# EFFICIENT AND PERCEPTUALLY PLAUSIBLE 3-D SOUND FOR VIRTUAL REALITY

Fabian Brinkmann and Hannes Gamper (Mentor)

Audio and Acoustics Group (Ivan Tashev)
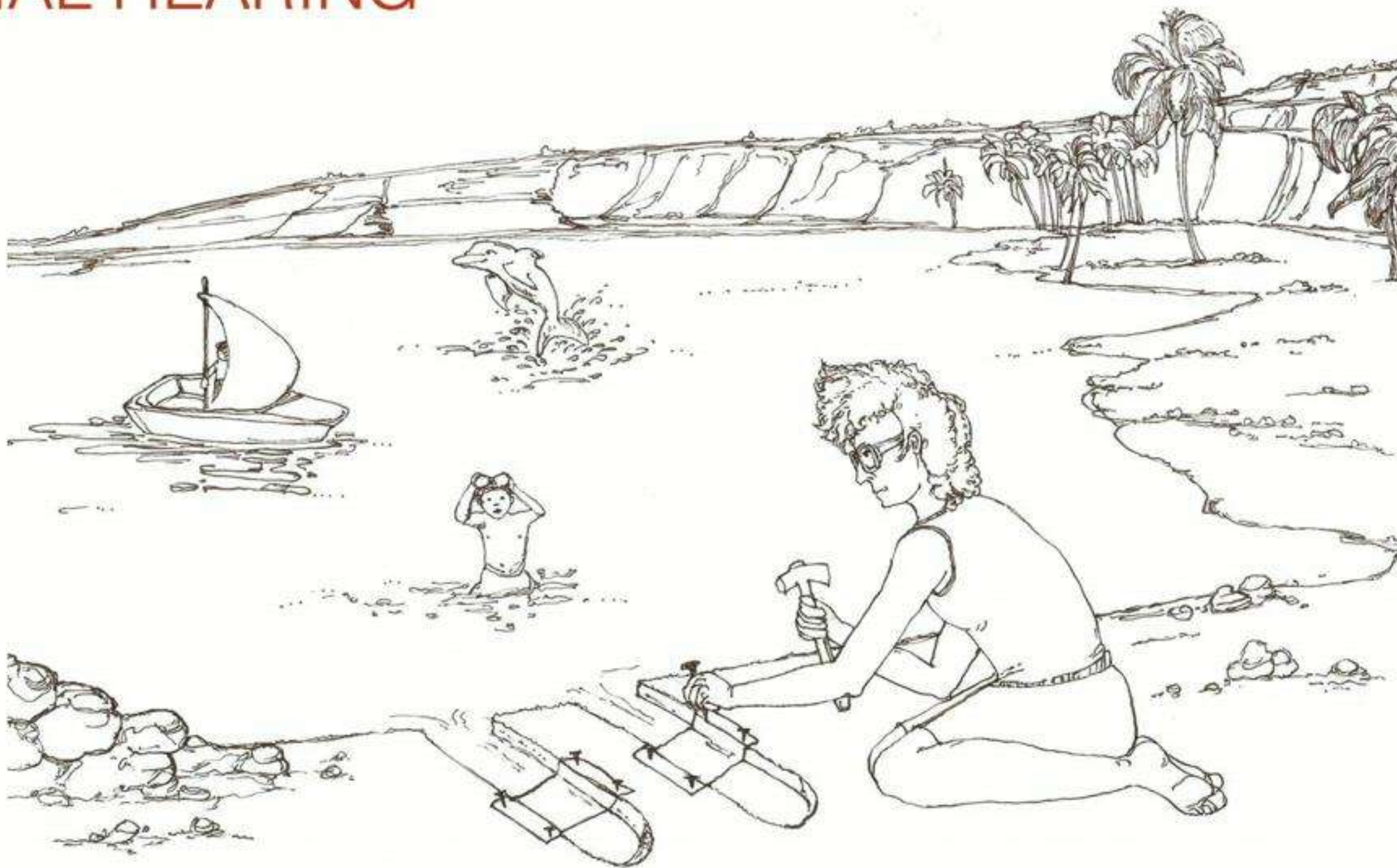
Microsoft Research, Redmond Lab

# MOTIVATION



"HoloLens Demo at Penn Museum" by pennlibtrl (CC BY-SA 2.0)

- 3-D Audio improves presence and immersion

- Graphics rendering consumes most of the compute

- Room acoustical simulation is expensive
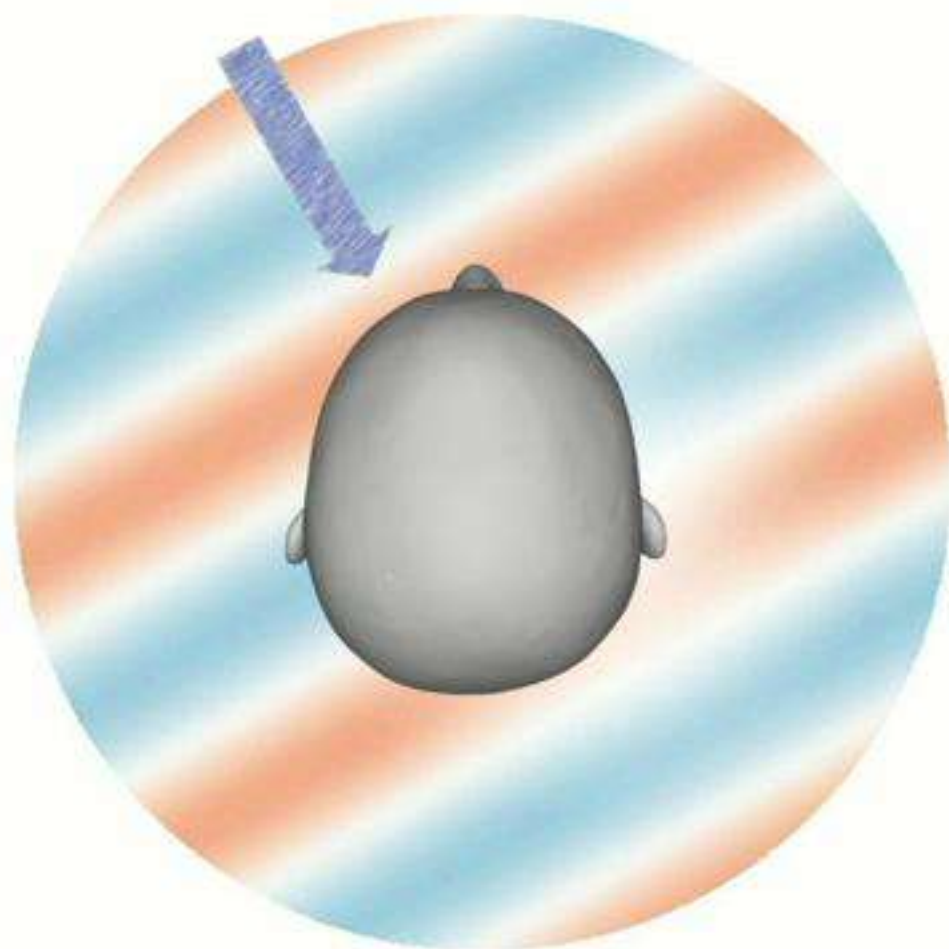
# SPATIAL HEARING



after A. S. Bregman (1994): Auditory scene analysis, MIT Press, Cambridge, USA, pp. 5.
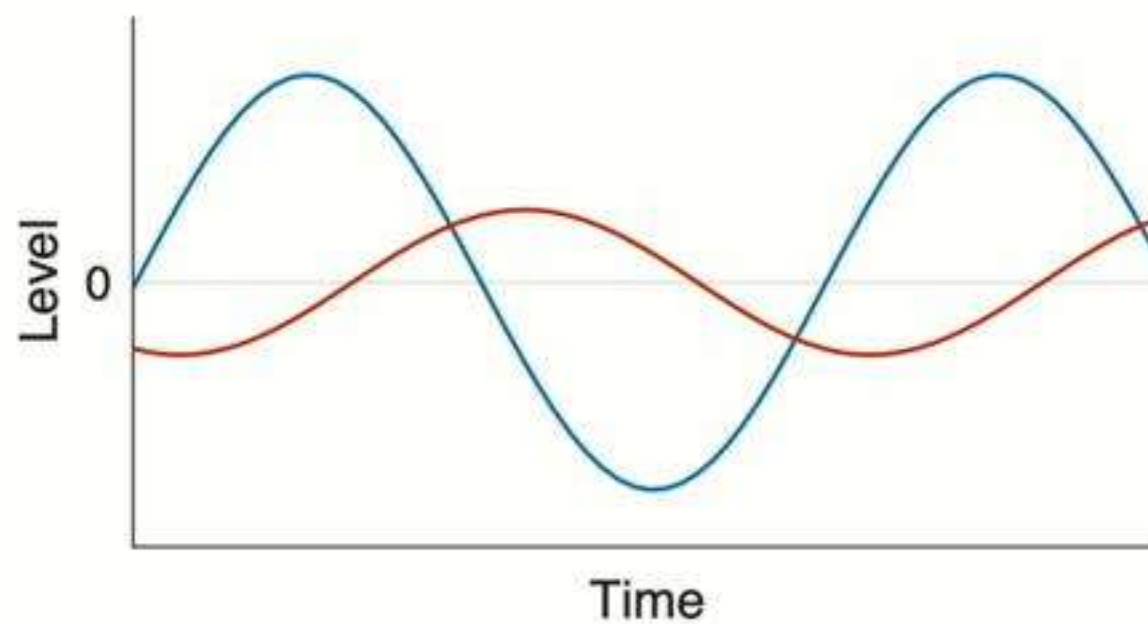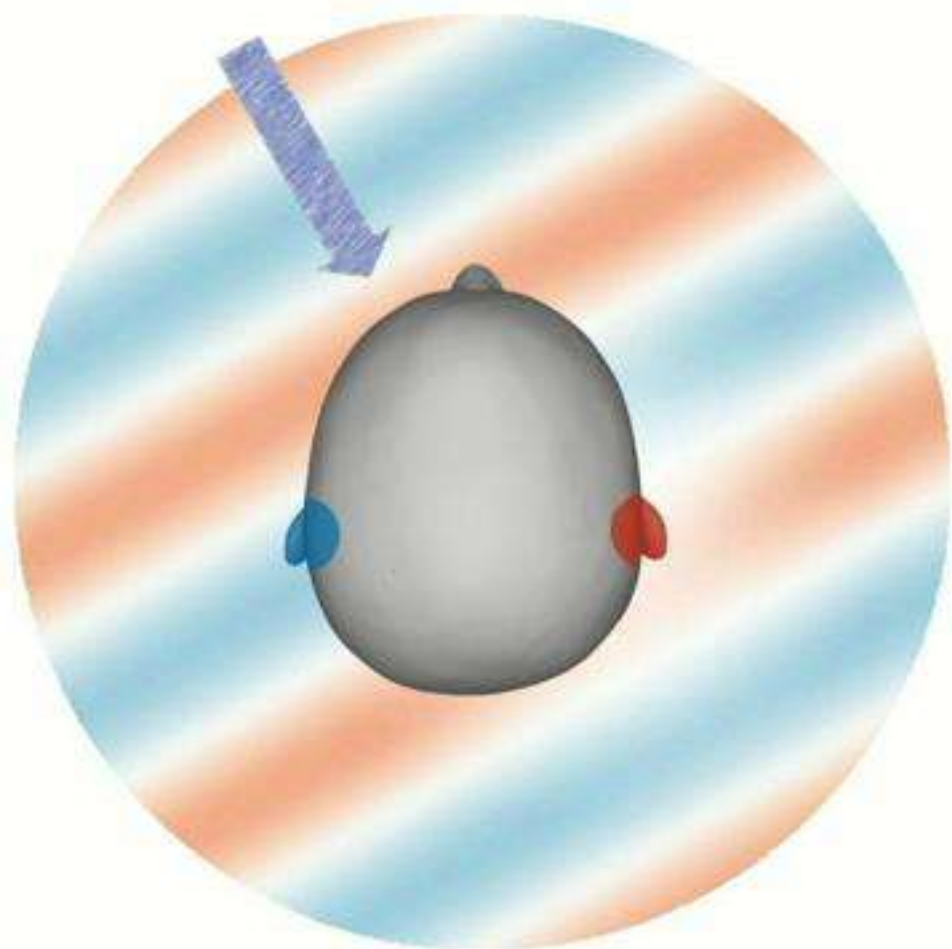Image: Fabian Brinkmann (CC-BY 4.0)

# SPATIAL HEARING



Blauert (**1997**). *Spatial hearing. The psychophysics of human sound localization*, (MIT Press, Cambridge, Massachusetts)
Images: Fabian Brinkmann (CC-BY 4.0)

# SPATIAL HEARING



Blauert (**1997**). *Spatial hearing. The psychophysics of human sound localization*, (MIT Press, Cambridge, Massachusetts)
Images: Fabian Brinkmann (CC-BY 4.0)

