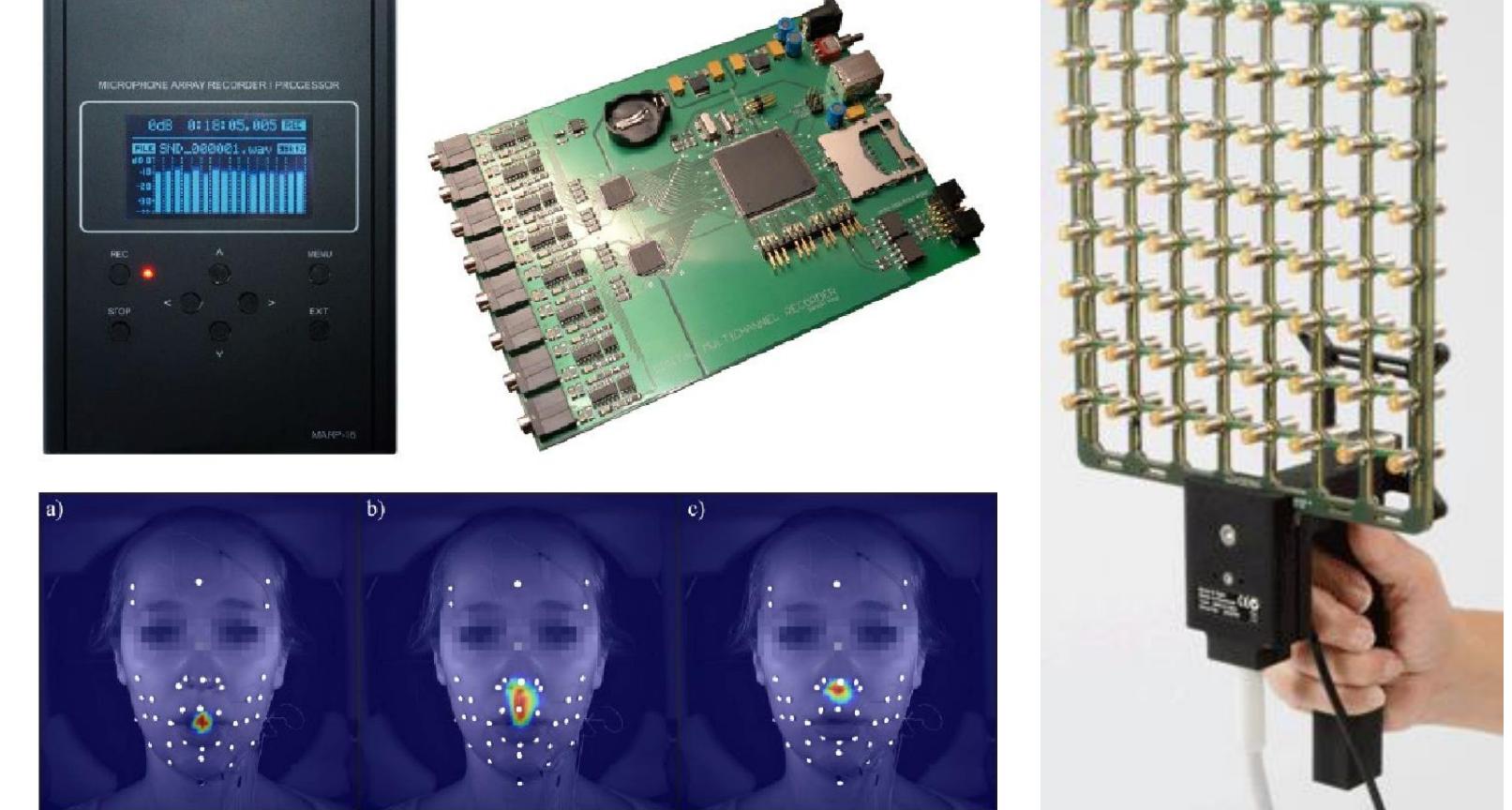


www.completespeech.com/smartpalate/.pl

Electromagnetic articulography



Acoustic holography



Acoustic Field Distribution in Speech with the use of
the Microphone Array

www.bksv.com

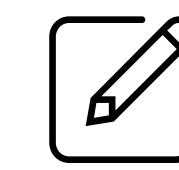


Measuring station

9



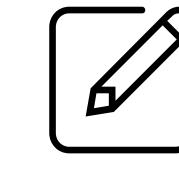
multi-channel, spatial and repeatable sound recording,



visually attractive for children



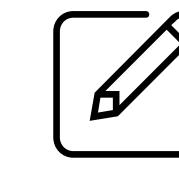
does not require a specialist measuring room



no impact on the articulation and phonation process



access to raw data

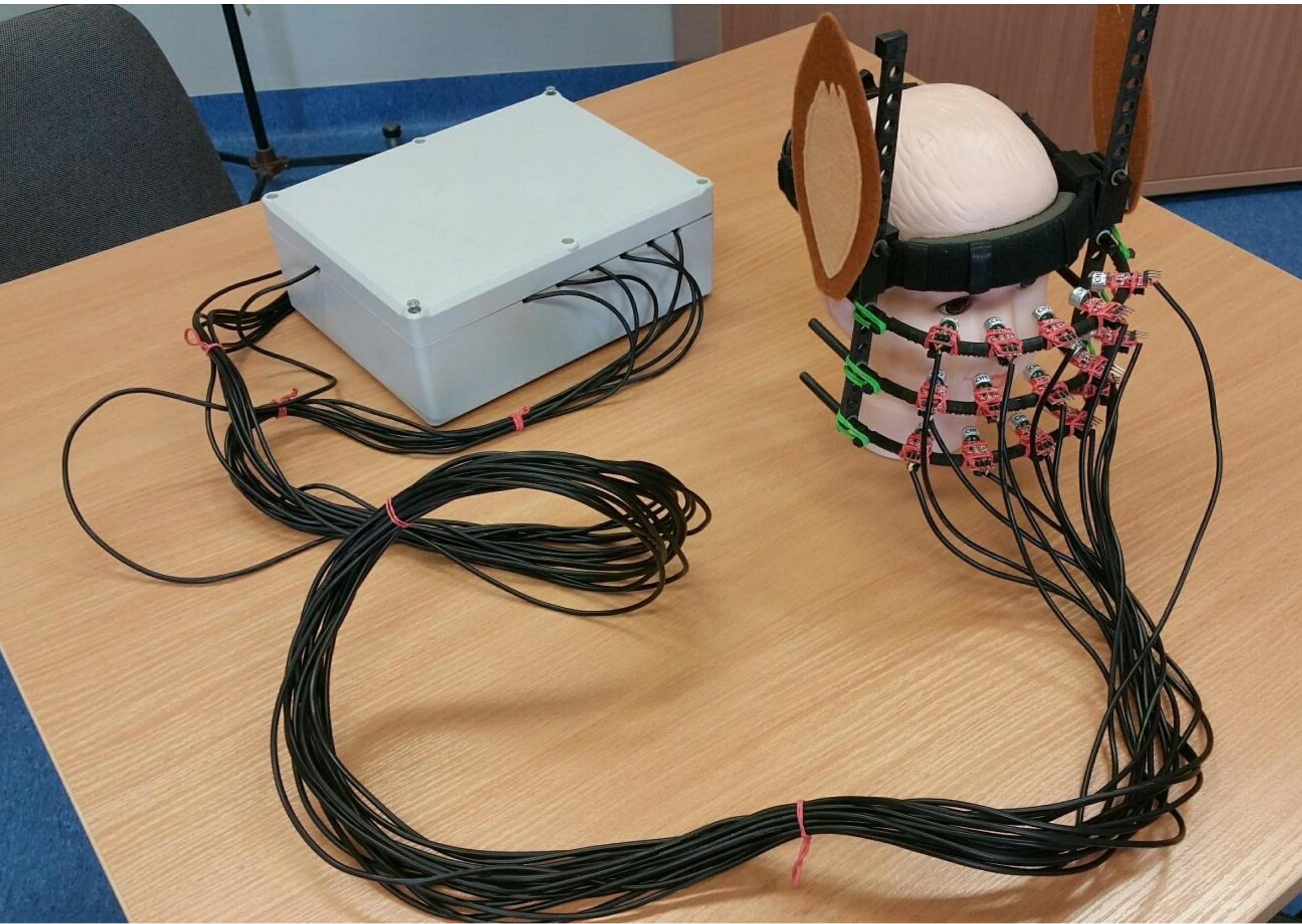


painless (non-invasive) data recording



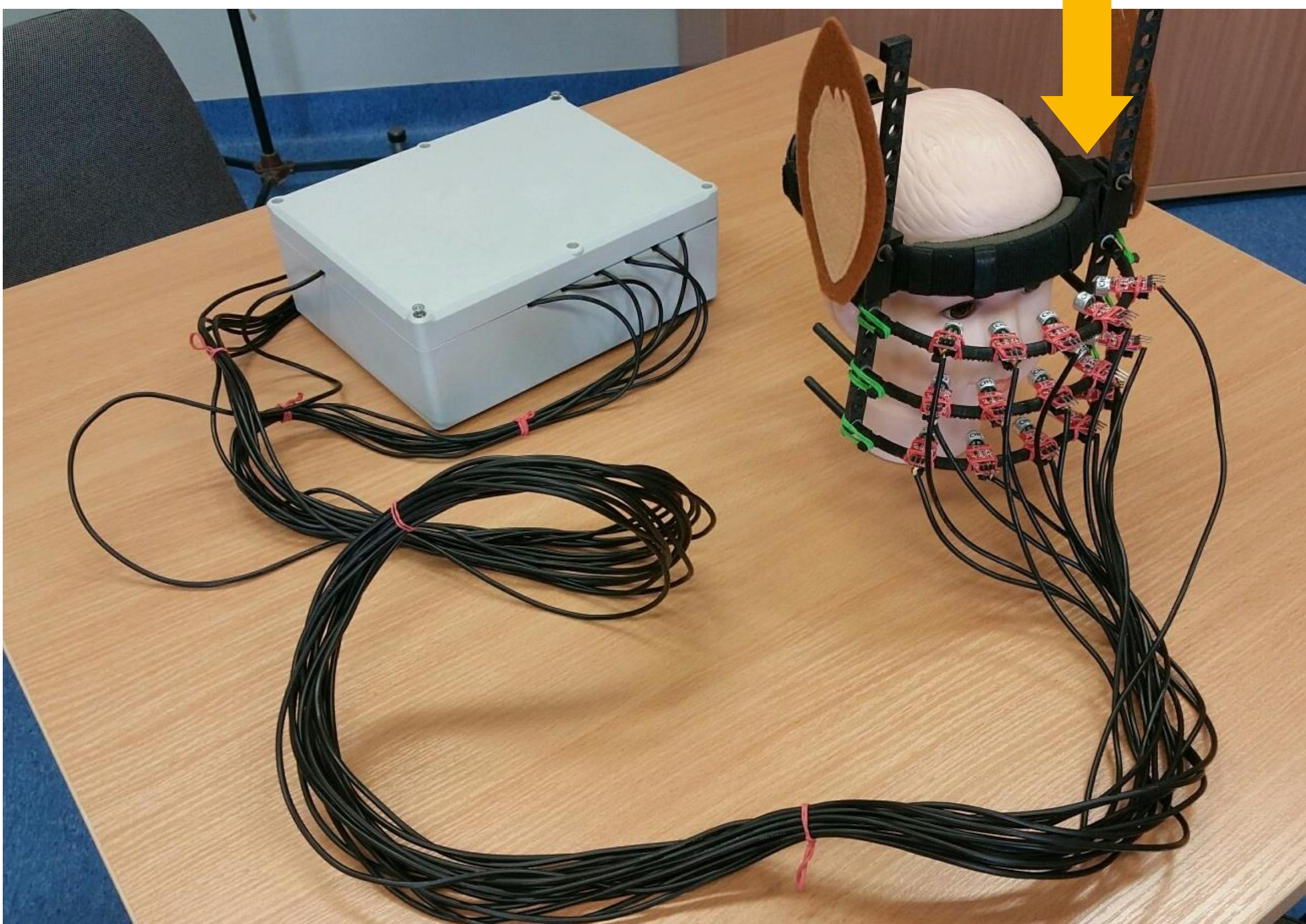
Measuring station

10



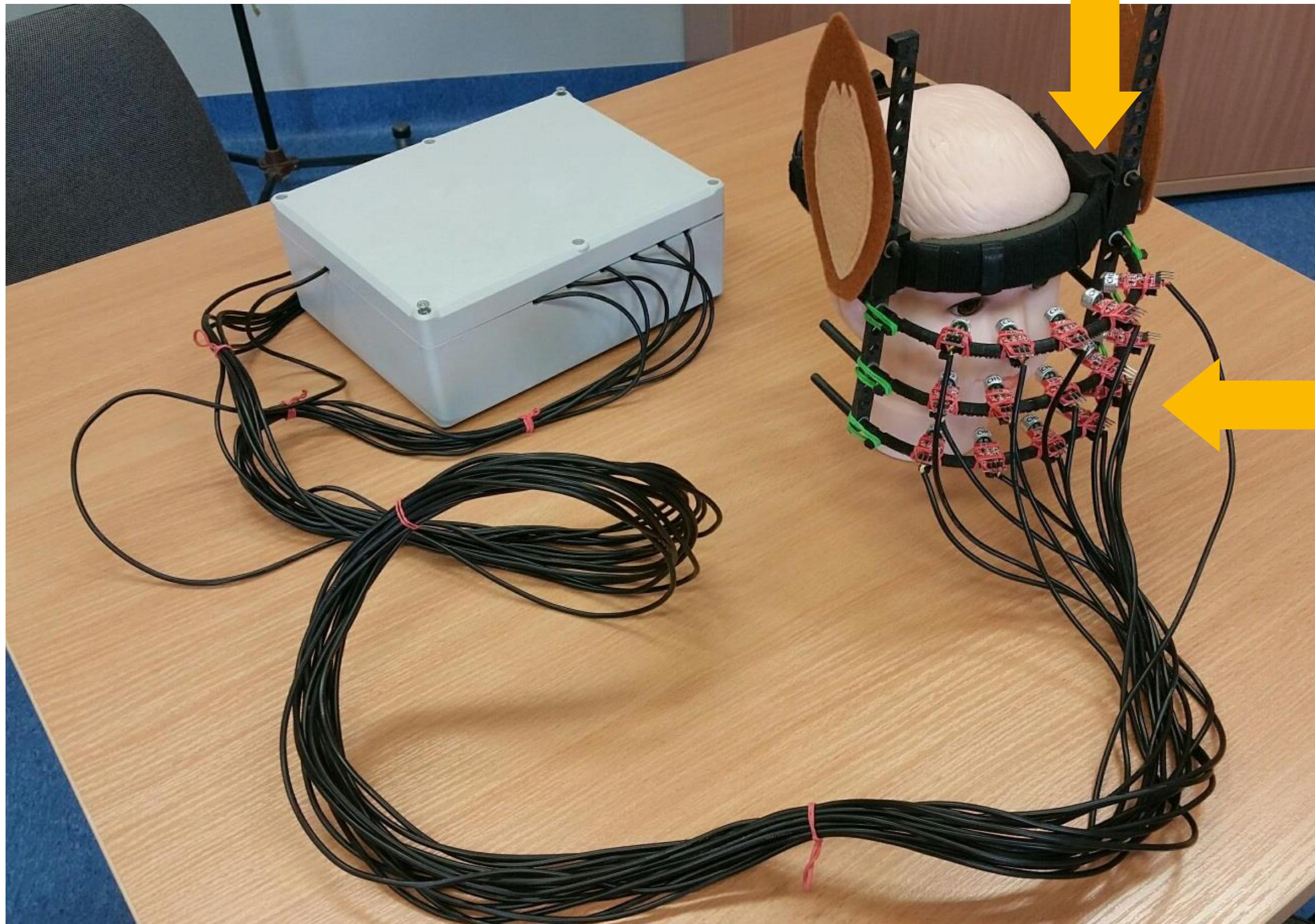
BUILDING A PROTOTYPE

Measuring station



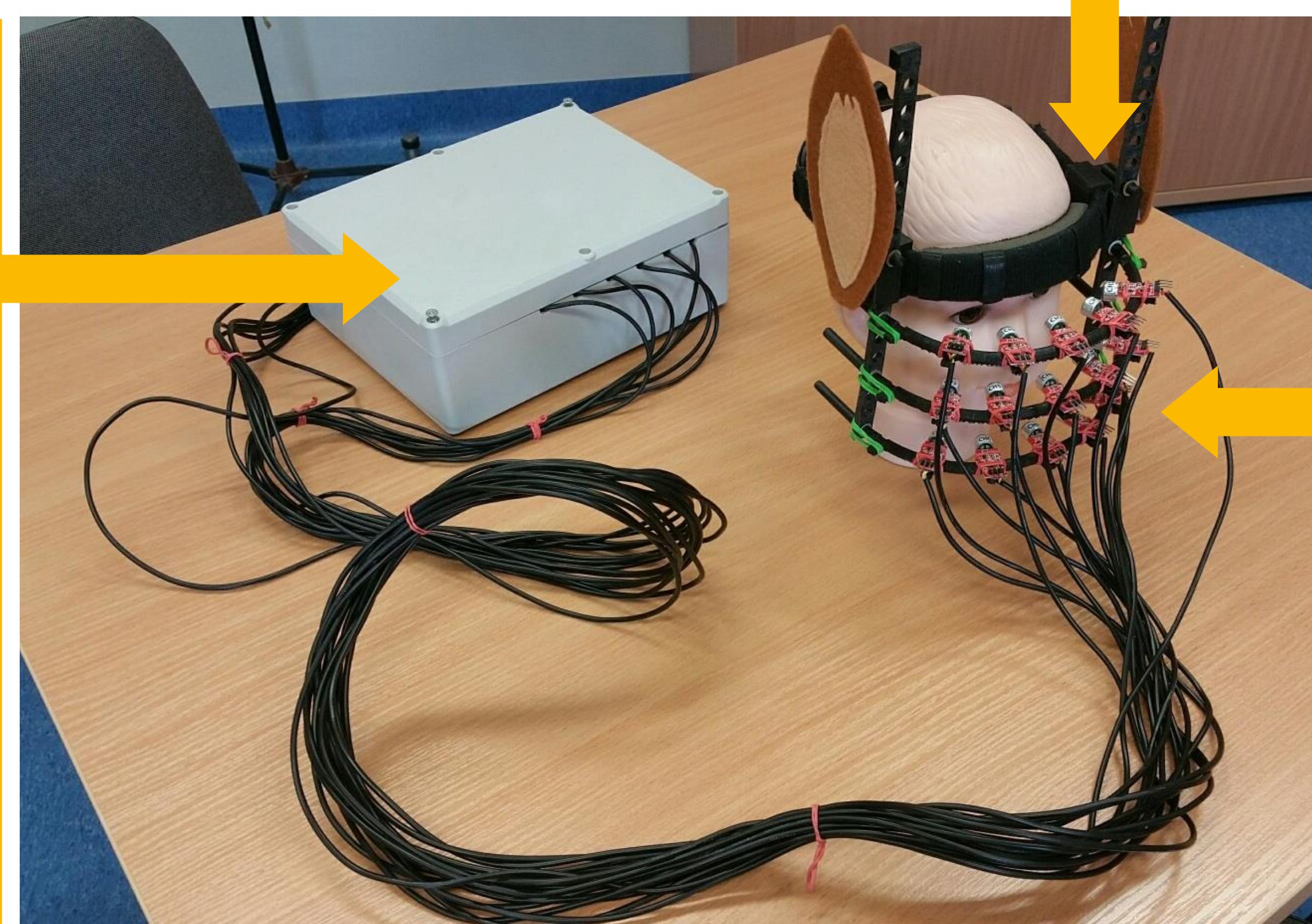
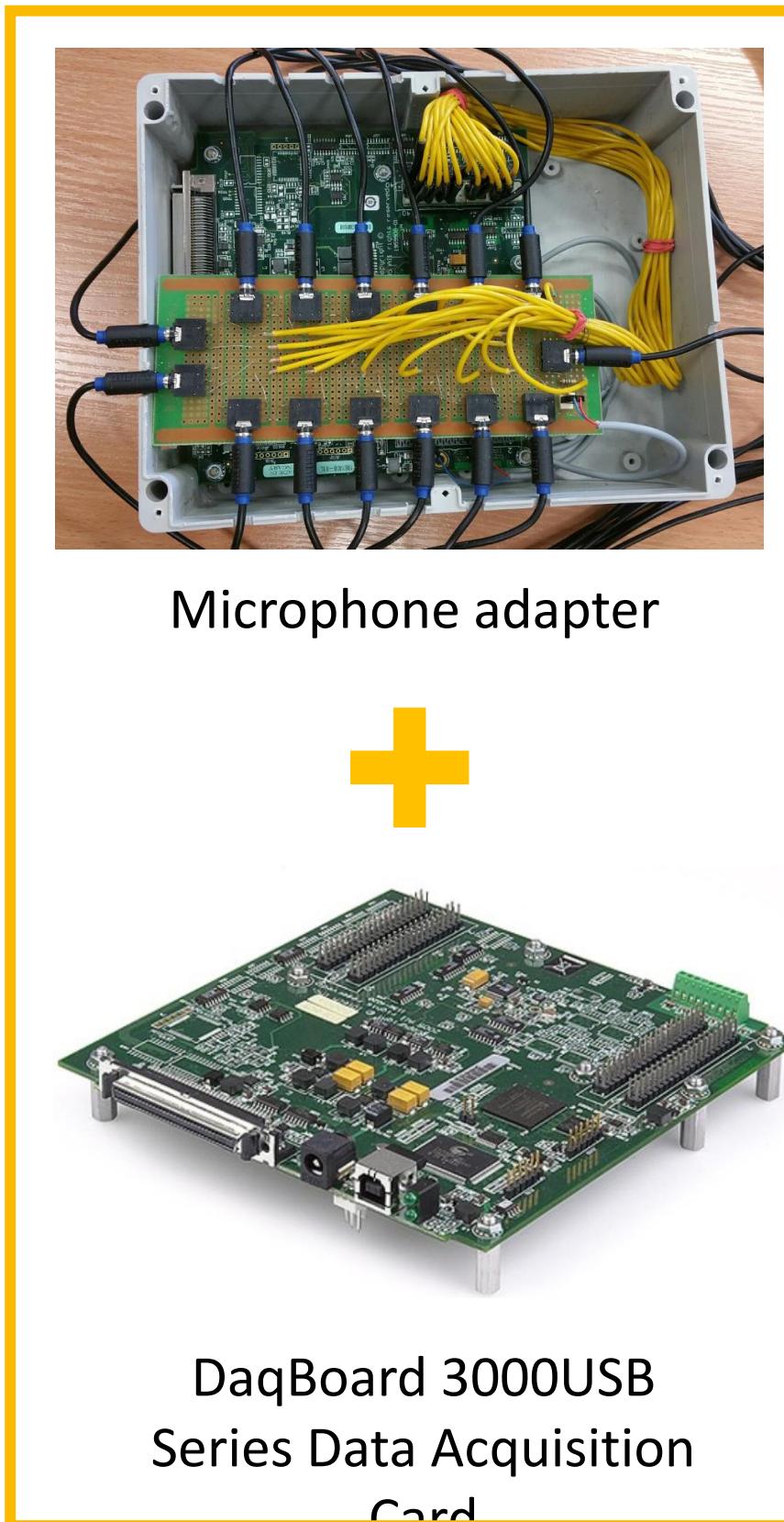
Measuring station

12

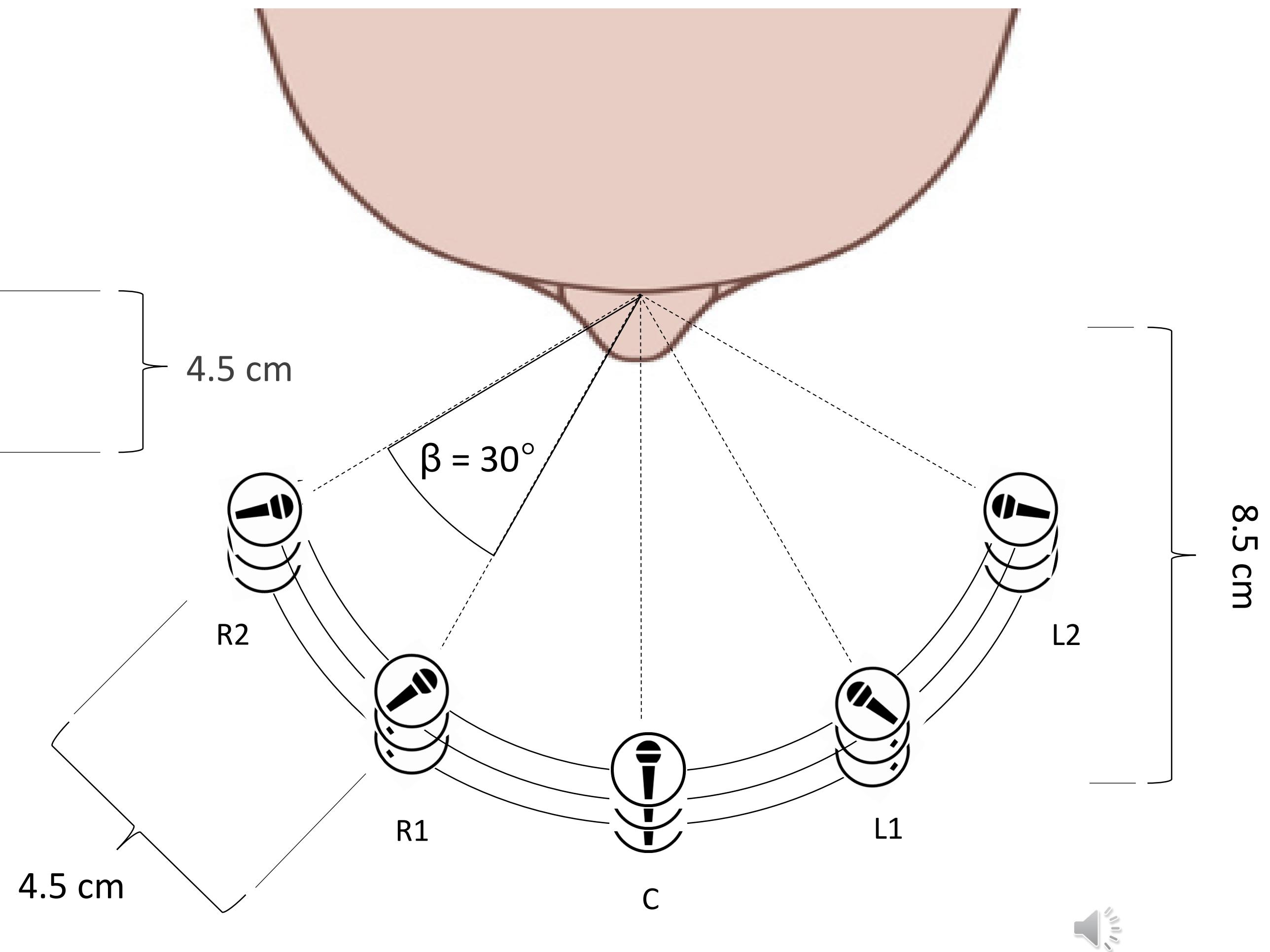
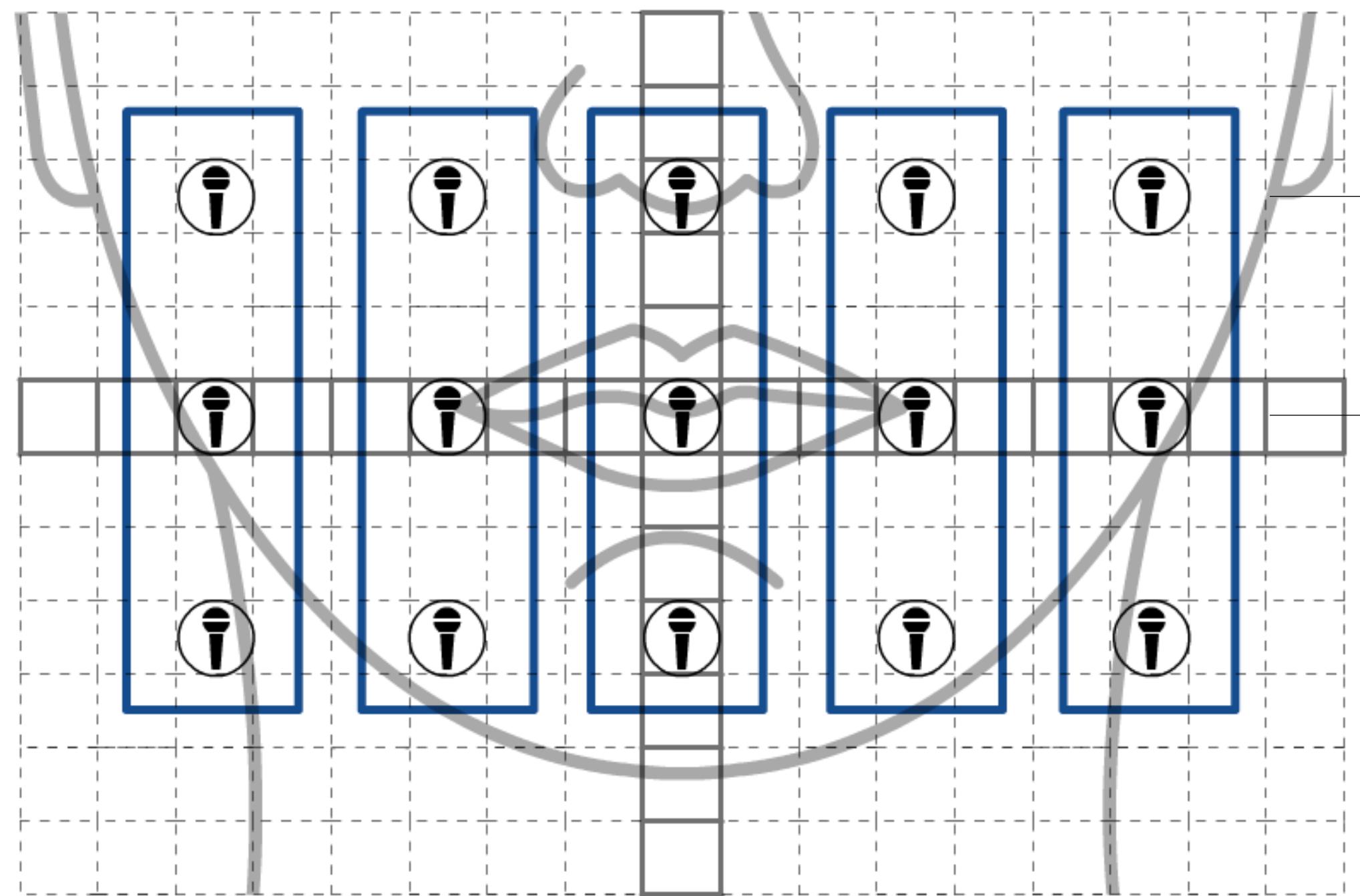


Measuring station

13



Measuring station

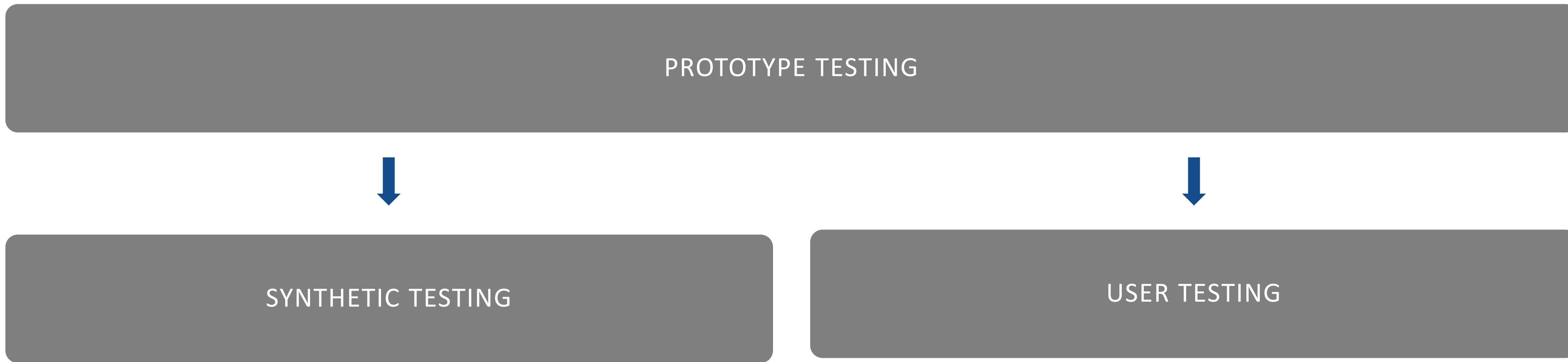


Measuring station

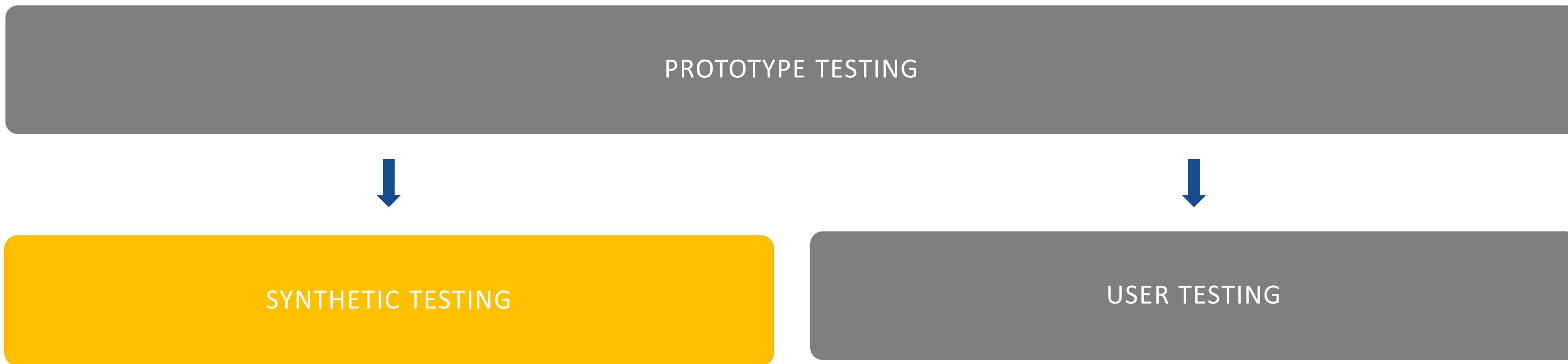
Number of channels	15
Number of groups	5
Distance from sound source	8,5 cm
Sampling rate	44,1 kHz
Microphone model	WM-61a
sensitivity	-35 dB
frequency range	20 Hz – 20 kHz
SNR	72 dB (65 dB)
THD	-72 dB (-65 dB)
power supply	2,3 V
Amplifier model	MAX9812
Gain	20 dB