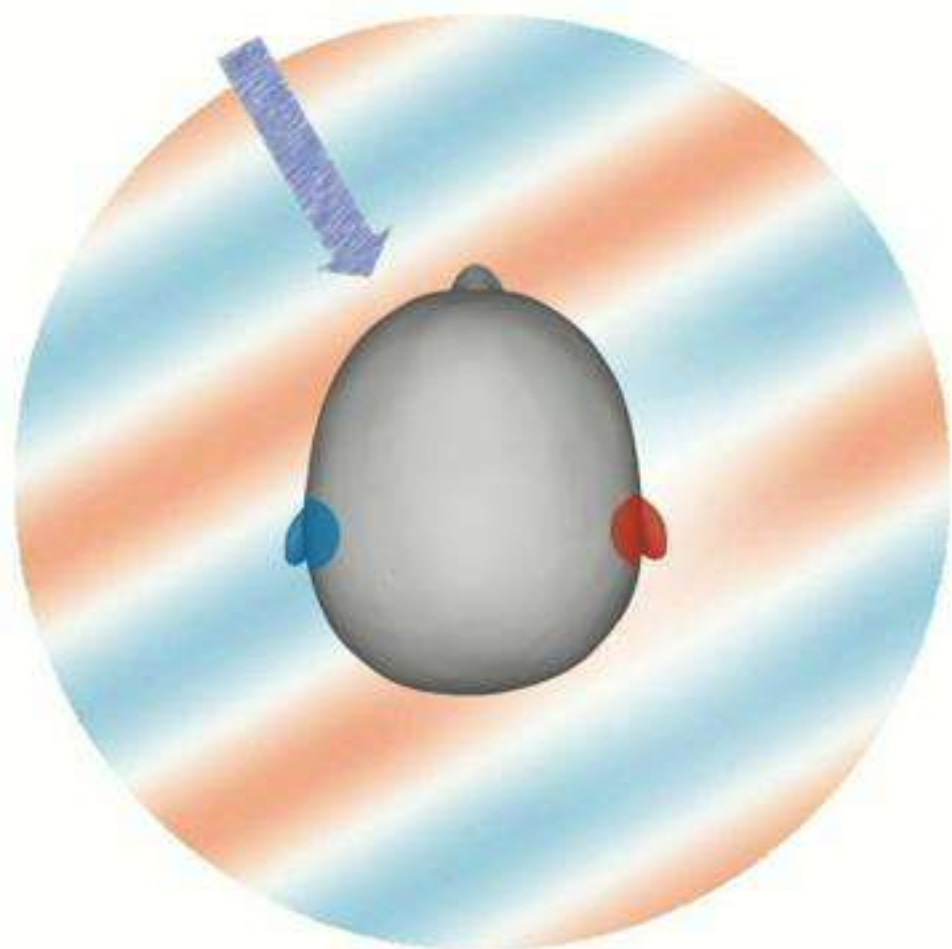# SPATIAL HEARING



Interaural time difference

# SPATIAL HEARING



Interaural time difference

Interaural level difference

# SPATIAL HEARING



Blauert (**1997**). *Spatial hearing. The psychophysics of human sound localization*, (MIT Press, Cambridge, Massachusetts)
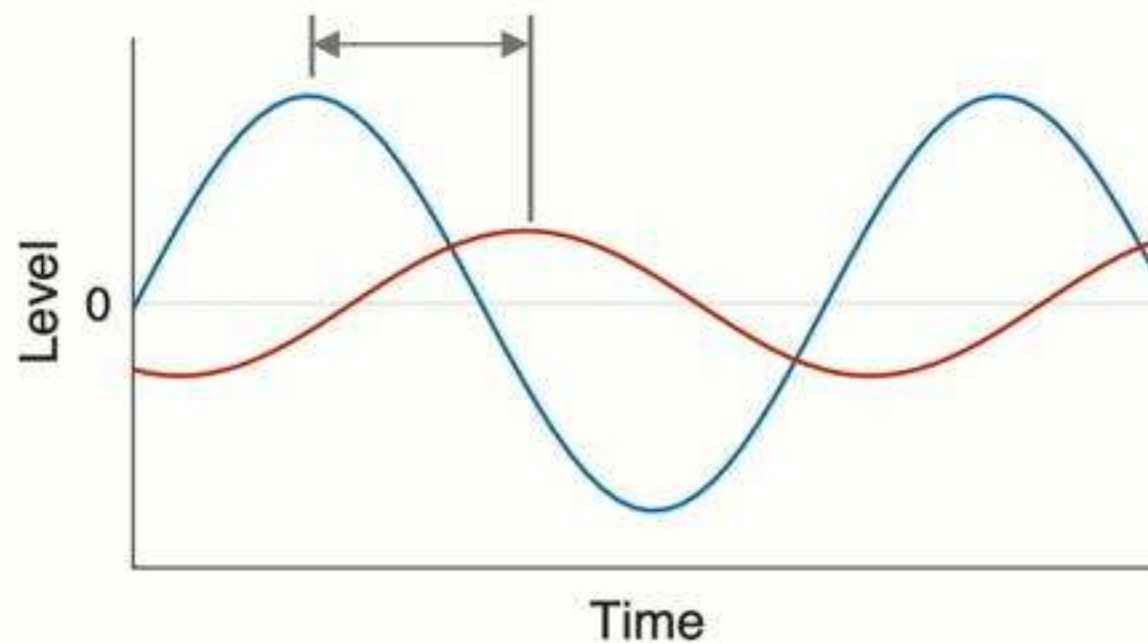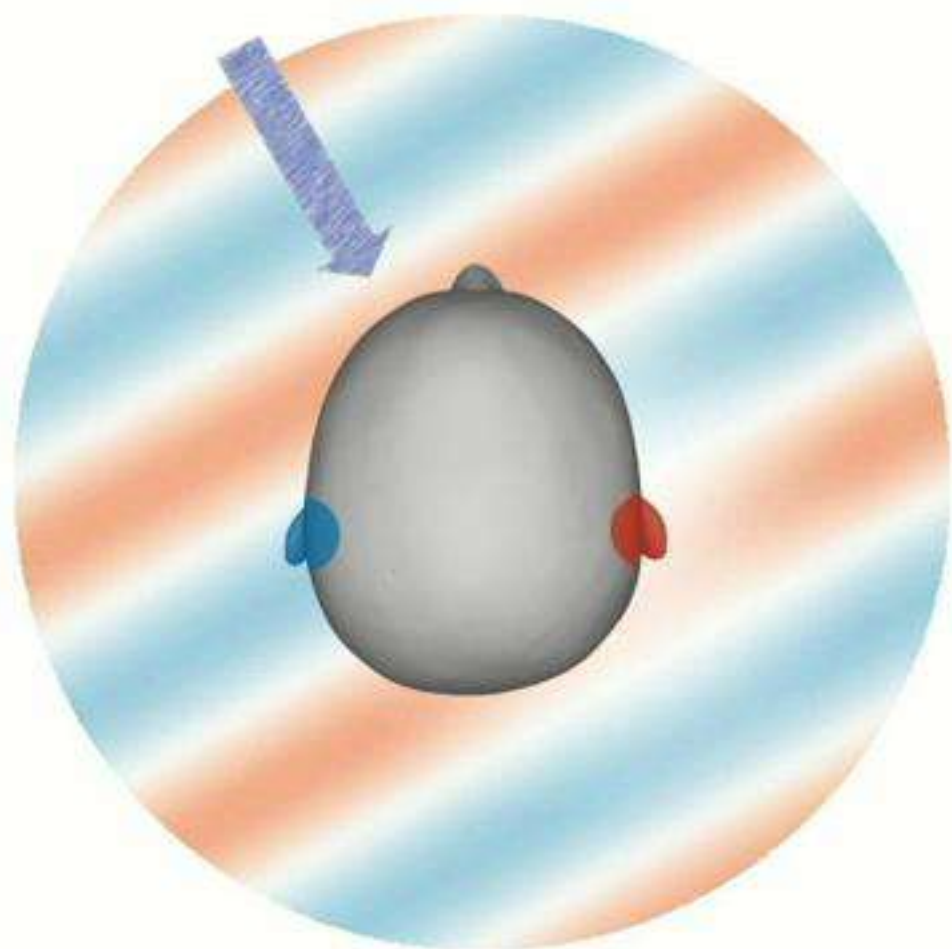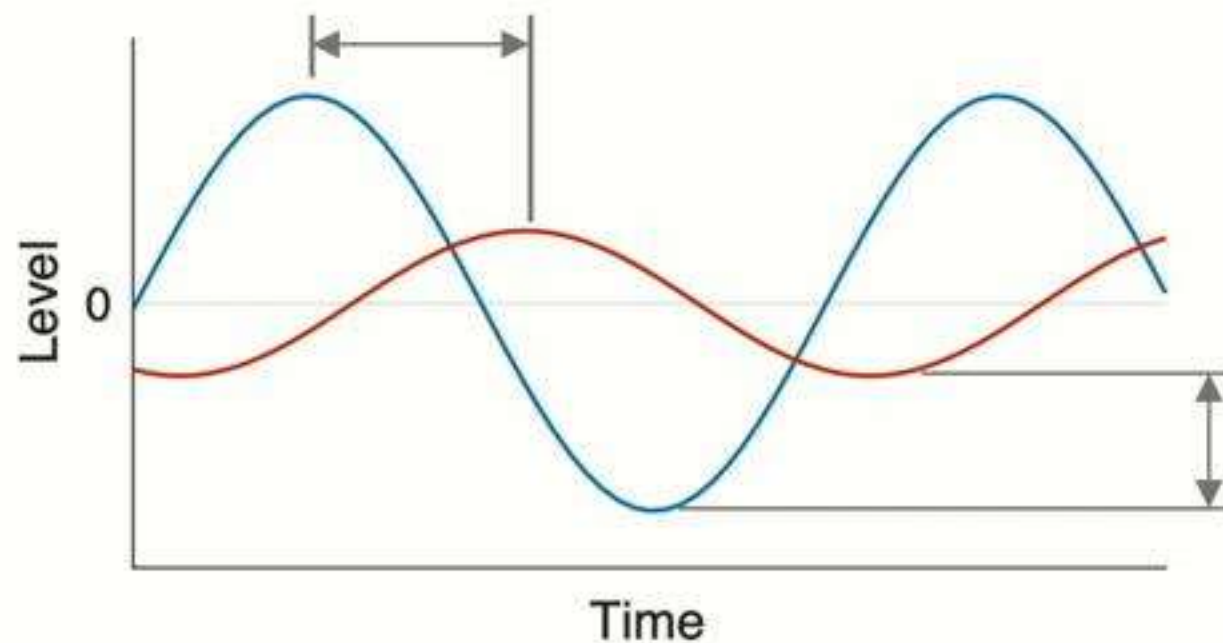Images: Fabian Brinkmann (CC-BY 4.0)
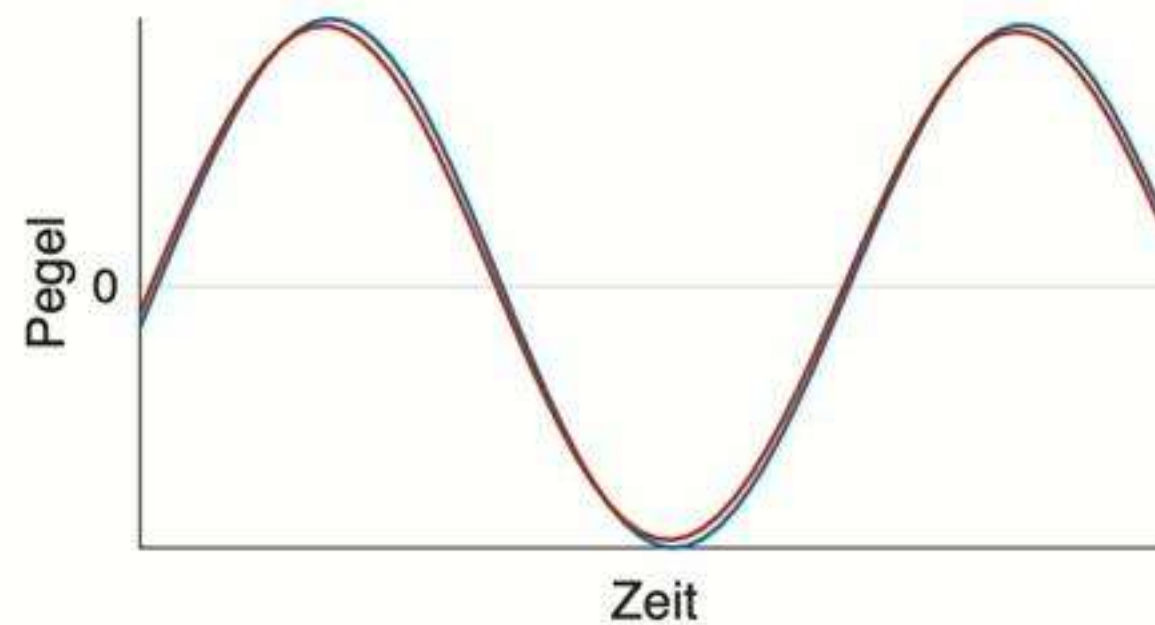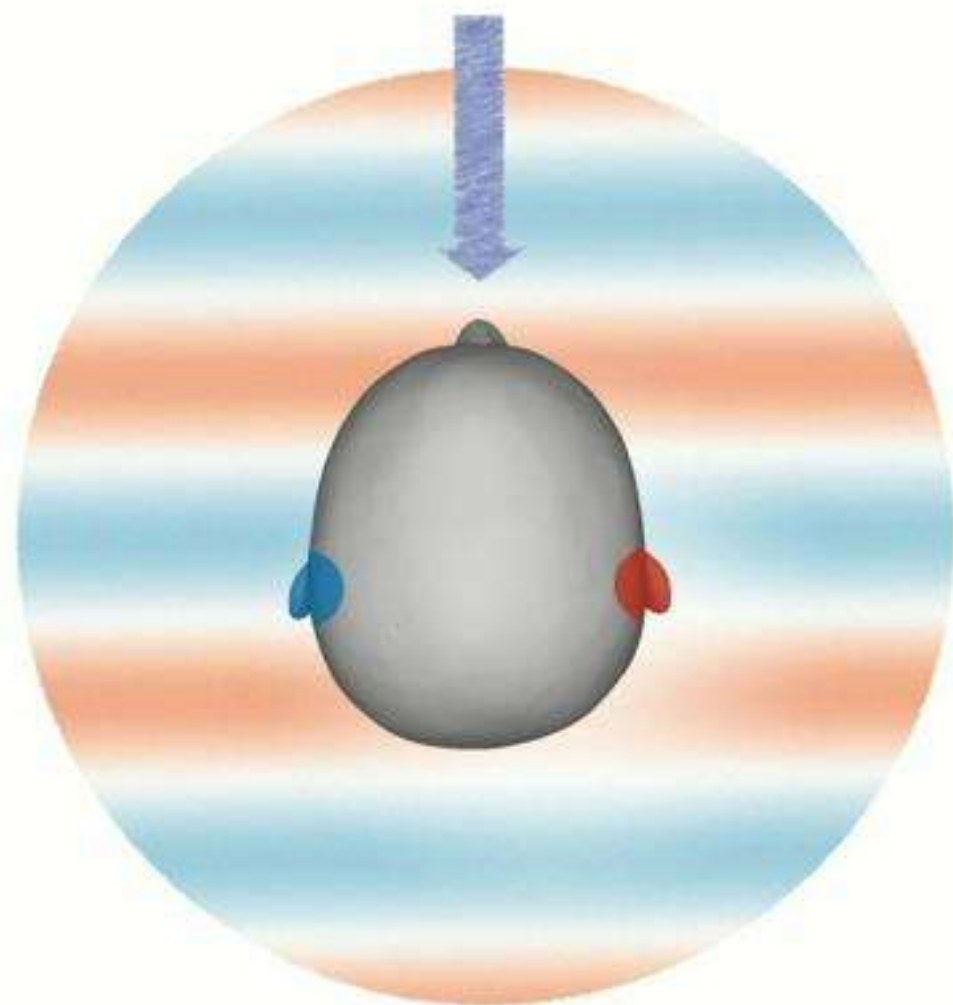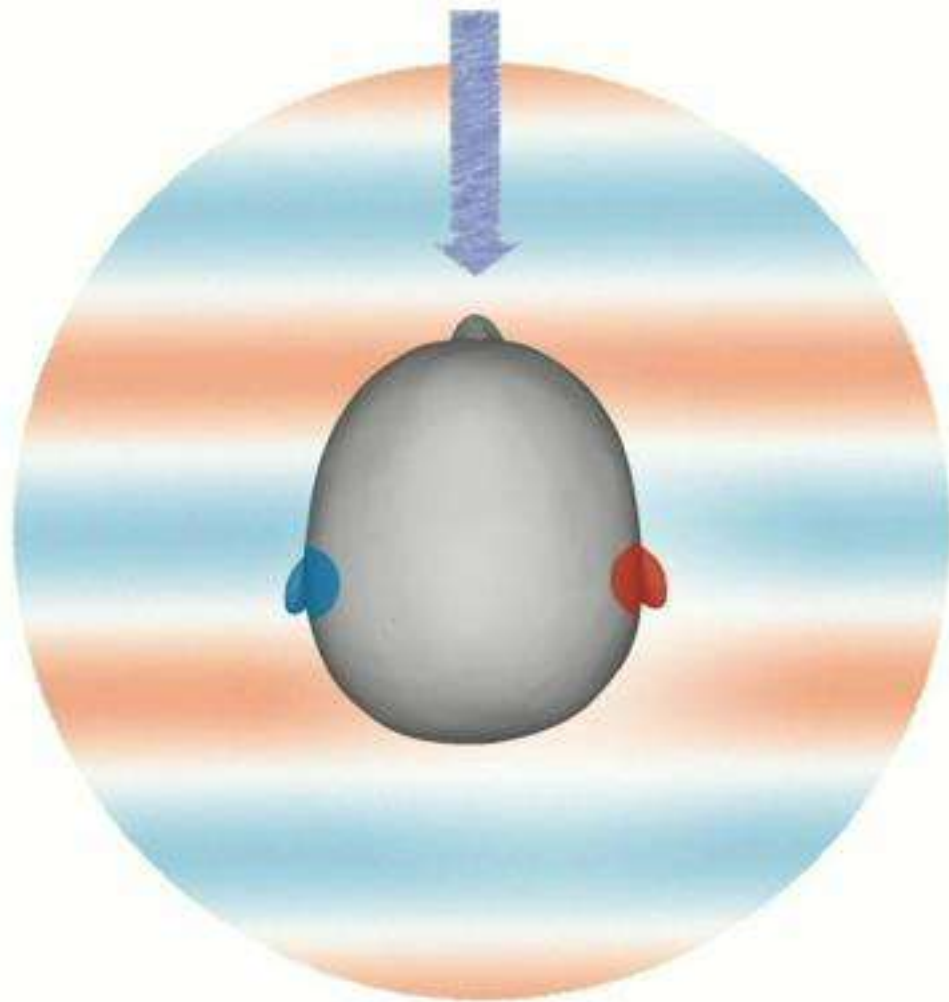
# SPATIAL HEARING



monaural spectral cues

# ACOUSTICAL SIMULATION



ROOM MODEL      SIMULATION      IMPULSE RESPONSE      AURALIZATION

Images: Fabian Brinkmann (CC-BY 4.0)

# PARAMETRIC SPATIAL AUDIO

Lindau *et al.* (2012) Perceptual evaluation of model- and signal-based predictors of the mixing time... JAES **60**(11): 887 – 898.

Godin *et al.* (2019) Aesthetic modification of room impulse responses... *AES Conf. Immersion and Interactive Audio*, York, UK.

Images: Fabian Brinkmann (CC-BY 4.0)

# PARAMETRIC SPATIAL AUDIO

Lindau *et al.* (2012) Perceptual evaluation of model- and signal-based predictors of the mixing time... JAES **60**(11): 887 – 898.

Godin *et al.* (2019) Aesthetic modification of room impulse responses... *AES Conf. Immersion and Interactive Audio*, York, UK.

# PARAMETRIC SPATIAL AUDIO

Lindau *et al.* (2012) Perceptual evaluation of model- and signal-based predictors of the mixing time... JAES **60**(11): 887 – 898.