Measuring station

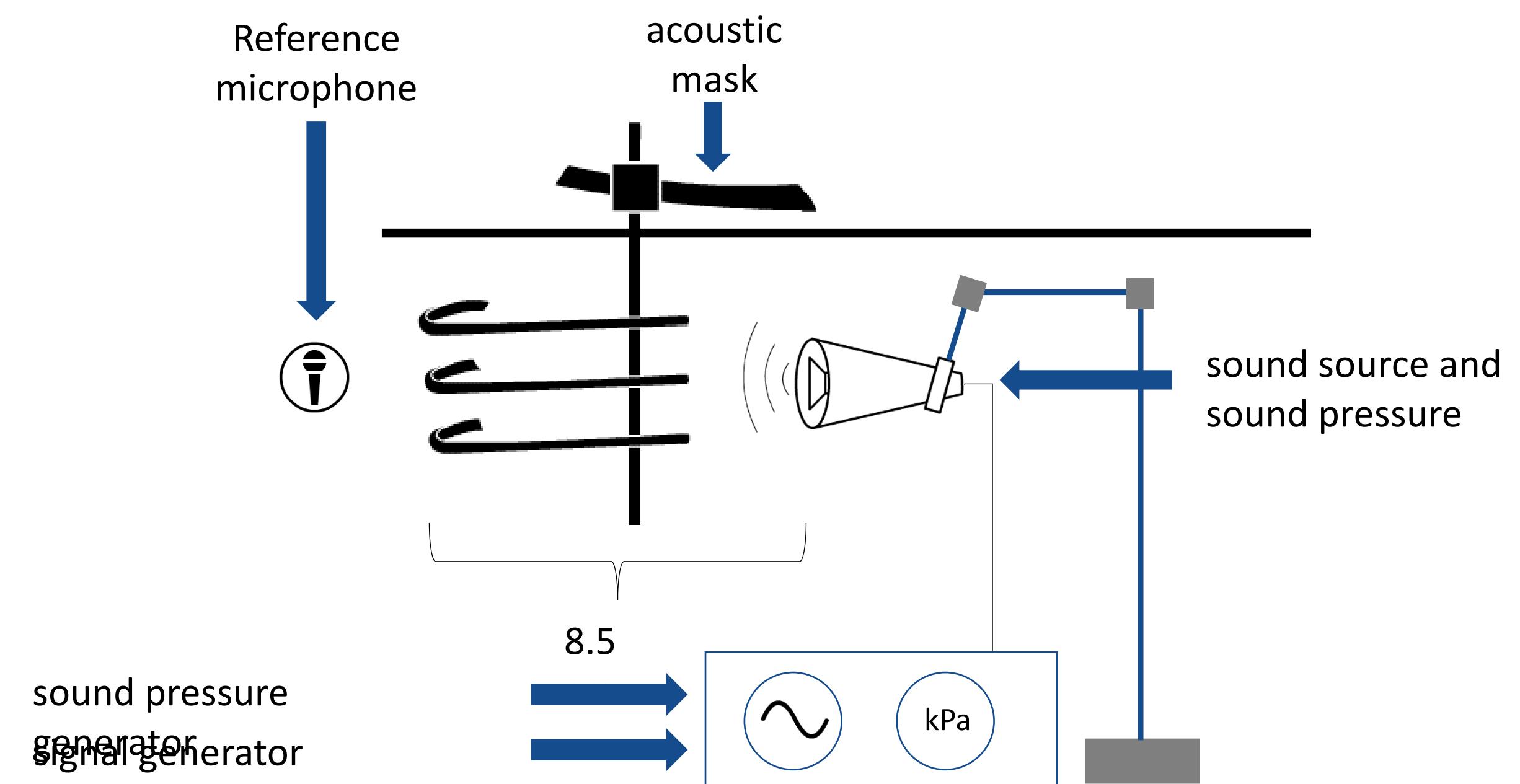
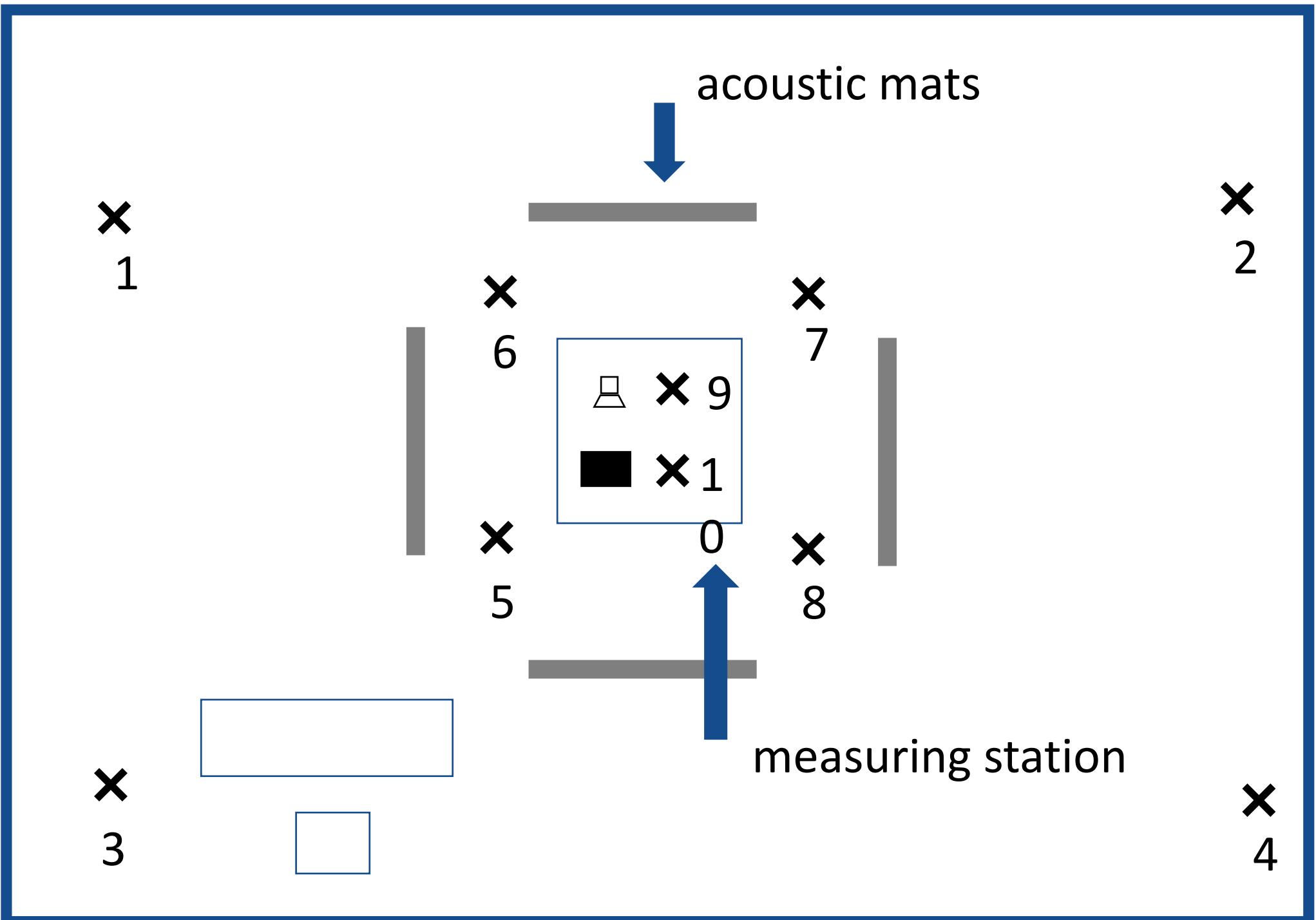


Measuring station



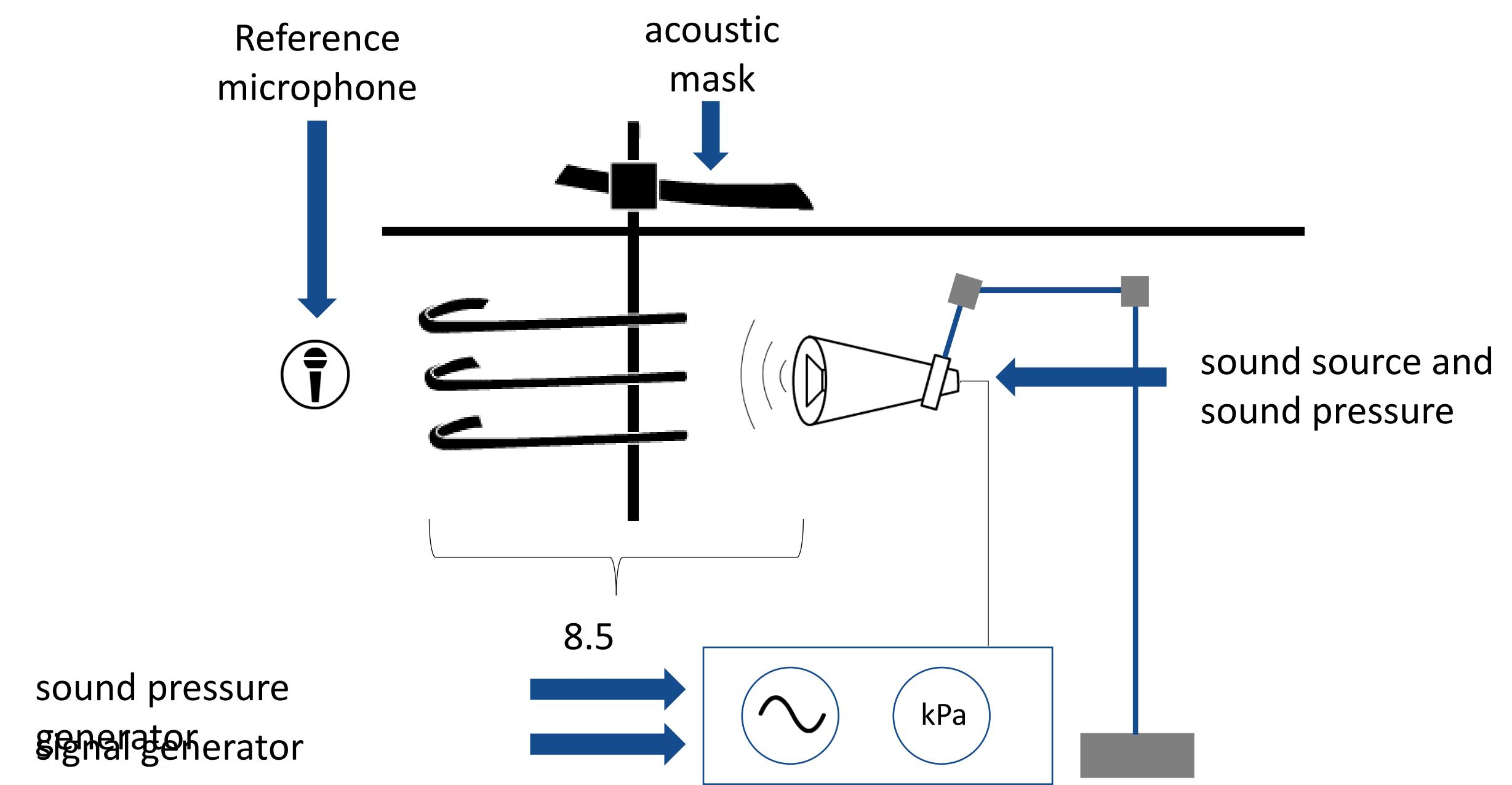
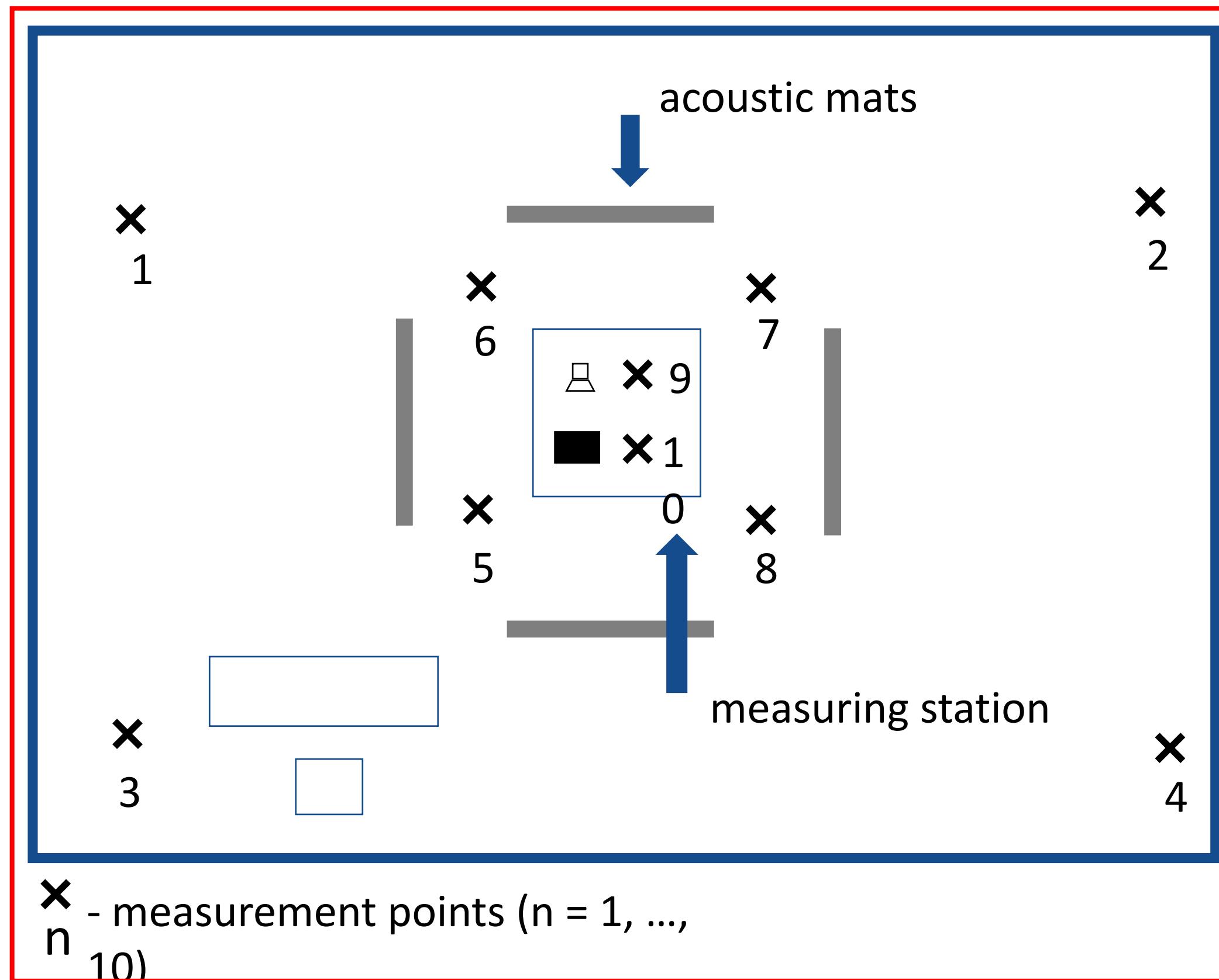
Measuring station

18

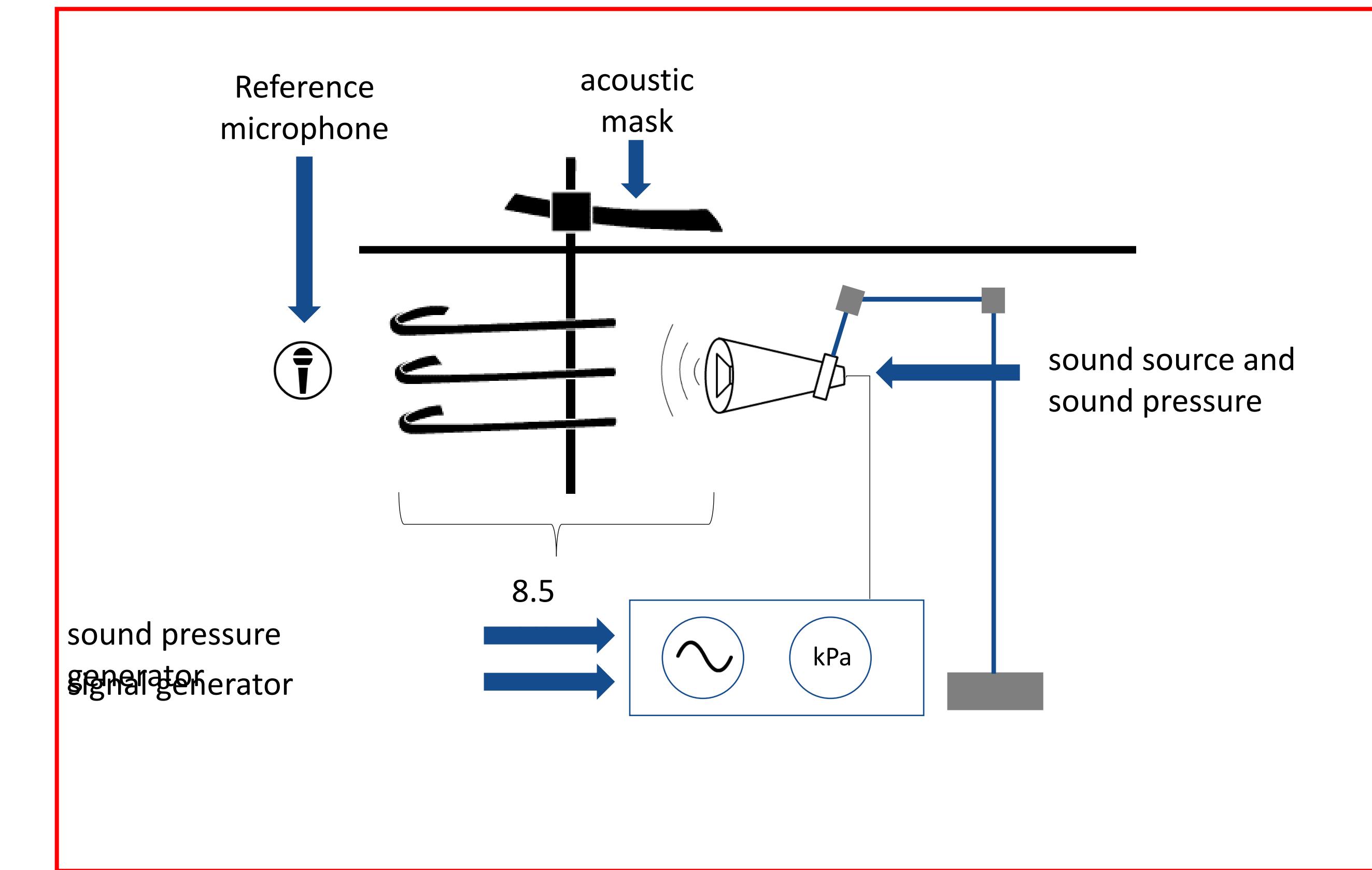
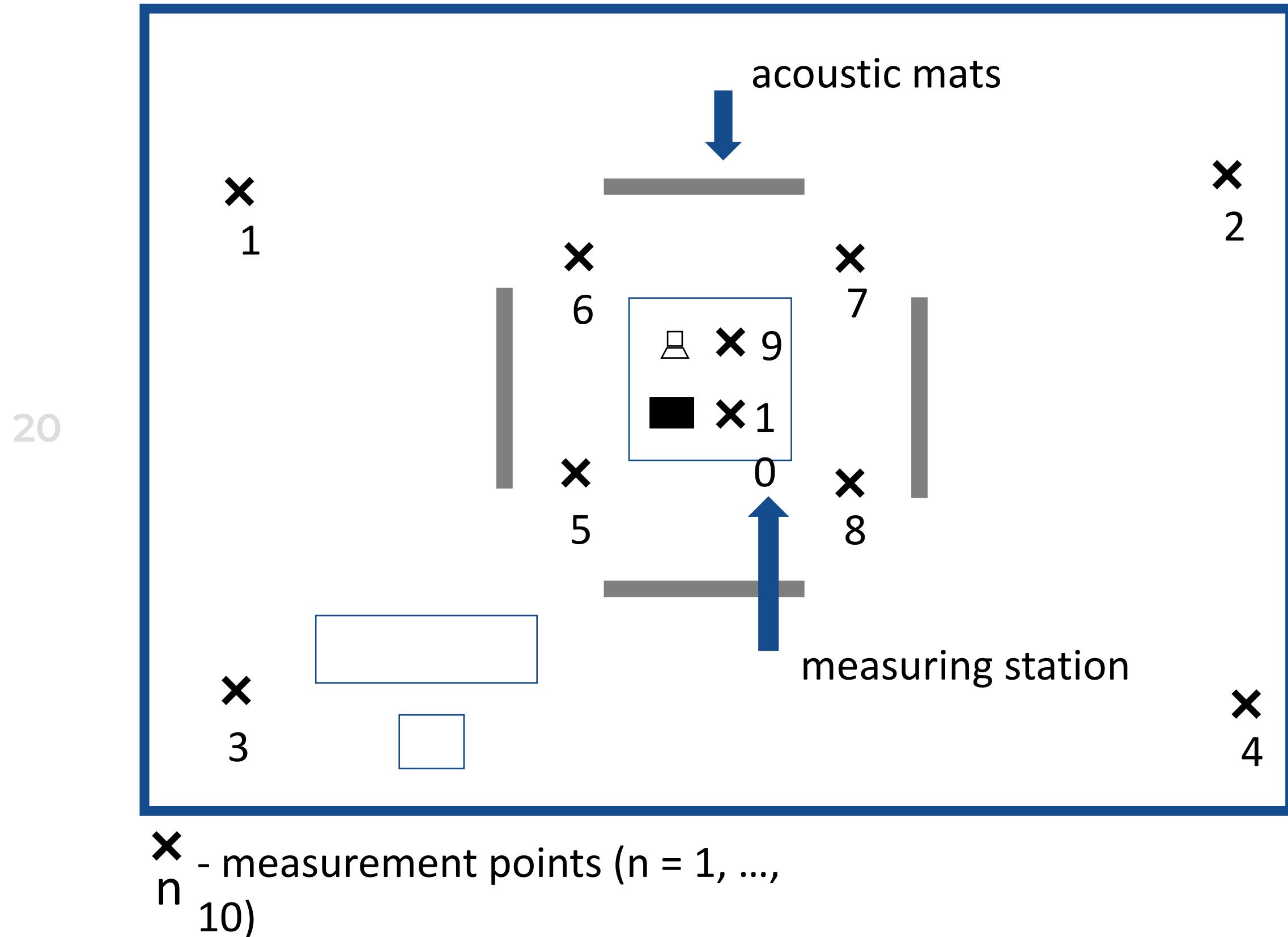


Measuring station

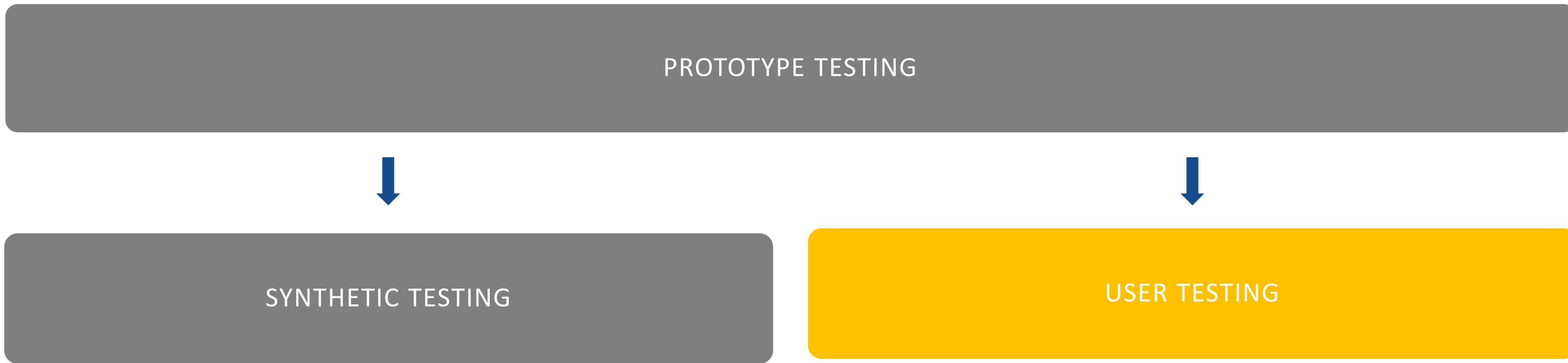
19



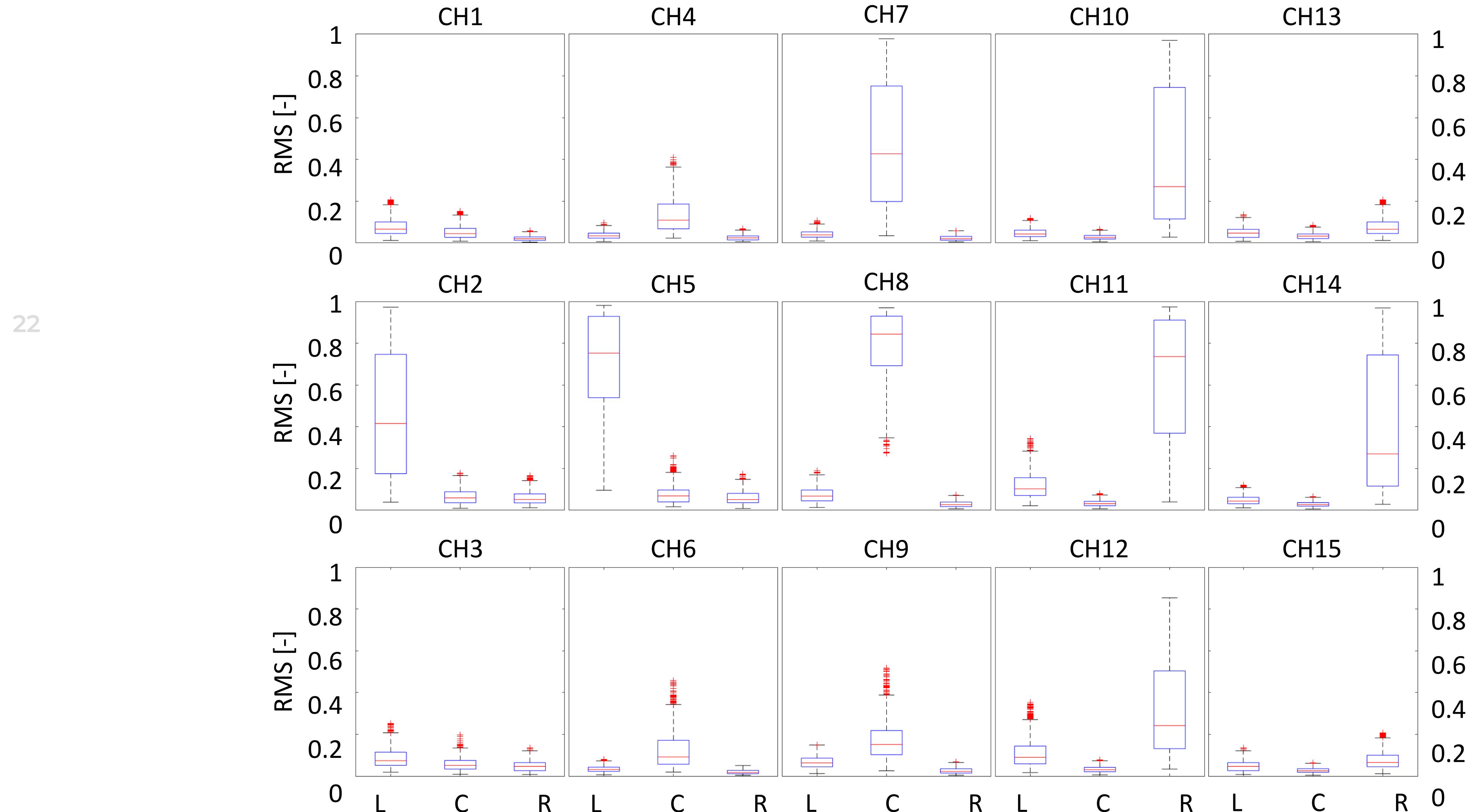
Measuring station



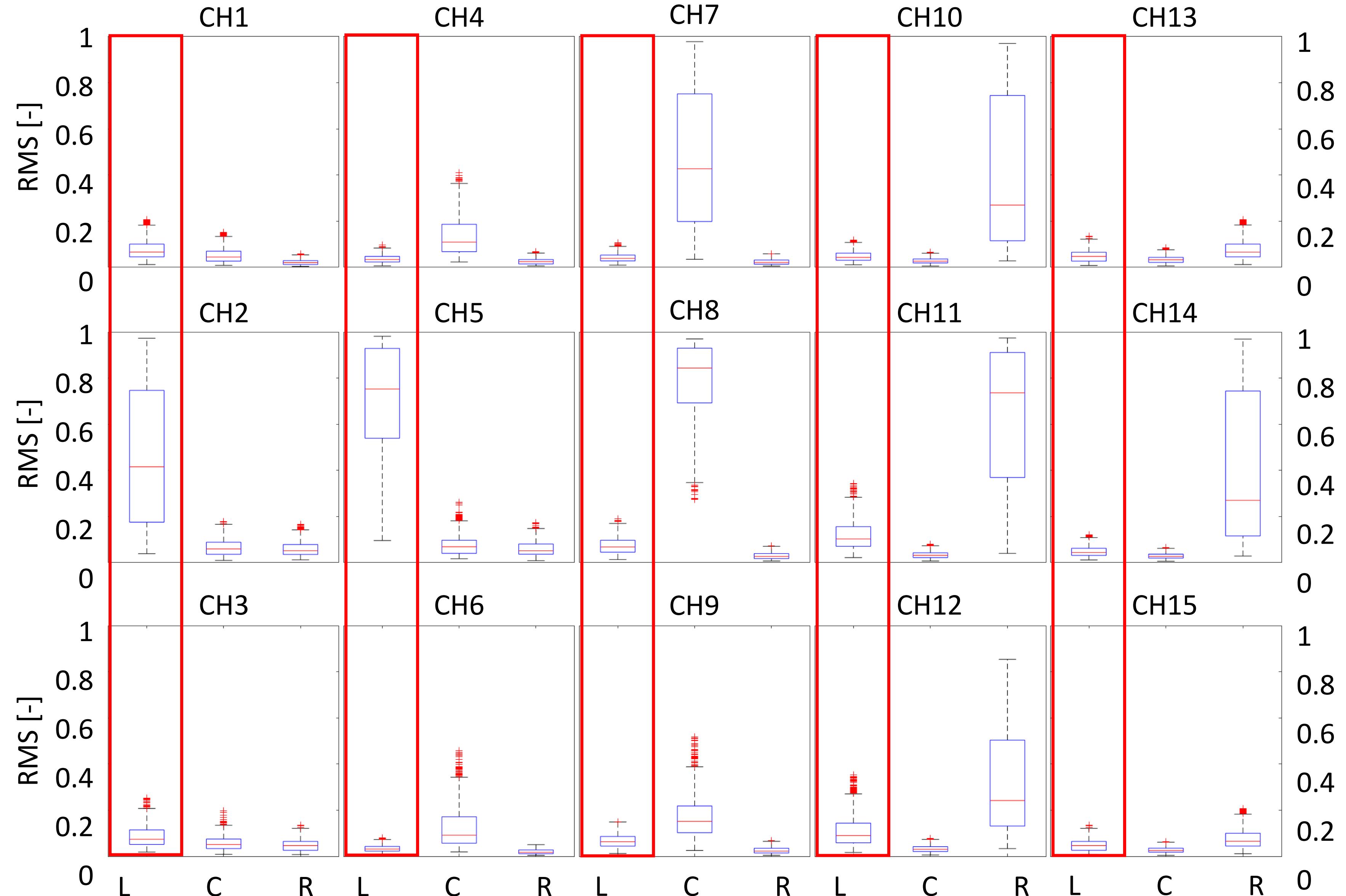
Measuring station



Measuring station

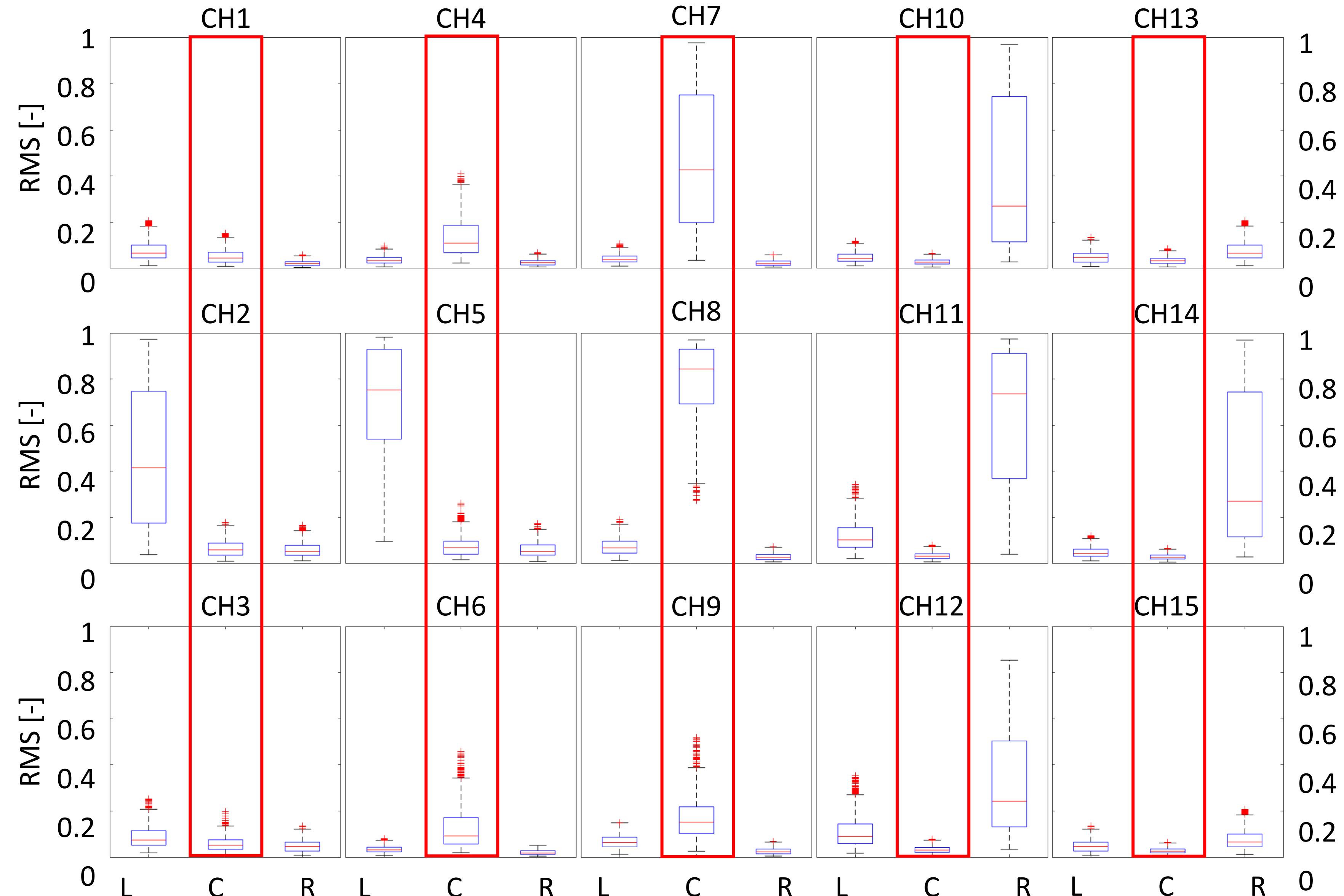


Measuring station

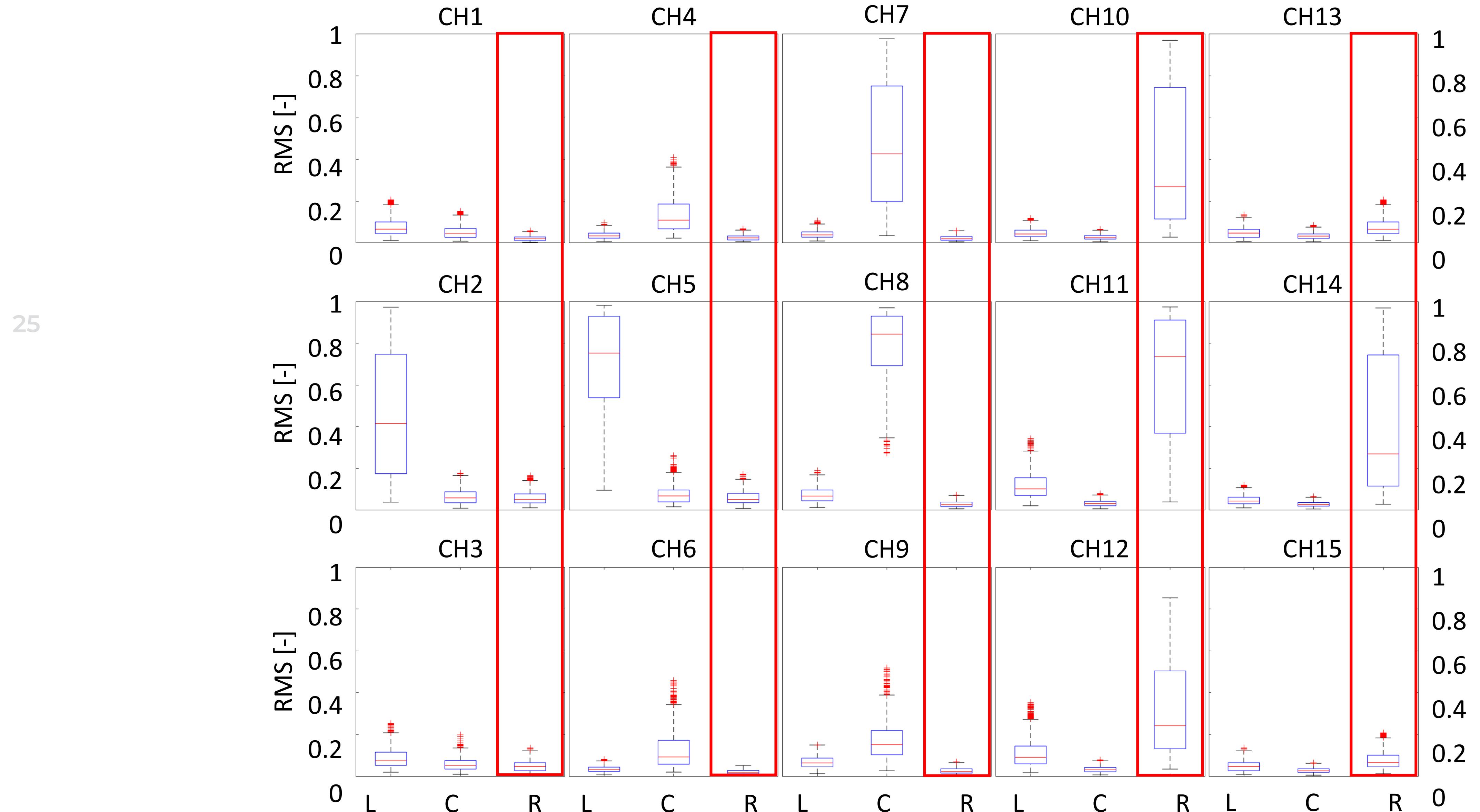


Measuring station

24



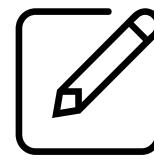
Measuring station



Measuring station



multi-channel, spatial and repeatable sound recording,



time synchronization,



visual adjustment,

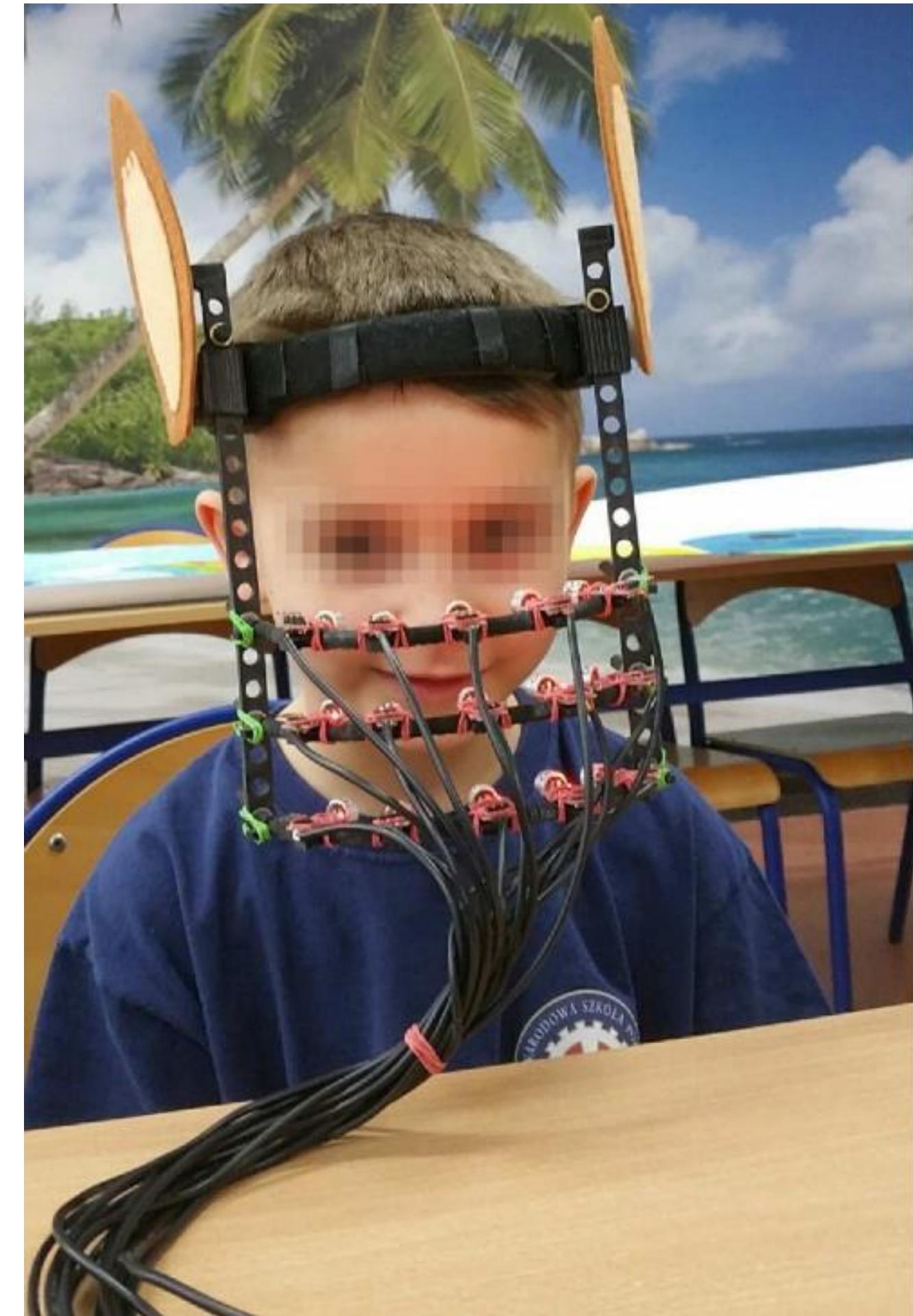


size adjustment,



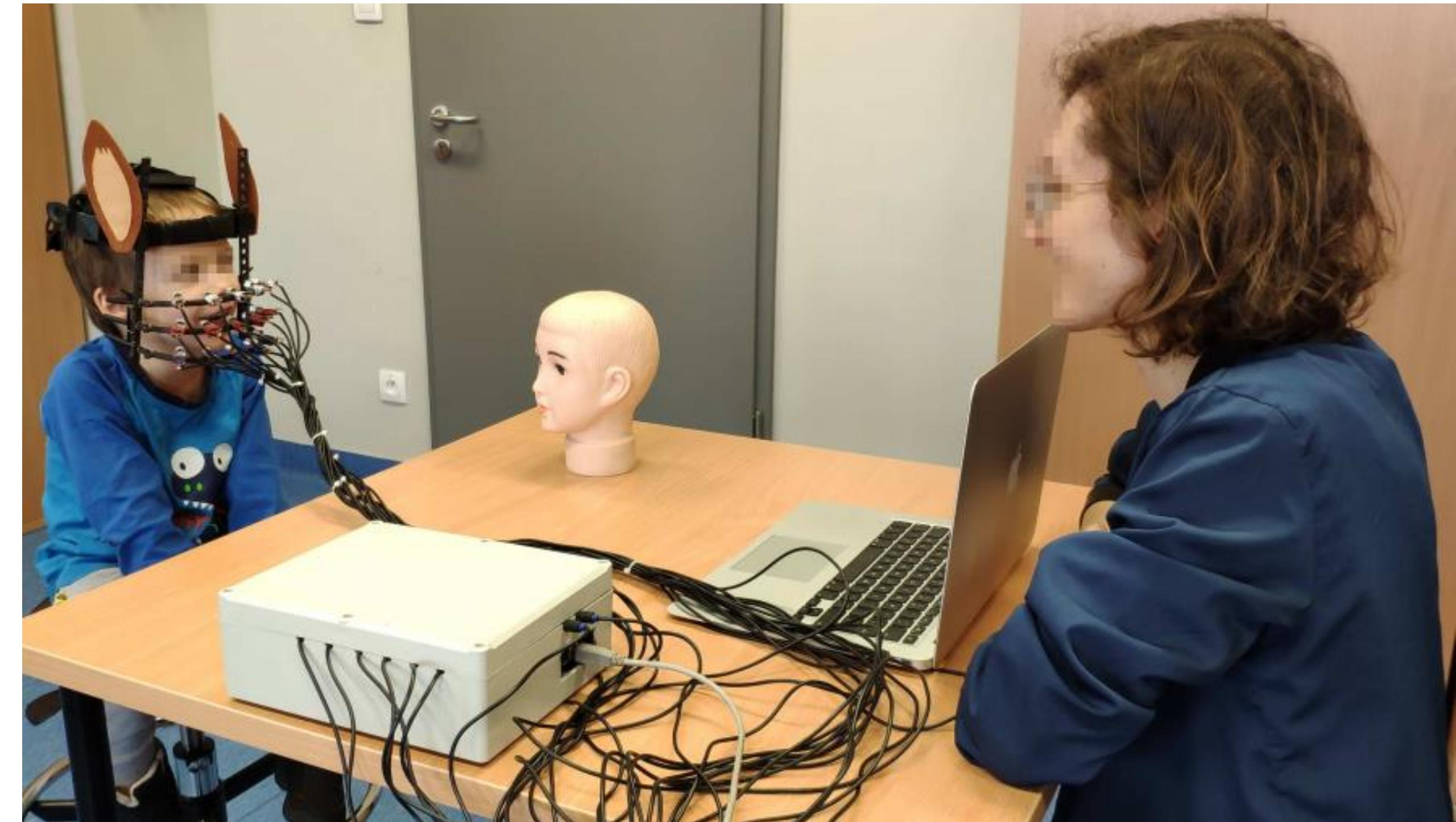
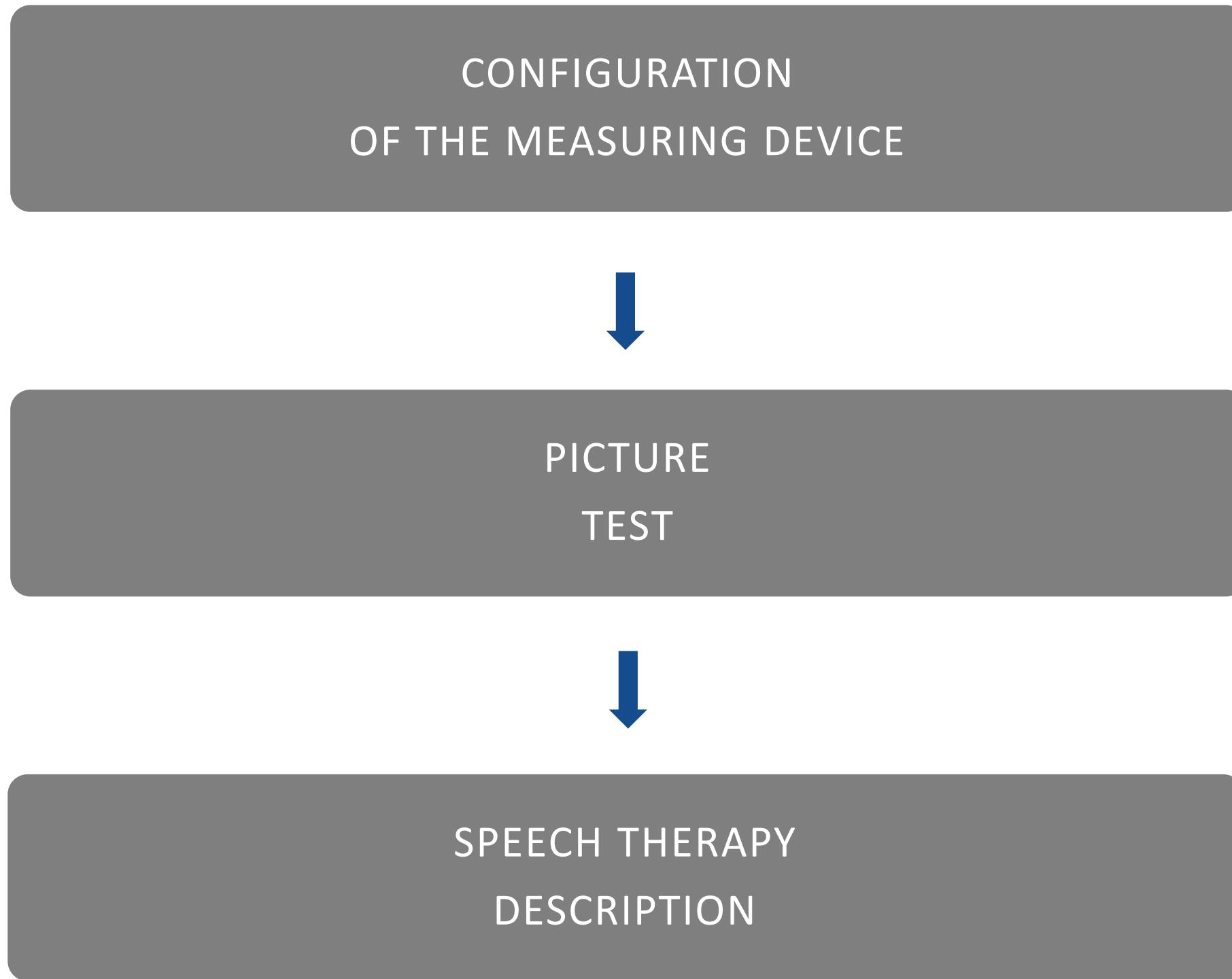
semi-cylindrical surface.

26



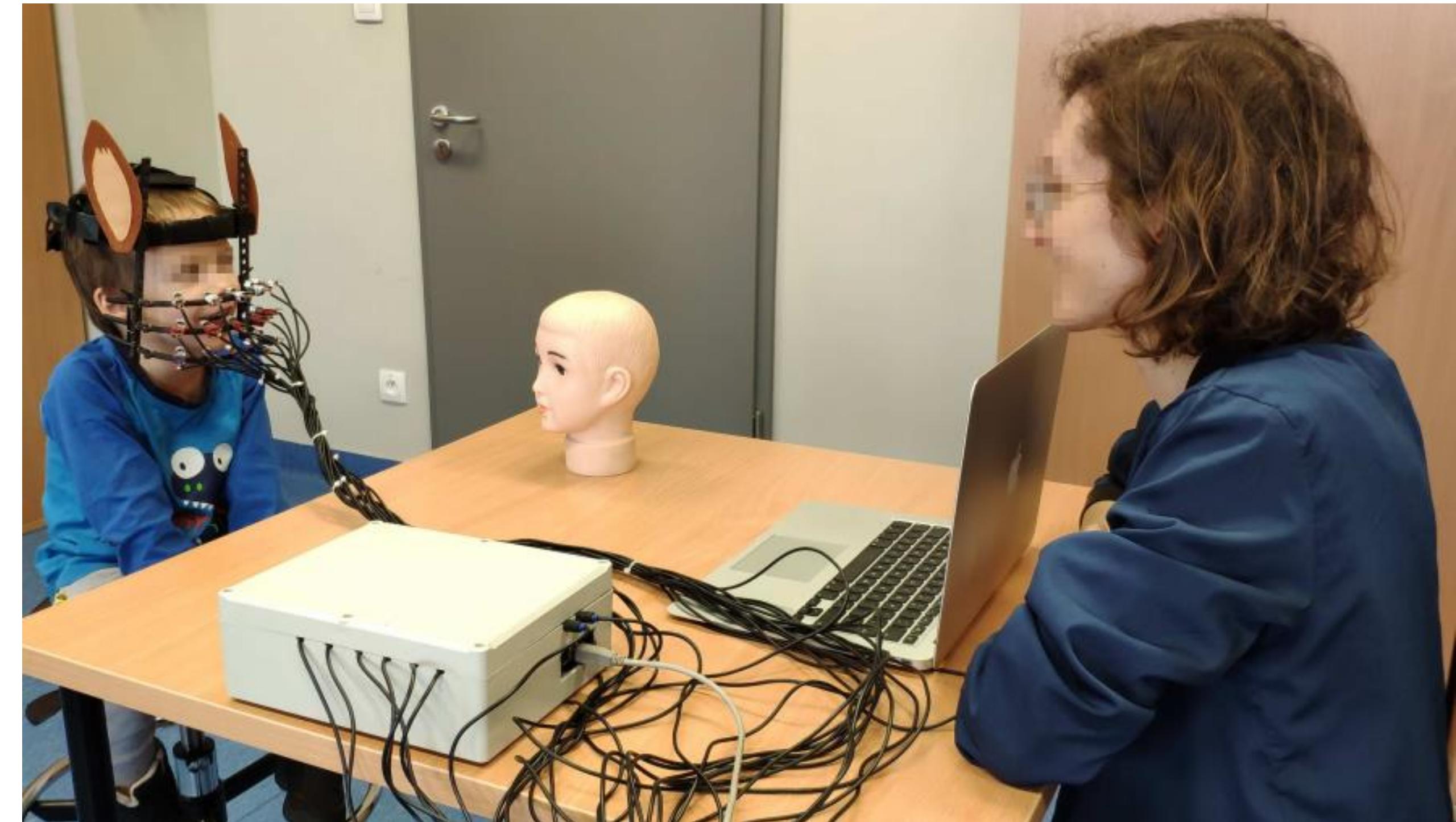
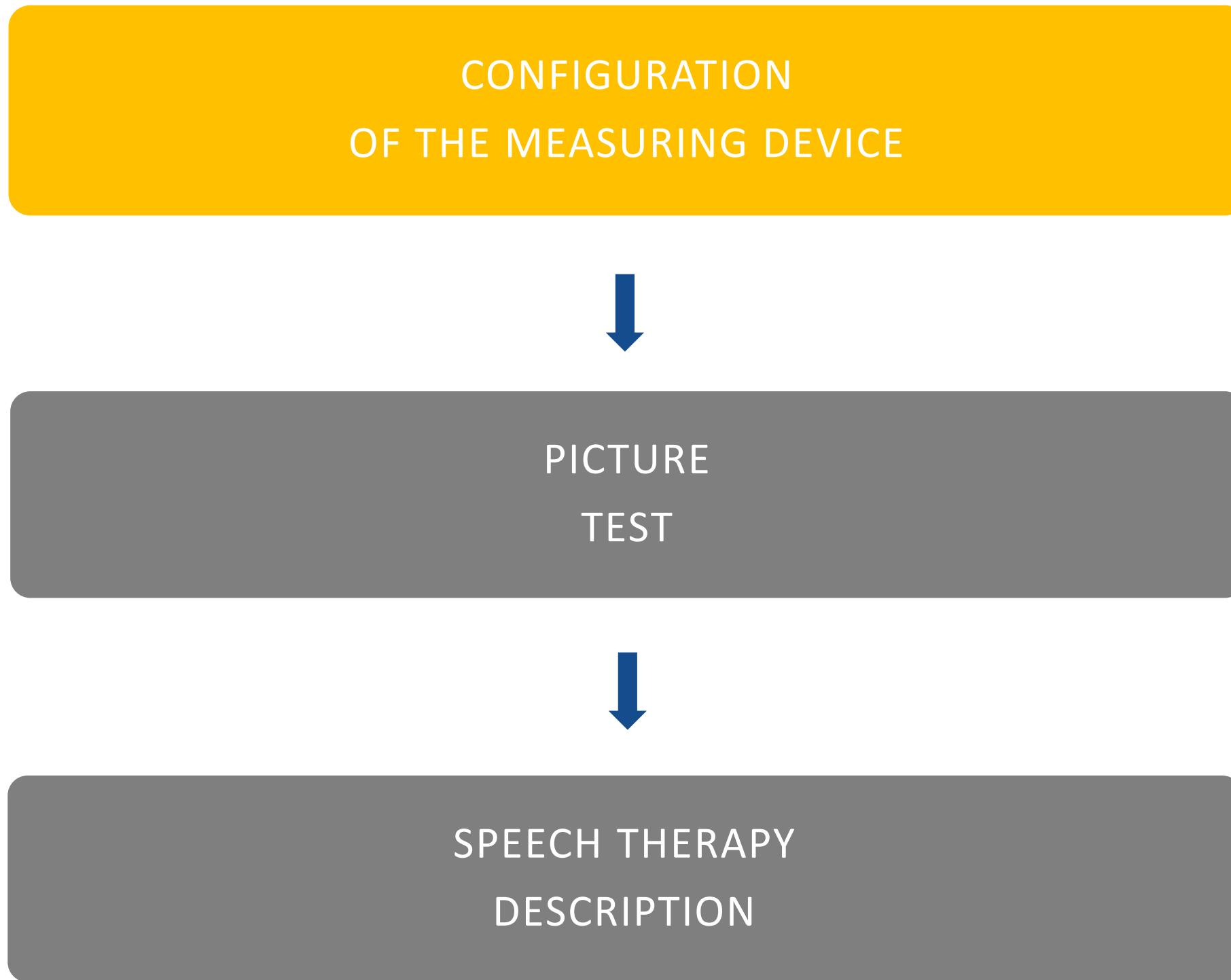
Speech Database

27



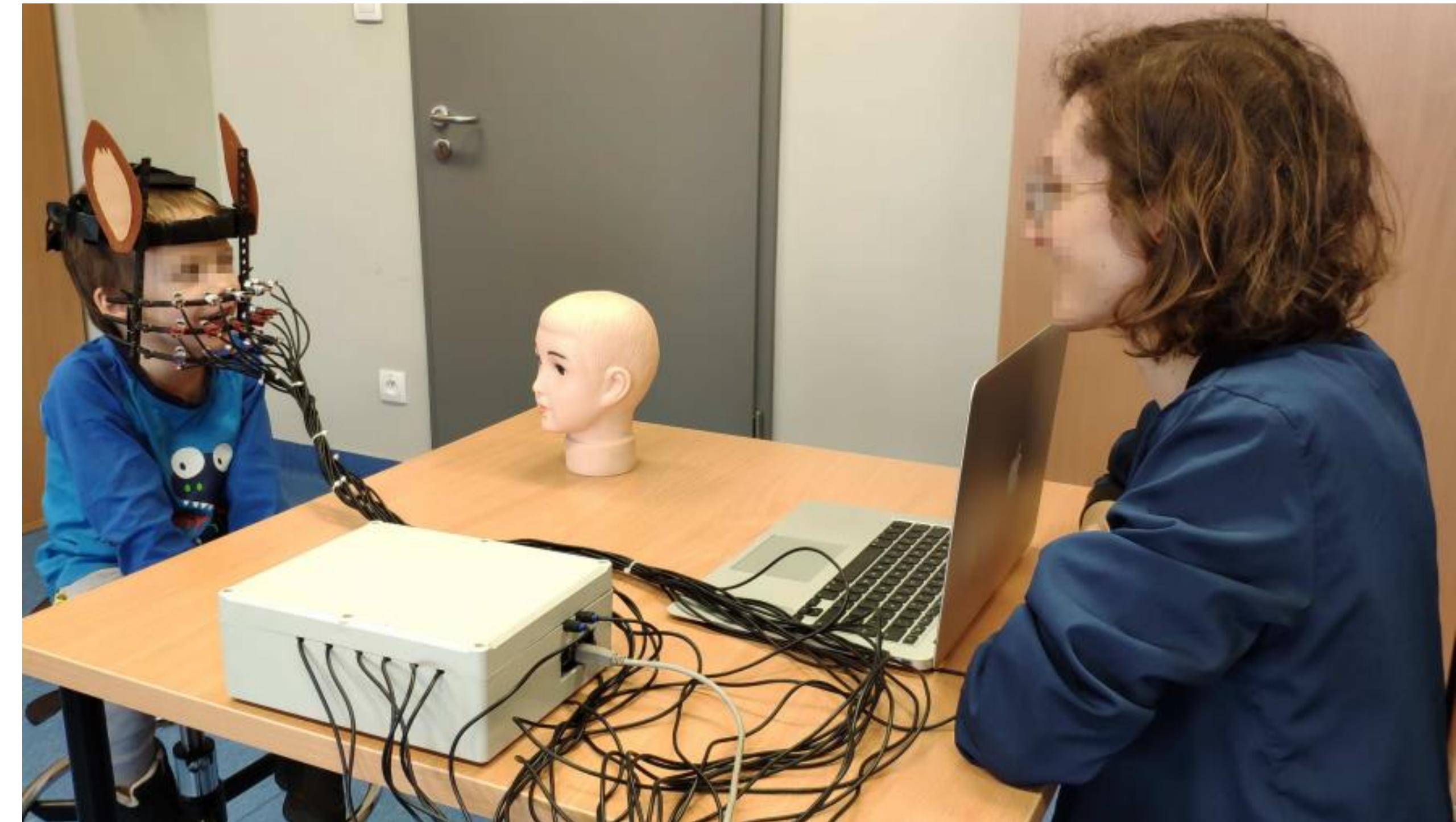
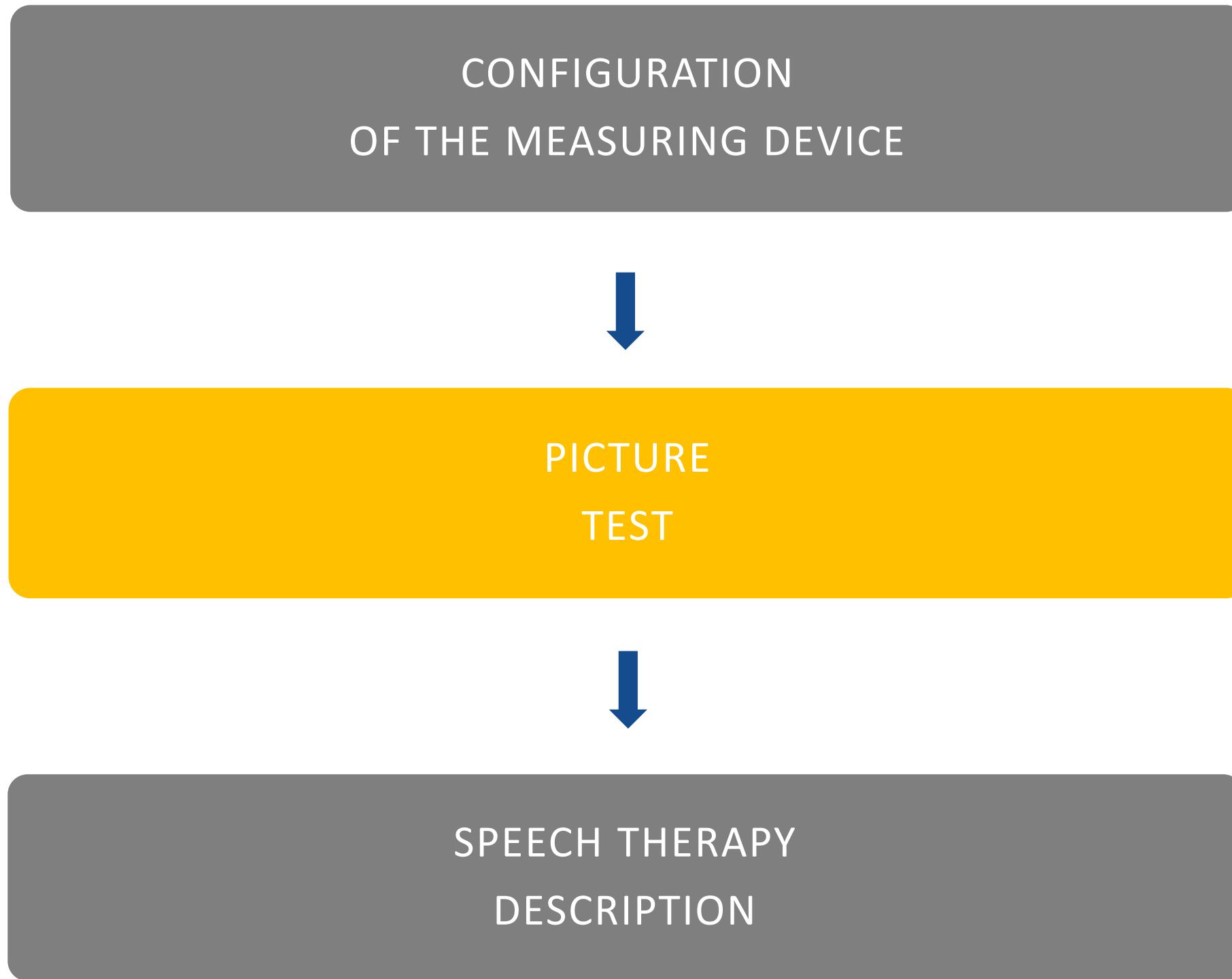
Speech Database

28



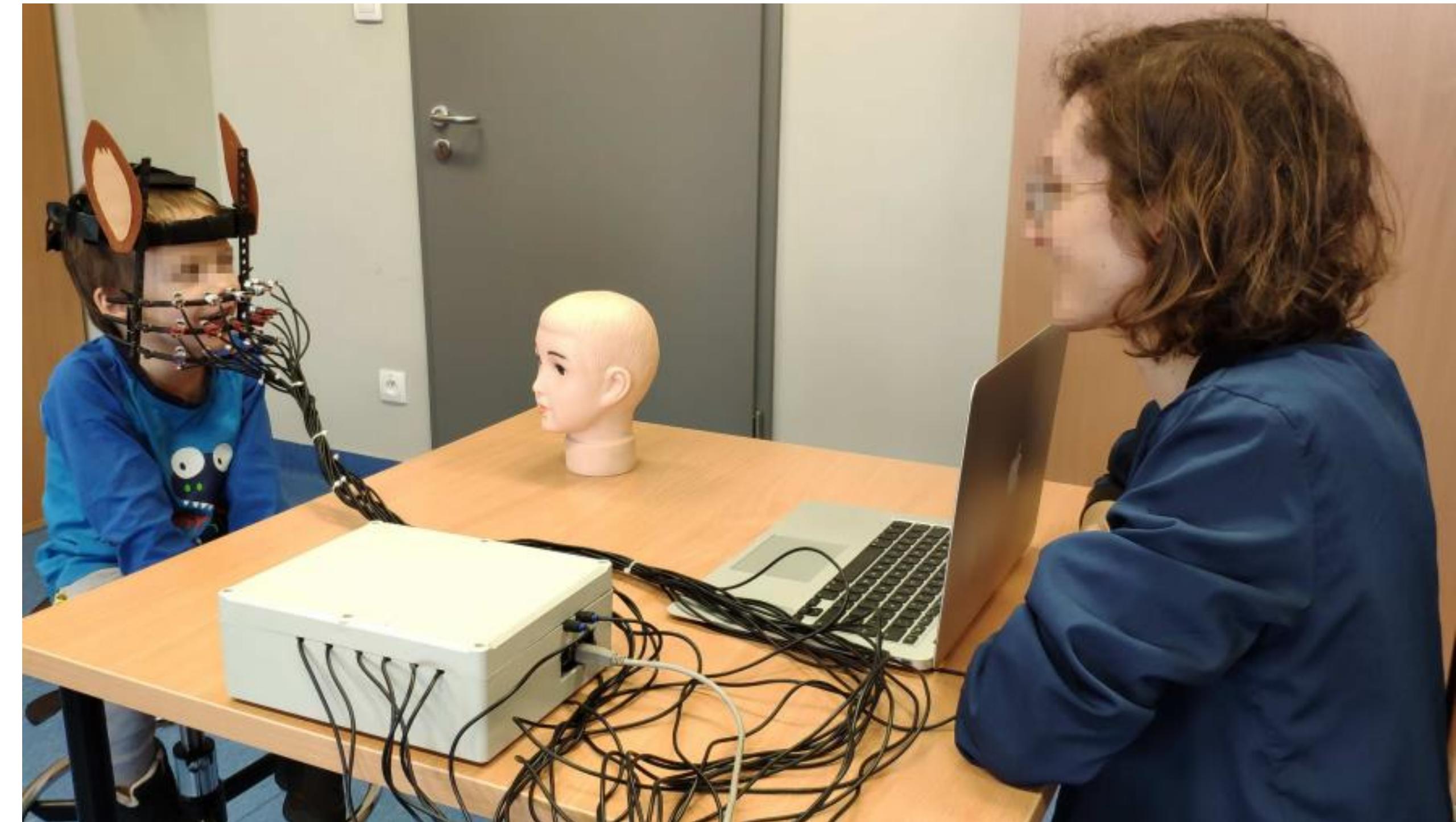
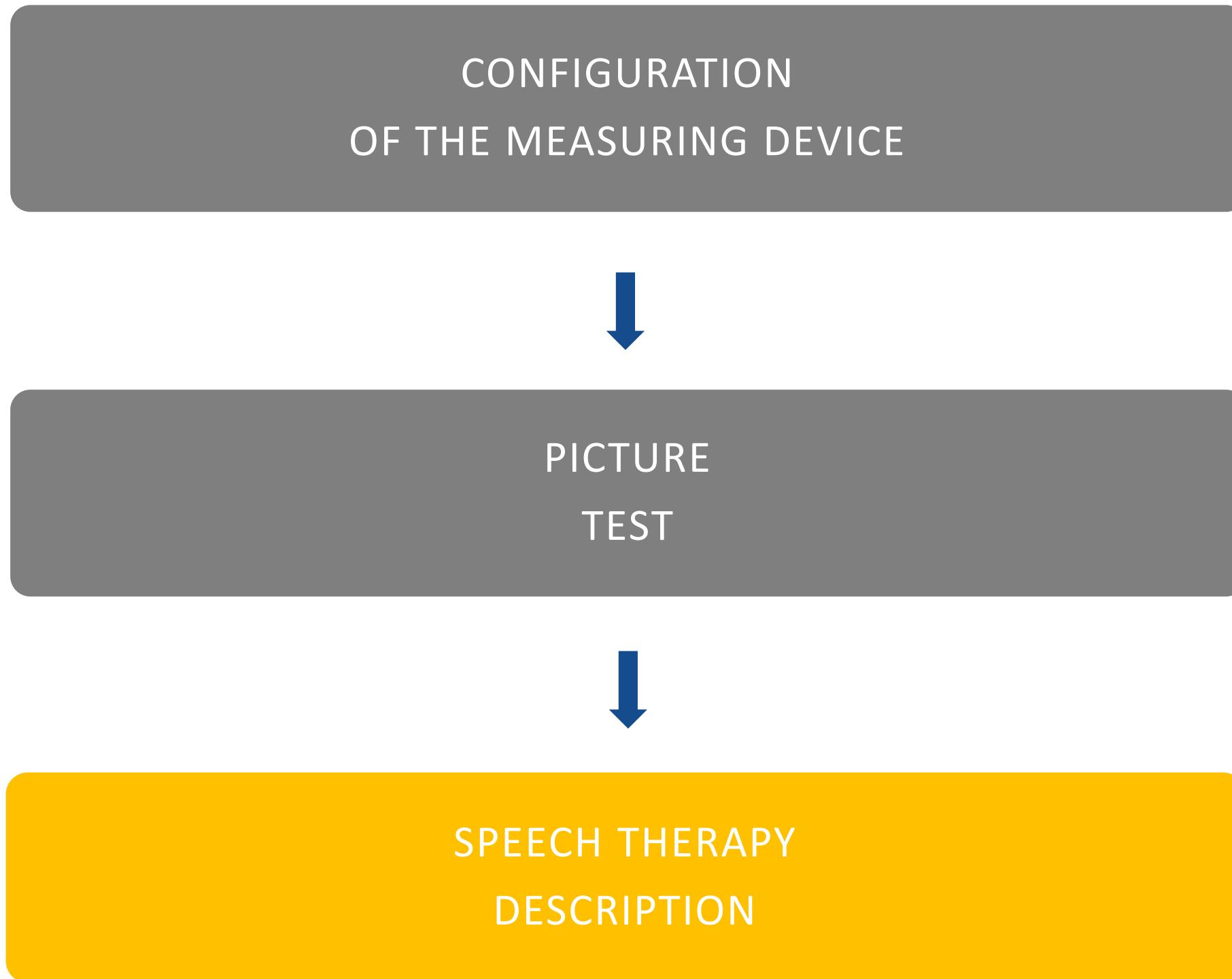
Speech Database

29



Speech Database

30



Speech Database

31



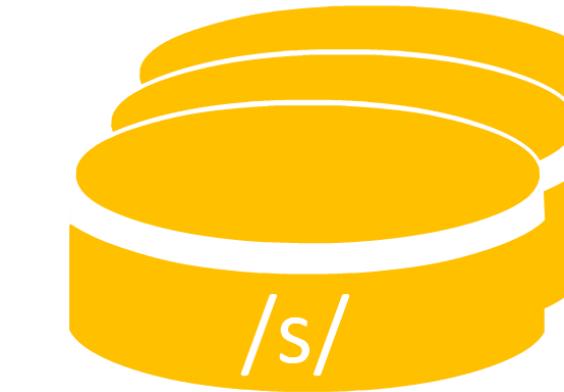
Dictionary
words with dentalized sounds
in different articulation
positions

98 speakers

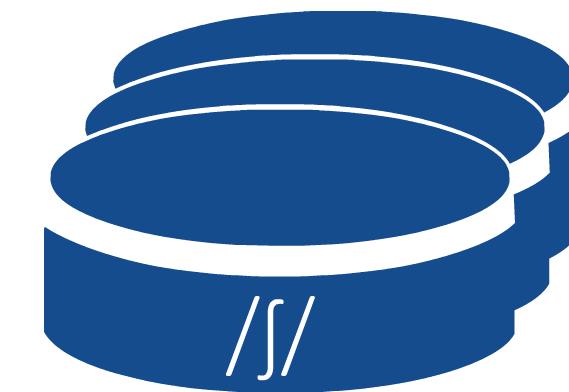
Age 6-7, 1362 recordings



Speech therapy
description
normative, interdental, addental
and dental realization



/s/



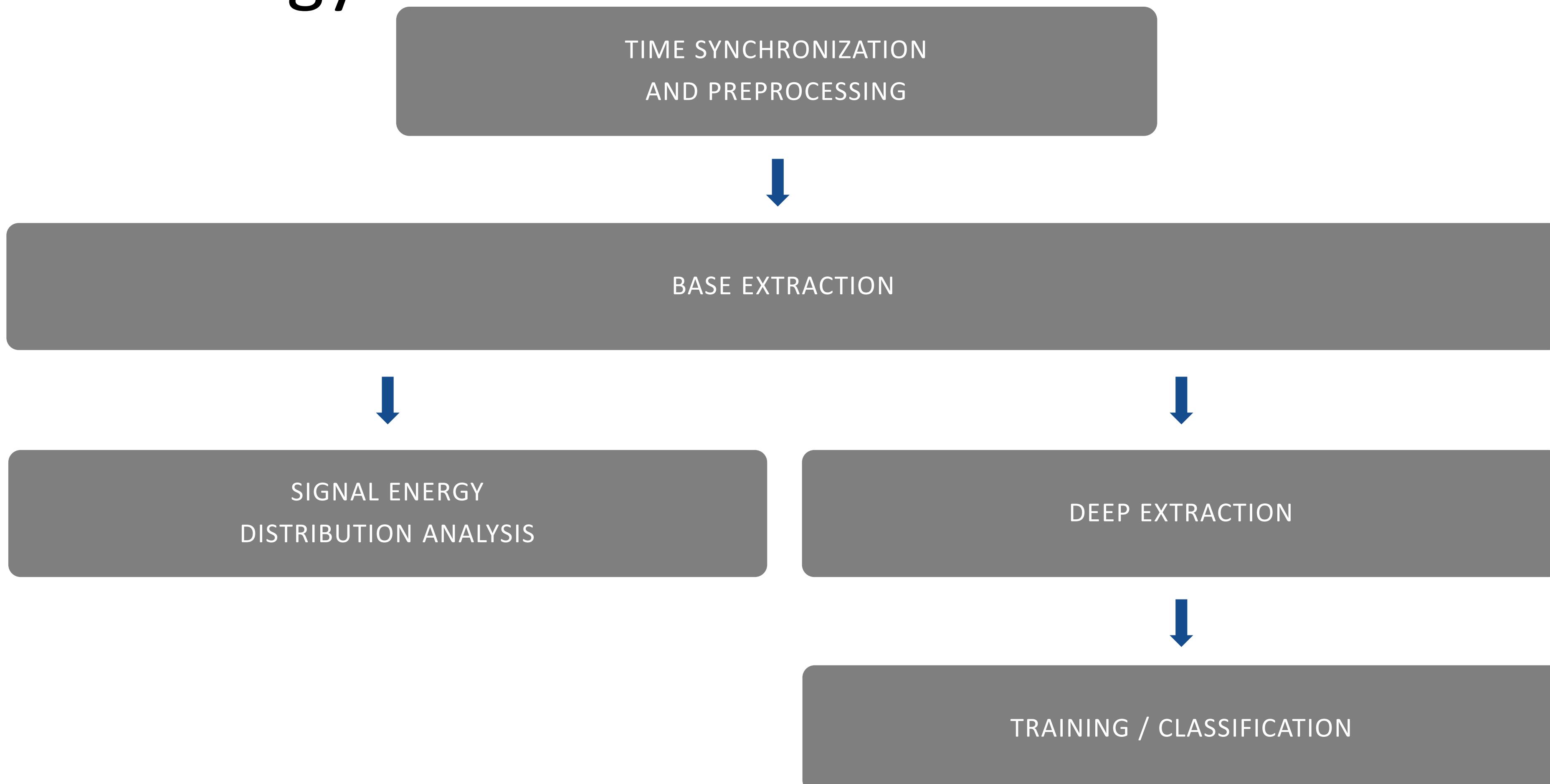
/ʃ/

Polish	English	number
samolot	plane	76
serce	heart	82
strażak	fireman	70
pasek	belt	81
parasol	umbrella	48
pies	dog	51
lis	fox	85