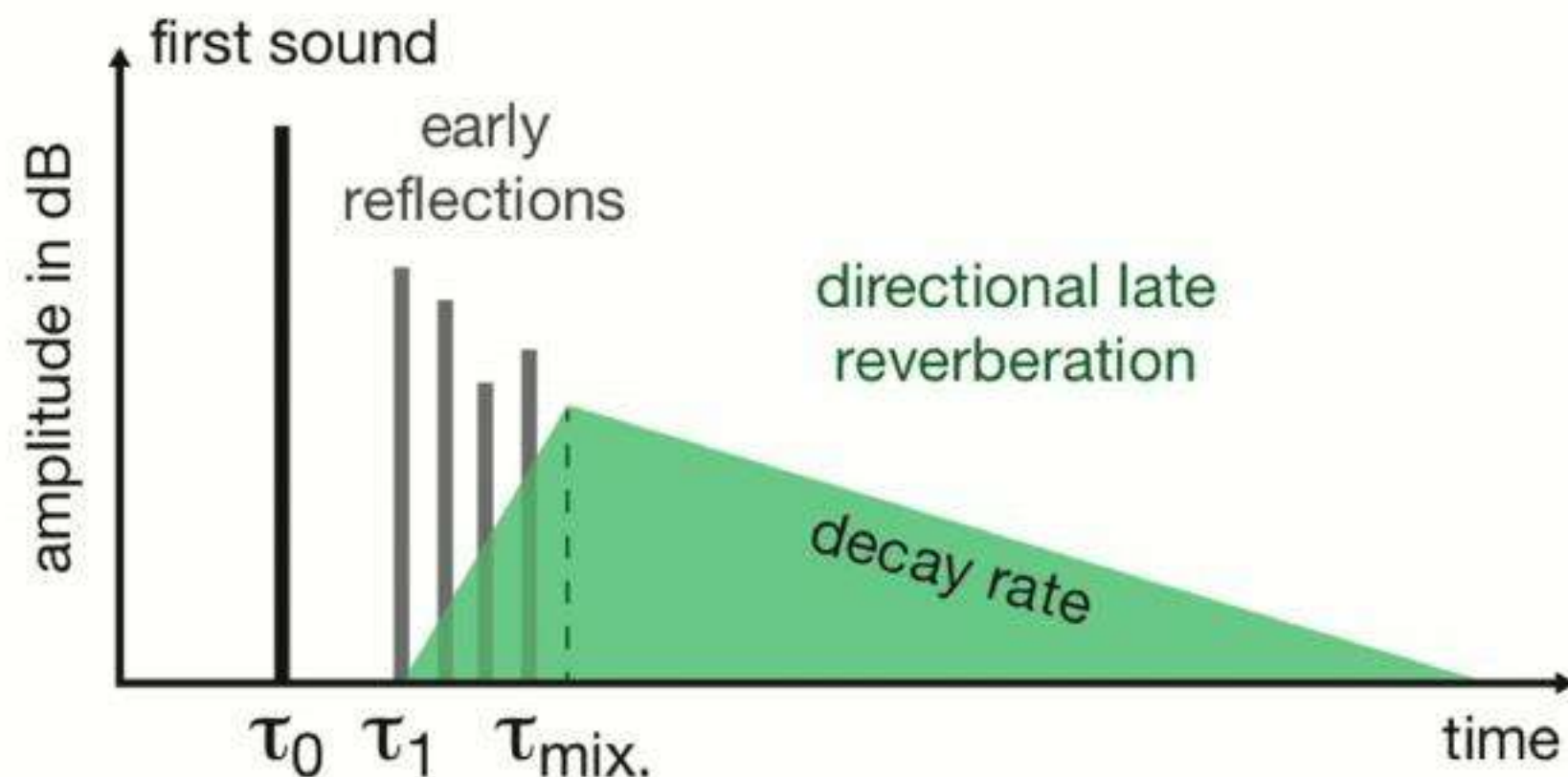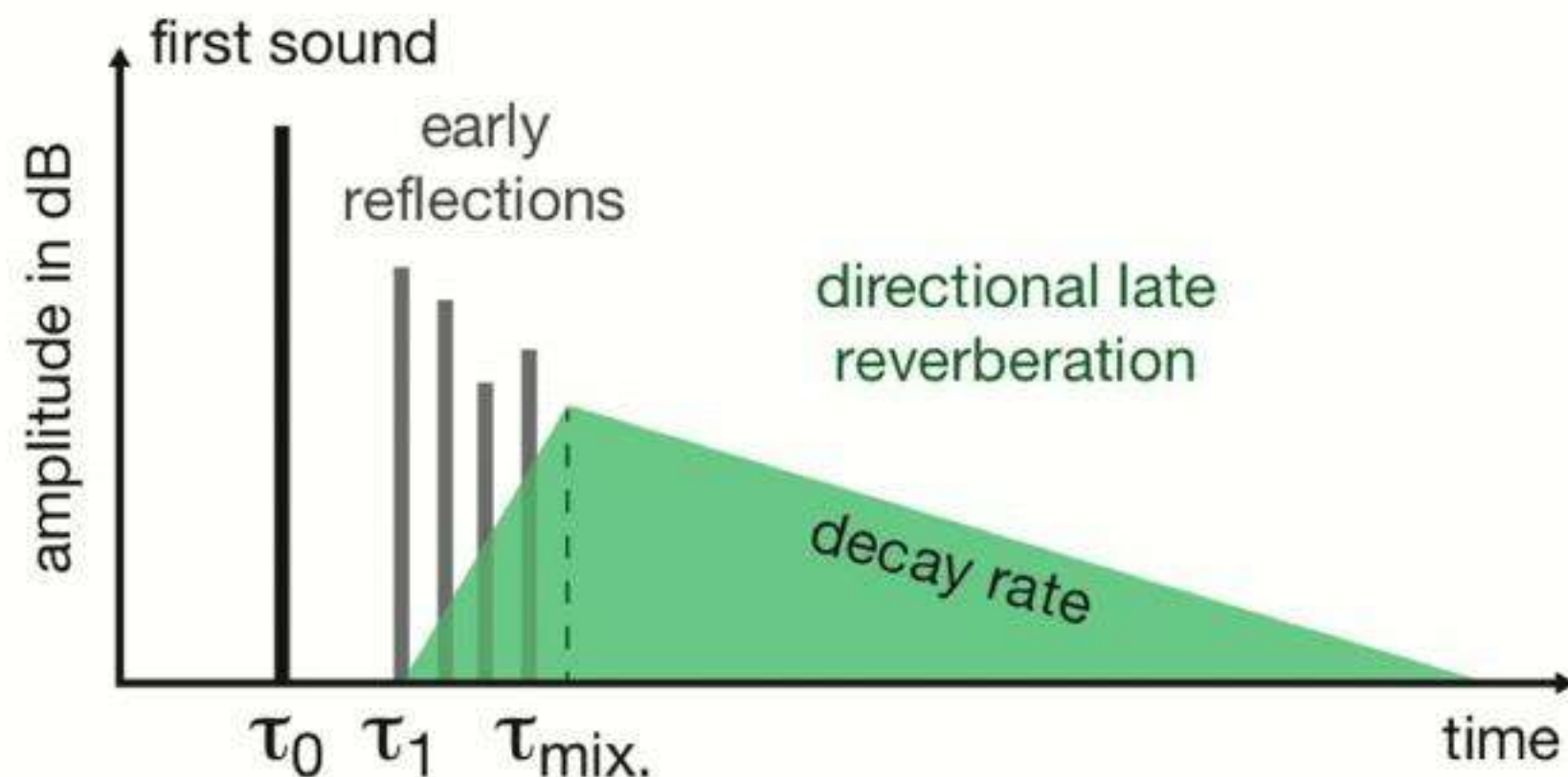Godin *et al.* (2019) Aesthetic modification of room impulse responses... *AES Conf. Immersion and Interactive Audio*, York, UK.

Images: Fabian Brinkmann (CC-BY 4.0)

# PARAMETRIC SPATIAL AUDIO



- Cheap
  rendering

- Low
  memory

- Aesthetic
  modification
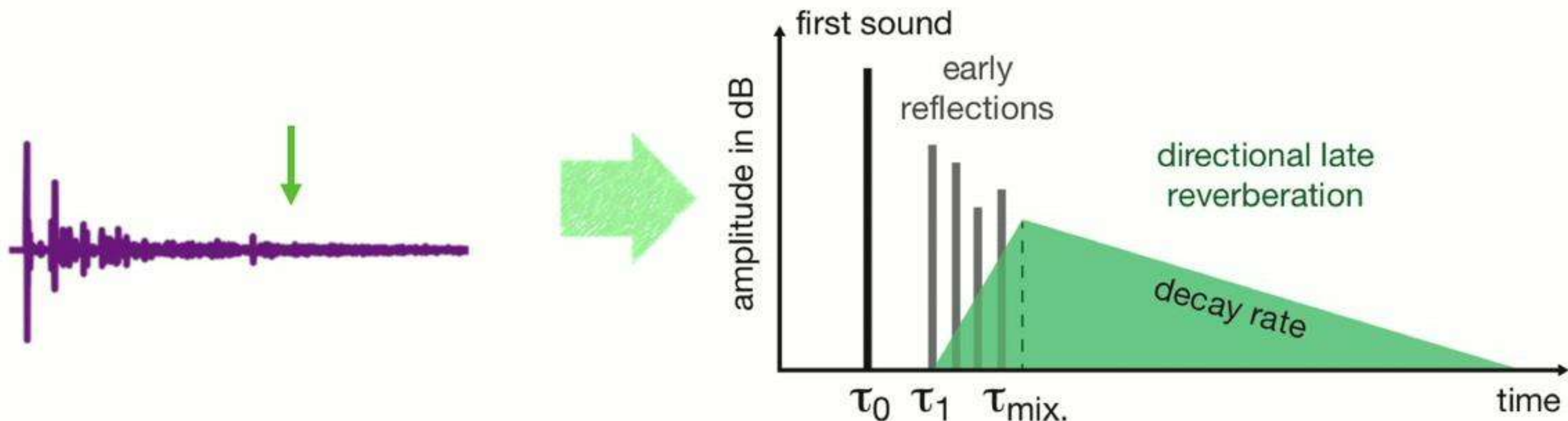
Lindau *et al.* (2012) Perceptual evaluation of model- and signal-based predictors of the mixing time... JAES **60**(11): 887 – 898.

Godin *et al.* (2019) Aesthetic modification of room impulse responses... *AES Conf. Immersion and Interactive Audio*, York, UK.

# PARAMETRIC DE- AND ENCODING IN CONSIDERATTION OF EARLY REFELCTIONS

- Encoding (offline)
  - Dependency on environment
  - Dependency on audio content
  - Scalable to match computational resources
  - Smooth spatial distribution
  - Low-cost and short memory

first sound

early reflections

directional late reverberation

decay rate

$\tau_0$  $\tau_1$  $\tau_{mix.}$

time

# PARAMETRIC DE- AND ENCODING IN CONSIDERATTION OF EARLY REFELCTIONS



- Encoding (offline)
  - Dependency on environment
  - Dependency on audio content
  - Scalable to match computational resources
  - Smooth spatial distribution
  - Low-cost and short memory

- Decoding (real-time)
  - Efficient
  - Perceptually plausible

Images: Fabian Brinkmann (CC-BY 4.0)

# PRECEDENCE IN ROOM ACOUSTICS – TEMPORAL ASPECTS

Litovsky *et al.* (1999): The precedence effect. *J. Acoust. Soc. Am.* **104**(4): 1633 – 1654.

Brown *et al.* (2015): The precedence effect in sound localization. *J. Assoc. Res. Otolaryng.* **1**(16):1 – 28.

# PRECEDENCE IN ROOM ACOUSTICS – TEMPORAL ASPECTS



Litovsky *et al.* (1999): The precedence effect. *J. Acoust. Soc. Am.* **104**(4): 1633 – 1654.

Brown *et al.* (2015): The precedence effect in sound localization. *J. Assoc. Res. Otolaryng.* **1**(16):1 – 28.

Images: Fabian Brinkmann (CC-BY 4.0)

# PRECEDENCE IN ROOM ACOUSTICS – TEMPORAL ASPECTS

Litovsky *et al.* (1999): The precedence effect. *J. Acoust. Soc. Am.* **104**(4): 1633 – 1654.

Brown *et al.* (2015): The precedence effect in sound localization. *J. Assoc. Res. Otolaryng.* **1**(16):1 – 28.