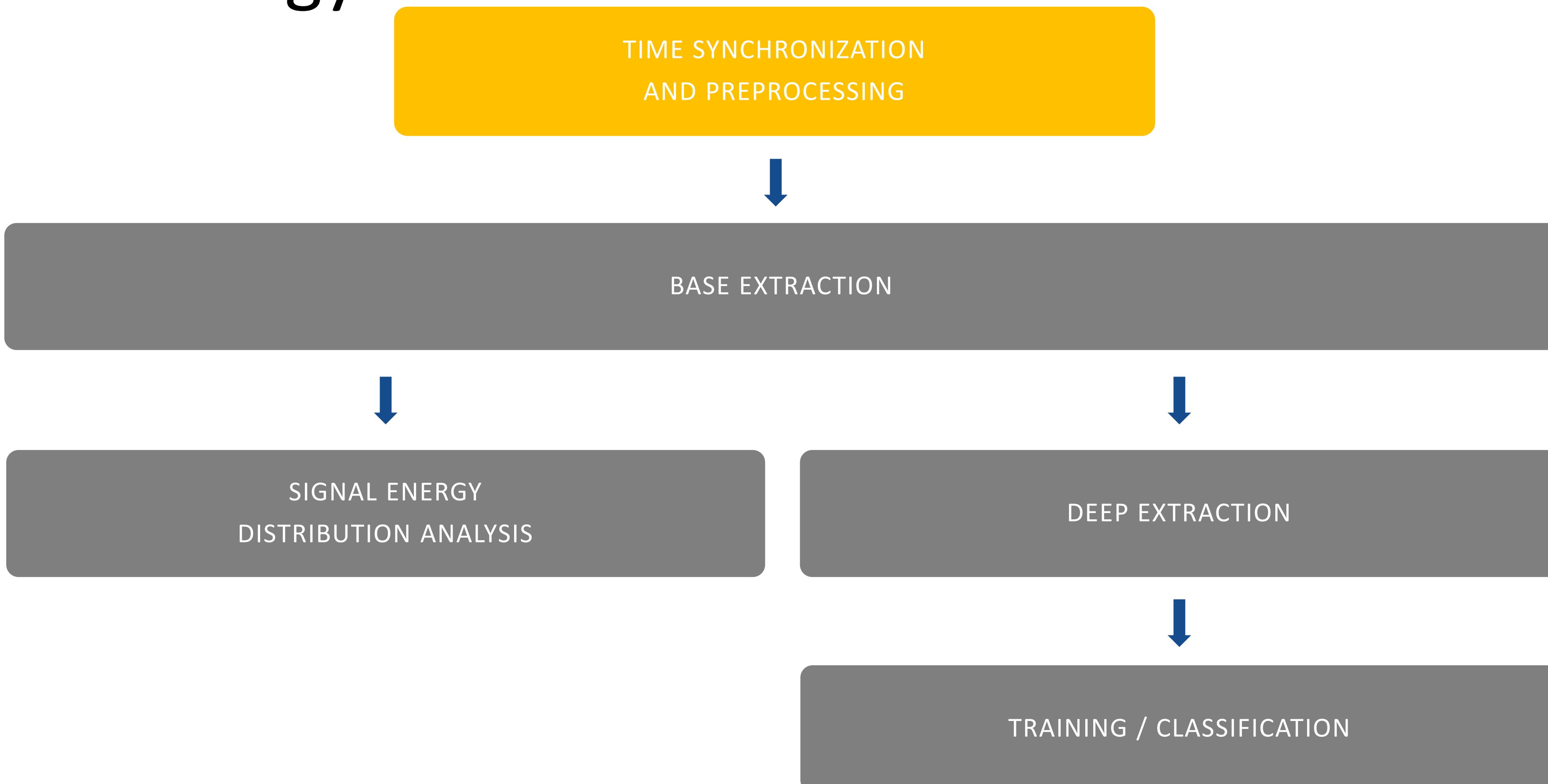
Polish	English	number
szafa	wardrobe	82
sznur	rope	11
szufelka	dustpan	79
kalosze	wellies	82
koszyk	basket	82
książka	book	11
lekarz	doctor	12
nóż	knife	78
wąż	snake	82



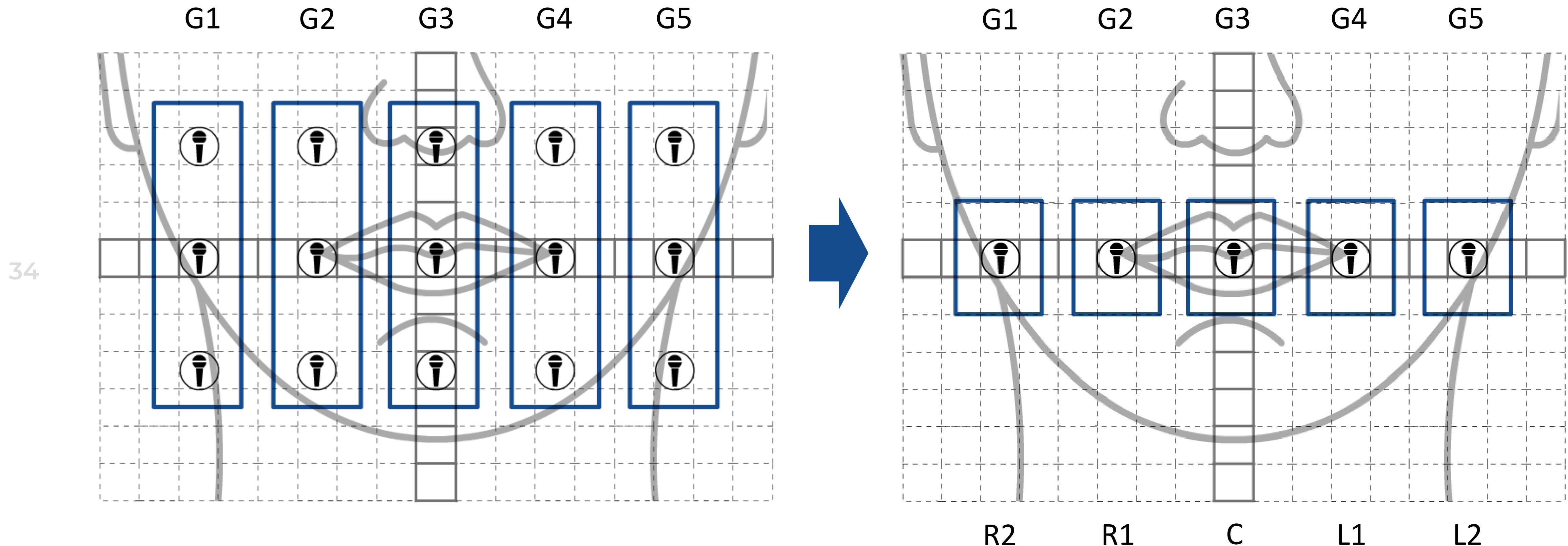
Methodology



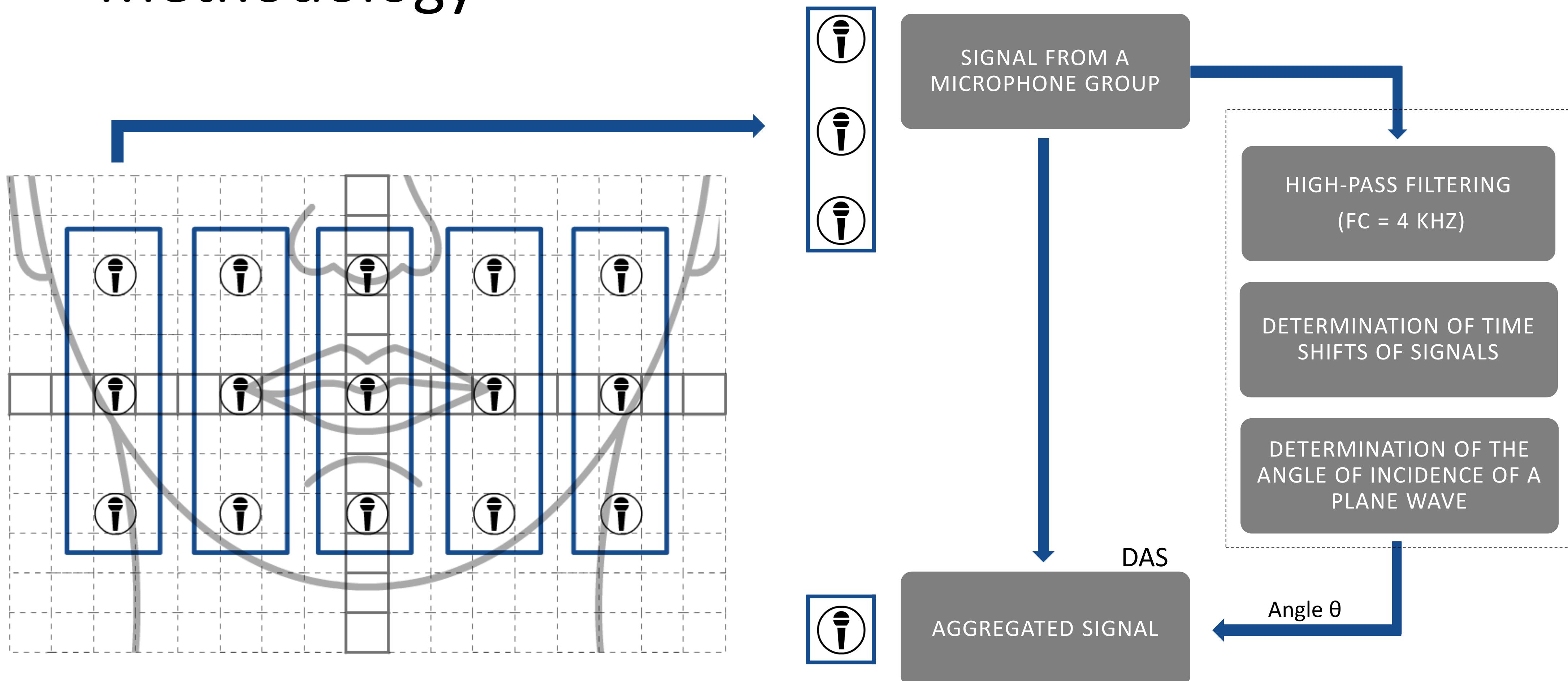
Methodology



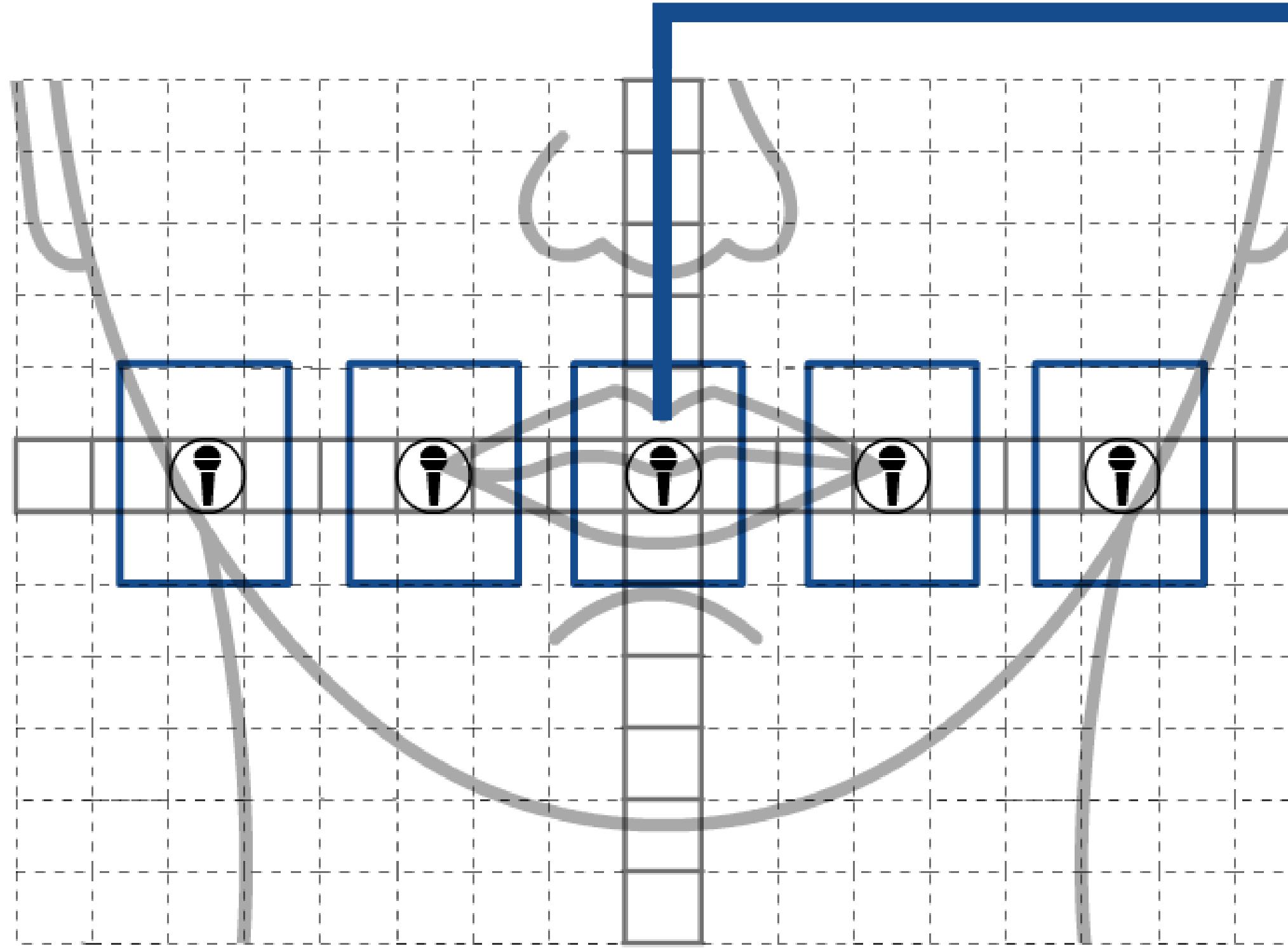
Methodology



Methodology

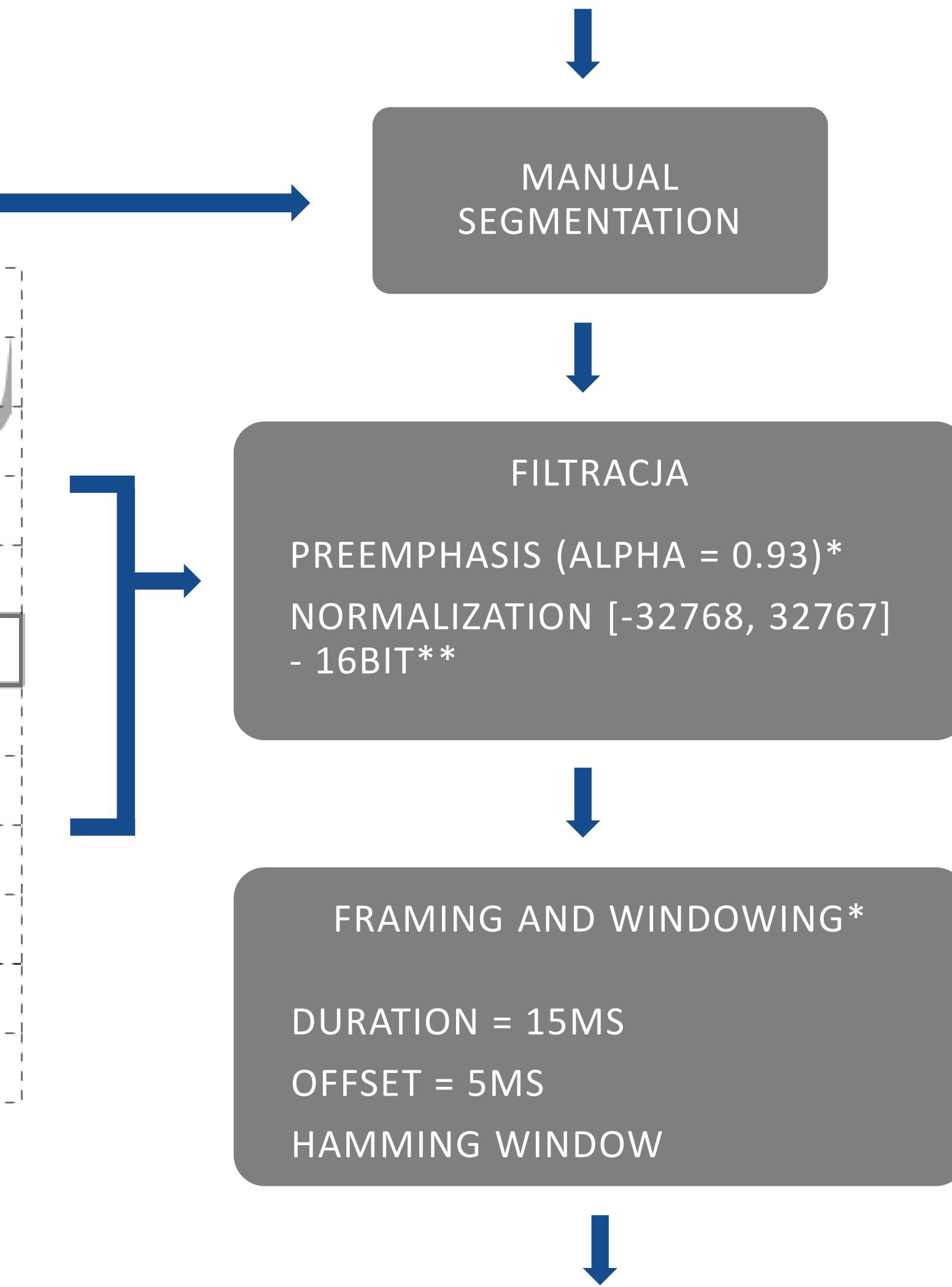


Methodology

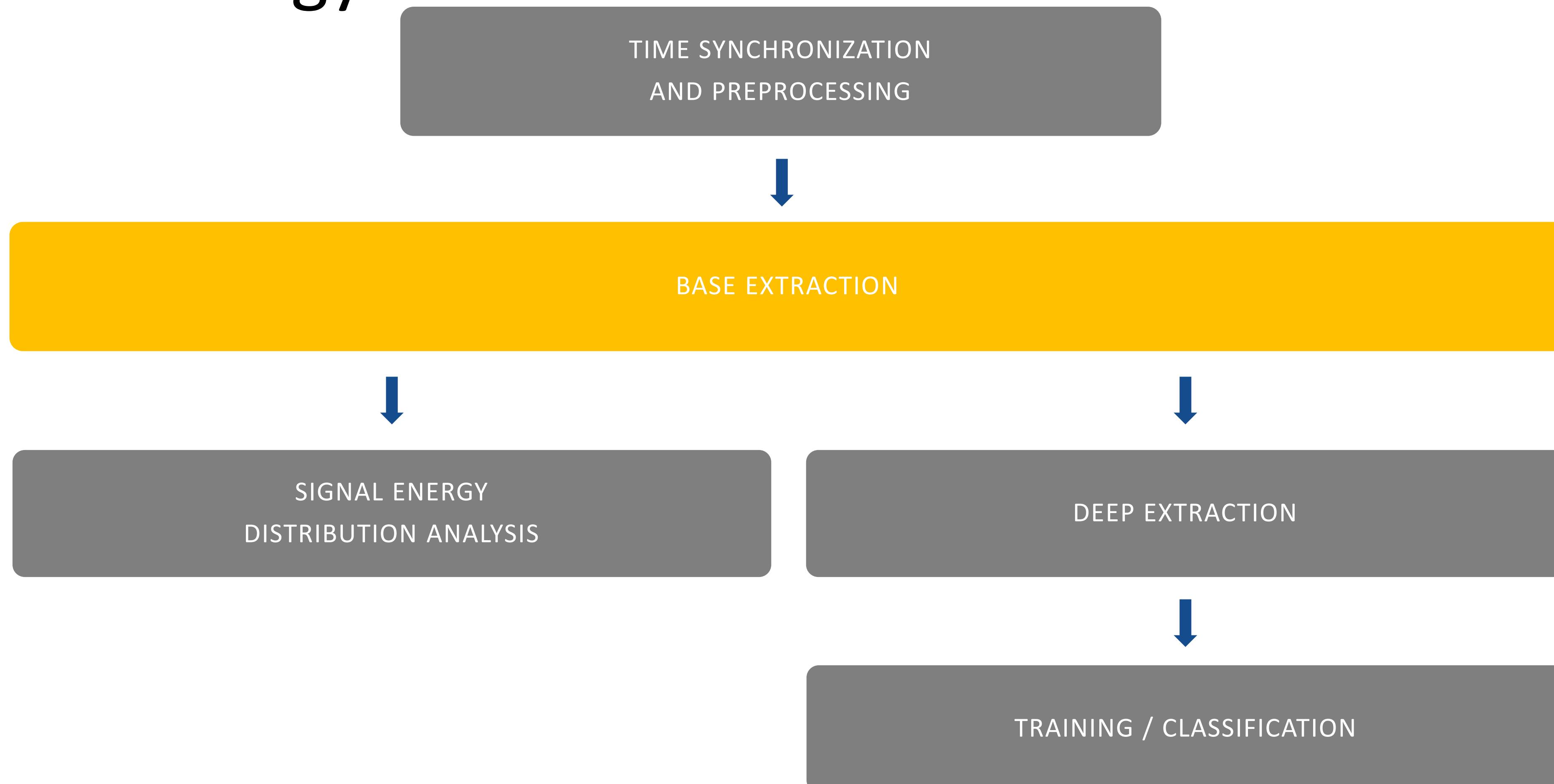


* - independently for each channel

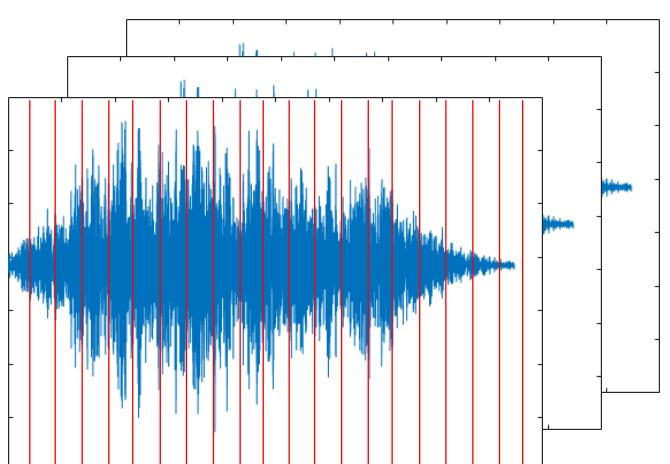
** - based on all channels



Methodology

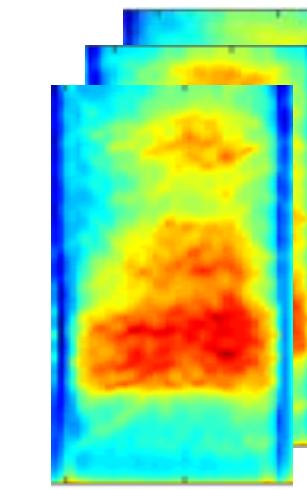


Methodology



DETERMINATION OF LPC COEFFICIENTS

PREDICTION RANK: 34



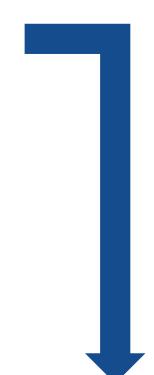
Filter Bank Energy (FBE)



Δ FBE

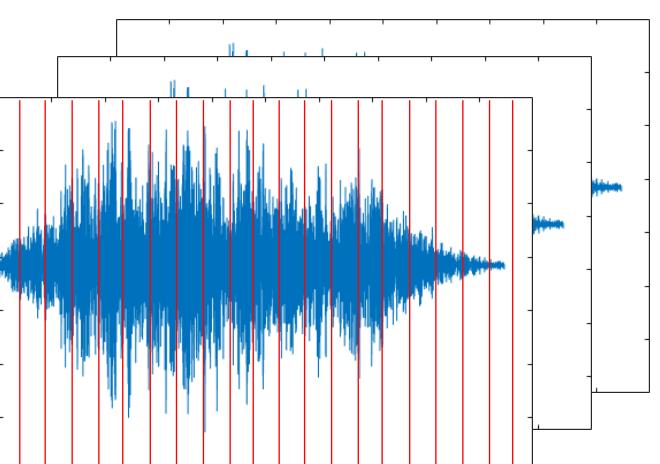


$\Delta\Delta$ FBE



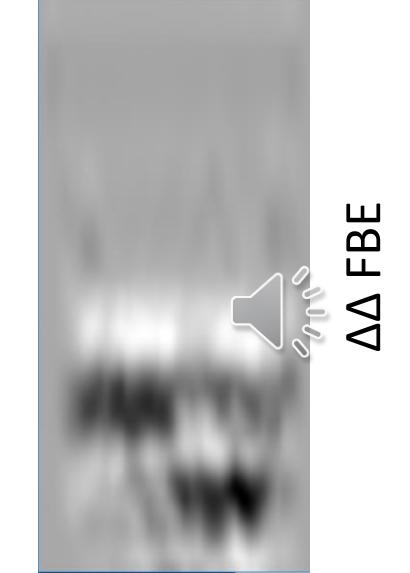
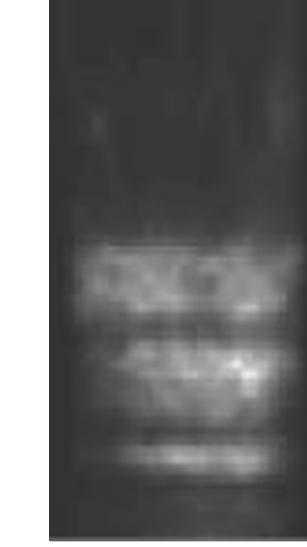
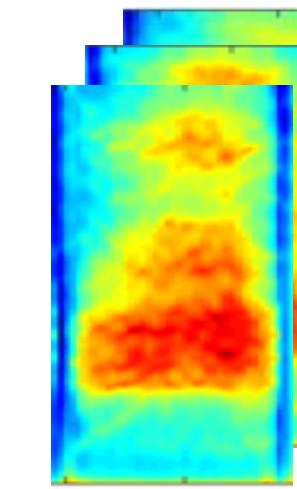
Methodology

39



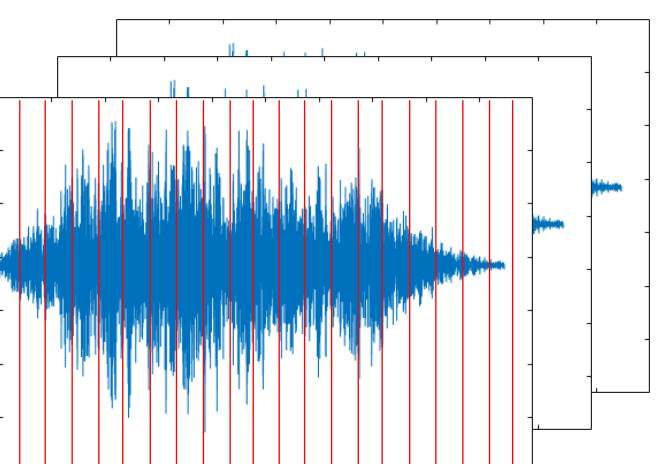
DETERMINATION OF LPC COEFFICIENTS

PREDICTION RANK: 34



Methodology

40

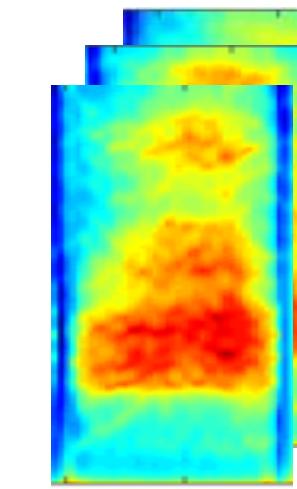


RMS



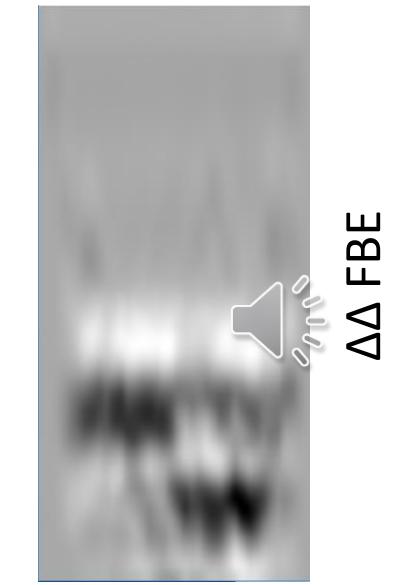
DETERMINATION OF LPC COEFFICIENTS

PREDICTION RANK: 34



FFT

FEATURE EXTRACTION

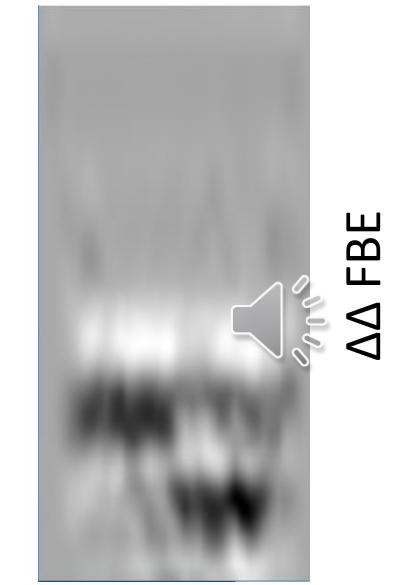


DETERMINATION OF LPC COEFFICIENTS

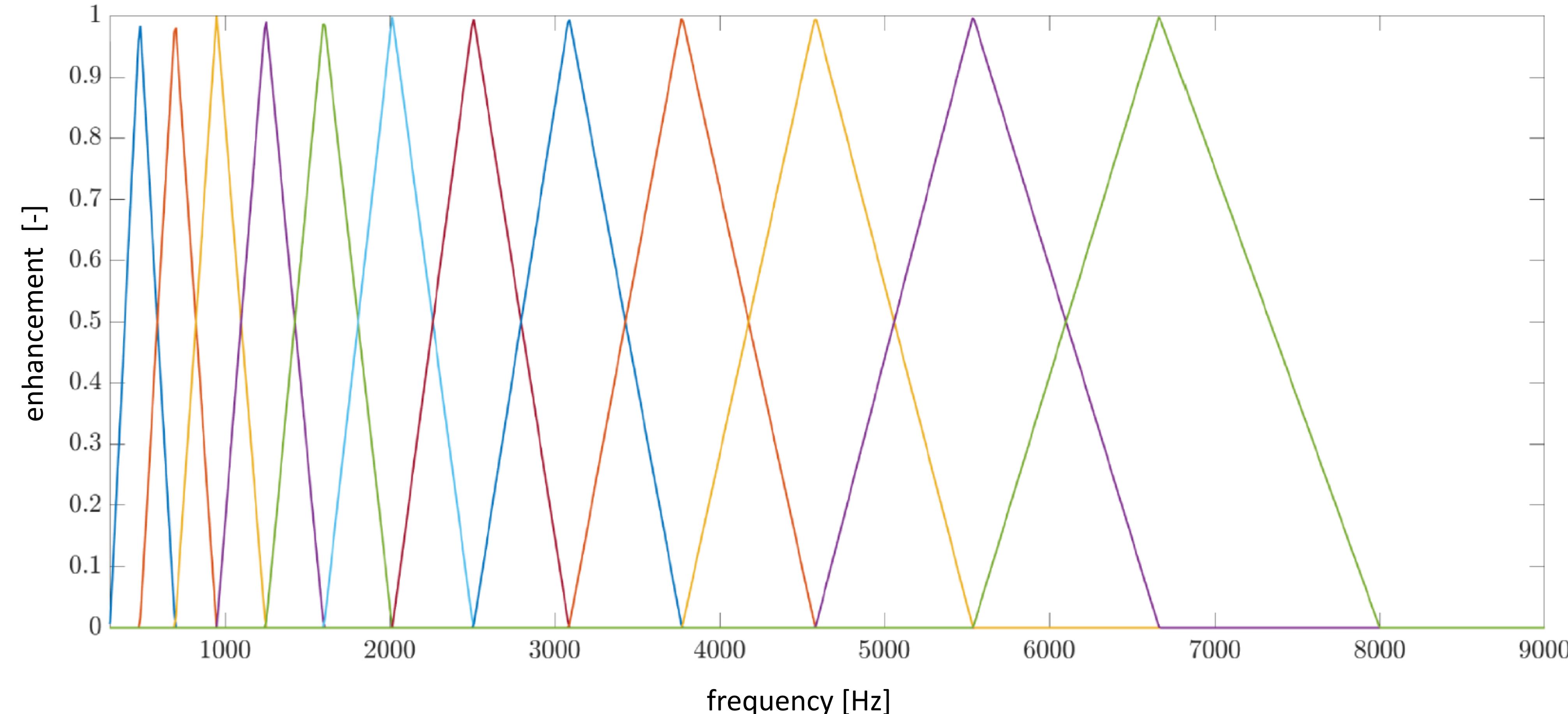
PREDICTION RANK: 34

FFT

FEATURE EXTRACTION

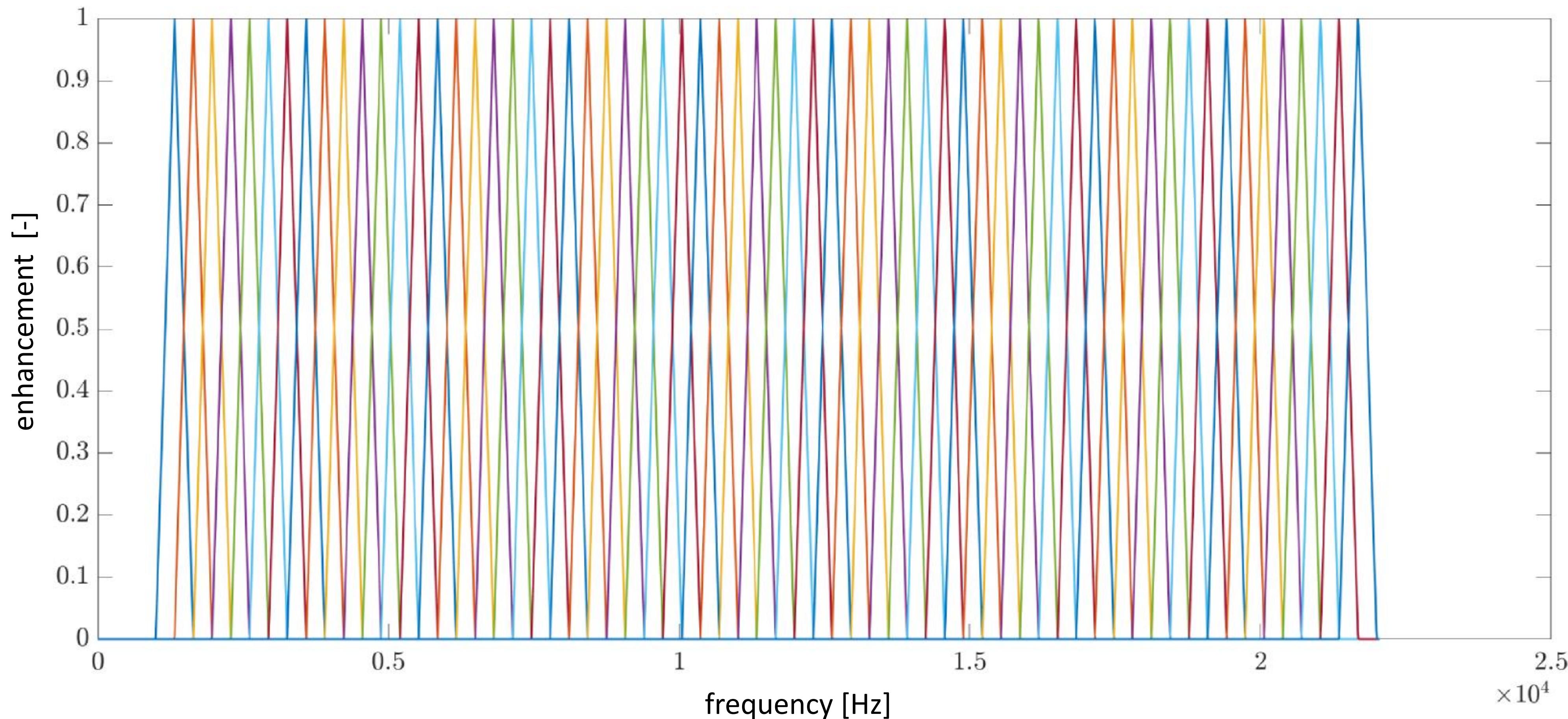


Methodology



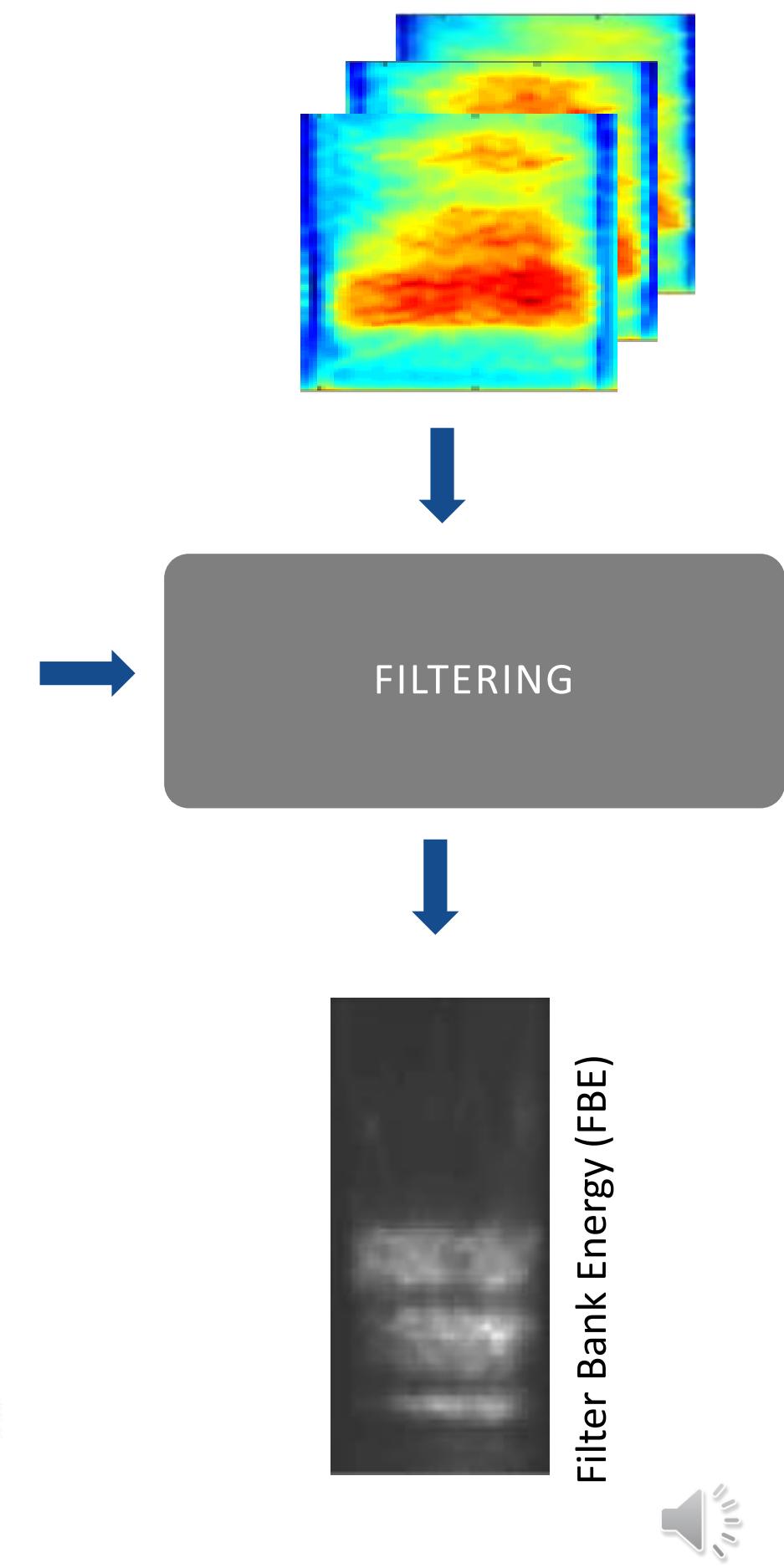
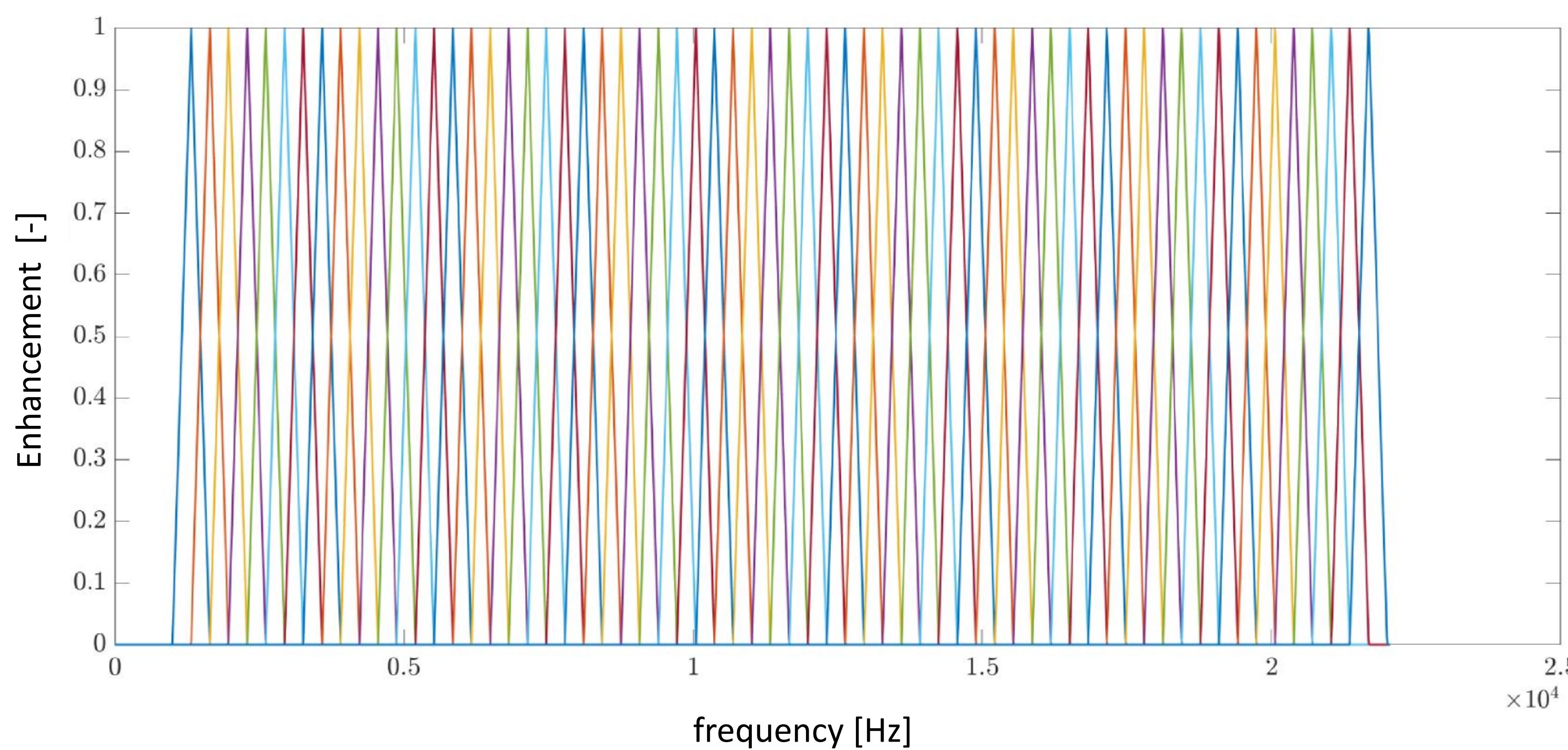
Methodology

42



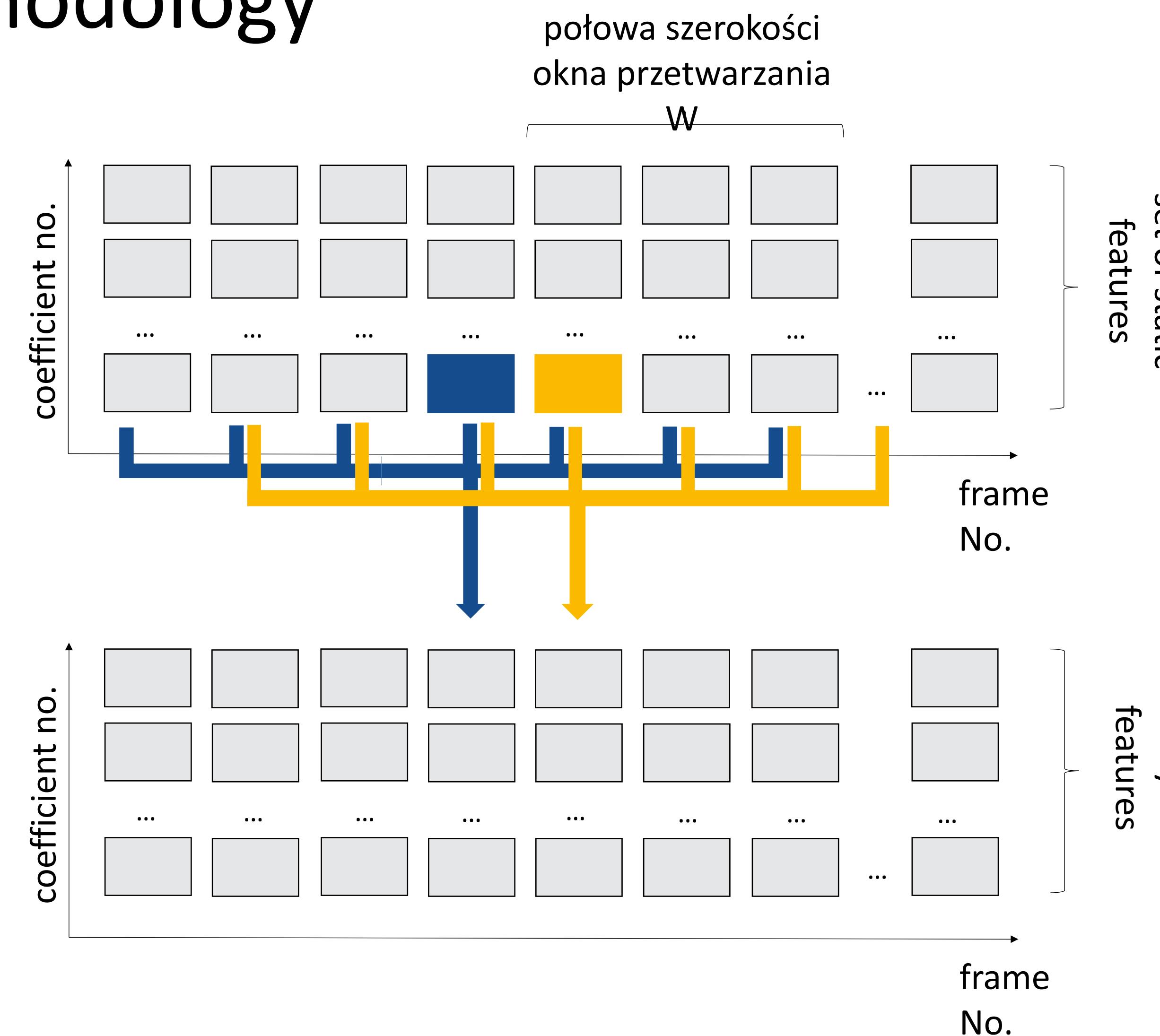
Methodology

43

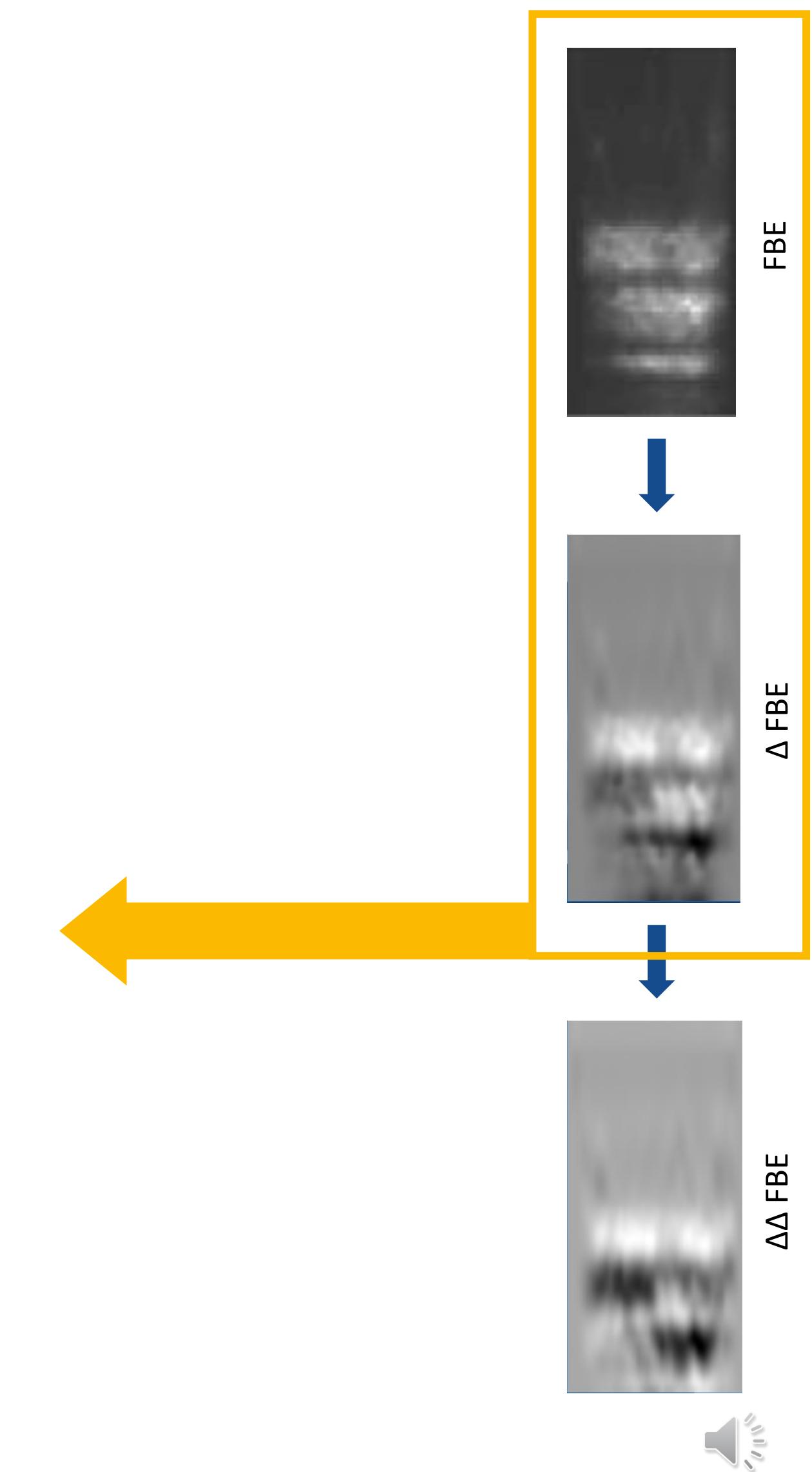


Methodology

np.: FBE 64 x 30



44
np.: FBE 64 x 30
+ ΔFBE 64 x 30

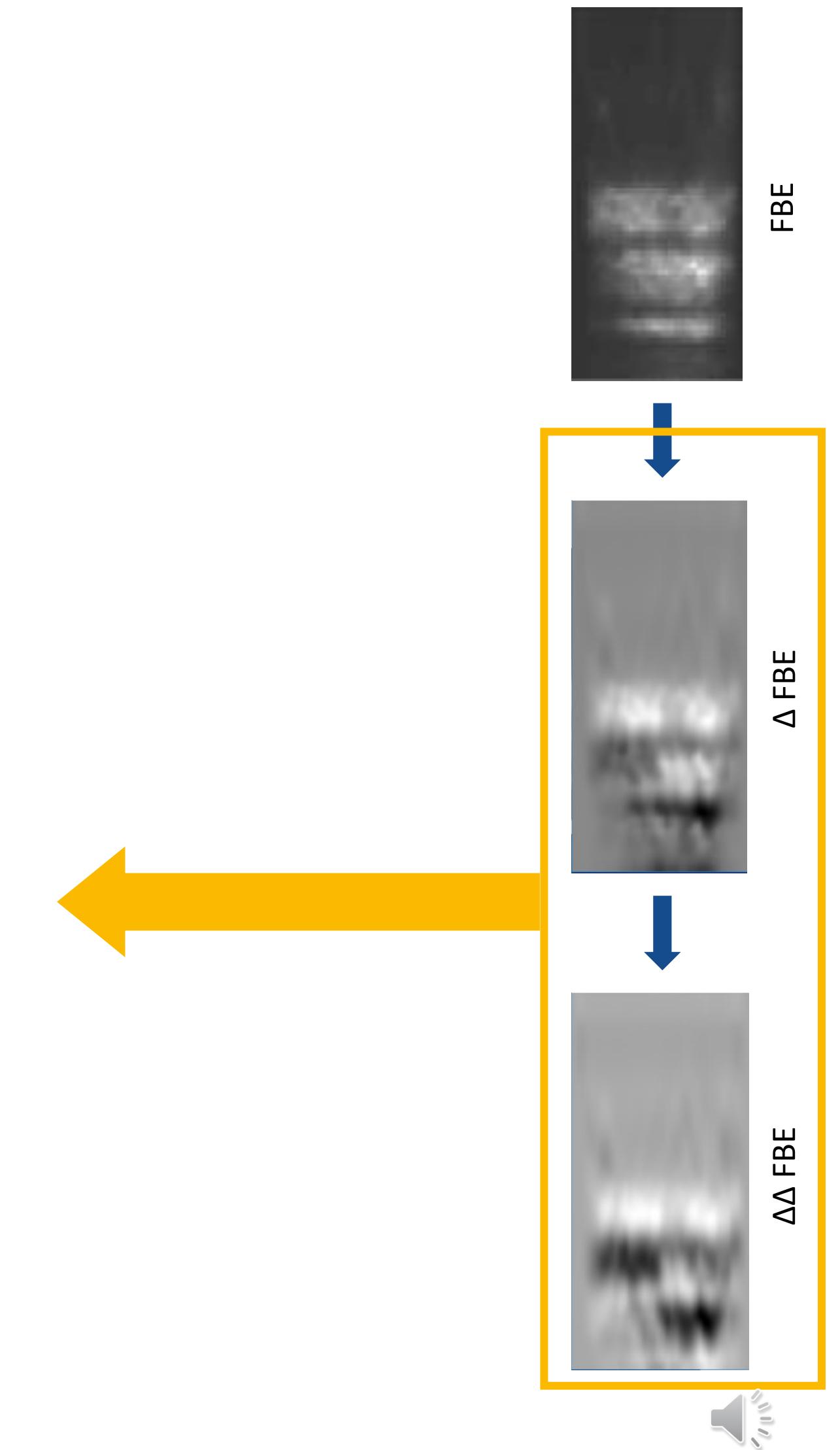
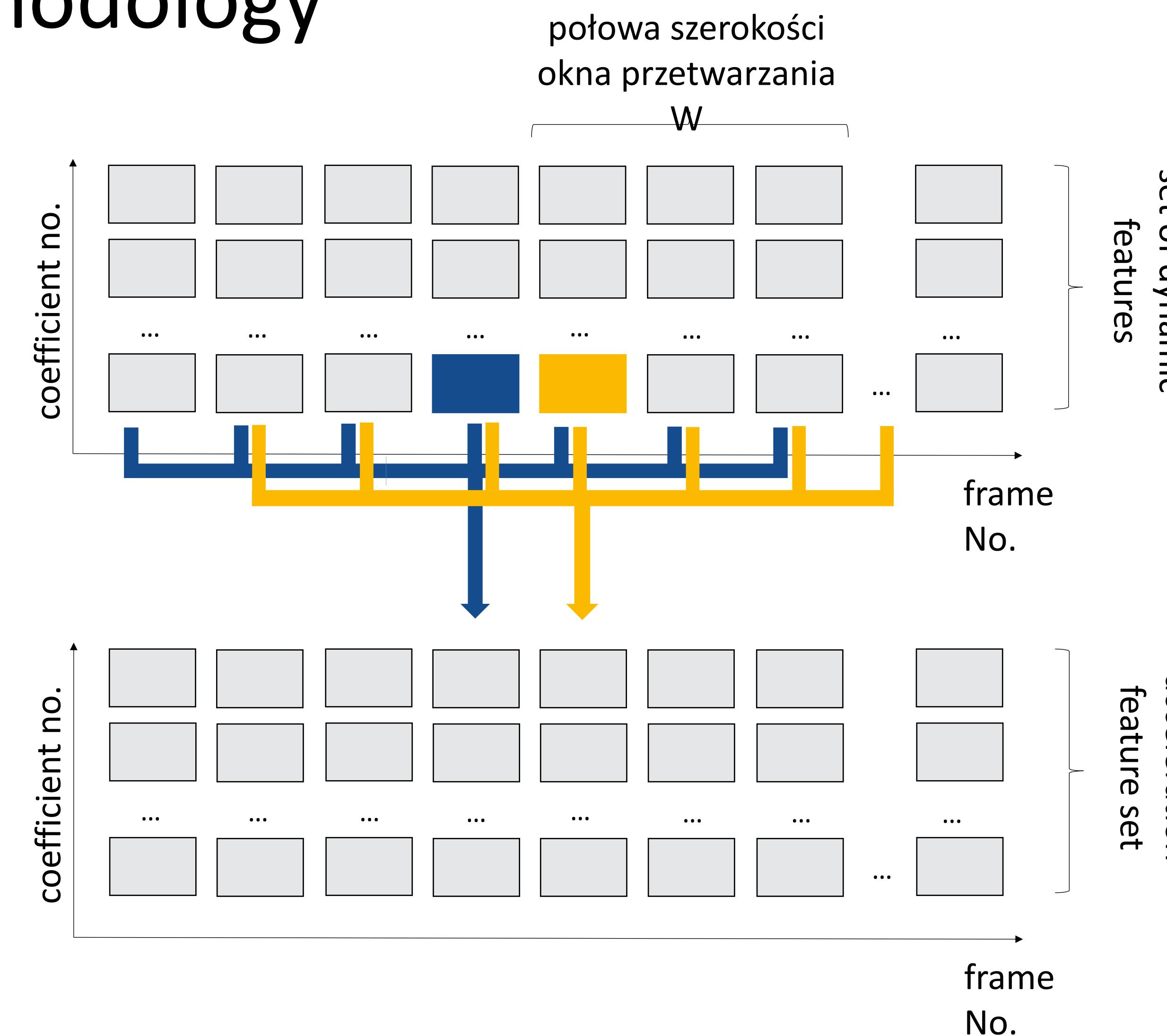


$W = 3$

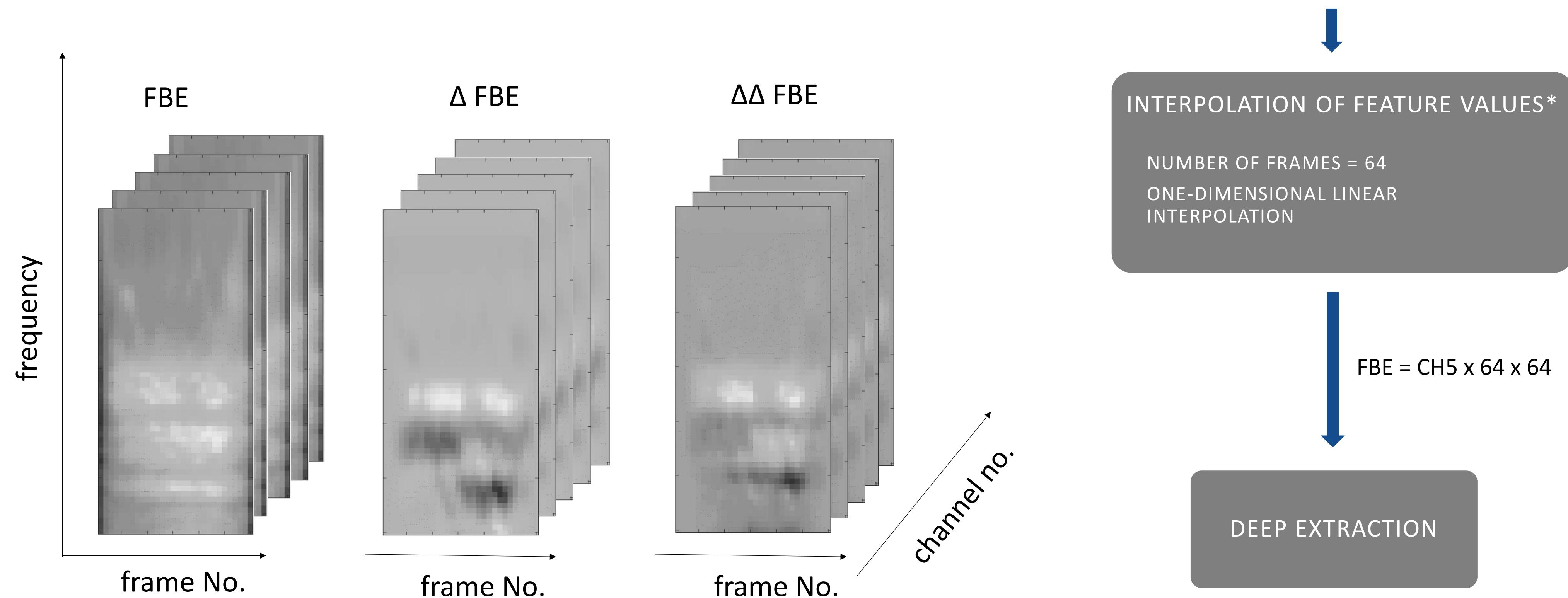


Methodology

np.: FBE 64 x 30
+ Δ FBE 64 x 30



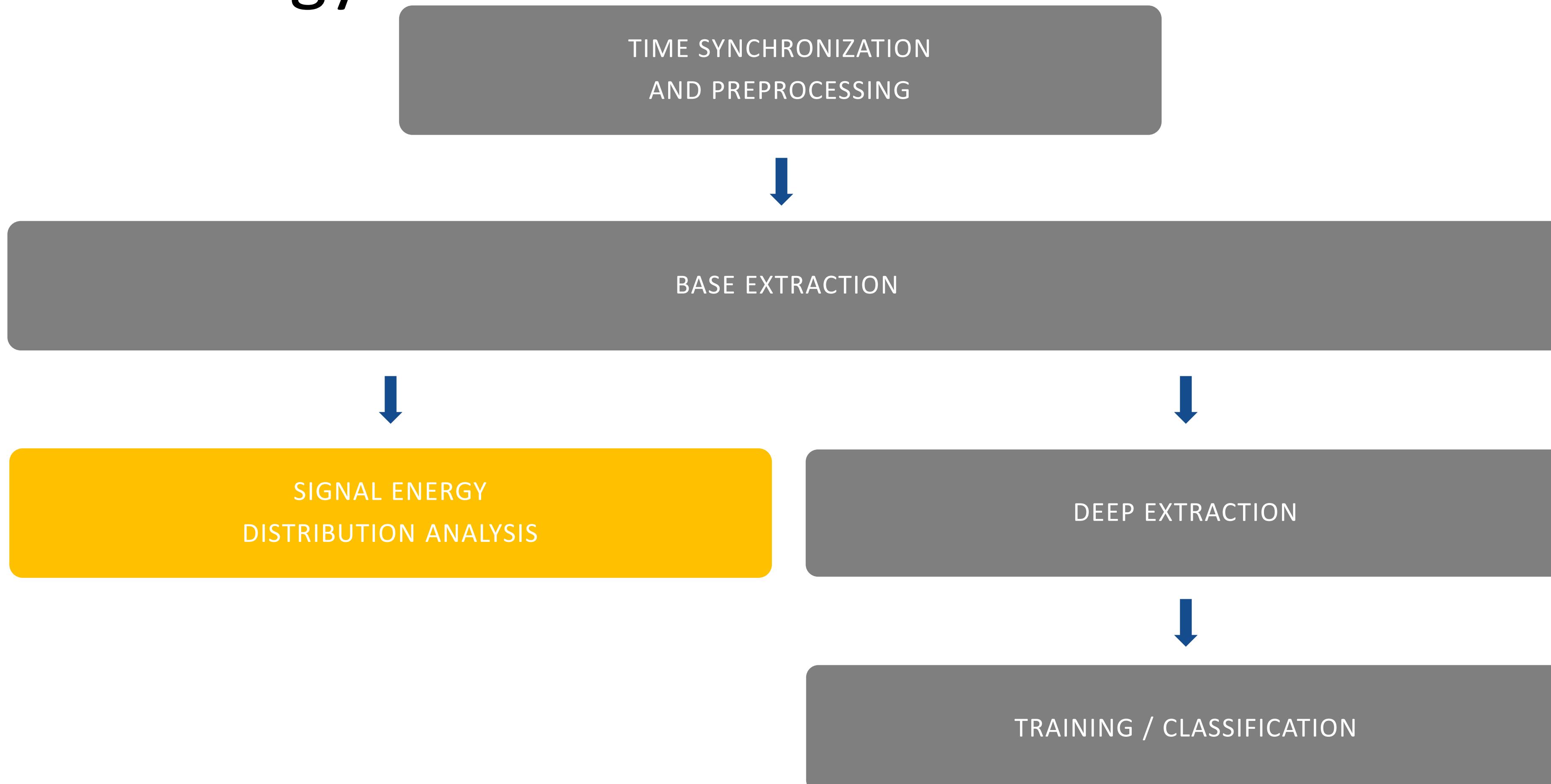
Methodology



* - independently for all features from a given channel



Methodology



Methodology

48

Class

normative realization for the sounds /s/ and /ʃ/

→ klasa s_{norm} i \int_{norm}

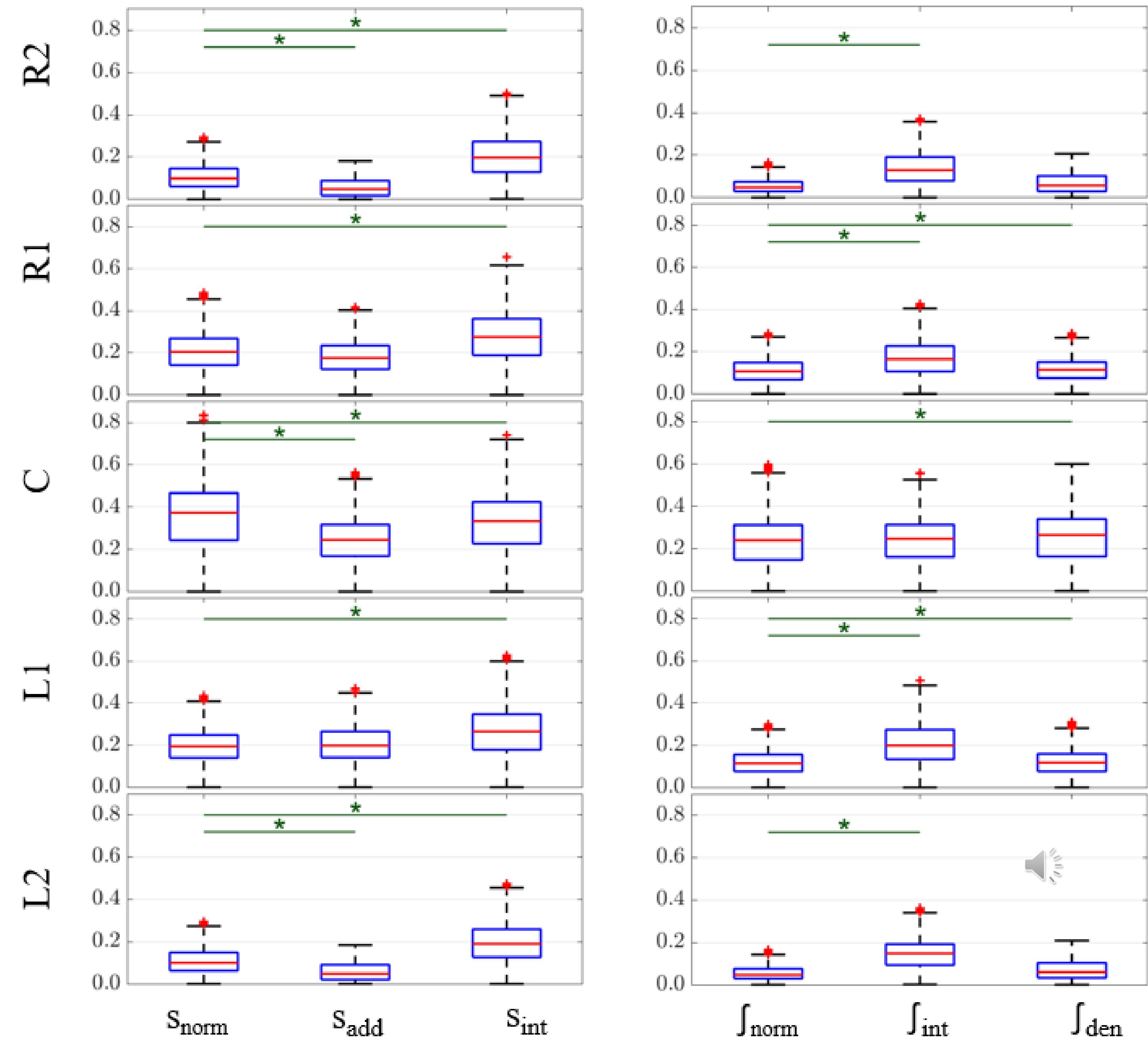
interdental and addental realization for the sound /s/

→ klasa s_{int} i s_{add}

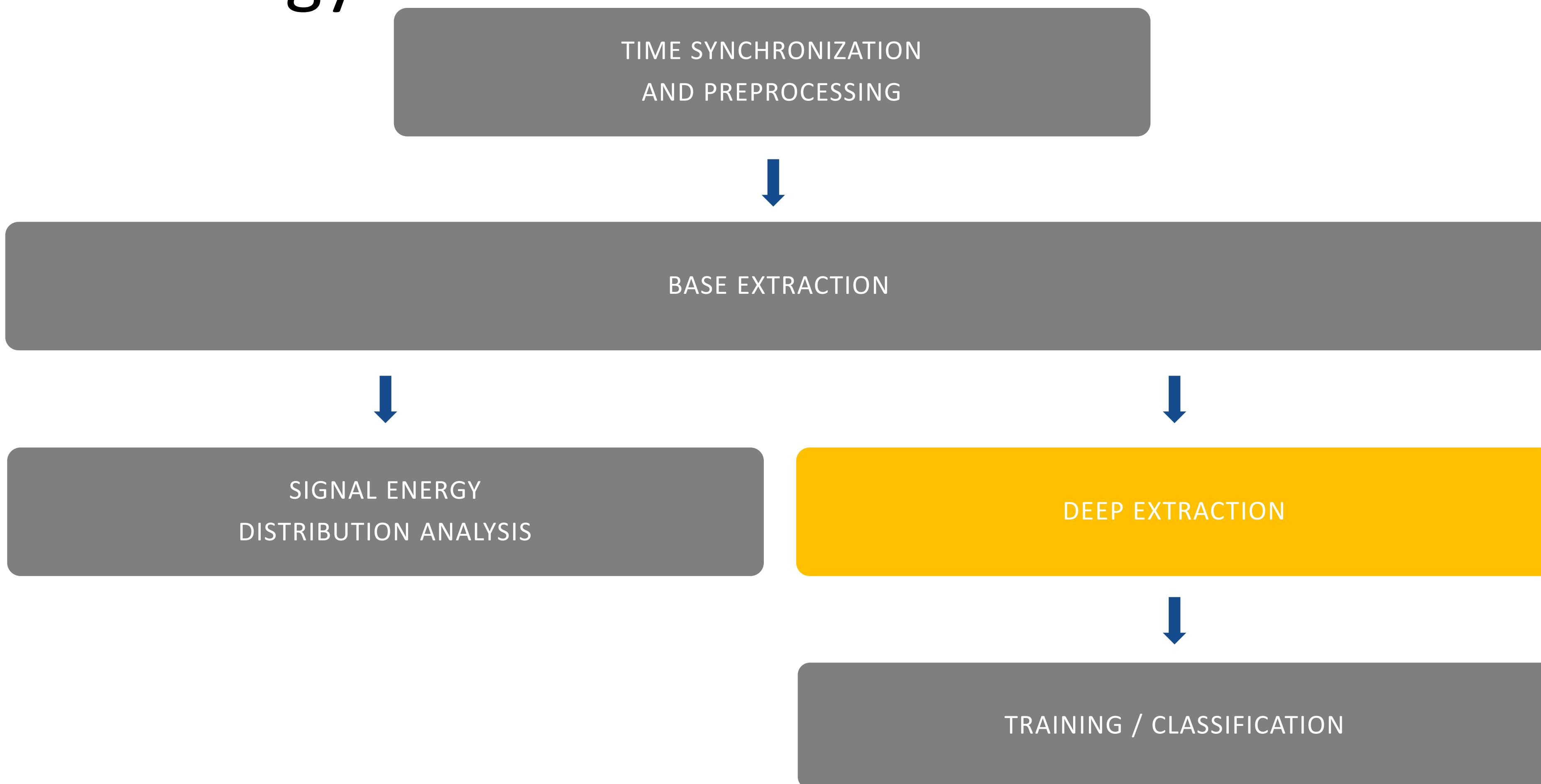
interdental and dental realization for the sound /ʃ/