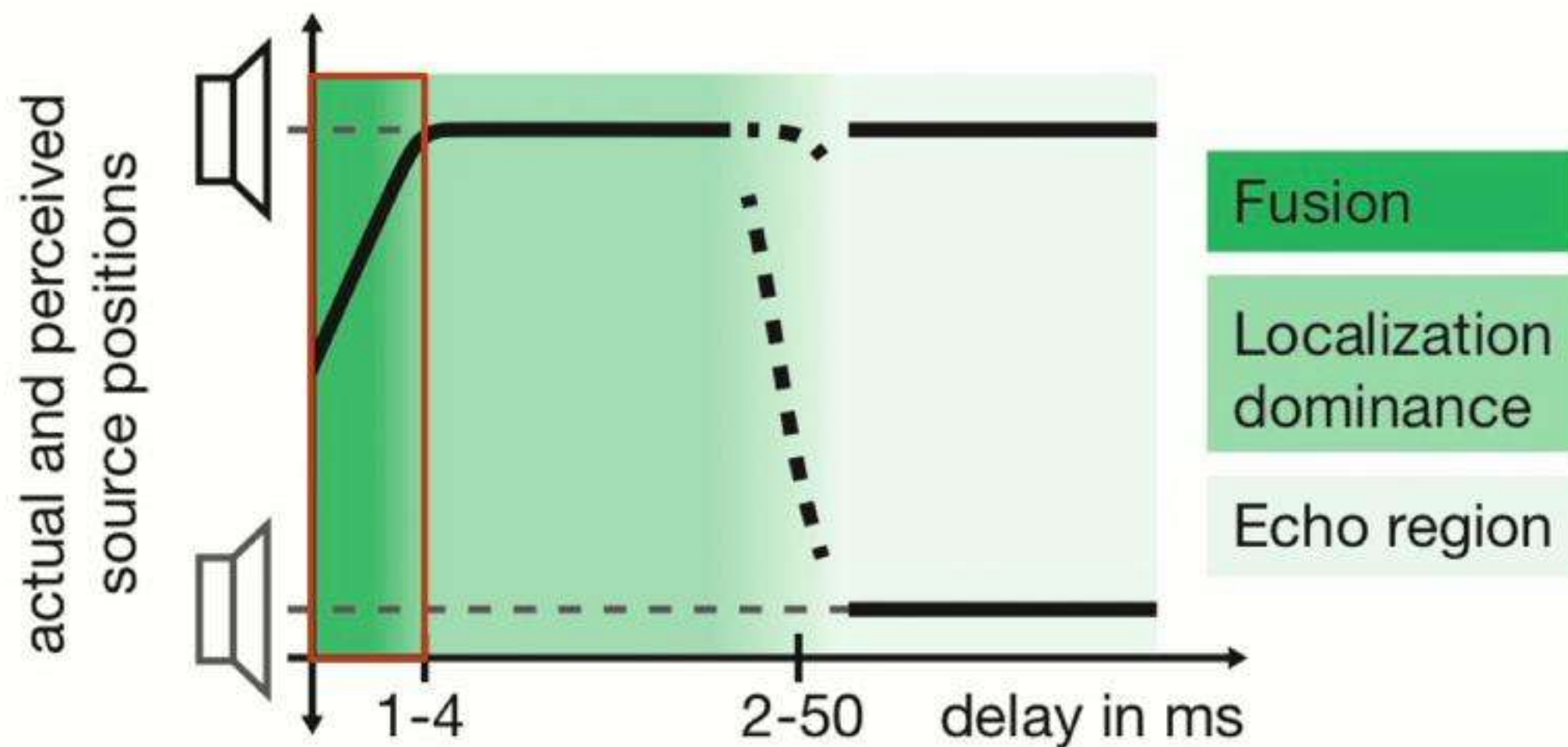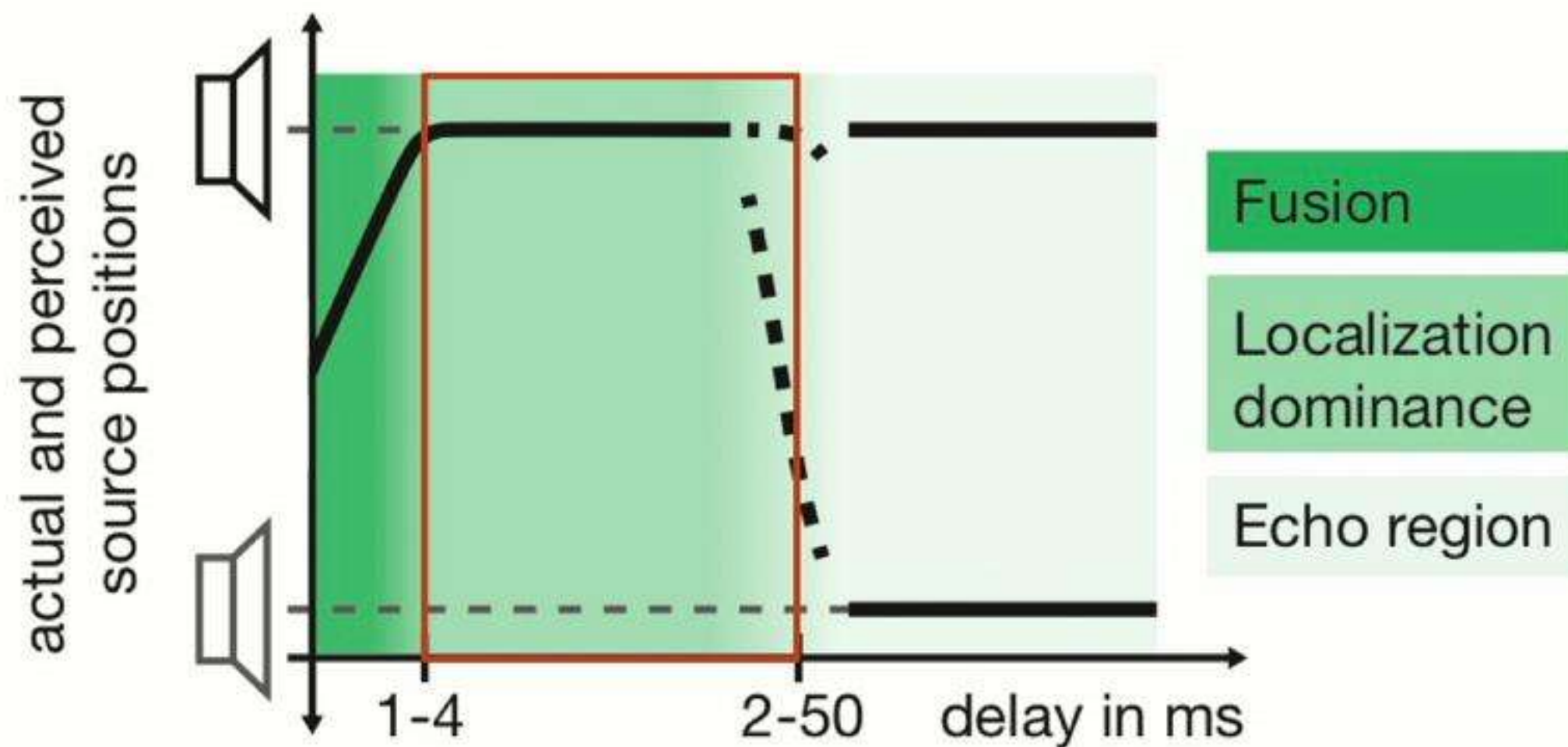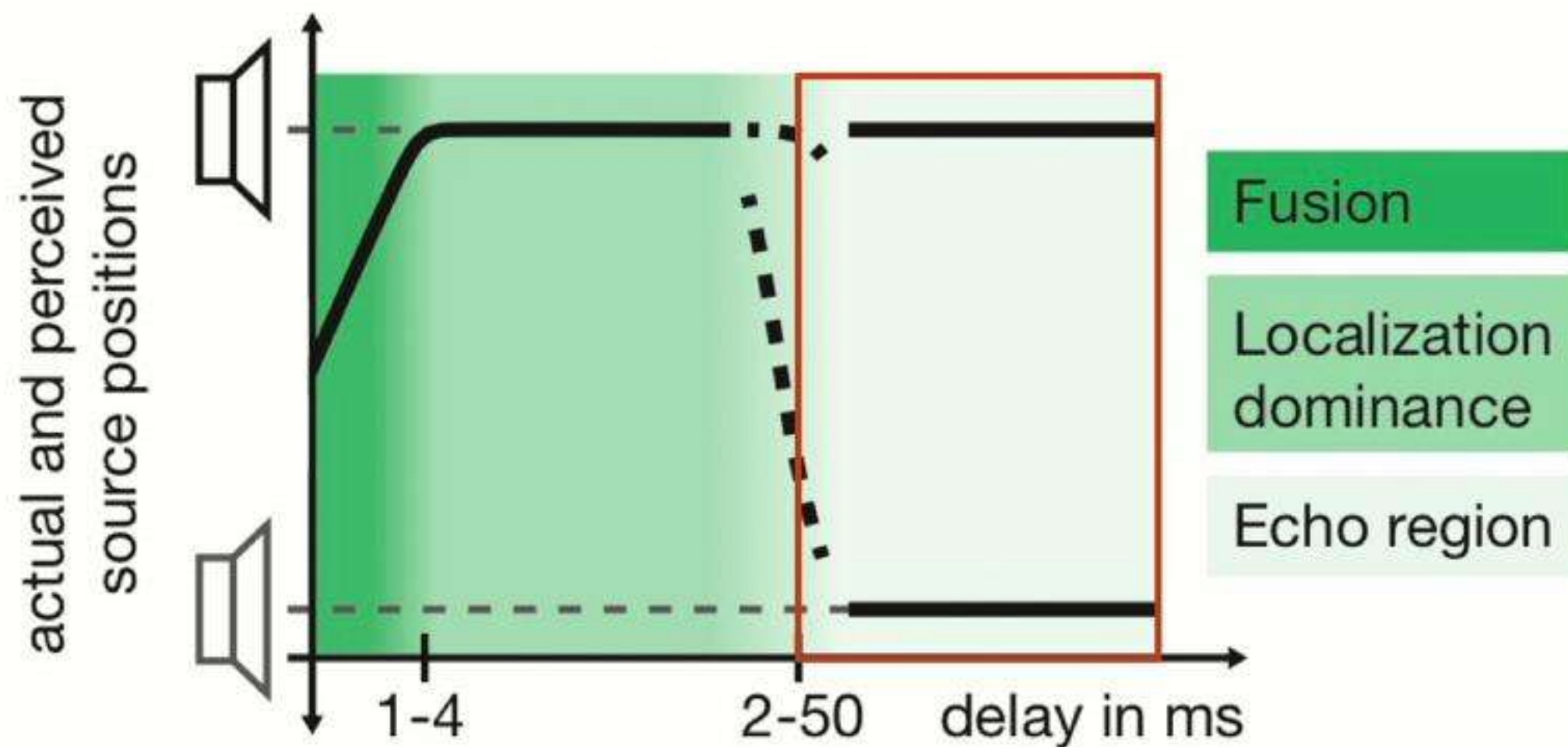# PRECEDENCE IN ROOM ACOUSTICS – TEMPORAL ASPECTS

Litovsky *et al.* (1999): The precedence effect. *J. Acoust. Soc. Am.* **104**(4): 1633 – 1654.

Brown *et al.* (2015): The precedence effect in sound localization. *J. Assoc. Res. Otolaryng.* **1**(16):1 – 28.

Images: Fabian Brinkmann (CC-BY 4.0)

# PRECEDENCE IN ROOM ACOUSTICS – ENERGETIC ASPECTS

Olive and Toole (1989): The detection of reflections in typical rooms. *J. Audio Eng. Soc.* **37**(7/8): 539–553.

Rakerd *et al.* (2000): Echo suppression in the horizontal and median sagittal plane. *J. Acoust. Soc. Am.* **107**(2):1061–1064.

Jensen and Welti (2003): The importance of reflections in a binaural room impulse response. *114th AES Convention.*

# PRECEDENCE IN ROOM ACOUSTICS – SPATIAL ASPECTS

- Decreased threshold with increasing spatial separation

- Differences of 10 - 15 dB

Bech (1995/1996): Timbral aspects of reproduced sound in small rooms. Part I/II. *J. Acoust. Soc. Am.* **97**(3) and **99**(6).

Litovsky and Cunningham (2001): Investigation of the relationship among three common measures of precedence. *J. Acoust. Soc. Am.*

Best *et al.* (2004): Separation of concurrent broadband sound sources... *J. Acoust. Soc. Am.* **115**(1):324–336.

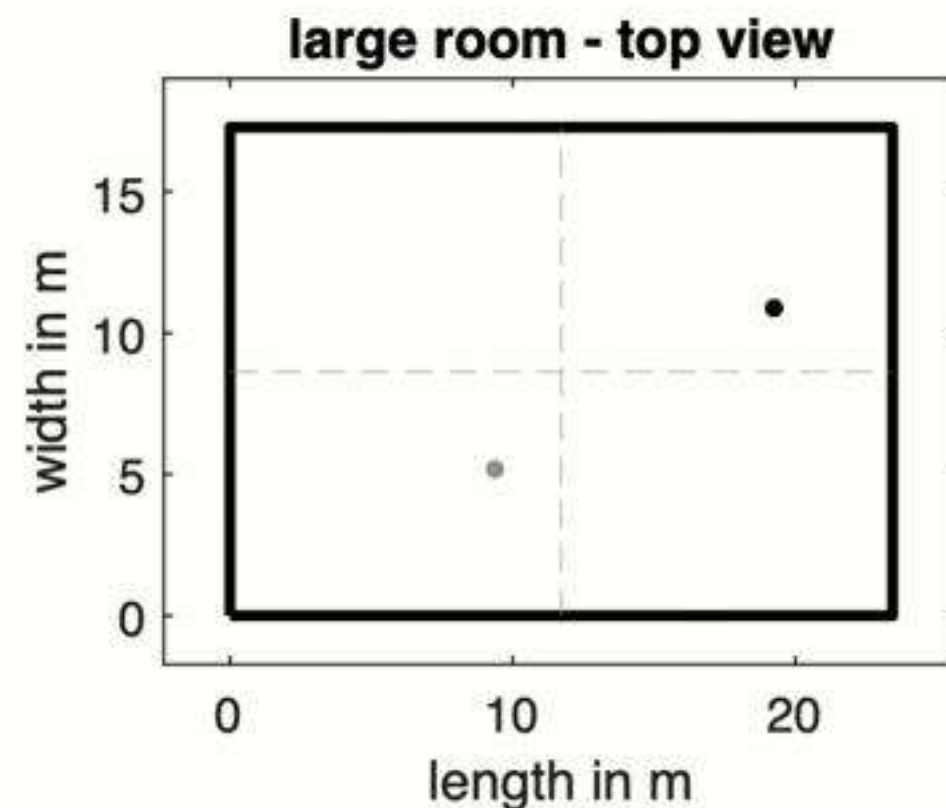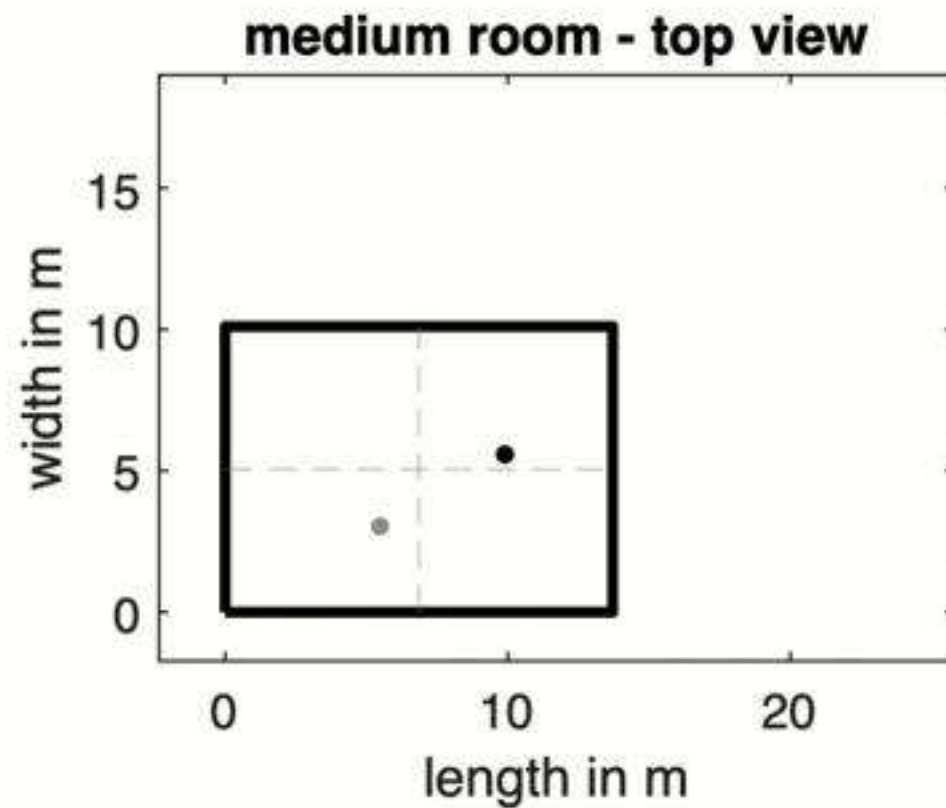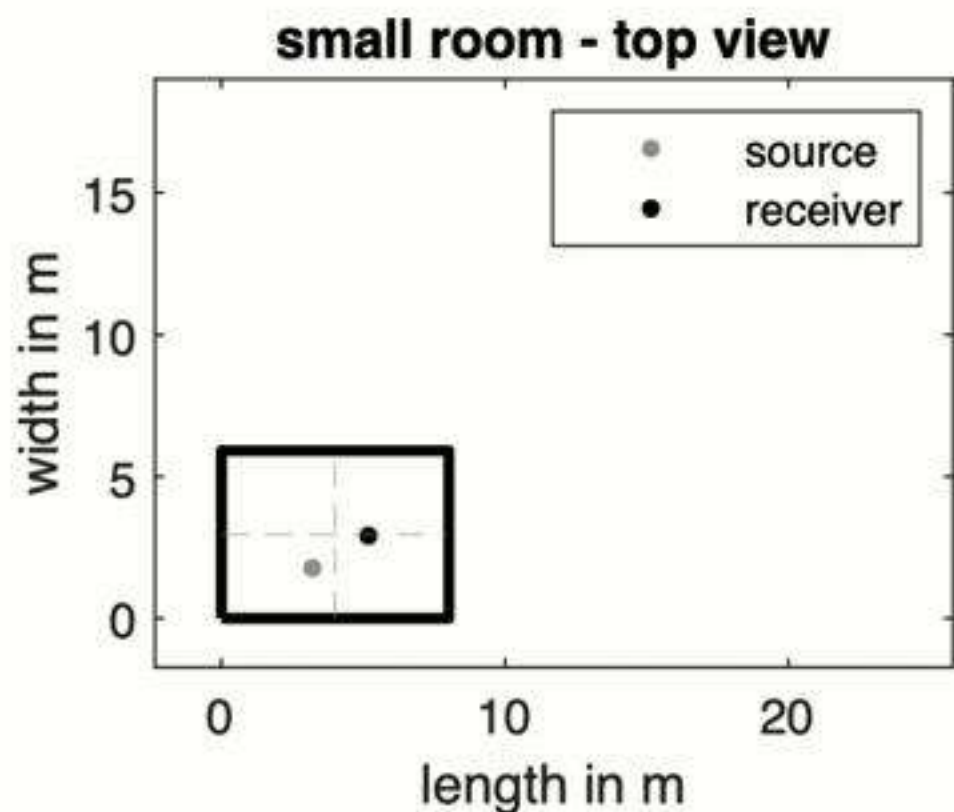INTRODUCTION    METHOD    RESULTS    DISCUSSION

# DATABASE

- 9 empty shoe box rooms
  - 3 Reverberation times: 0.5 s, 1 s, 2 s
  - 3 Volumes: 200 m³, 1000 m³, 5000 m³

# DATABASE – IMAGE SOURCE MODEL



Brinkmann and Weinzierl (2017): AKtools. 142nd AES Convention, Berlin, Germany

# DATABASE – IMAGE SOURCE MODEL



Brinkmann and Weinzierl (2017): AKtools. 142nd AES Convention, Berlin, Germany

# DATABASE - TRITON

- Wave based simulation up to 8 kHz

- Direction from sound intensity $\boldsymbol{I}$

$$\boldsymbol{I} = p\boldsymbol{v}, \qquad \boldsymbol{v} = -\frac{1}{\rho_0}\int \nabla p \, \mathrm{d}t$$

- Pressure gradient $\nabla p$ from neighboring points in x/y/z-direction

- $10^{\text{th}}$ order, zero phase low-pass @ 2 kHz

Raghuvanshi and Snyder (2018): Parametric directional coding for precomputed sound propagation. *ACM Trans. Graph.*

# DATABASE - TRITON



Raghuvanshi and Snyder (2018): Parametric directional coding for precomputed sound propagation. *ACM Trans. Graph.*

# DATABASE - TRITON



Raghuvanshi and Snyder (2018): Parametric directional coding for precomputed sound propagation. *ACM Trans. Graph.*

# DATABASE - TRITON



Raghuvanshi and Snyder (2018): Parametric directional coding for precomputed sound propagation. *ACM Trans. Graph.*