→ klasa \int_{int} i \int_{den}

* - significant difference at the level of $p < 0.05$

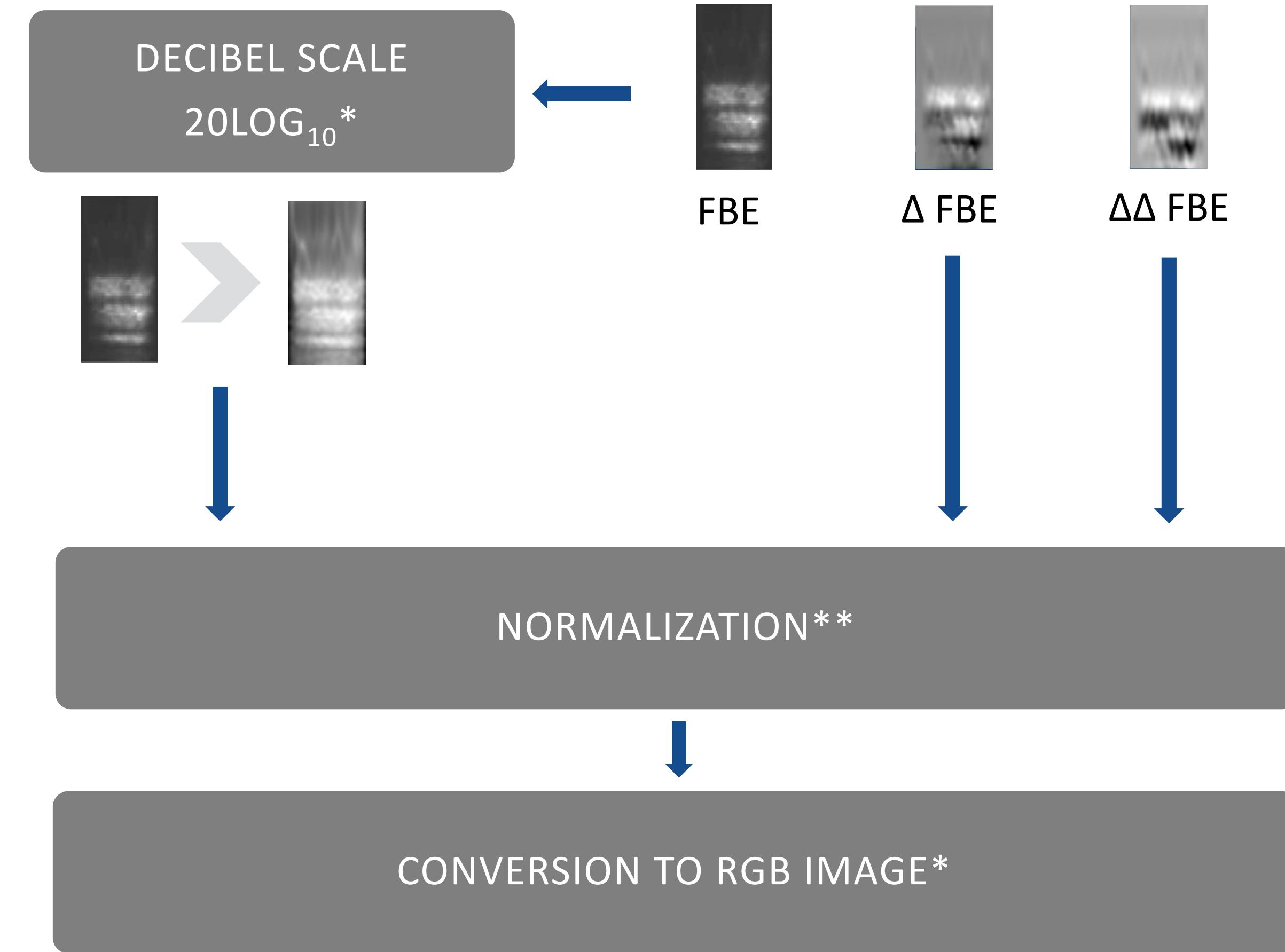


Methodology



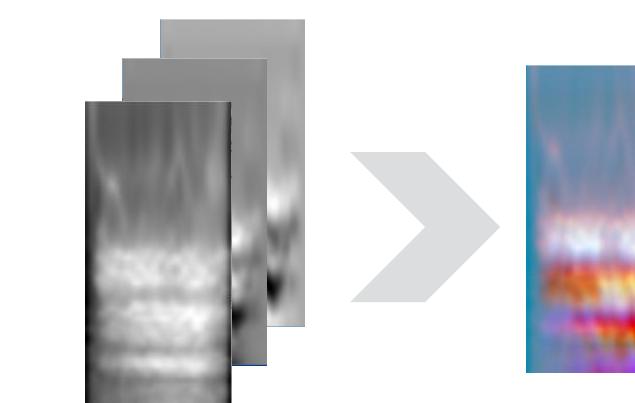
Methodology

50

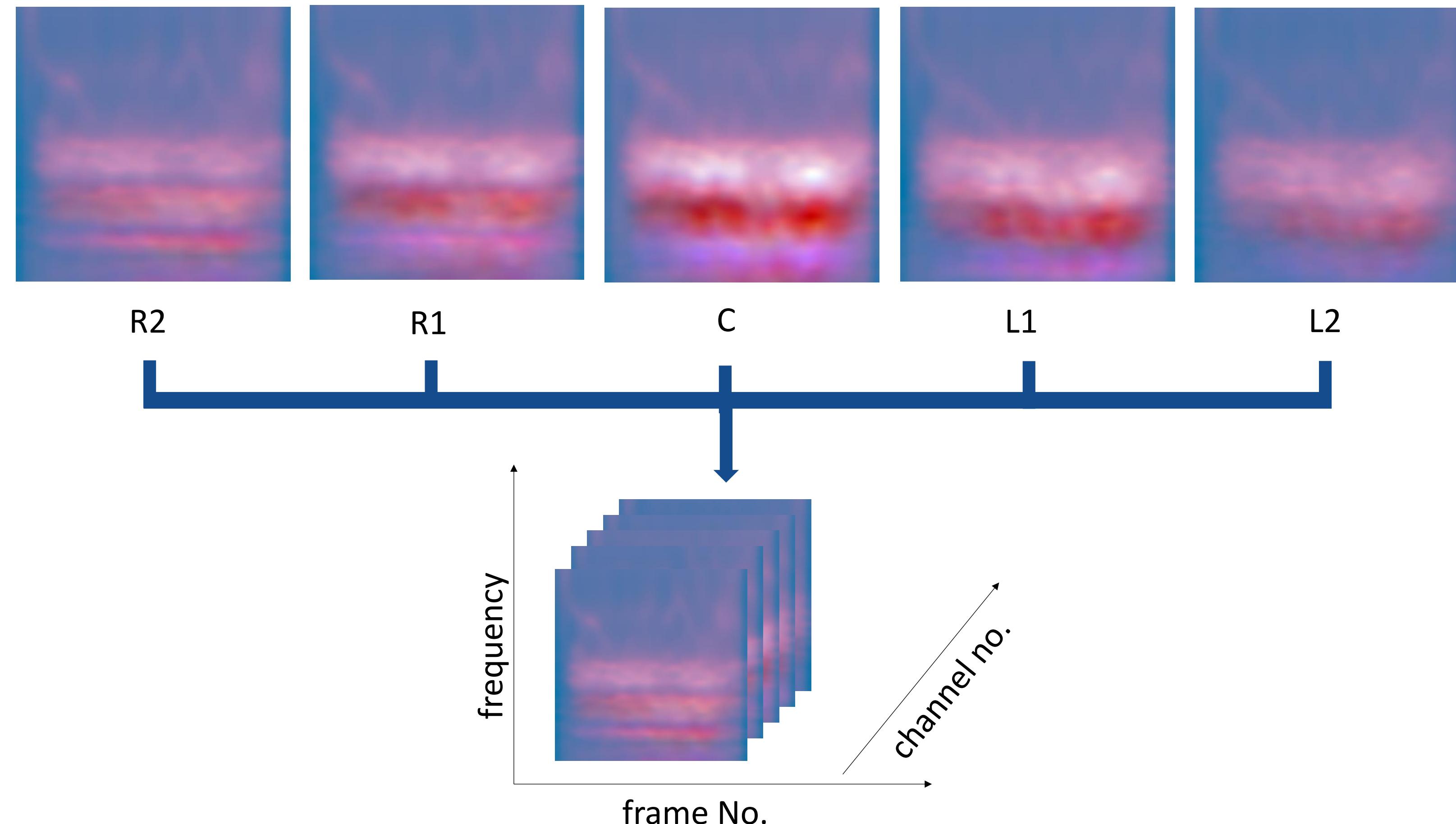


* - independently for each channel

** - based on all channels

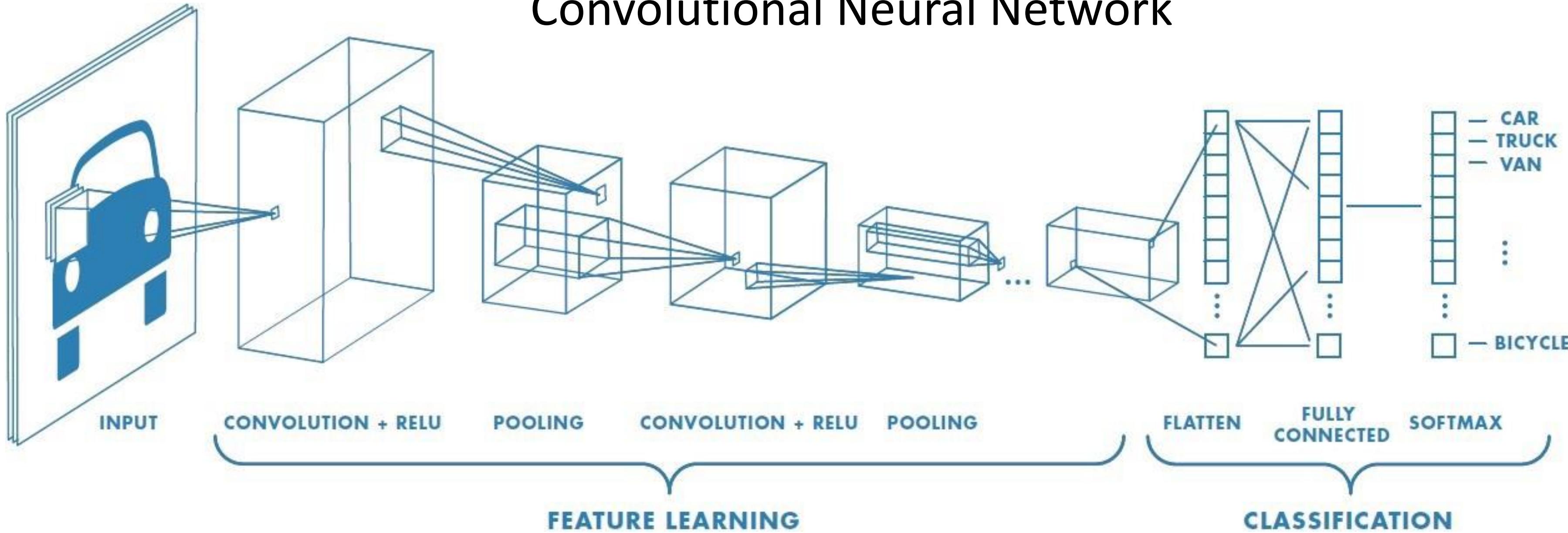


Methodology



Methodology

Convolutional Neural Network

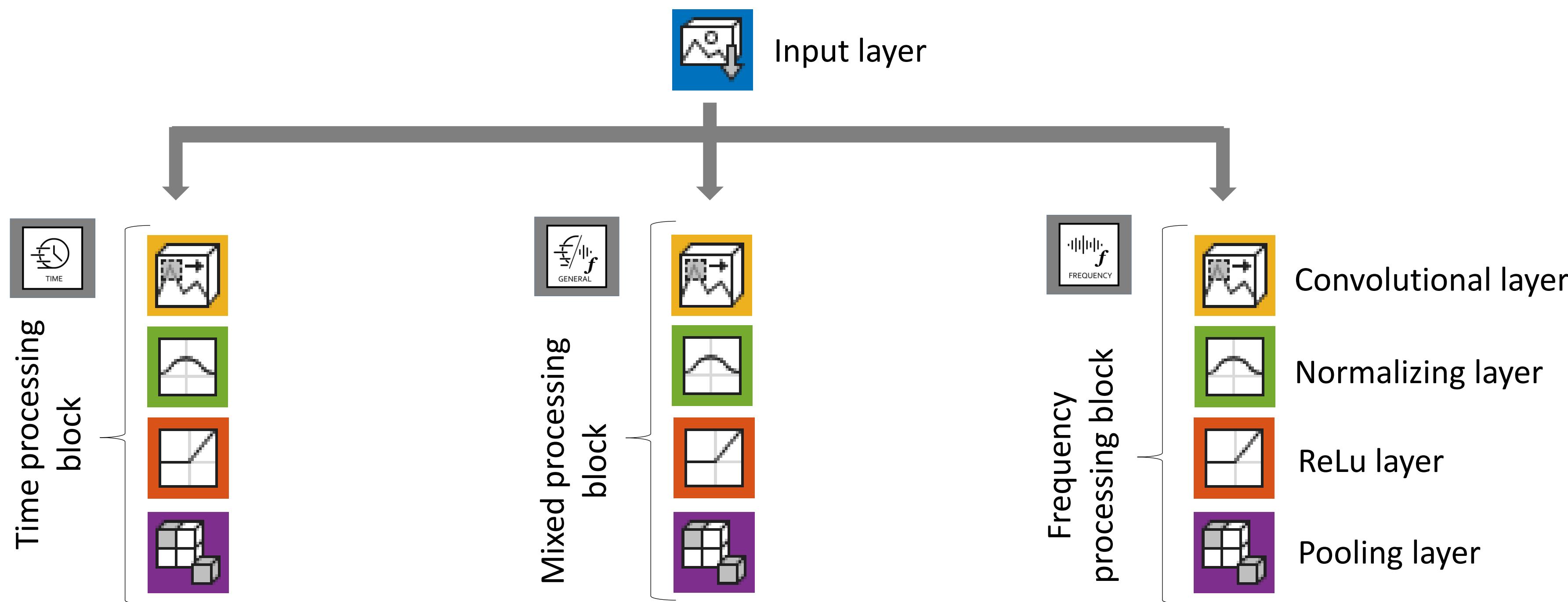


<https://www.mathworks.com/>



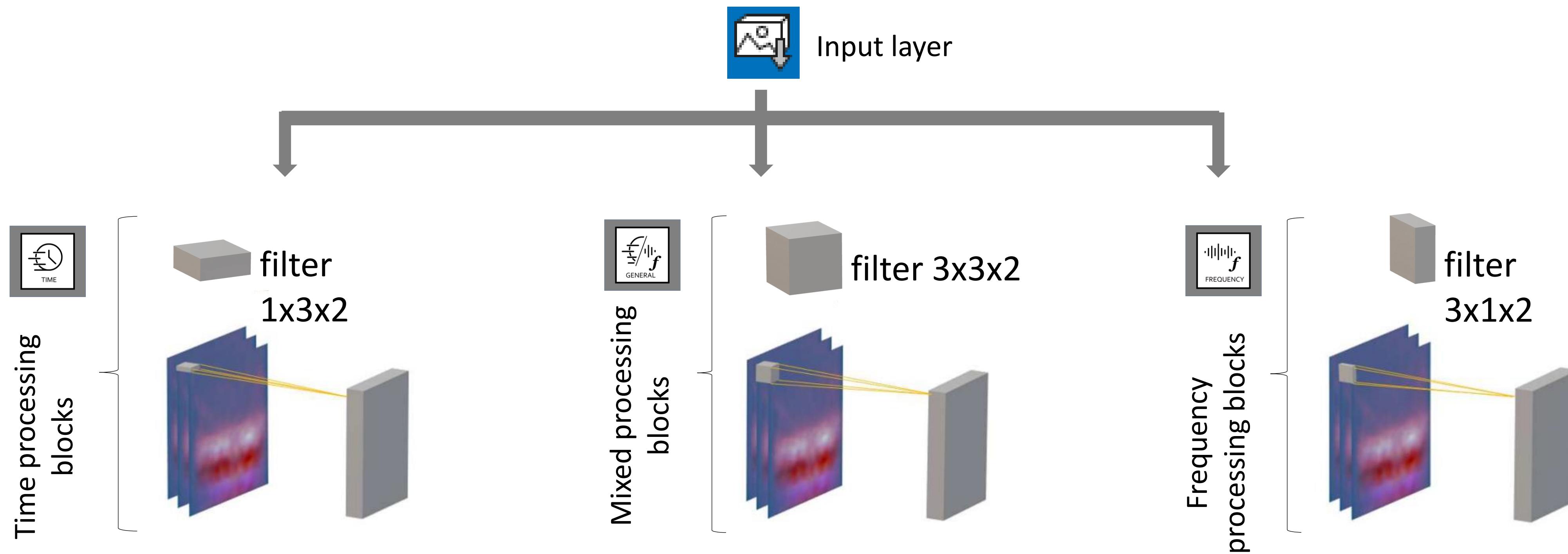
Methodology

53



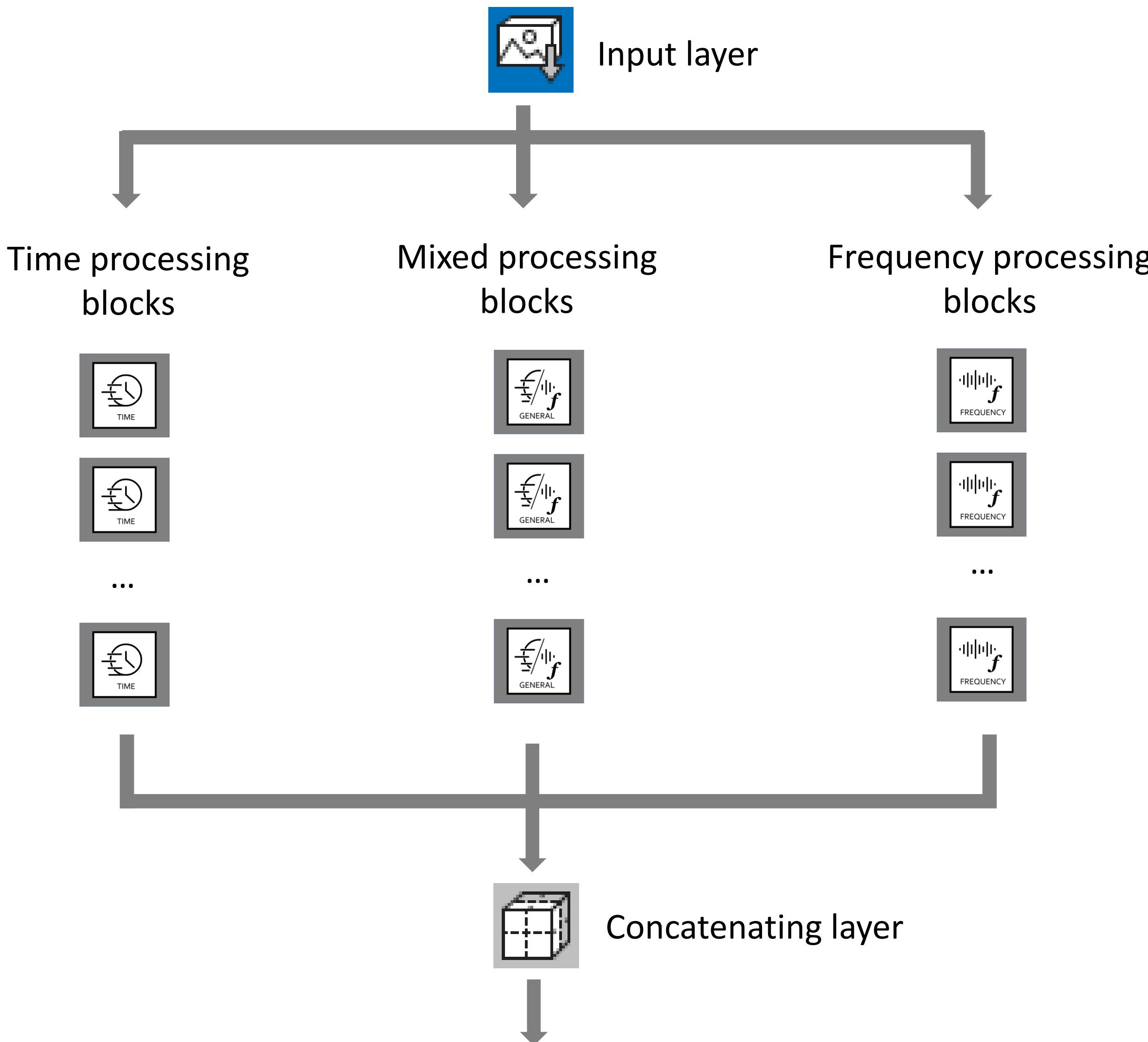
Methodology

54



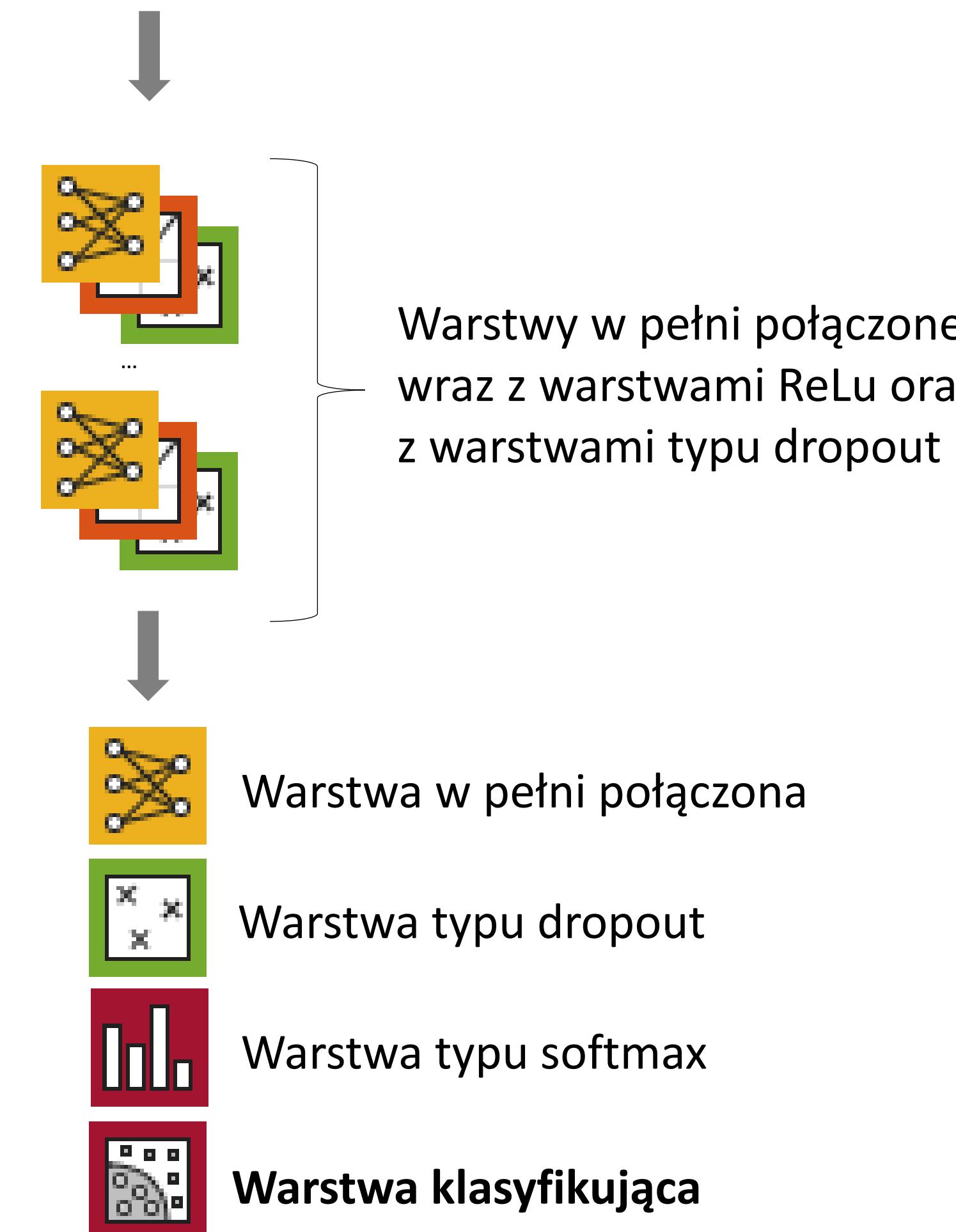
Methodology

55

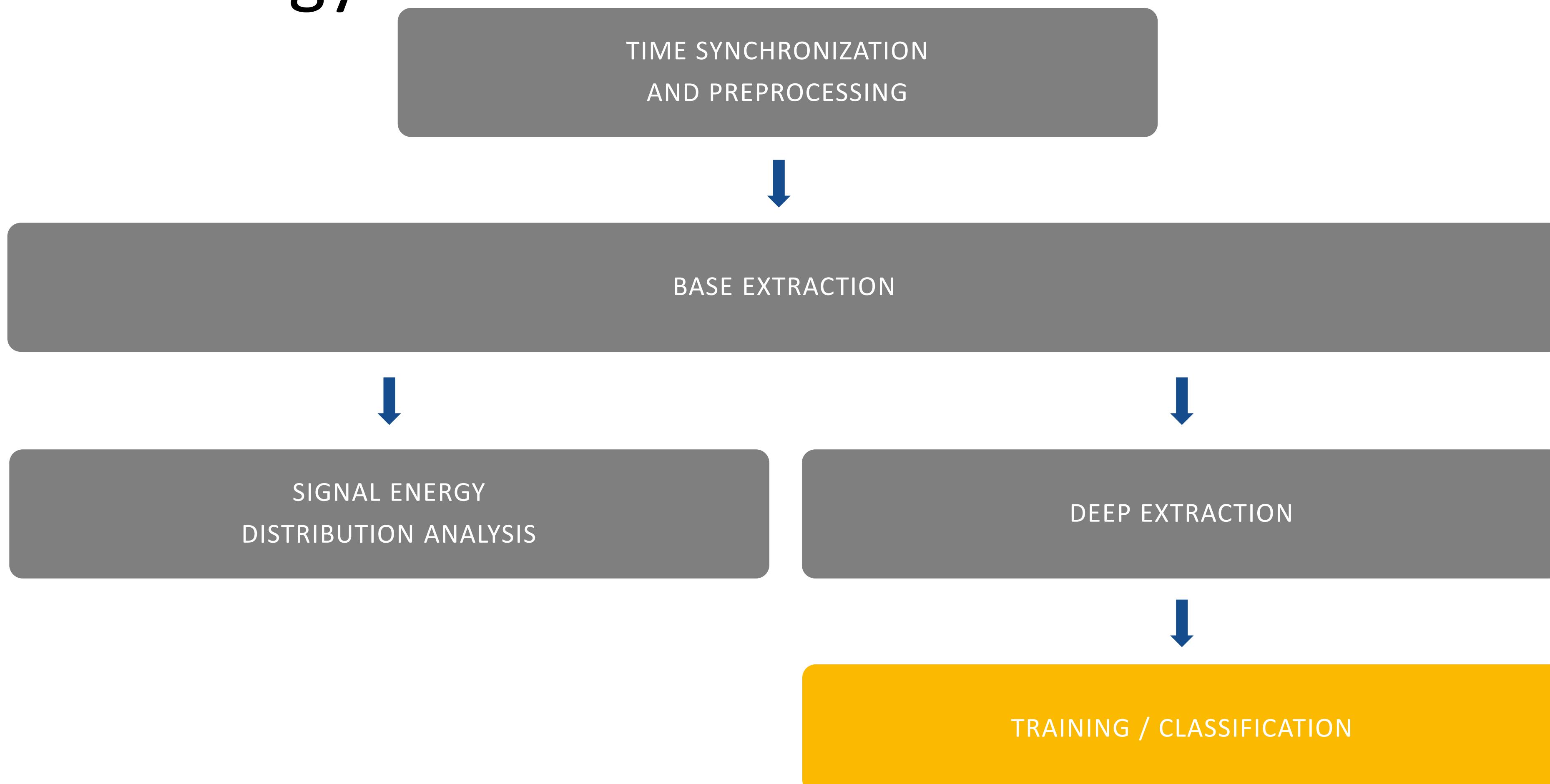


Methodology

56

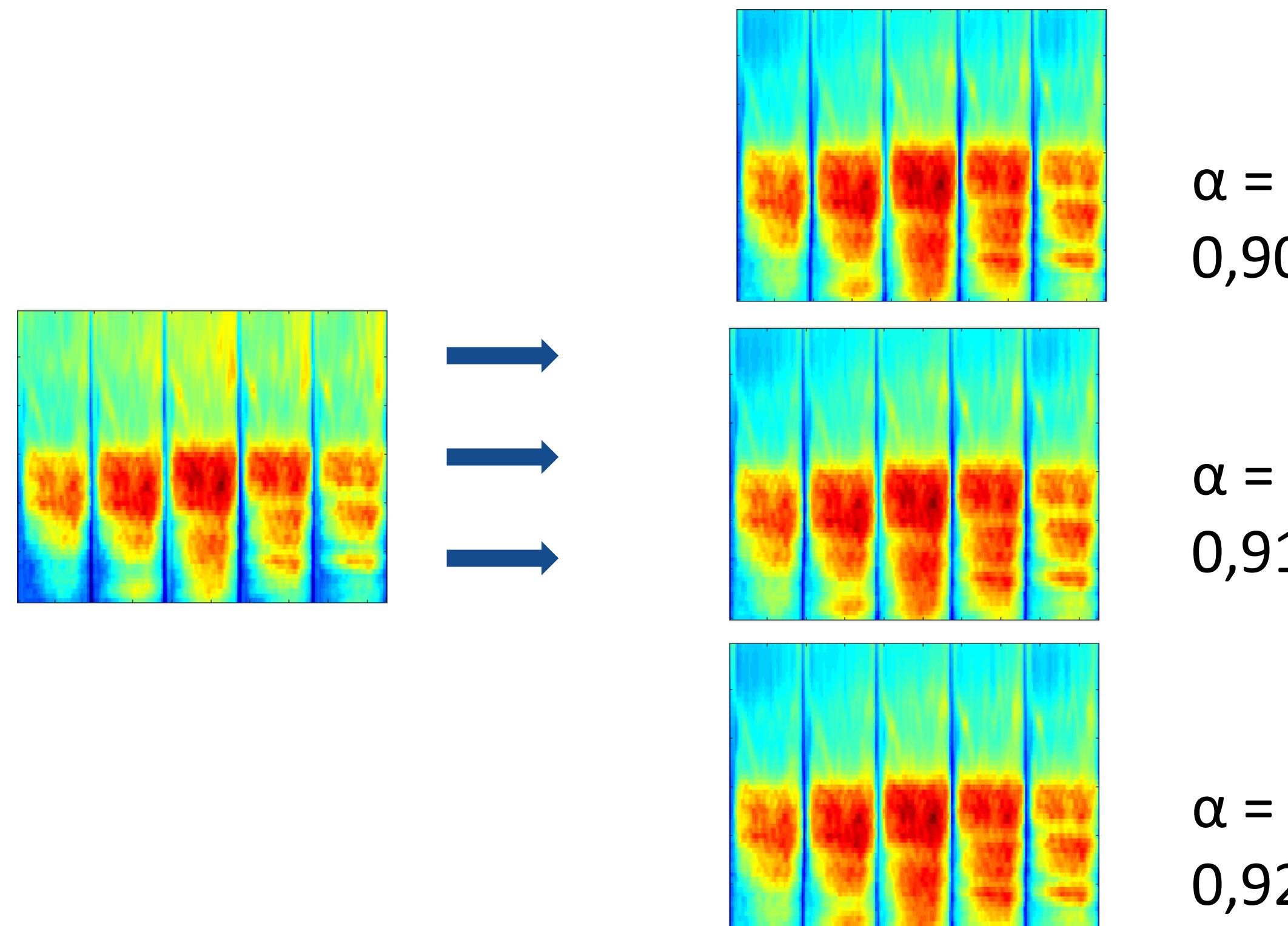


Methodology



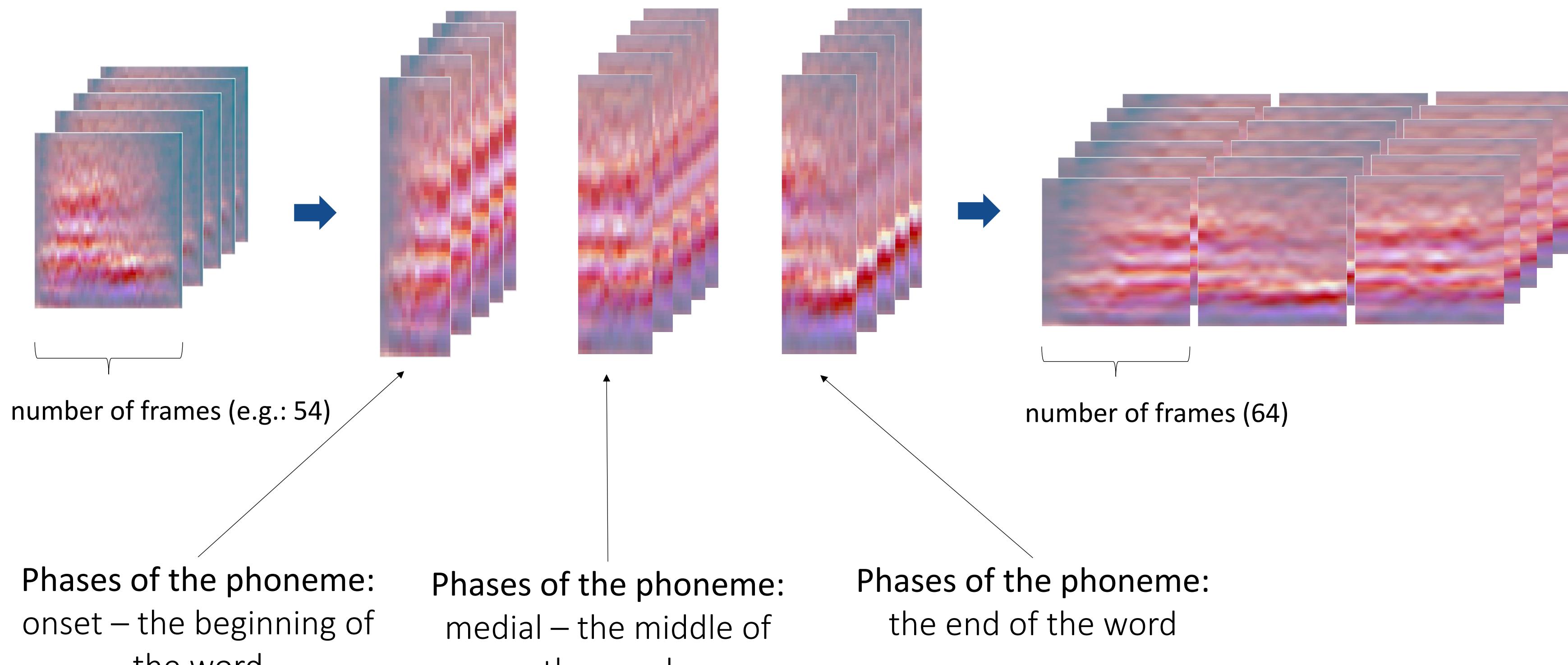
Methodology

58



Methodology

59



Methodology

Network parameter name	Parameter value
	$3 \times 1 \times 2$
filter size	$1 \times 3 \times 2$
	$3 \times 3 \times 2$
input data size	$64 \times 64 \times 5 (\times 3)$
model depth - number of processing blocks	5
model width - number of filters	16-64
training batch size	128

Summary of the parameters of the proposed model used during the experiments



Methodology

61

Network parameter name	Parameter value	Learning process configuration	Hyperparameter value
filter size	3 x 1 x 2	optimization algorithms	SGDM
	1 x 3 x 2		
	3 x 3 x 2		
input data size	64 x 64 x 5 (x 3)	number of epochs	25
model depth - number of processing blocks	5	number of epochs to stop the learning process early	5
model width - number of filters	16-64	learning rate coefficient	0,001
training batch size	128	learning rate update coefficient	0,1
		number of epochs to change the learning rate	10
		L2	0,01
		momentum	0,95

Summary of parameters and hyperparameters of the proposed model used during the experiments



Experiments

ANALYSIS OF SIBILANTS MODELS

62

PHONEME
CLASSIFICATION

SENSITIVITY ANALYSIS OF THE CNN MODEL



Experiments

ANALYSIS OF
SIBILANTS MODELS

63

PHONEME
CLASSIFICATION

SENSITIVITY ANALYSIS OF THE CNN MODEL



Experiments

PHONE ME CLASSIFICATION

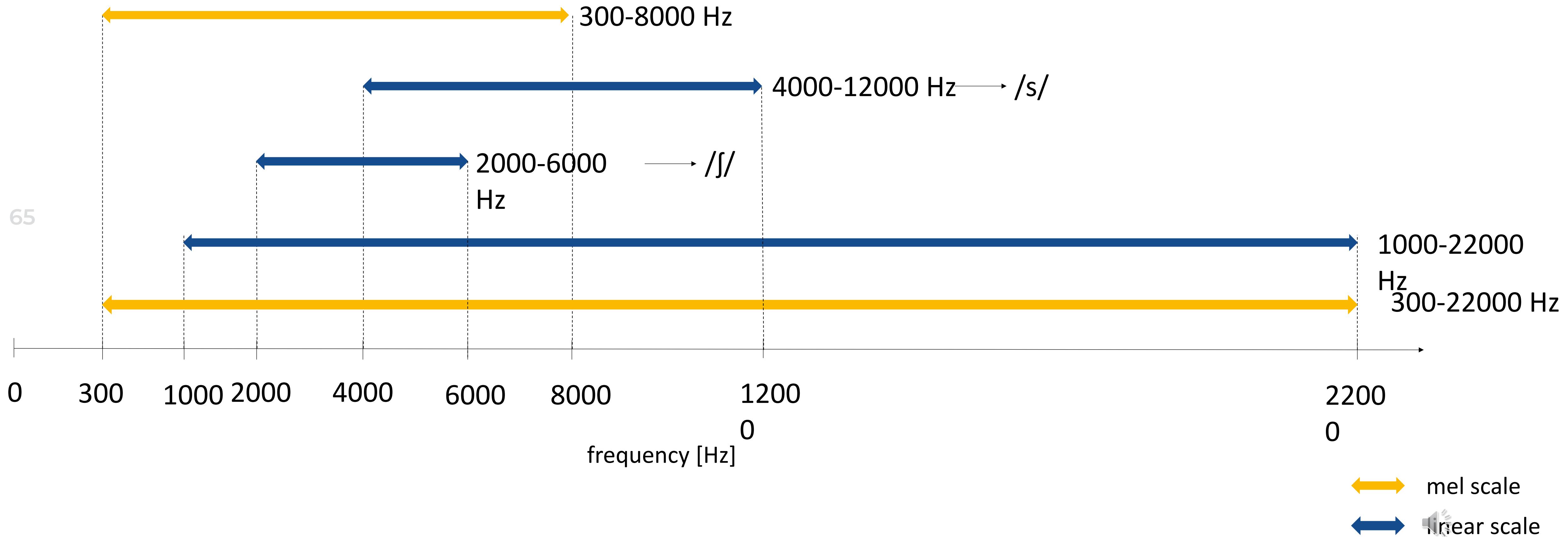
64

- frequency range of the filter set,
- number of channels (1-CH vs 5-CH vs 15-CH),
- different variants of processing blocks (PB1 vs PB2).