# DATABASE - TRITON



Raghuvanshi and Snyder (2018): Parametric directional coding for precomputed sound propagation. *ACM Trans. Graph.*

# PROPOSED ENCODING

a)   estimate $\tau_0$, $a_0$, and $\measuredangle_0 = \{\phi_0, \theta_0\}$

# PROPOSED ENCODING

a) estimate $\tau_0$, $a_0$, and $\measuredangle_0 = \{\phi_0, \theta_0\}$

b) estimate $\tau_i$, $a_i$, and $\measuredangle_i$ based on masking threshold

# PROPOSED ENCODING

a) estimate $\tau_0$, $a_0$, and $\measuredangle_0 = \{\phi_0, \theta_0\}$

b) estimate $\tau_i$, $a_i$, and $\measuredangle_i$ based on masking threshold

c) select $N$ reflections

# PROPOSED ENCODING

a)   estimate $\tau_0$, $a_0$, and $\sphericalangle_0 = \{\phi_0, \theta_0\}$

b)   estimate $\tau_i$, $a_i$, and $\sphericalangle_i$ based on masking threshold

c)   select $N$ reflections

d)   estimate late reverberation based on residual energy
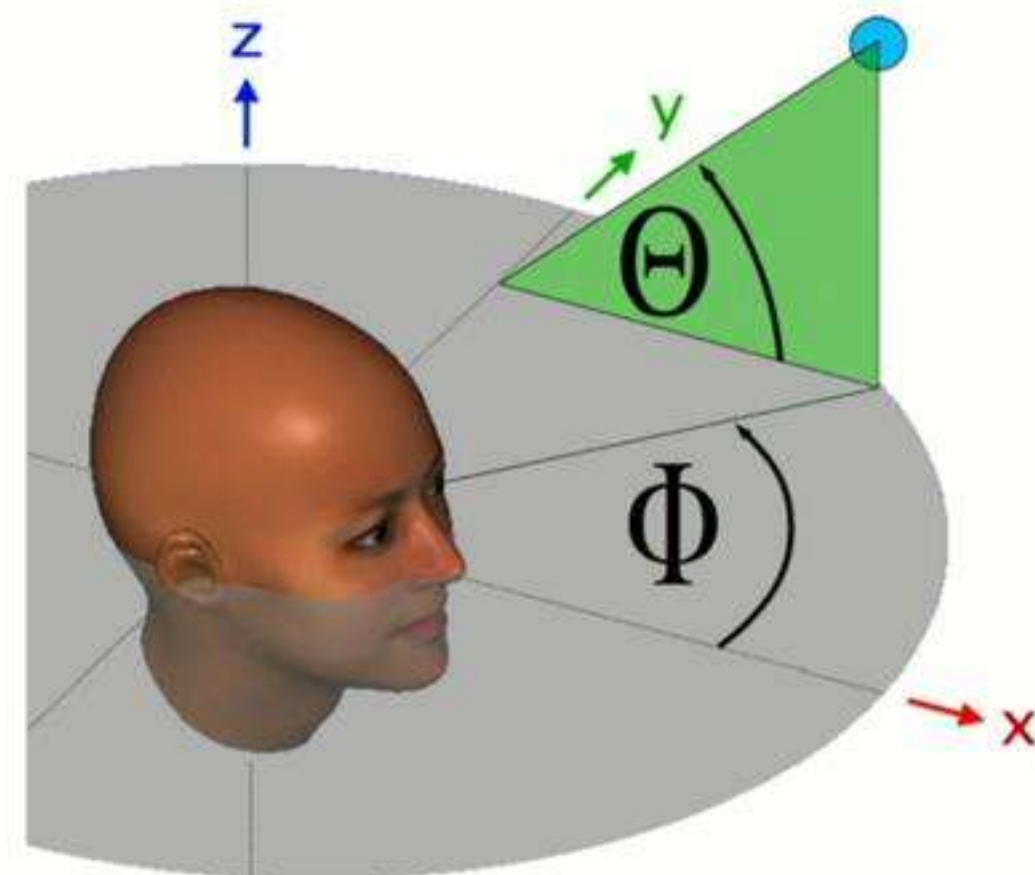
# PROPOSED ENCODING: $\tau_0, a_0$

$\tau_0$: Onset estimator according to [1]

$$a_0 = \left[ \frac{1}{1.5 \text{ ms}} \int_{\tau_0 - 0.5 \text{ ms}}^{\tau_0 + 1 \text{ ms}} p(t)^2 \mathrm{d}t \right]^{1/2}$$

[1] Raghuvanshi and Snyder (2018): Parametric directional coding for precomputed sound propagation. *ACM Trans. Graph.*

# PROPOSED ENCODING: $\sphericalangle_0$

$$\phi_0 = \frac{1}{\int_{\tau_0-0.5\text{ ms}}^{\tau_0+1\text{ ms}} p^2 \mathrm{d}t} \int\limits_{\tau_0-0.5\text{ ms}}^{\tau_0+1\text{ ms}} p(t)^2 \phi(t) \mathrm{d}t$$



Ziegelwanger and Majdak (2014): Modeling the direction continuous time-of-arival. *J. Acust. Soc. Am.* **135**(3): 1278 – 1293.

# PROPOSED ENCODING: $\sphericalangle_0$

$$\phi_0 = \frac{1}{\int_{\tau_0-0.5\,\text{ms}}^{\tau_0+1\,\text{ms}} p^2 \mathrm{d}t} \int_{\tau_0-0.5\,\text{ms}}^{\tau_0+1\,\text{ms}} p(t)^2 \phi(t) \mathrm{d}t$$

$$\theta_0 = \angle \left( \int_{\tau_0-0.5\,\text{ms}}^{\tau_0+1\,\text{ms}} p(t)^2 e^{-i\theta(t)} \mathrm{d}t \right)$$
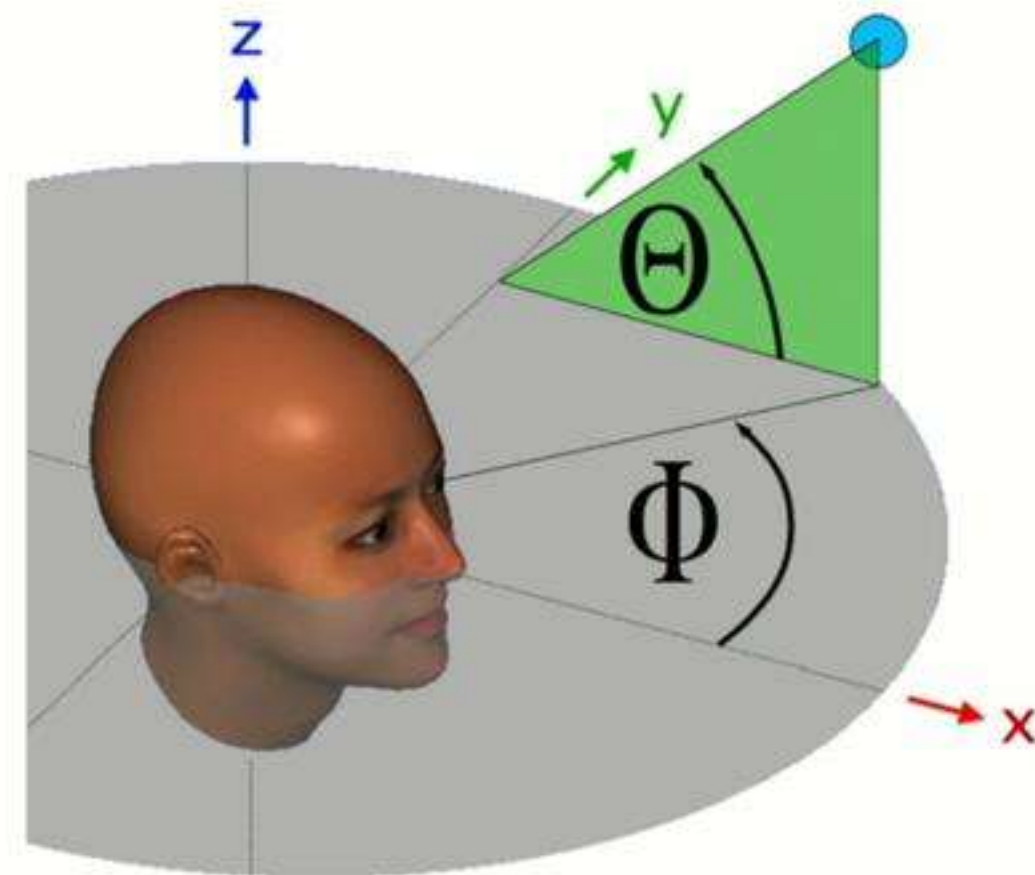


Ziegelwanger and Majdak (2014): Modeling the direction continuous time-of-arival. *J. Acust. Soc. Am.* **135**(3): 1278 – 1293.

# PROPOSED ENCODING: $\sphericalangle_0$

$$\phi_0 = \frac{1}{\int_{\tau_0 - 0.5 \text{ ms}}^{\tau_0 + 1 \text{ ms}} p^2 \mathrm{d}t} \int_{\tau_0 - 0.5 \text{ ms}}^{\tau_0 + 1 \text{ ms}} p(t)^2 \phi(t) \mathrm{d}t$$
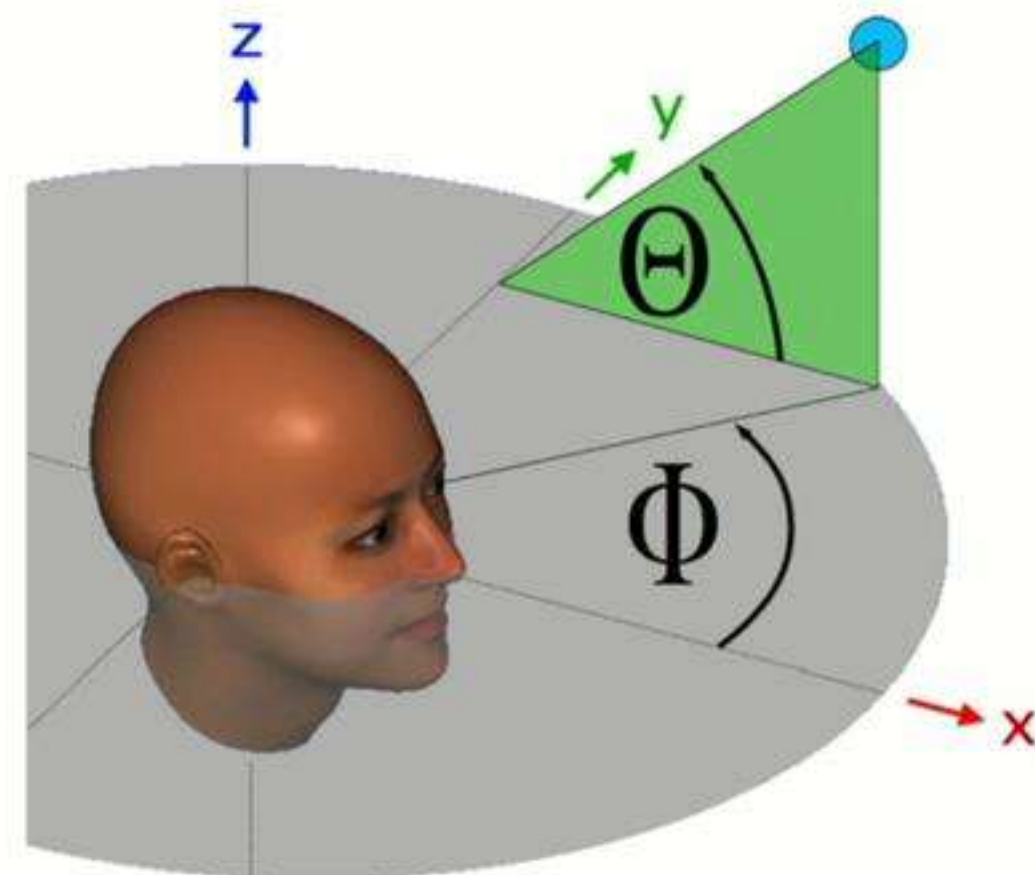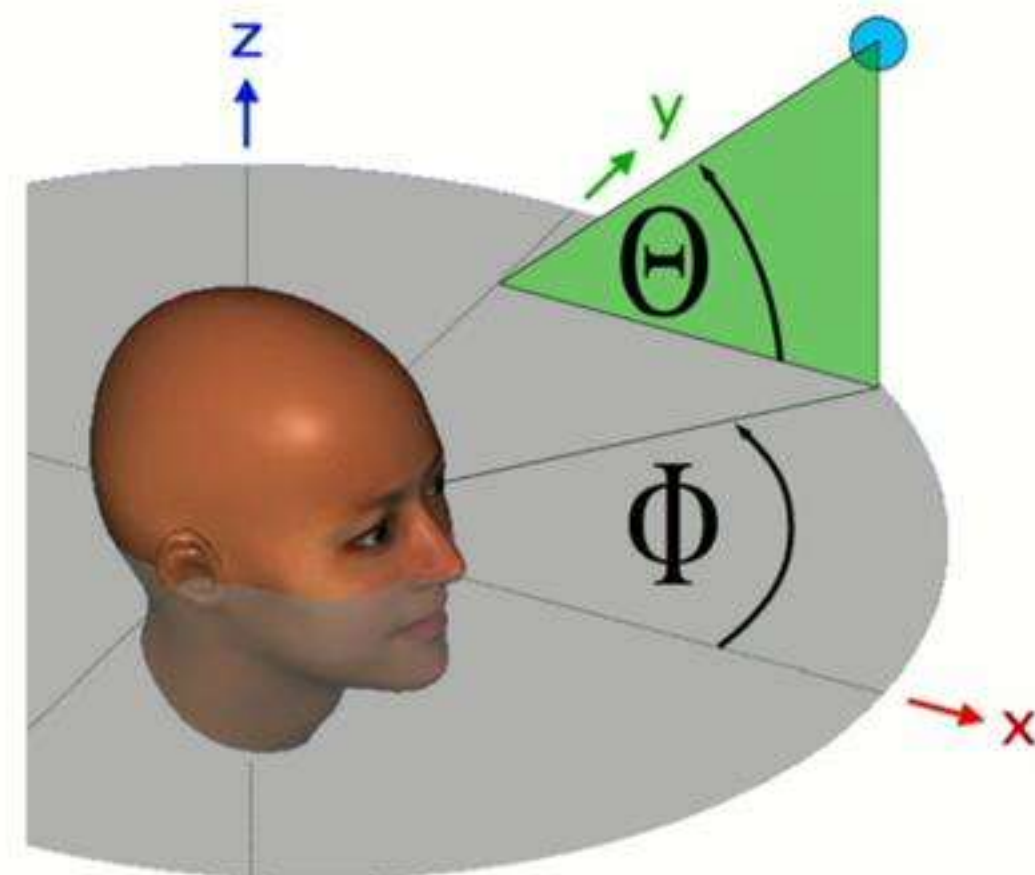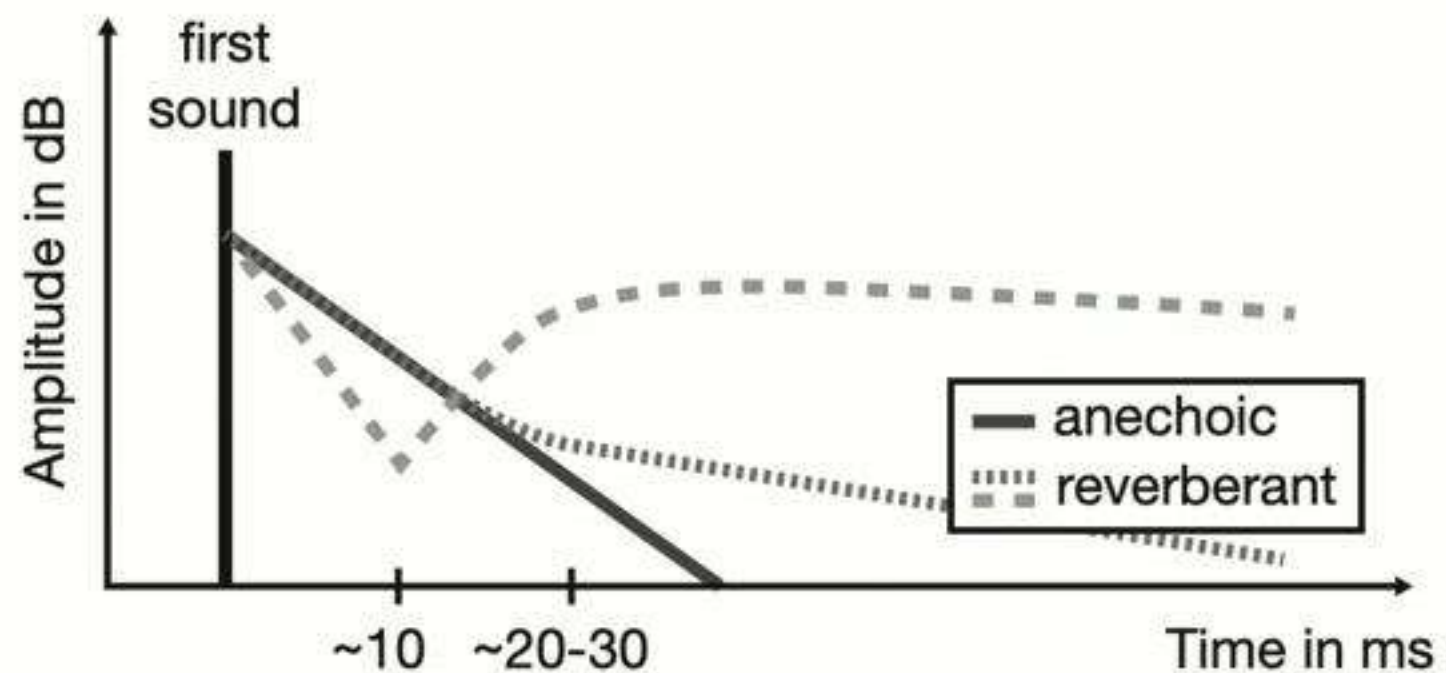


Ziegelwanger and Majdak (2014): Modeling the direction continuous time-of-arival. *J. Acust. Soc. Am.* **135**(3): 1278 – 1293.

# PROPOSED ENCODING: $\sphericalangle_0$

$$\phi_0 = \frac{1}{\int_{\tau_0 - 0.5 \text{ ms}}^{\tau_0 + 1 \text{ ms}} p^2 \mathrm{d}t} \int\limits_{\tau_0 - 0.5 \text{ ms}}^{\tau_0 + 1 \text{ ms}} p(t)^2 \phi(t) \mathrm{d}t$$

$$\theta_0 = \angle \left( \int\limits_{\tau_0 - 0.5 \text{ ms}}^{\tau_0 + 1 \text{ ms}} p(t)^2 e^{-i\theta(t)} \, \mathrm{d}t \right)$$



Ziegelwanger and Majdak (2014): Modeling the direction continuous time-of-arival. *J. Acust. Soc. Am.* **135**(3): 1278 – 1293.

# PROPOSED ENCODING: MASKING THRESHOLD



slope: -1 dB/ms

offset: -10 dB

v-shape: add 35% reflection energy to threshold

# PROPOSED ENCODING: MASKING THRESHOLD



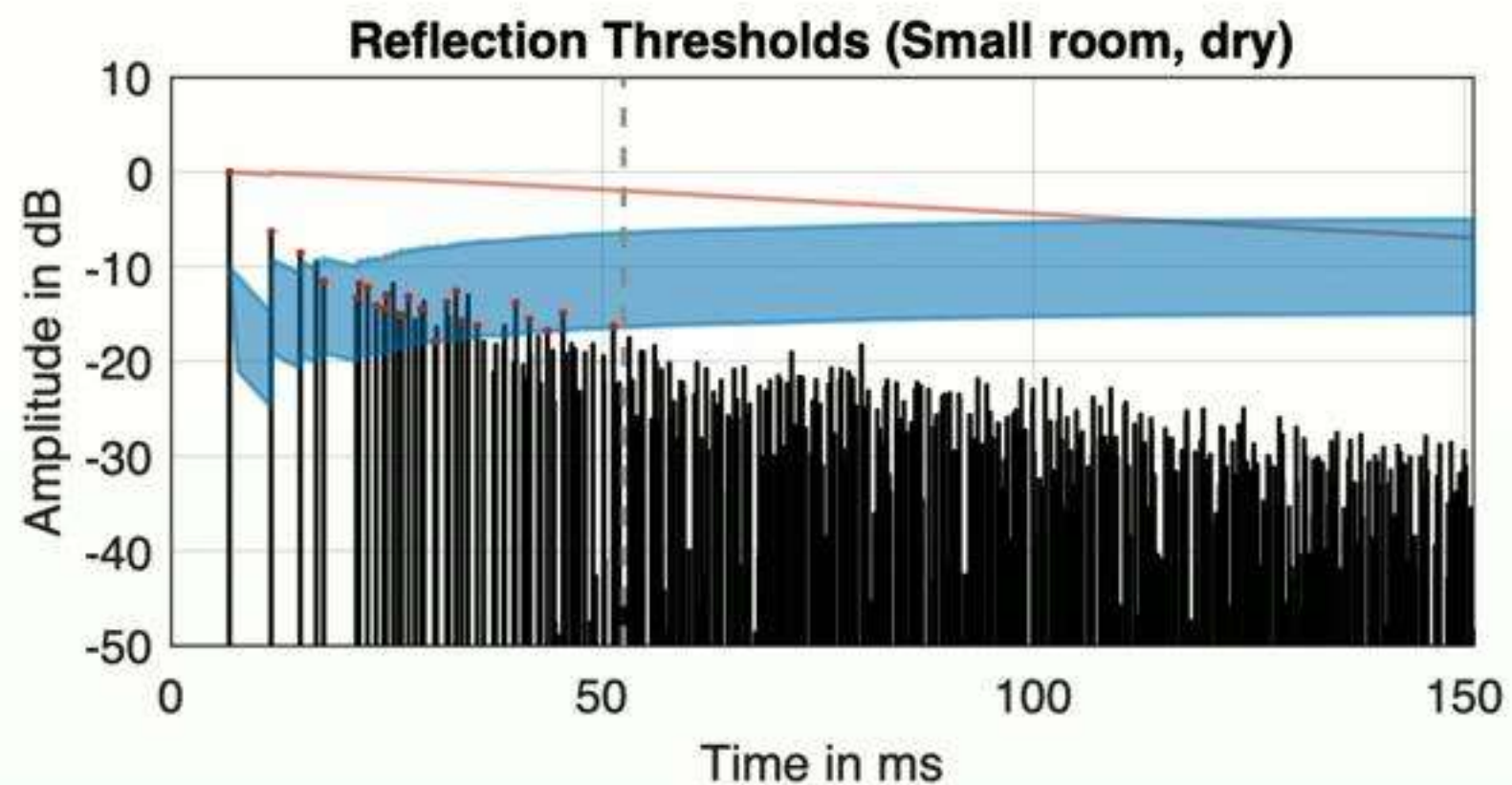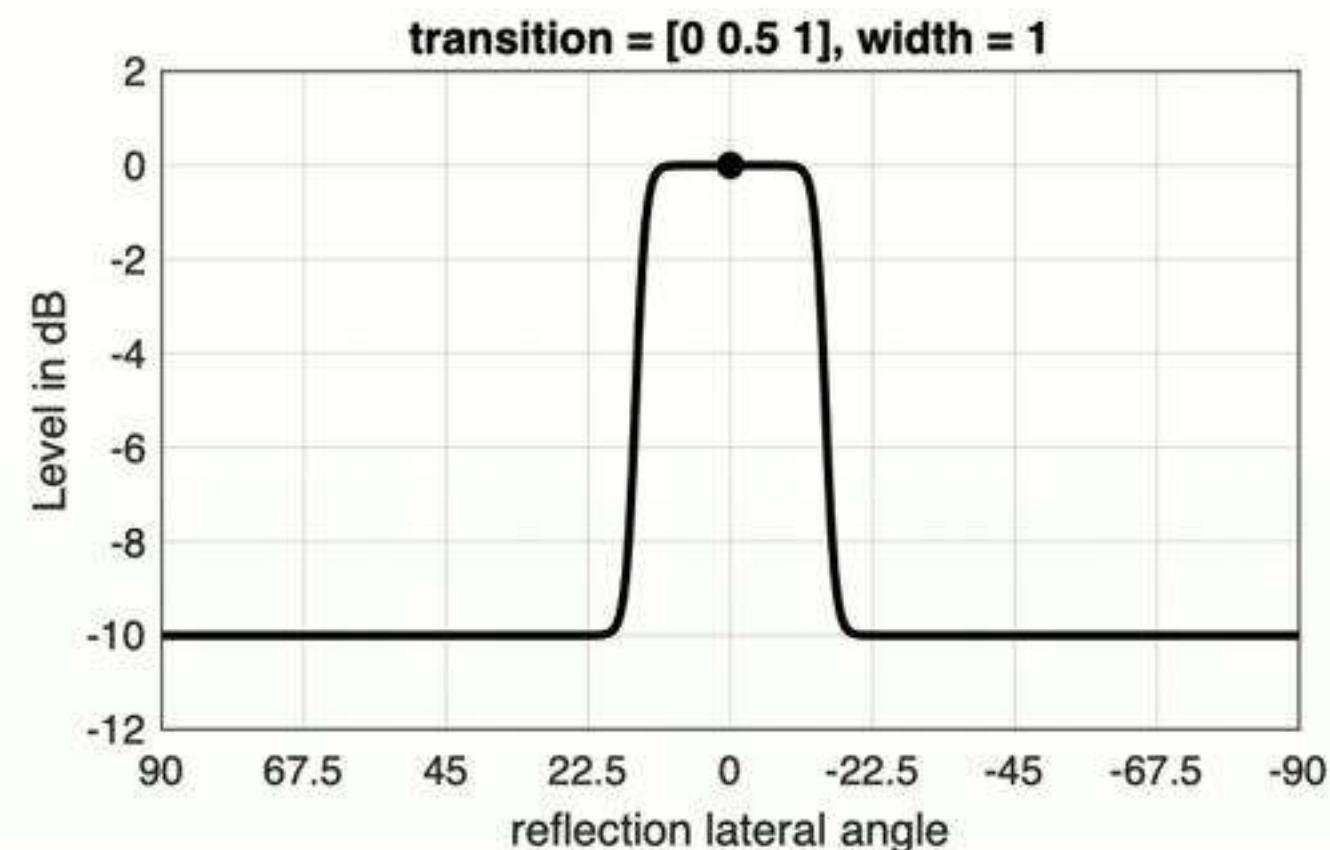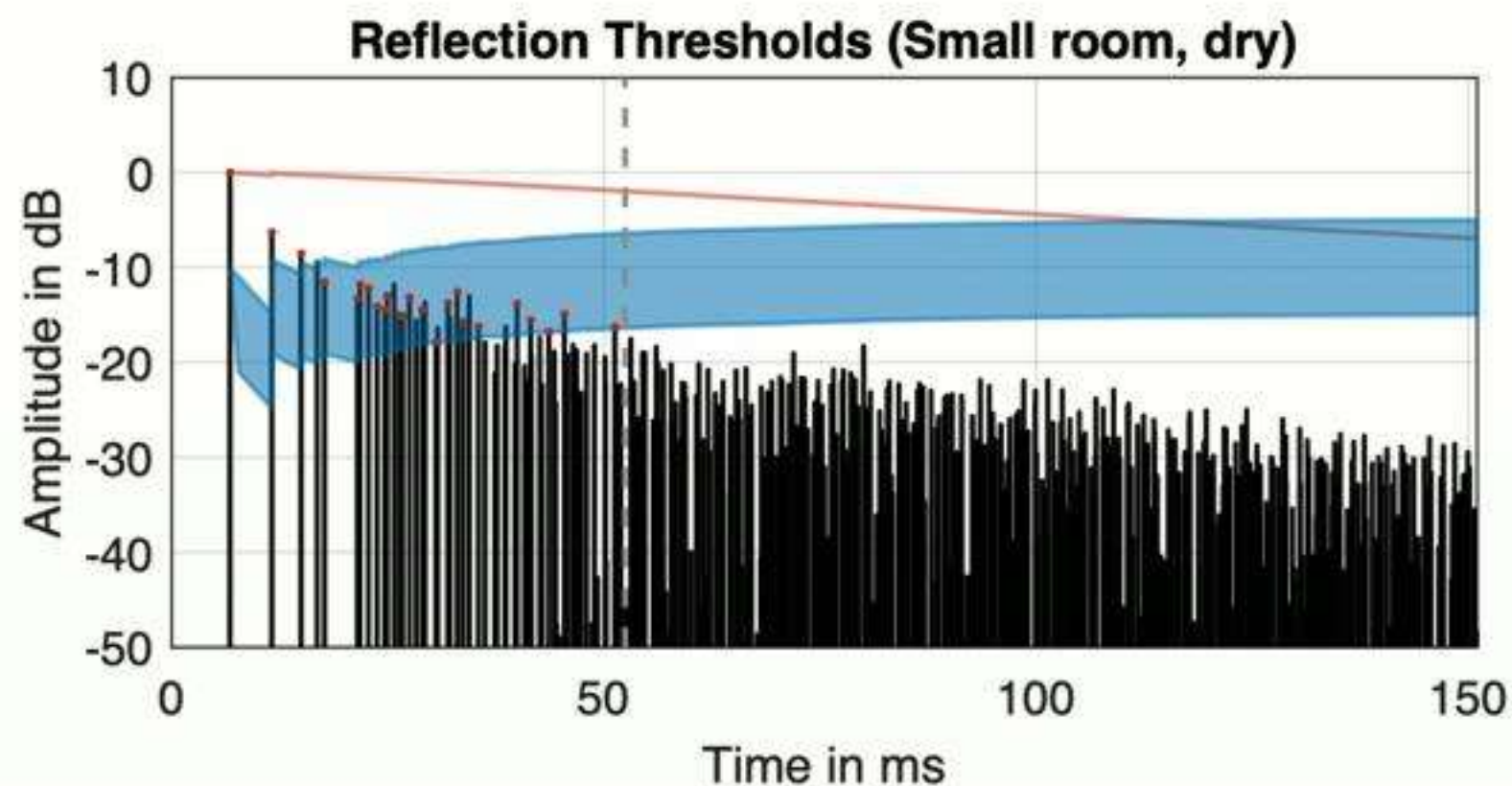Reflection Thresholds (Small room, dry)

slope: -1 dB/ms

offset: -10 dB

v-shape: add 35% reflection energy to threshold

# PROPOSED ENCODING: MASKING THRESHOLD



slope: -1 dB/ms

offset: -10 dB

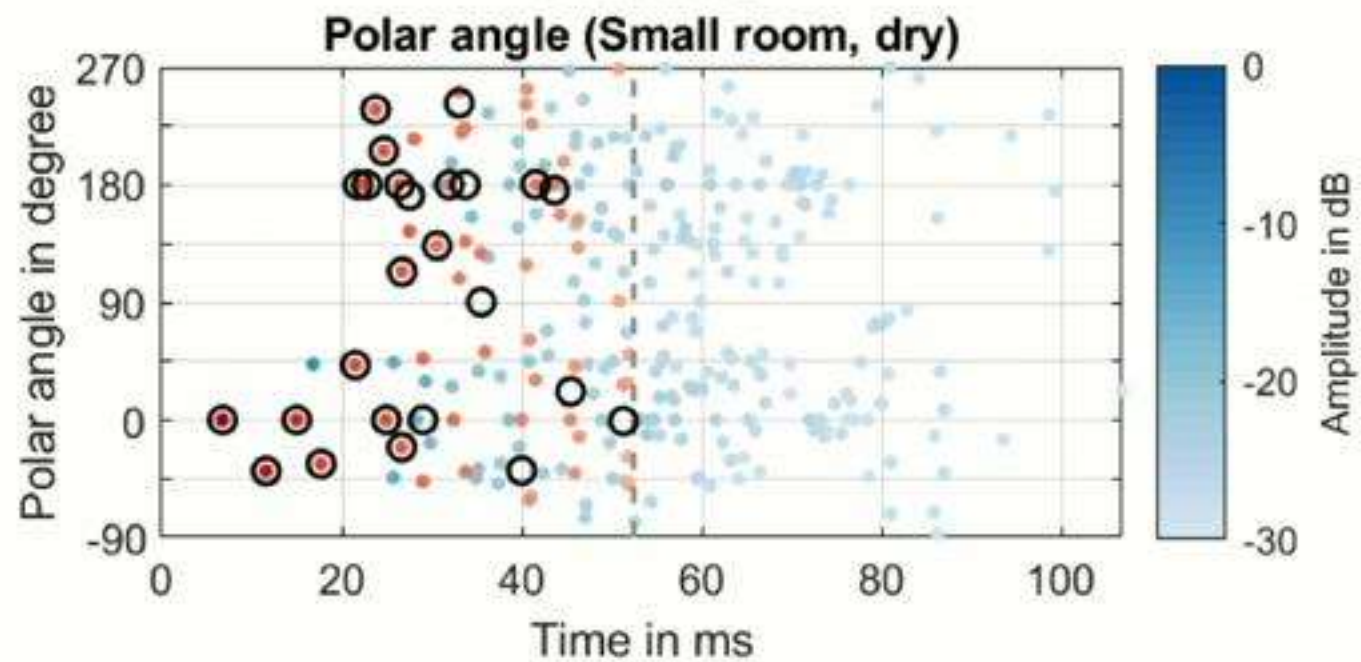v-shape: add 35% reflection energy to threshold