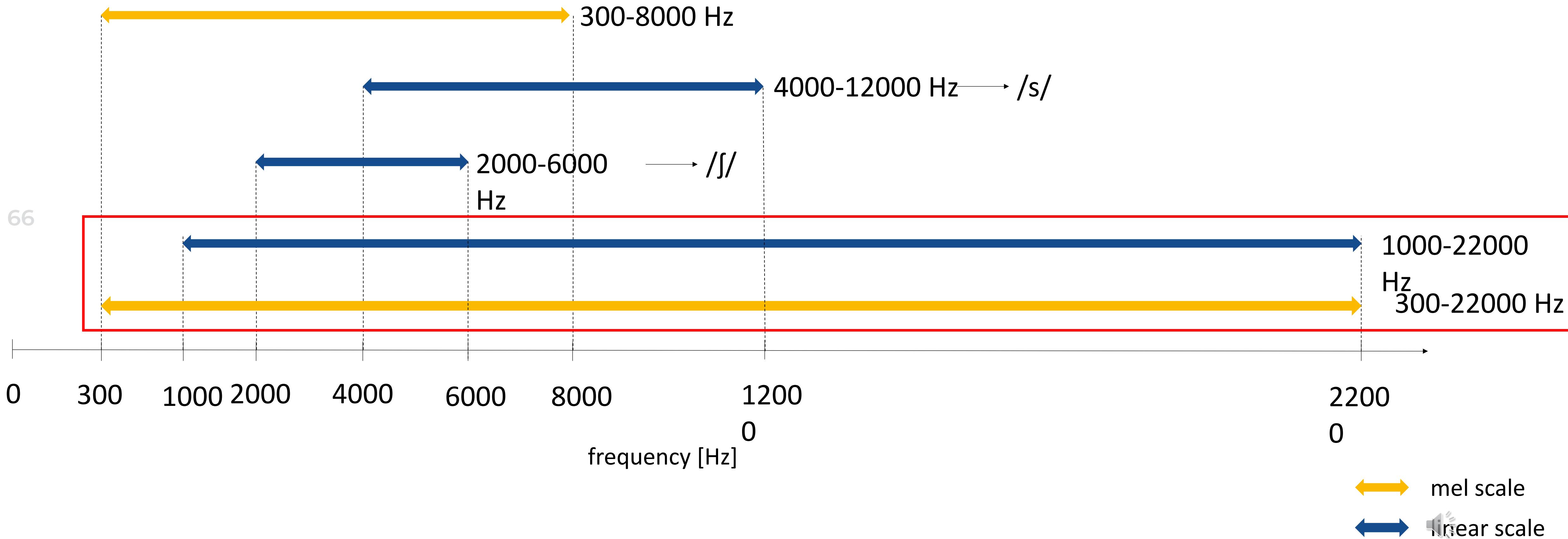


Experiments





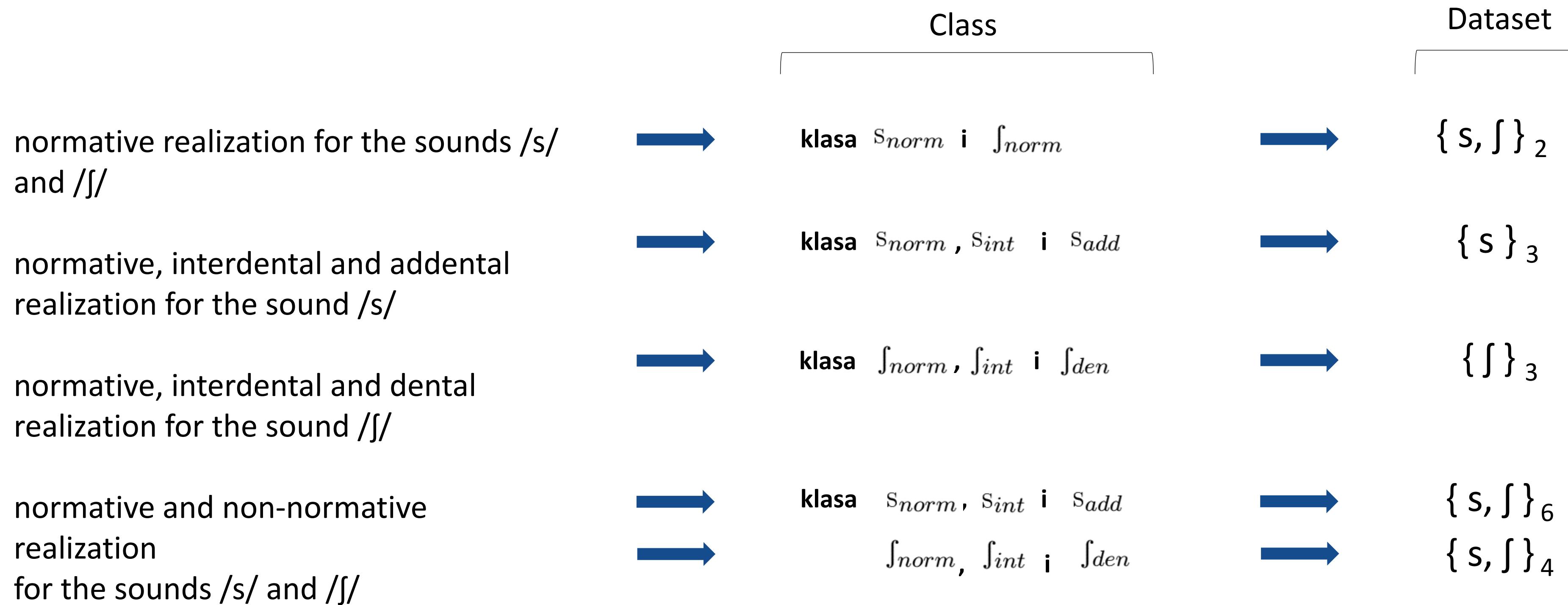
Experiments





Experiments

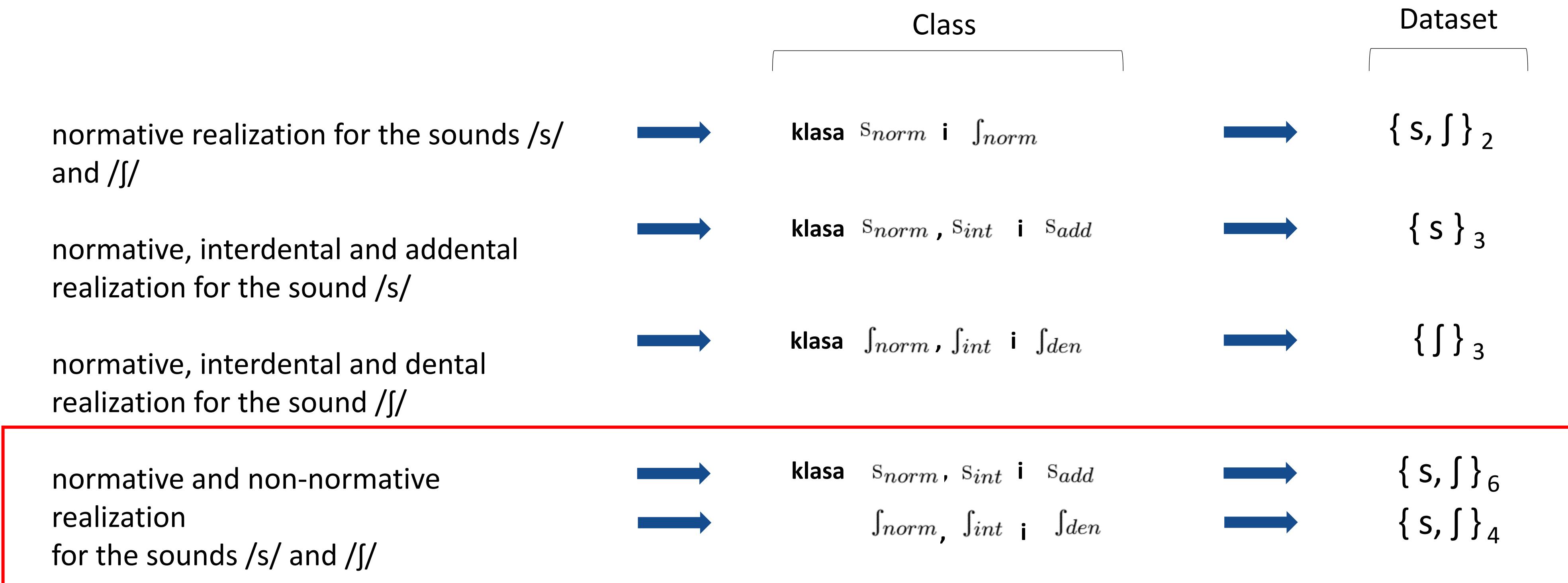
67





Experiments

68





Experiments

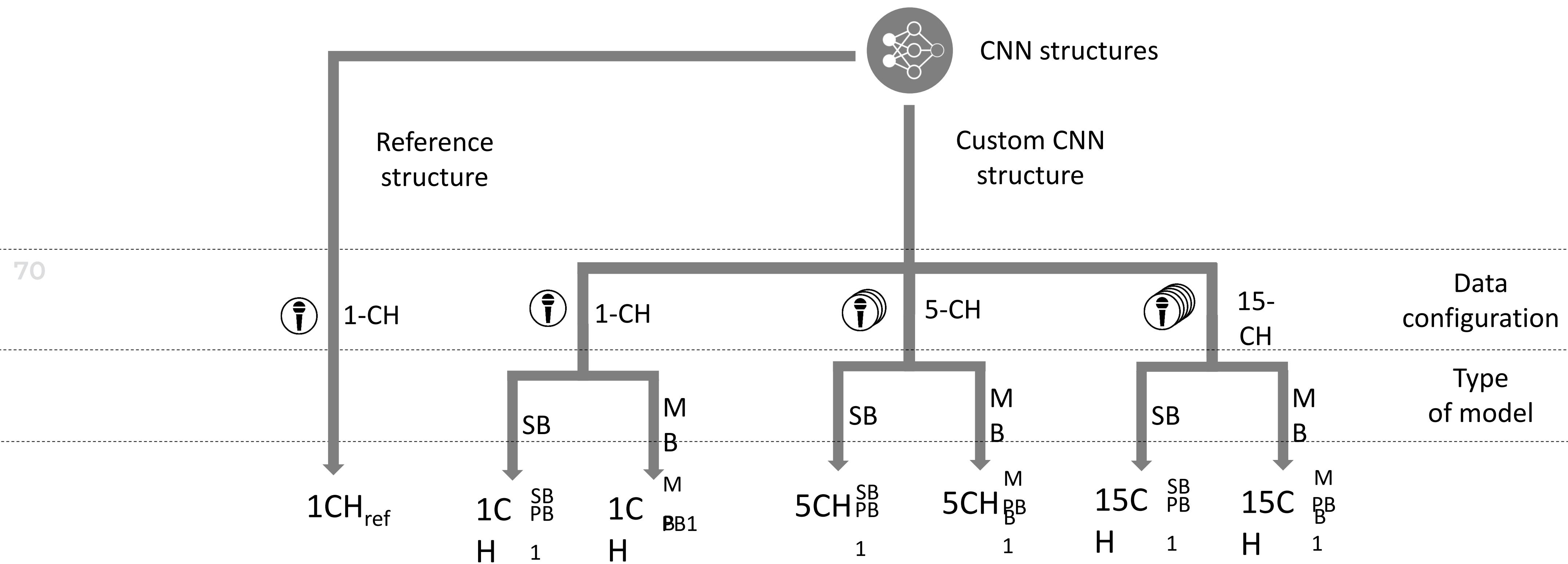
	Type of realisation	/s/			/ʃ/		
		normative	addentalna	interdentalna	normative	dentalna	interdentalna
Type of set	training 70%	151 (906)	64 (384)	116 (696)	163 (978)	54 (324)	84 (504)
	validation 15%	33	14	26	35	12	18
	test 15%	33	14	26	35	12	18

The size of the dataset (before and after augmentation) in recognizing different realizations of the sounds /s/ and /ʃ/



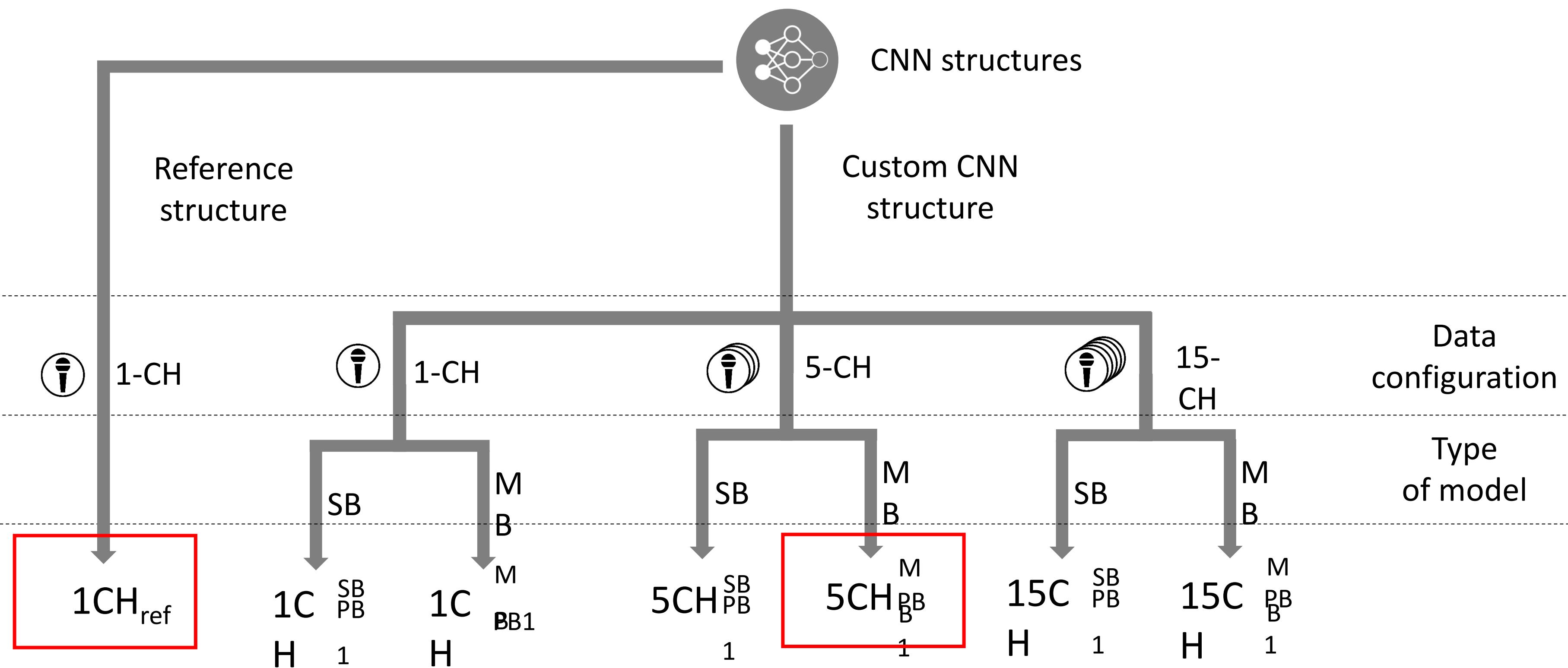


Experiments



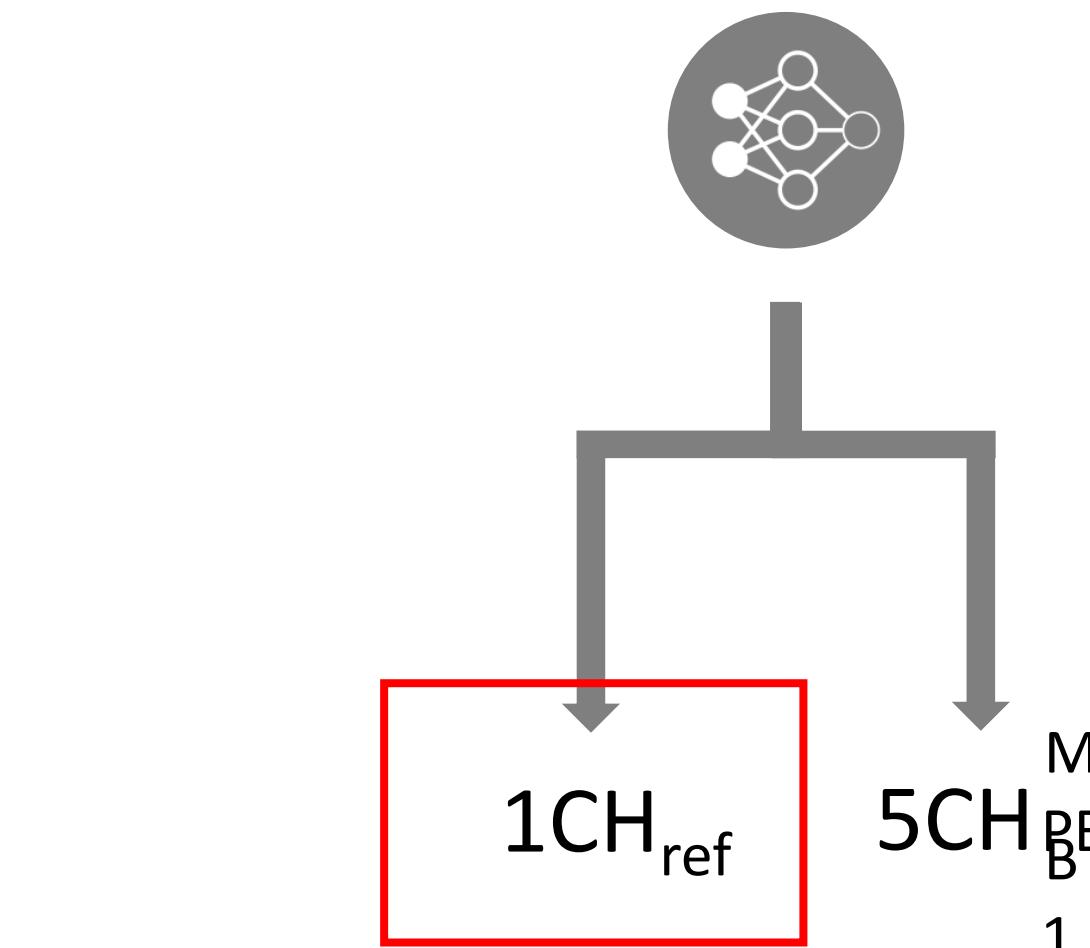


Experiments

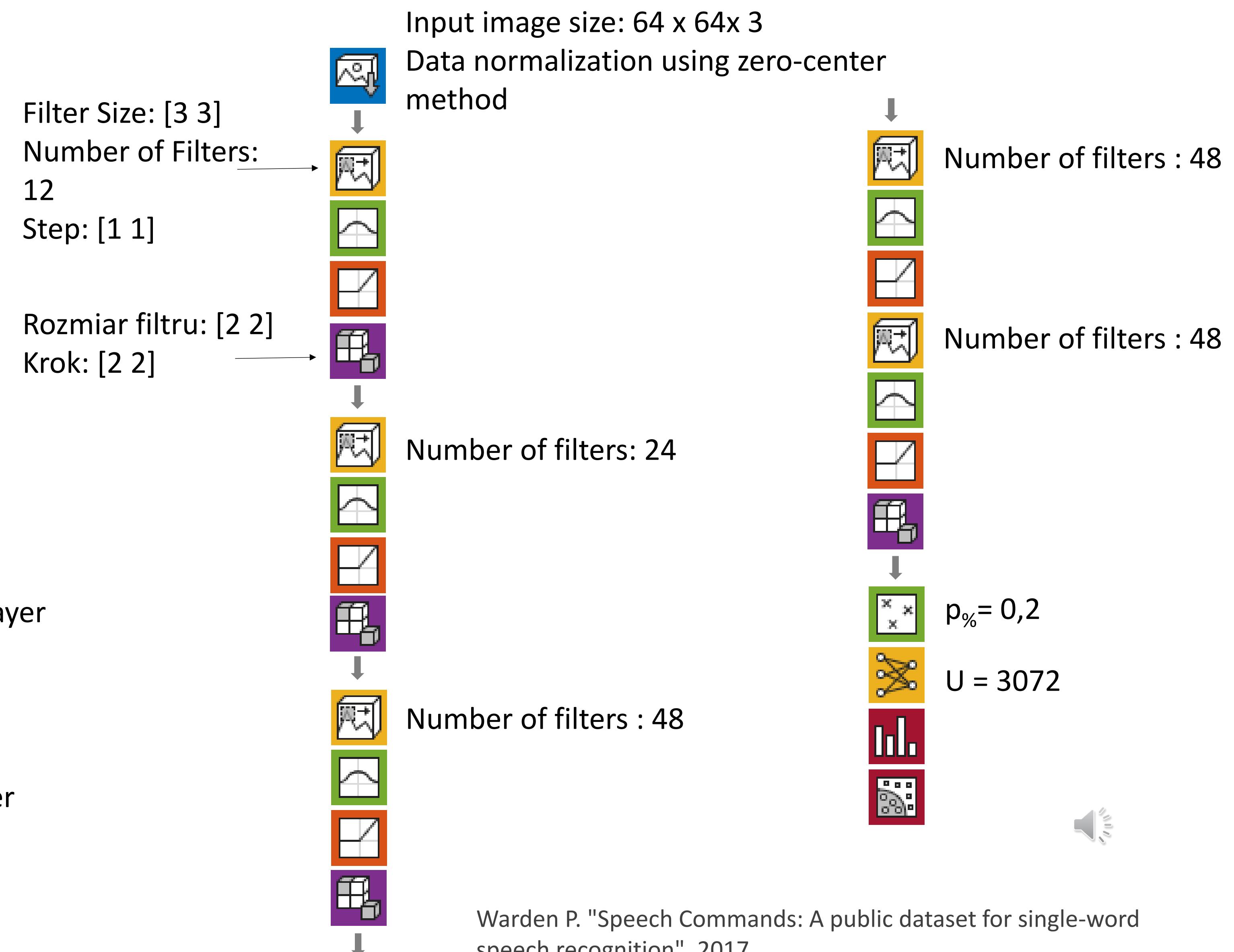


Experiments

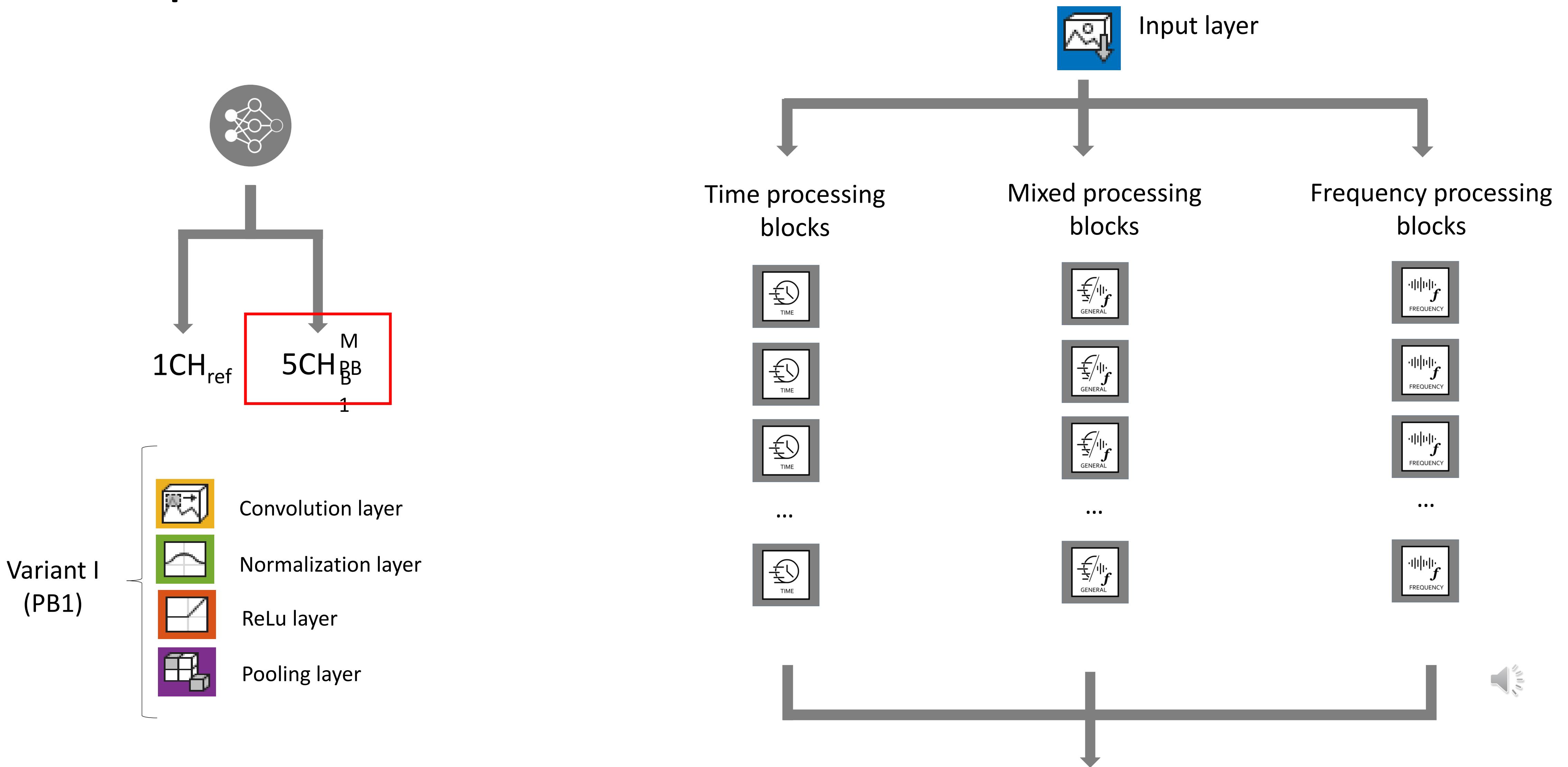
72



- | | |
|--|-----------------------|
| | Convolution layer |
| | Normalization layer |
| | ReLU layer |
| | Pooling layer |
| | Fully connected layer |
| | Dropout layer |
| | Softmax layer |
| | Classification layer |

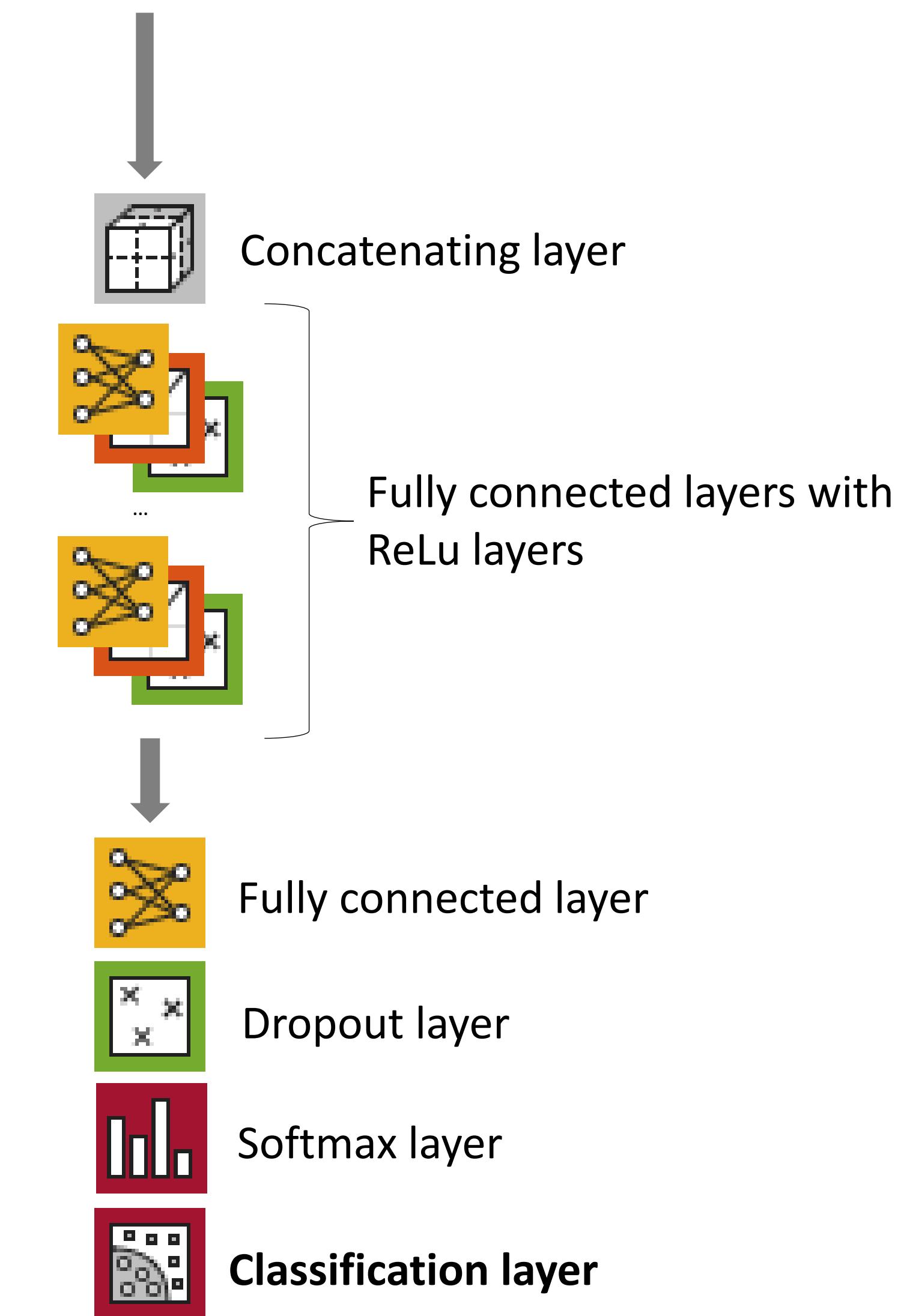
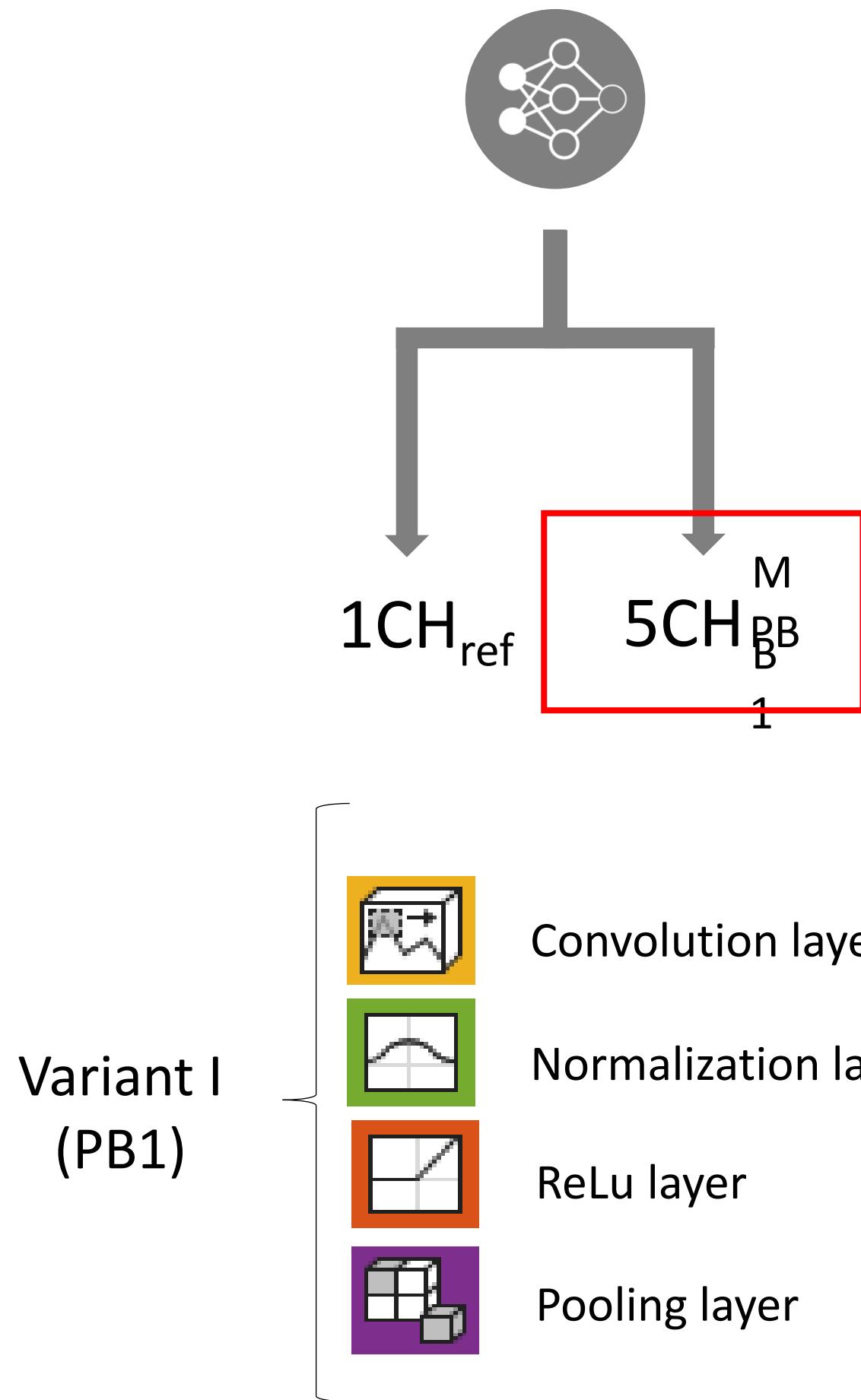


Experiments



Experiments

74





Experiments

1CH_{ref}

[1000-22000 Hz]

Reference	S _{norm}	S _{add}	S _{int}	ʃ _{norm}	ʃ _{int}	ʃ _{den}	PPV
S _{norm}	289	46	140	5	7	173	43.8%
S _{add}	92	23	55	8	5	97	8.2%
S _{int}	112	24	135	35	37	176	26.0%
ʃ _{norm}	1		19	516	67	92	74.2%
ʃ _{int}	9	6	33	115	22	52	9.3%
ʃ _{den}	61	5	54	63	20	157	43.6%
TPR	51.2%	22.1%	31.0%	69.5%	13.9%	21.0%	
	S _{norm}	S _{add}	S _{int}	ʃ _{norm}	ʃ _{int}	ʃ _{den}	
	Prediction						

ACC: 41.5%
TPR: 34.1%
TNR: 88.1%

5CH_{BB}^M

[1000-22000 Hz]

Reference	S _{norm}	S _{add}	S _{int}	ʃ _{norm}	ʃ _{int}	ʃ _{den}	PPV
S _{norm}	414	47	59	19	1	120	62.7%
S _{add}	69	159	8	8	2	34	56.8%
S _{int}	36	5	374	13	57	34	72.1%
ʃ _{norm}	10	5	21	551	46	62	79.3%
ʃ _{int}		1	95	50	78	13	32.9%
ʃ _{den}	88	12	43	92	24	101	28.1%
TPR	67.1%	69.4%	62.3%	75.2%	37.5%	27.7%	
	S _{norm}	S _{add}	S _{int}	ʃ _{norm}	ʃ _{int}	ʃ _{den}	
	Prediction						

ACC: 60.9%
TPR: 55.3%
TNR: 92.0%





Experiments



1CH_{ref}

[1000-22000 Hz]

Reference	PPV					
	S _{norm}	S _{add}	S _{int}	ʃ _{norm}	ʃ _{int}	ʃ _{den}
S _{norm}	289	46	140	5	7	173
S _{add}	92	23	55	8	5	97
S _{int}	112	24	135	35	37	176
ʃ _{norm}	1		19	516	67	92
ʃ _{int}	9	6	33	115	22	52
ʃ _{den}	61	5	54	63	20	157
TPR	51.2%	22.1%	31.0%	69.5%	13.9%	21.0%
	S _{norm}	S _{add}	S _{int}	ʃ _{norm}	ʃ _{int}	ʃ _{den}

Prediction



5CH_{BB}^M

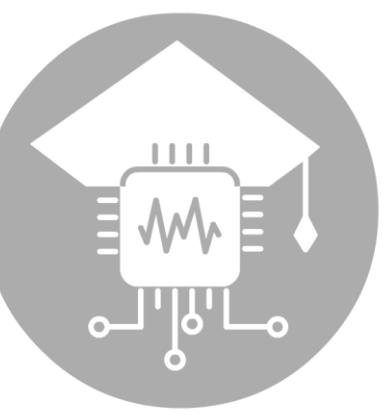
1

[1000-22000 Hz]

Reference	PPV					
	S _{norm}	S _{add}	S _{int}	ʃ _{norm}	ʃ _{int}	ʃ _{den}
S _{norm}	414	47	59	19	1	120
S _{add}	69	159	8	8	2	34
S _{int}	36	5	374	13	57	34
ʃ _{norm}	10	5	21	551	46	62
ʃ _{int}		1	95	50	78	13
ʃ _{den}	88	12	43	92	24	101
TPR	67.1%	69.4%	62.3%	75.2%	37.5%	27.7%
	S _{norm}	S _{add}	S _{int}	ʃ _{norm}	ʃ _{int}	ʃ _{den}