width according to [1]
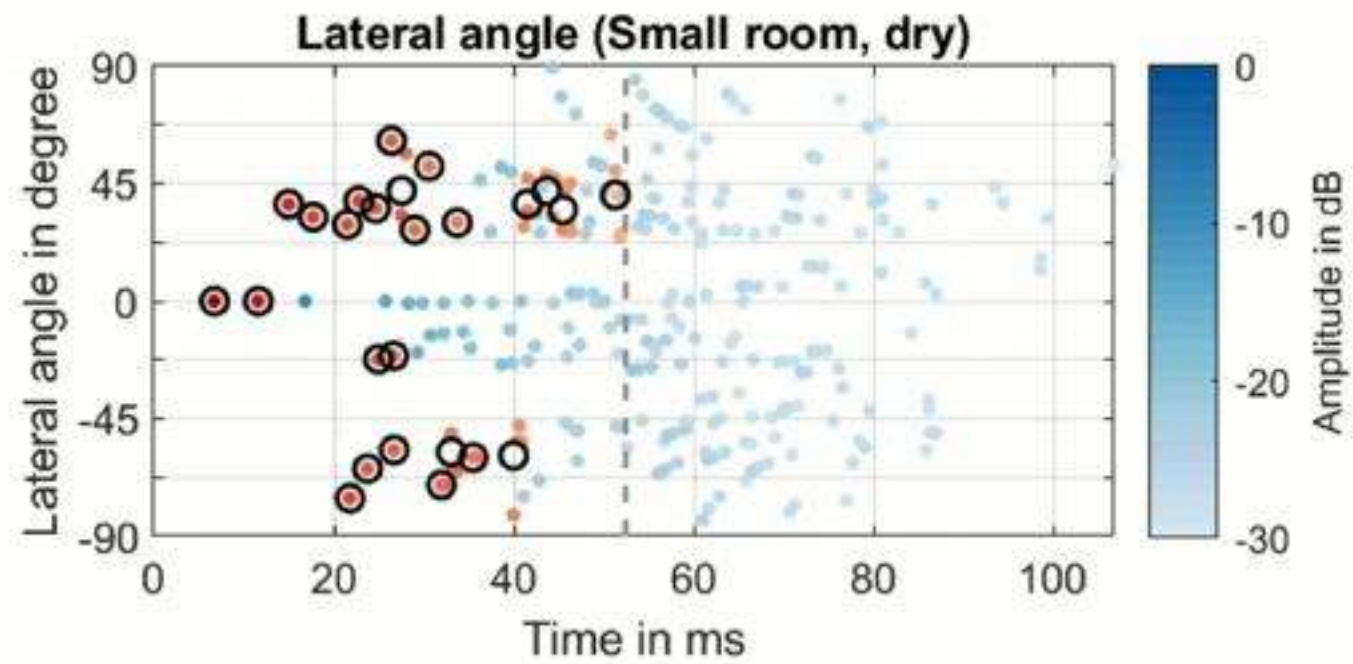
depth: -10 dB

[1] Best *et al.* (2004): Separation of concurrent broadband sound sources... *J. Acoust. Soc. Am.* **115**(1):324–336.
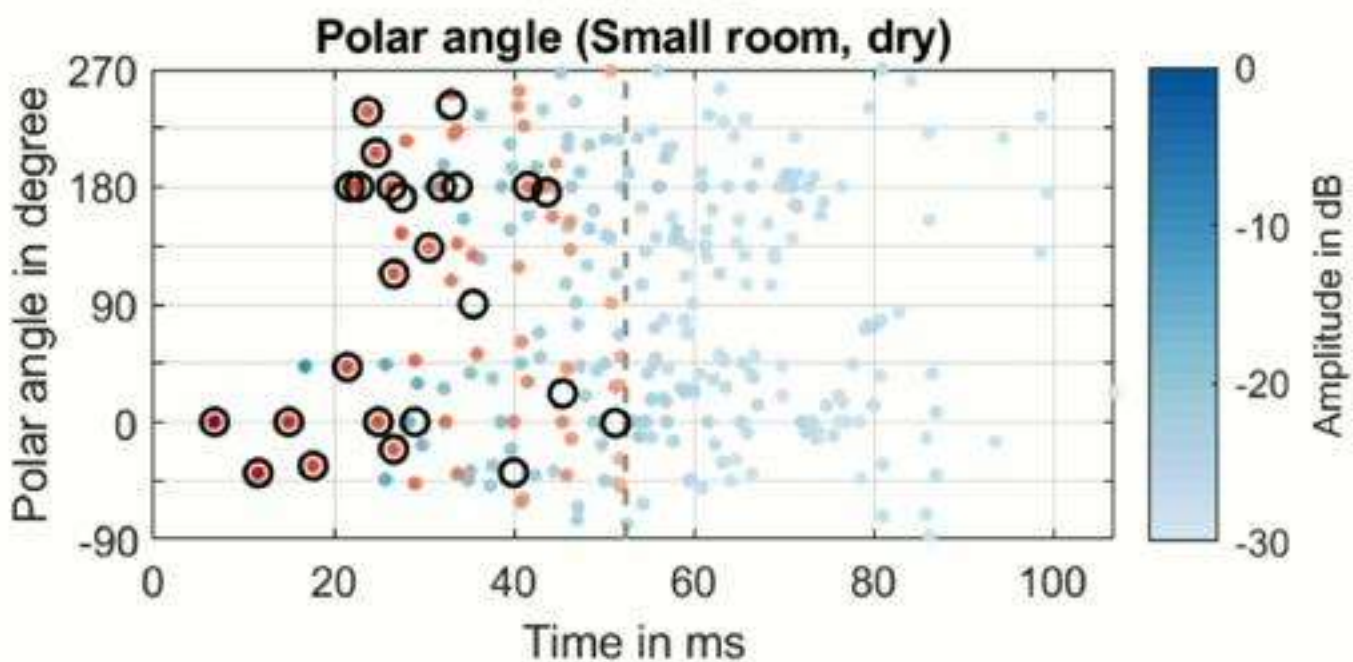
## Image source model
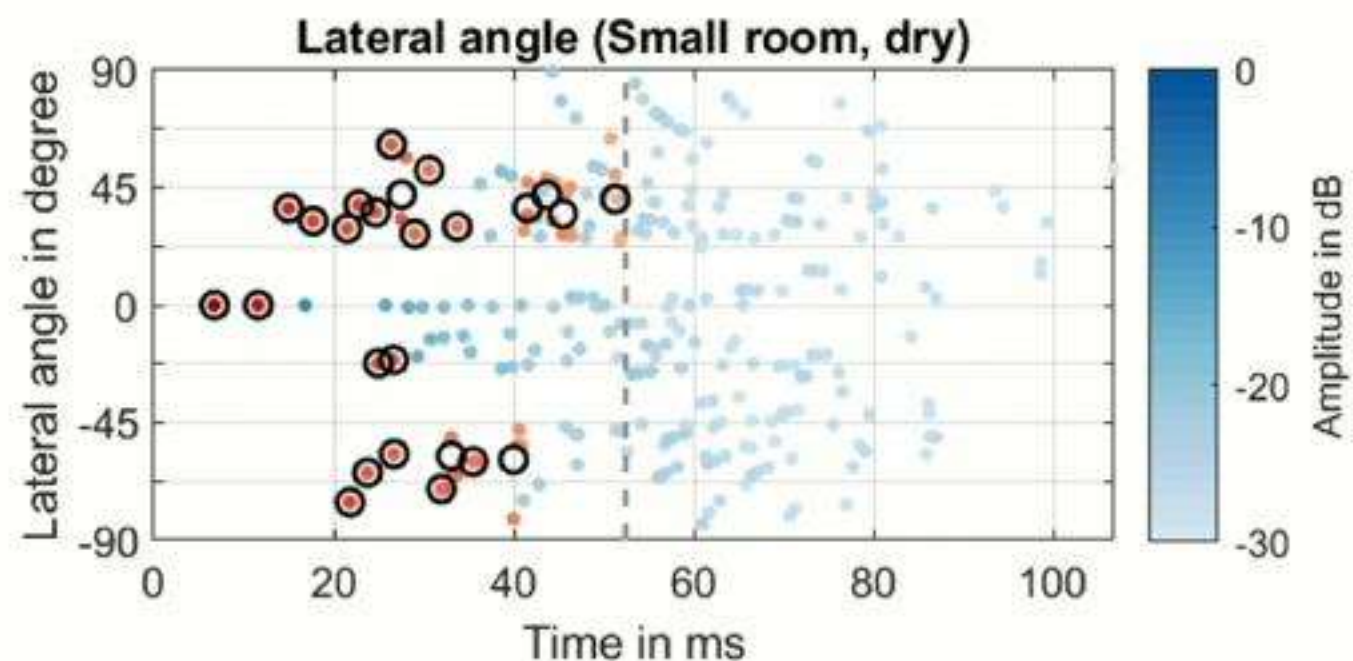
# PROPOSED ENCODING: DETECTED REFLECTIONS



## Image source model

## Triton

# PROPOSED ENCODING: SELECTION



EDC (Small room dry)



Lateral angle (Small room dry)

**First**

# PROPOSED ENCODING: SELECTION



**First**

**Exceed**

# PROPOSED ENCODING: SELECTION



## First       Exceed       Loudest

# PROPOSED ENCODING: LATE REVERBERATION

# PROPOSED ENCODING: LATE REVERBERATION

INTRODUCTION

METHOD

RESULTS

DISCUSSION

# PHYSICAL EVALUATION

- 0, 1, 2, 4, 6, 8, 10 and 15 reflections

- Three selection methods (first, exceed, loudest)

- Two late reverberations (single, double ramp)

- Comparison against reference

# PHYSICAL EVALUATION

**Max. absolute gain missmatch**

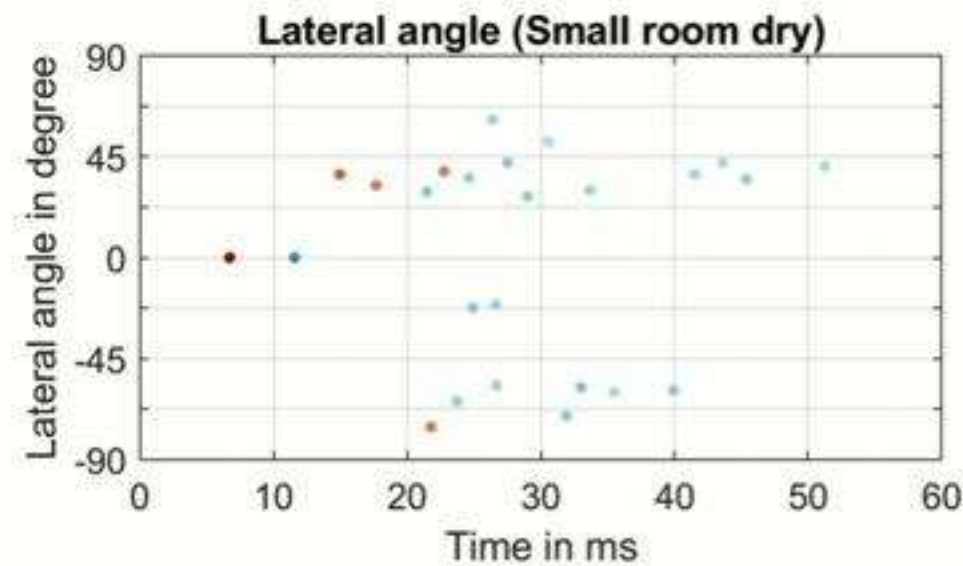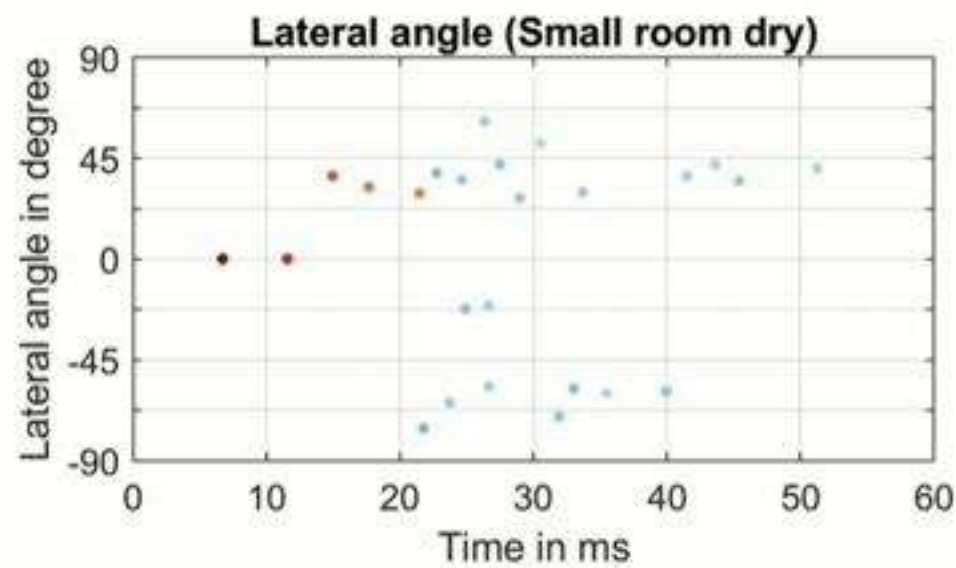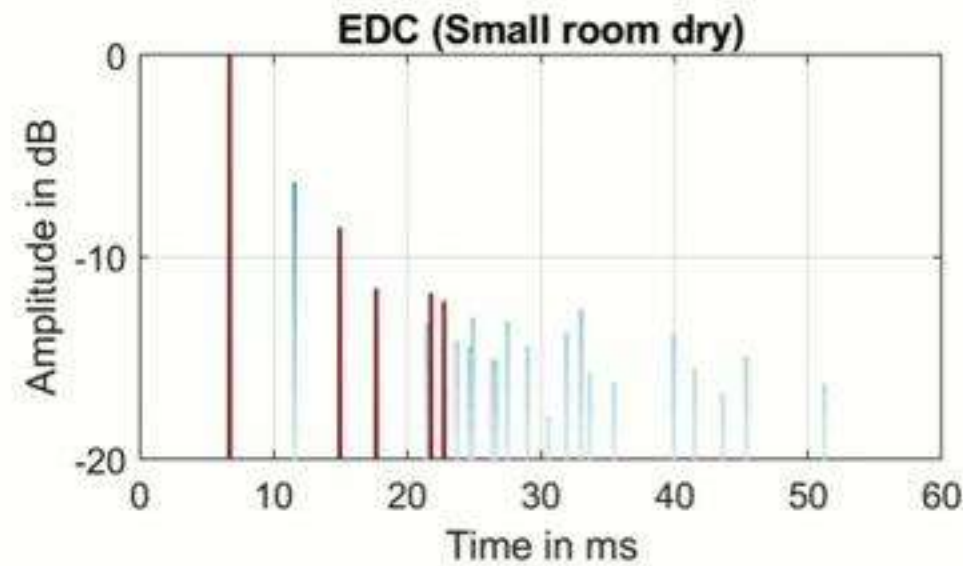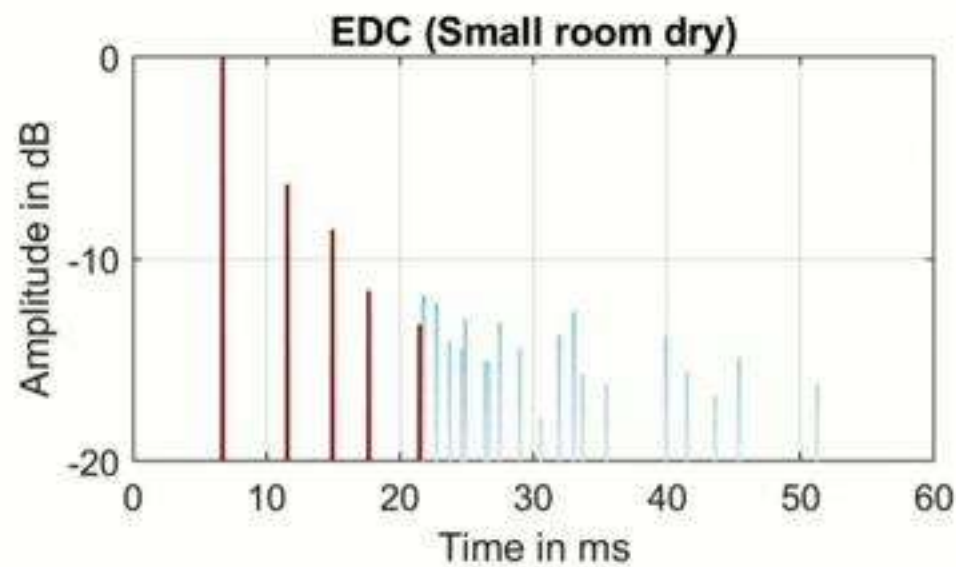| | first a | first c | loudest a | loudest c | exceed a | exceed c | mean |
|---|---|---|---|---|---|---|---|
| Small room dry | 0.6 | 0.4 | 0.9 | 0.4 | 0.9 | 0.3 | 0.6 |
| Medium room dry | 0.5 | 0.5 | 0.5 | 0.5 | 0.8 | 0.3 | 0.5 |
| Large room dry | 1 | 1 | 0.6 | 0.3 | 0.6 | 0.3 | 0.7 |
| Small room medium | 0.6 | 0.6 | 1.1 | 0.6 | 0.9 | 0.5 | 0.7 |
| Medium room medium | 0.4 | 0.4 | 0.6 | 0.4 | 0.7 | 0.5 | 0.5 |
| Large room medium | 0.7 | 0.8 | 0.7 | 0.8 | 0.5 | 0.8 | 0.7 |
| Small room wet | 1.3 | 1.3 | 1.5 | 1.3 | 1.4 | 1.2 | 1.3 |
| Medium room wet | 0.2 | 0.2 | 0.5 | 0.2 | 0.2 | 0.2 | 0.2 |
| Large room wet | 0.6 | 0.6 | 0.6 | 0.6 | 0.5 | 0.6 | 0.6 |
| mean | 0.7 | 0.7 | 0.8 | 0.6 | 0.7 | 0.5 | 0.7 |