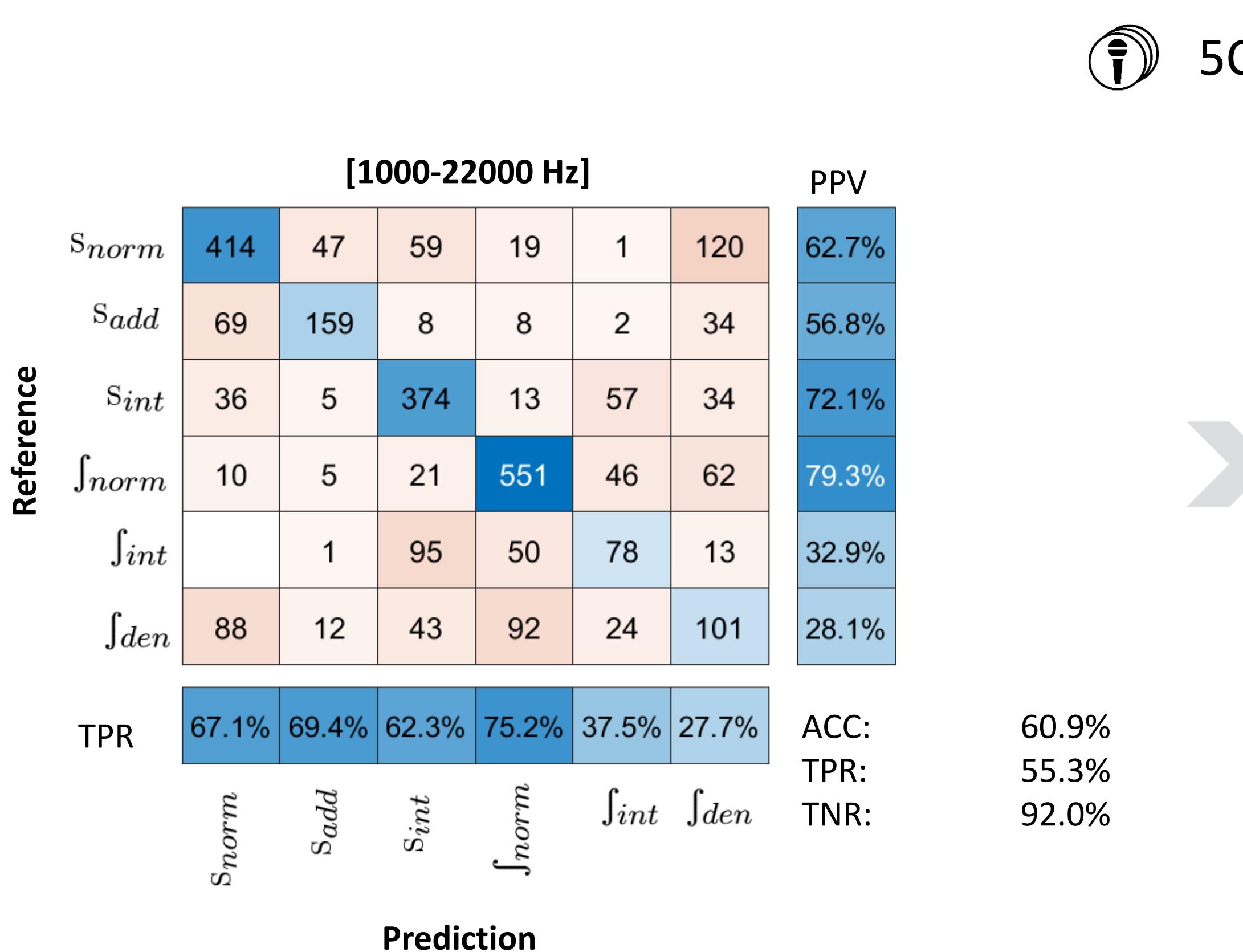
Prediction



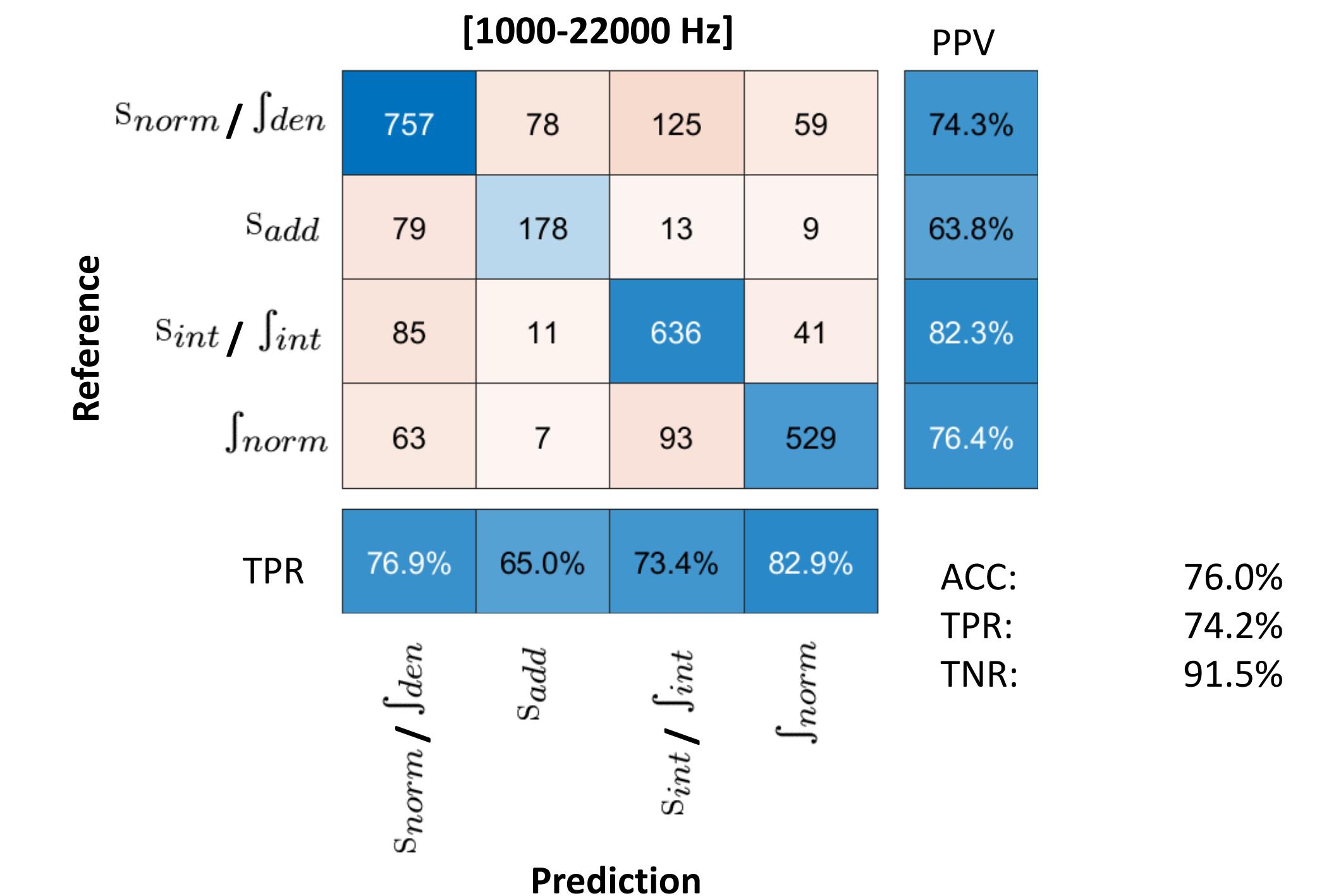


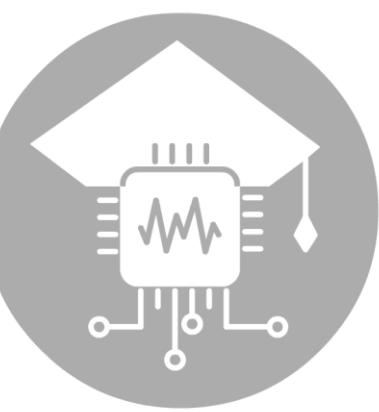
Experiments

77

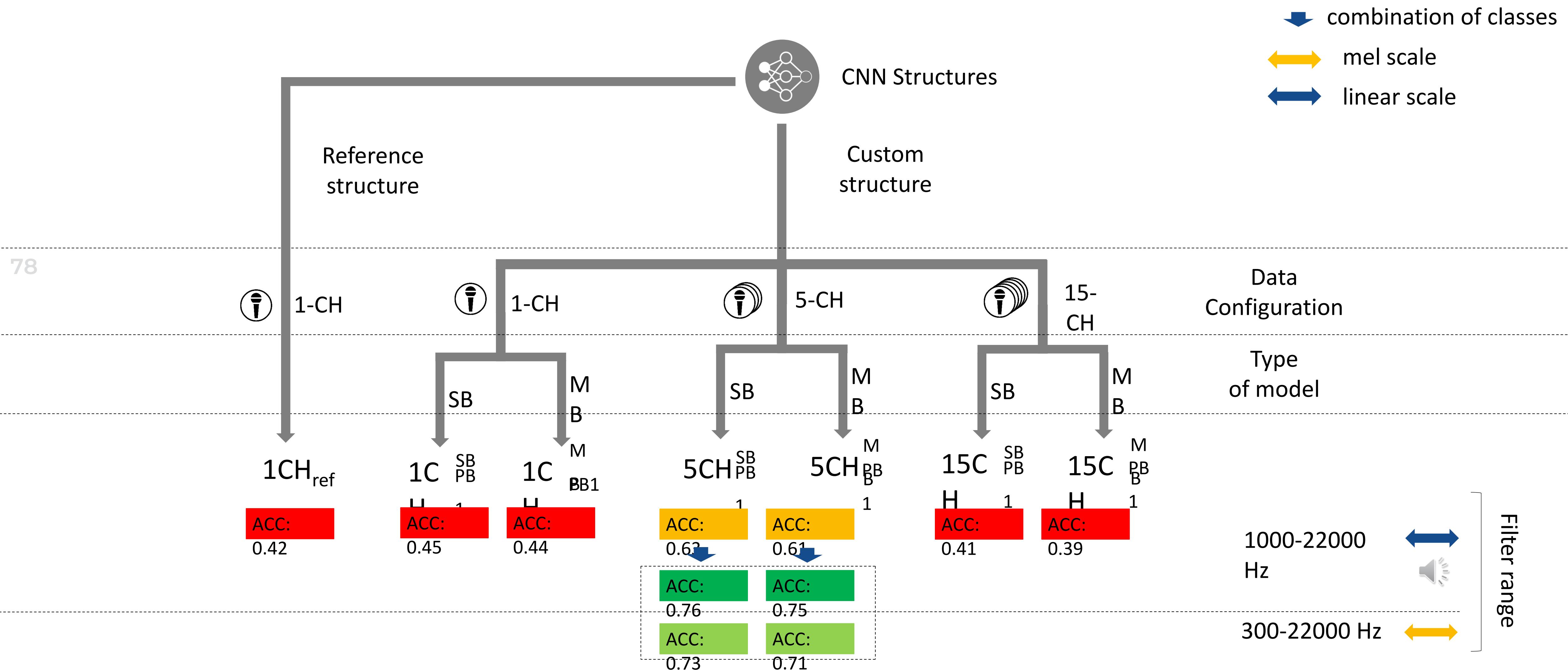


5CH_B^M
1





Experiments

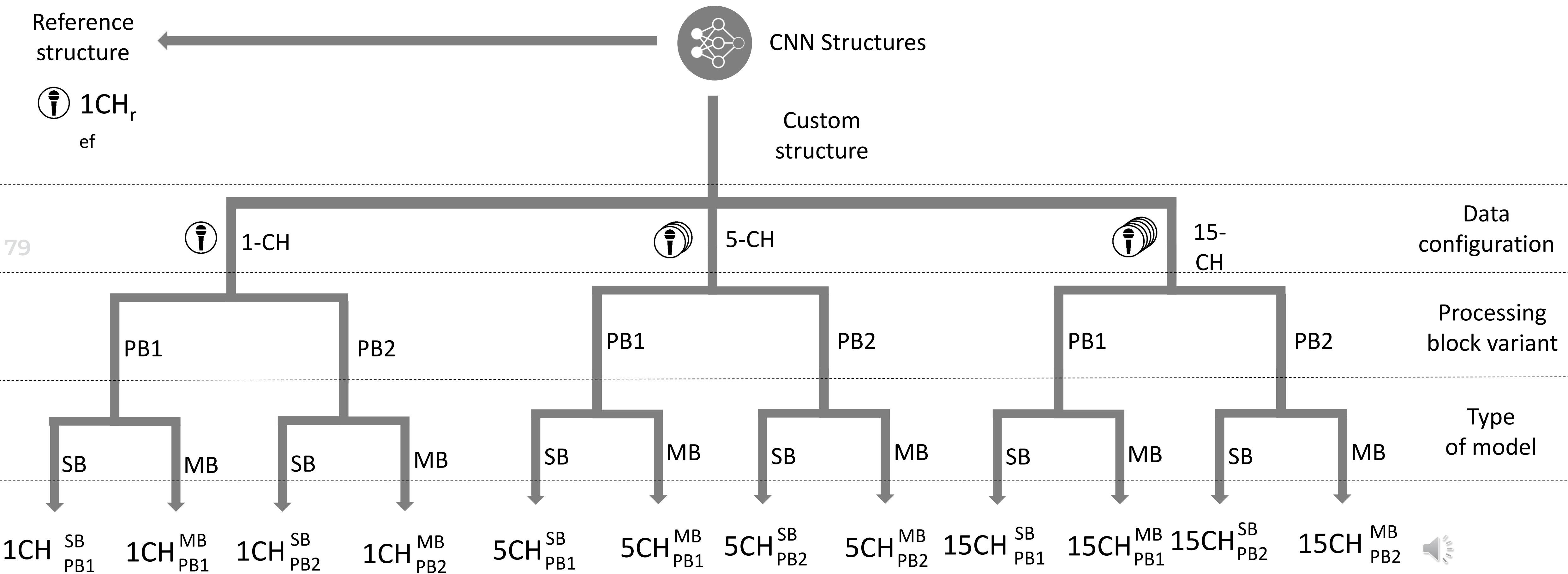


Experiments

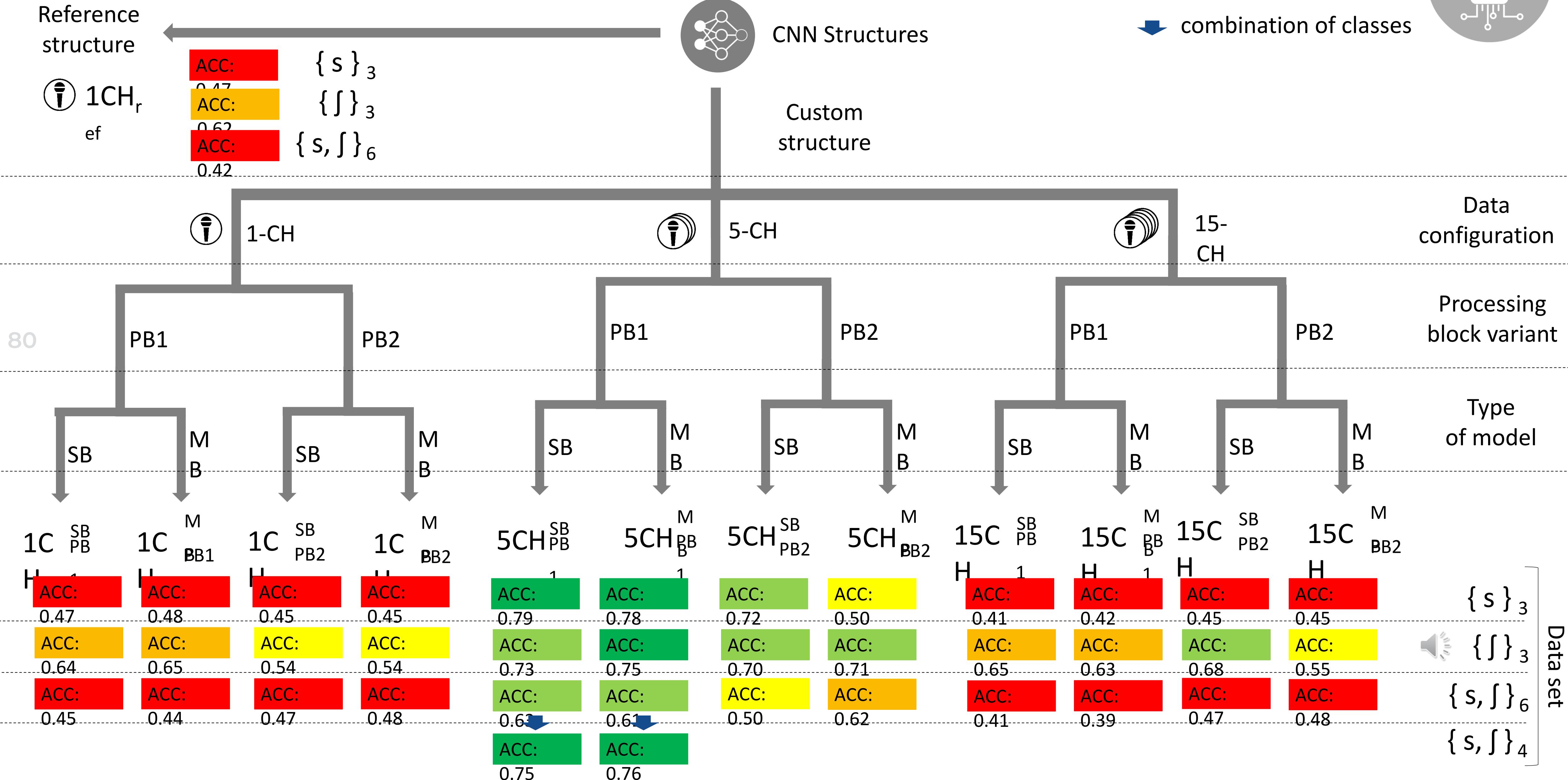


linear scale filter range: 1000-
22000 Hz

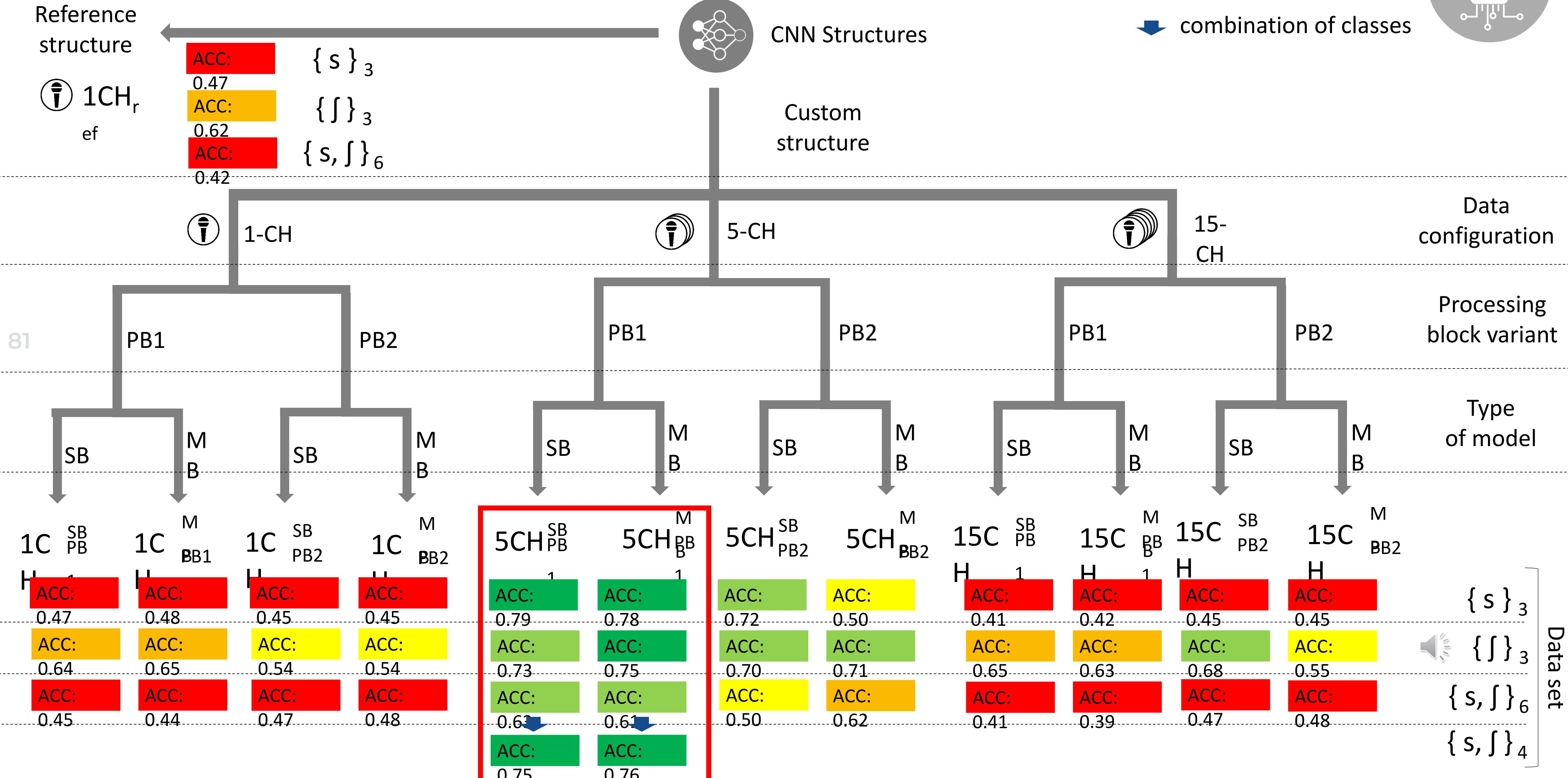
 combination of classes



PHONEM CLASSIFICATION – RESULTS OF TESTED CNN STRUCTURES



PHONEM CLASSIFICATION – RESULTS OF TESTED CNN STRUCTURES



Experiments

ANALYSIS OF SIBILANTS MODELS

82

PHONEME
CLASSIFICATION

SENSITIVITY ANALYSIS OF THE CNN MODEL



Experiments

ANALYSIS OF SIBILANTS MODELS

83

PHONEME
CLASSIFICATION

SENSITIVITY ANALYSIS OF THE CNN MODEL



Experiments

SENSITIVITY ANALYSIS OF THE CNN MODEL

84

- the influence of selected hyperparameters and network structure on classification results,
- verification of the validity of network structure parallelization.





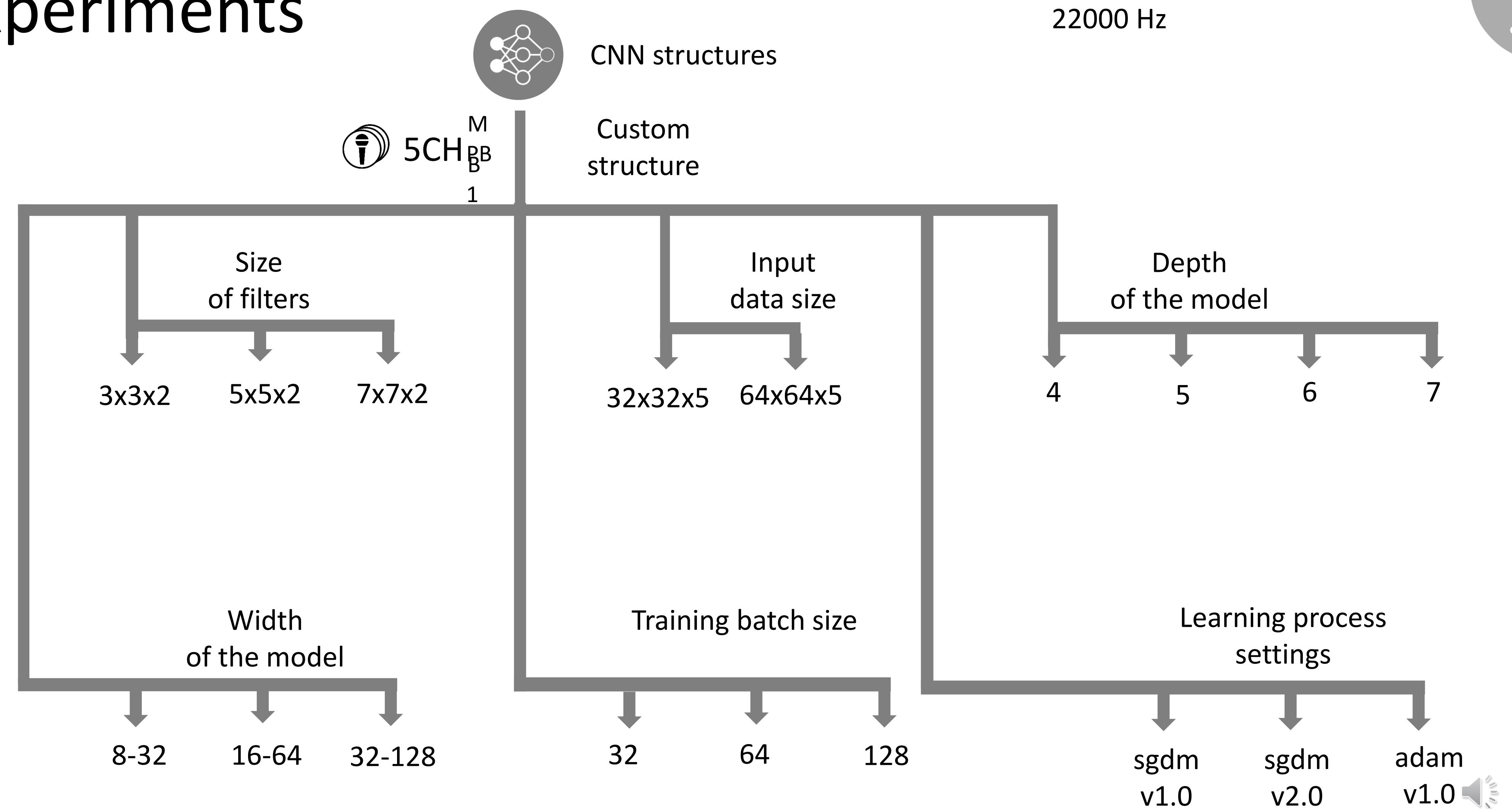
Experiments

	Type of realisation	/s/ normative	/ʃ/ normative	/s/ + /ʃ/ interdentalna	/s/ addentalna
Type of set	training 70%	235 (1410)	163 (978)	118 (708)	116 (696)
	validation 15%	33	35	14	26
	test 15%	33	35	14	26

The size of data sets (before and after augmentation) in recognizing different realizations of the sounds /s/ and /ʃ/



Experiments



Experiments



CNN structures



$5CH_{BB}^M$

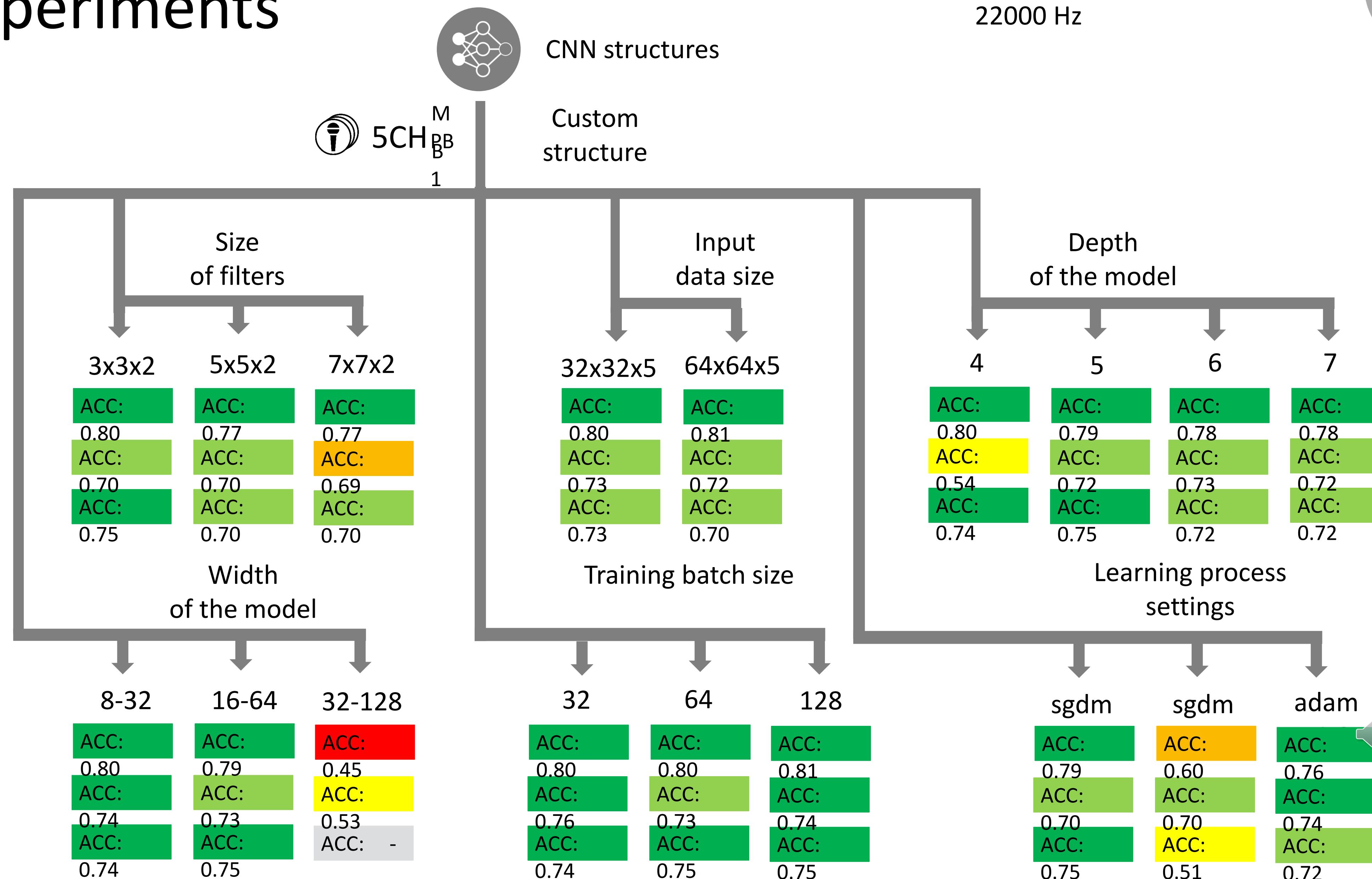


Custom structure

linear scale filter range: 1000-
22000 Hz



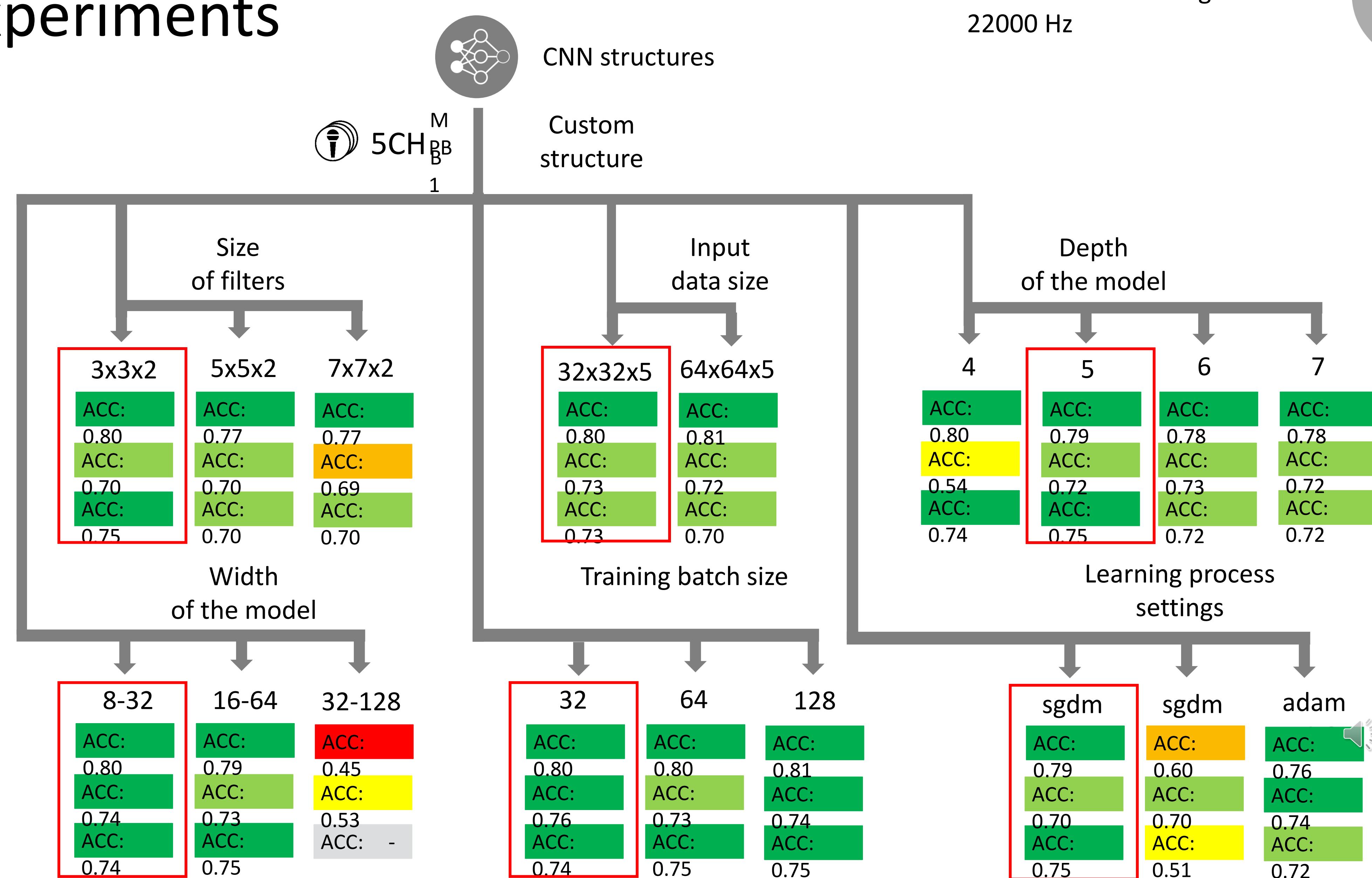
87



Experiments

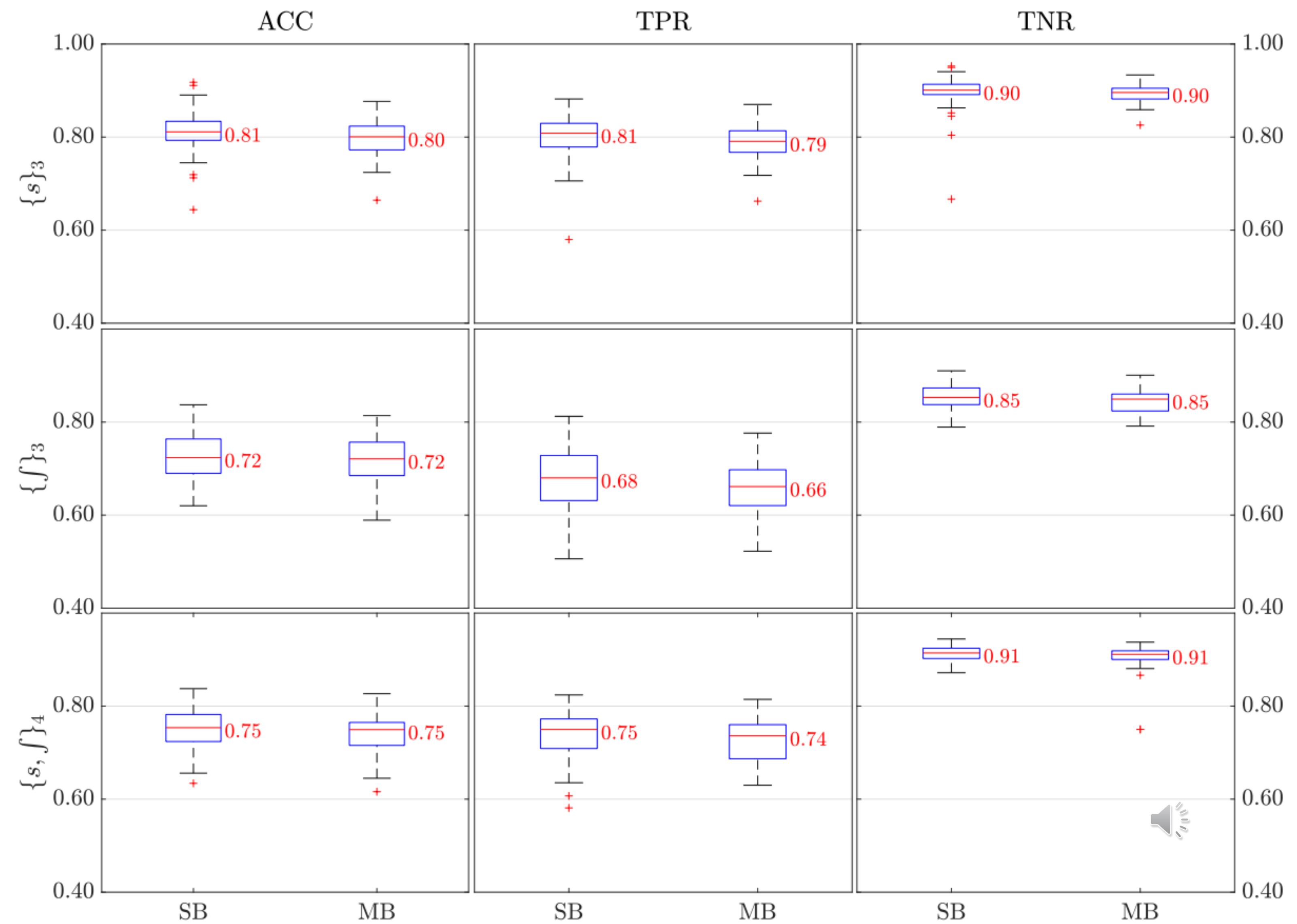
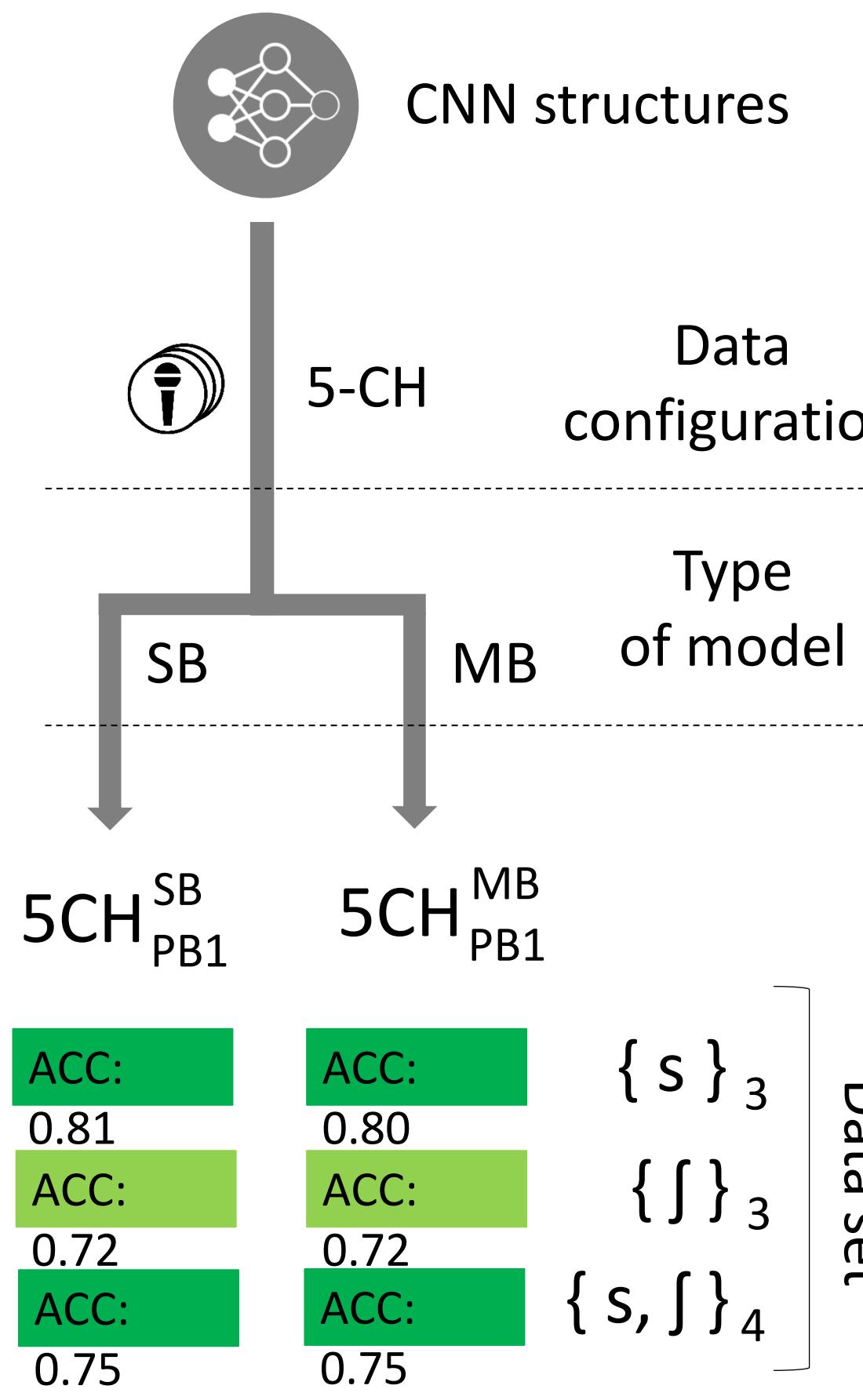


88



Experiments

89



Summary

90



Design, construction and verification of a measuring device.



Development of a methodology for spatial processing of the speech signal dedicated to dentalized sounds.



Development of the measuring device along with the expansion of the speech database, extension of the analysis to include additional dentalized sounds.



Reference

- **Automated detection of sigmatism using deep learning applied to multichannel speech signal**
Krecichwost, M., Mocko, N., Badura, P.
Biomedical Signal Processing and Control, 68, 102612. Elsevier, 2021.
- **Hybridization of Acoustic and Visual Features of Polish Sibilants Produced by Children for Computer Speech Diagnosis**
Sage, A., Miodonska, Z., Kręcichwost, M., Badura, P.
Sensors, 24(16), 5360. Multidisciplinary Digital Publishing Institute, 2024.
- **Multichannel acoustic analysis of phoneme/s/mispronunciation for lateral sigmatism detection**
Krecichwost, M., Miodonska, Z., Badura, P., Trzaskalik, J., Mocko, N.
Biocybernetics and Biomedical Engineering, 39(1), 246-255. Elsevier, 2019.
- **Multichannel Speech Acquisition and Analysis for Computer-Aided Sigmatism Diagnosis in Children**
Krecichwost, M., Miodonska, Z., Trzaskalik, J., Badura, P.
IEEE Access. IEEE, 2020.