AE in dB

- Good level preservation

- Slightly audible differences

- Minor differences across algorithms

# PHYSICAL EVALUATION



Max. absolute spectral difference

- Slightly audible coloration
- Minor differences across algorithms

# PHYSICAL EVALUATION

**Max. absolute ILD difference**

| | first a | first c | loudest a | loudest c | exceed a | exceed c | mean |
|---|---|---|---|---|---|---|---|
| Small room dry | 1.3 | 1.2 | 0.9 | 0.8 | 1.1 | 0.9 | 1 |
| Medium room dry | 1.1 | 1.1 | 0.8 | 0.8 | 0.8 | 0.6 | 0.8 |
| Large room dry | 0.8 | 0.7 | 0.8 | 0.7 | 0.8 | 0.7 | 0.7 |
| Small room medium | 1.1 | 1 | 0.8 | 0.8 | 0.9 | 0.9 | 0.9 |
| Medium room medium | 0.7 | 0.7 | 0.7 | 0.6 | 0.6 | 0.6 | 0.7 |
| Large room medium | 0.9 | 0.9 | 0.9 | 0.9 | 0.8 | 0.7 | 0.8 |
| Small room wet | 1 | 1 | 0.8 | 0.8 | 0.9 | 0.9 | 0.9 |
| Medium room wet | 0.5 | 0.5 | 0.6 | 0.5 | 0.4 | 0.4 | 0.5 |
| Large room wet | 0.5 | 0.5 | 0.5 | 0.5 | 0.5 | 0.4 | 0.5 |
| mean | 0.9 | 0.8 | 0.7 | 0.7 | 0.7 | 0.7 | 0.8 |

AE in dB

- Slight mismatches in source position

- Minor differences across algorithms

# PHYSICAL EVALUATION

**Max. absolute IACC difference**

| | first a | first c | loudest a | loudest c | exceed a | exceed c | ref. |
|---|---|---|---|---|---|---|---|
| Small room dry | 0.1 | 0.1 | 0 | 0.1 | 0.1 | 0.1 | 0.4 |
| Medium room dry | 0.1 | 0.1 | 0 | 0.1 | 0.1 | 0.1 | 0.5 |
| Large room dry | 0.1 | 0.1 | 0 | 0.1 | 0 | 0.1 | 0.8 |
| Small room medium | 0.1 | 0.1 | 0 | 0.1 | 0.1 | 0.1 | 0.2 |
| Medium room medium | 0.1 | 0.1 | 0 | 0 | 0.1 | 0.1 | 0.2 |
| Large room medium | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.3 |
| Small room wet | 0 | 0 | 0 | 0.1 | 0.1 | 0.1 | 0.2 |
| Medium room wet | 0 | 0 | 0 | 0 | 0 | 0 | 0.1 |
| Large room wet | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 | 0.1 |
| mean | 0.1 | 0.1 | 0 | 0.1 | 0.1 | 0.1 | 0.1 |

AE in dB

- Slight mismatches in source width & envelopment

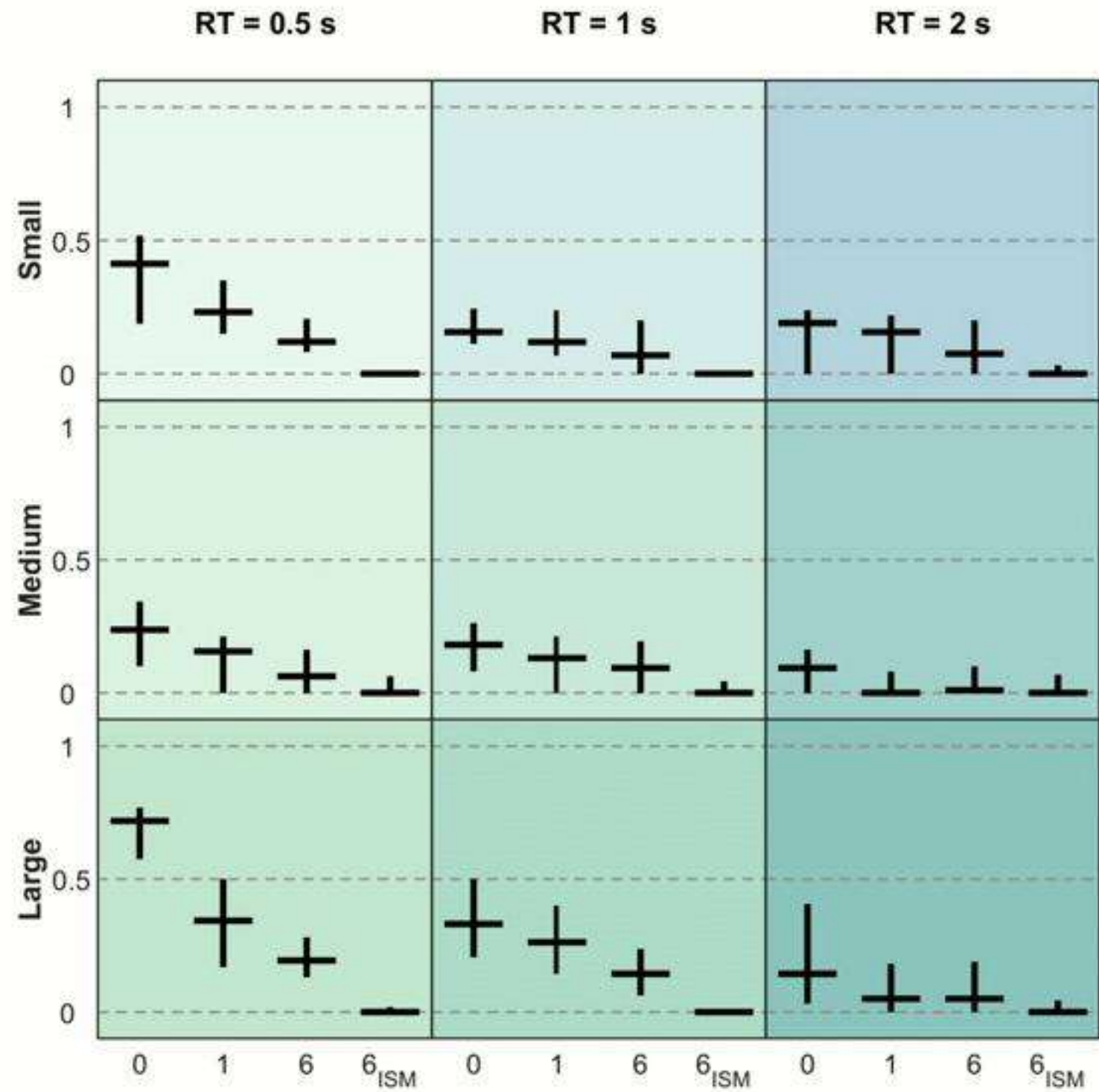- Minor differences across algorithms

# PERCEPTUAL EVALUATION

- 27 participants (5 female, 22 male, 38 yrs.)

- 0, 1, and 6 reflections

- First order image source model as anchor ($6_{ISM}$)

- 'Loudest' selection method

- 'Double sloped' late reverberation

- Running anechoic speech

- Rated against reference

# PERCEPTUAL EVALUATION

# PERCEPTUAL EVALUATION



- Differences decrease with **N**

# PERCEPTUAL EVALUATION



- Differences decrease with **N**

- $6_{ISM}$ is always rated 0

# PERCEPTUAL EVALUATION



- Differences decrease with **N**

- $6_{ISM}$ is always rated 0

- Differences decrease with **RT**

# PERCEPTUAL EVALUATION



- Differences decrease with **N**

- $6_{ISM}$ is always rated 0

- Differences decrease with **RT**
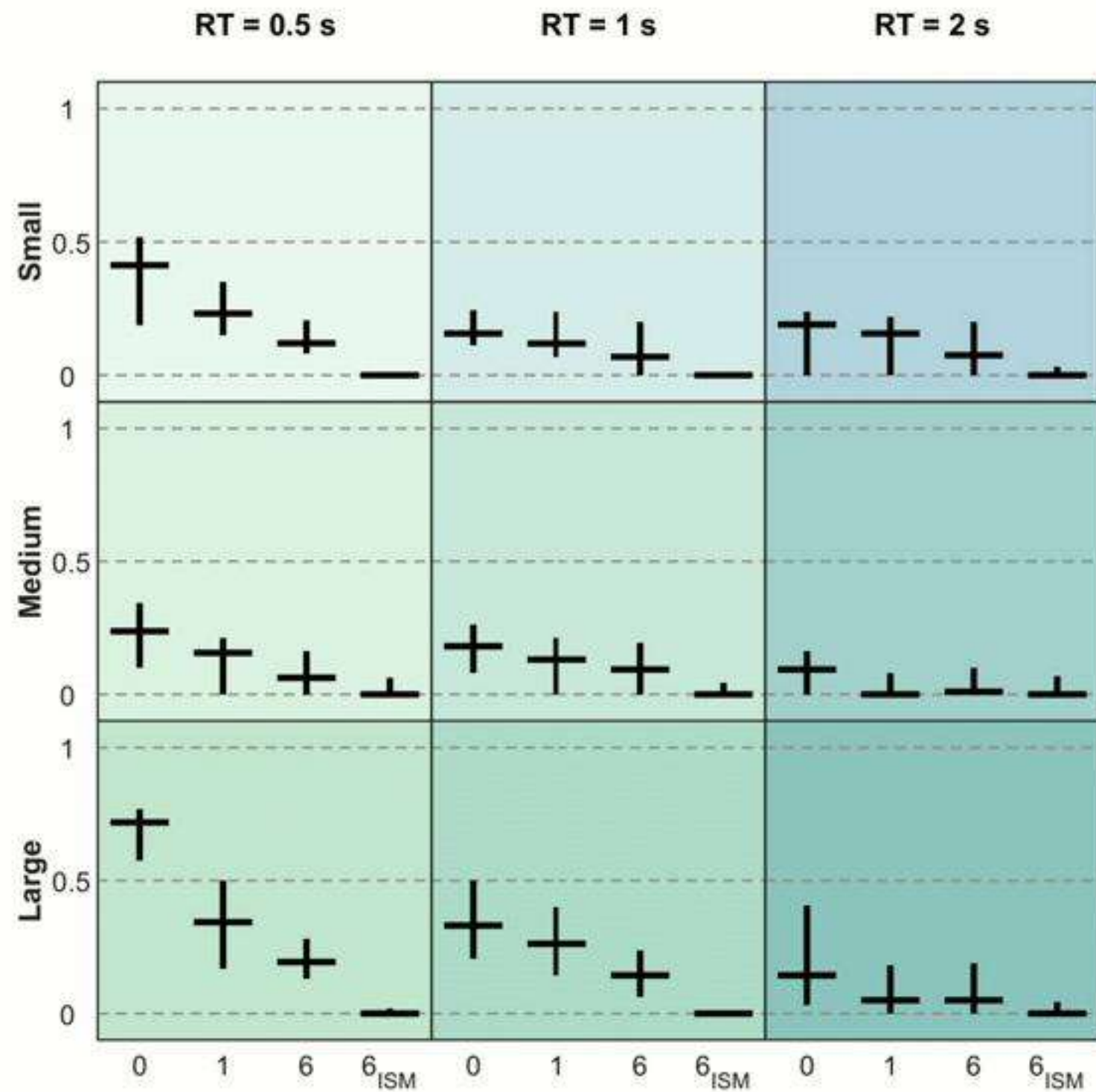
- Effect of **V** (?)
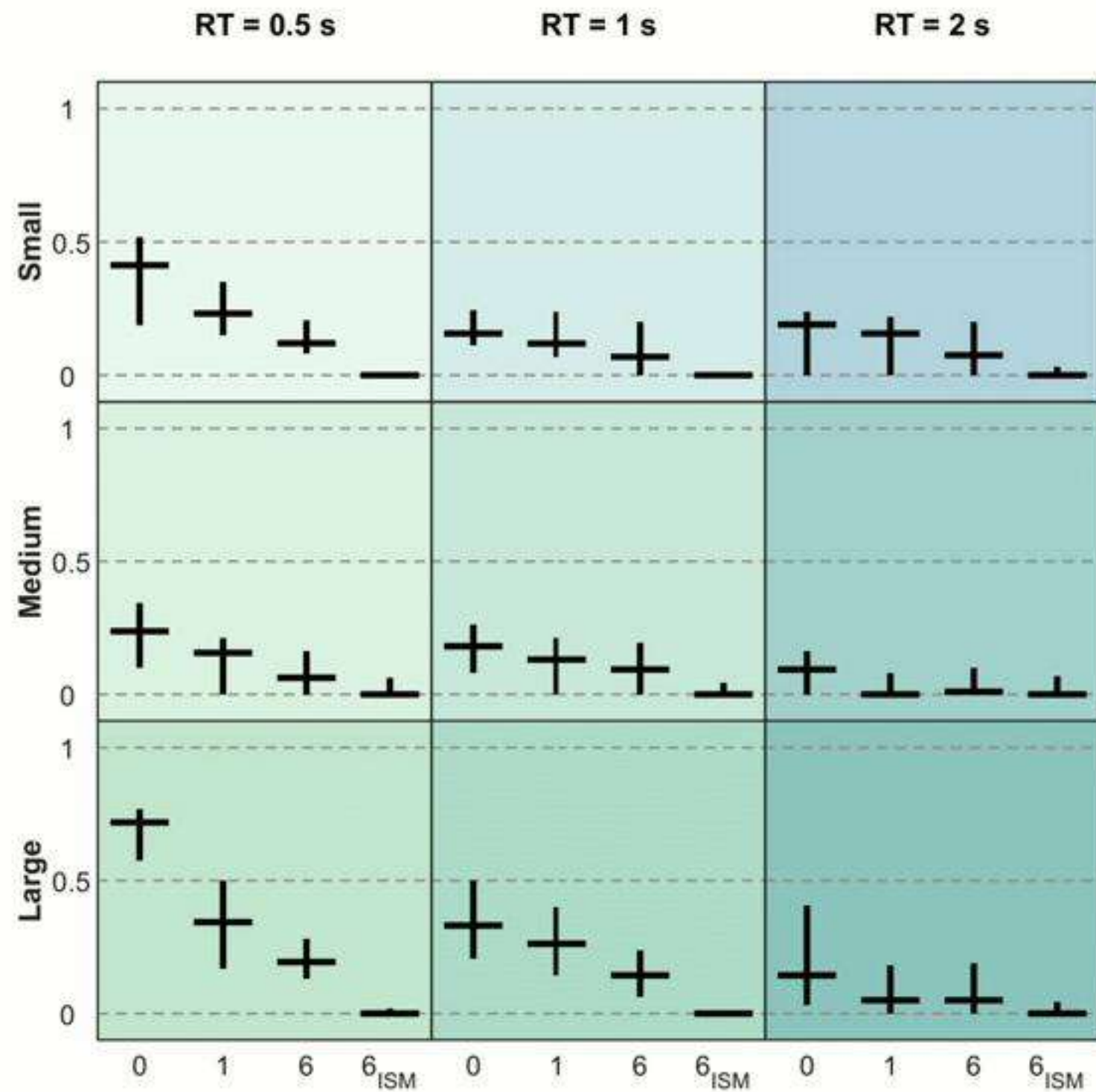
# PERCEPTUAL EVALUATION

# PERCEPTUAL EVALUATION
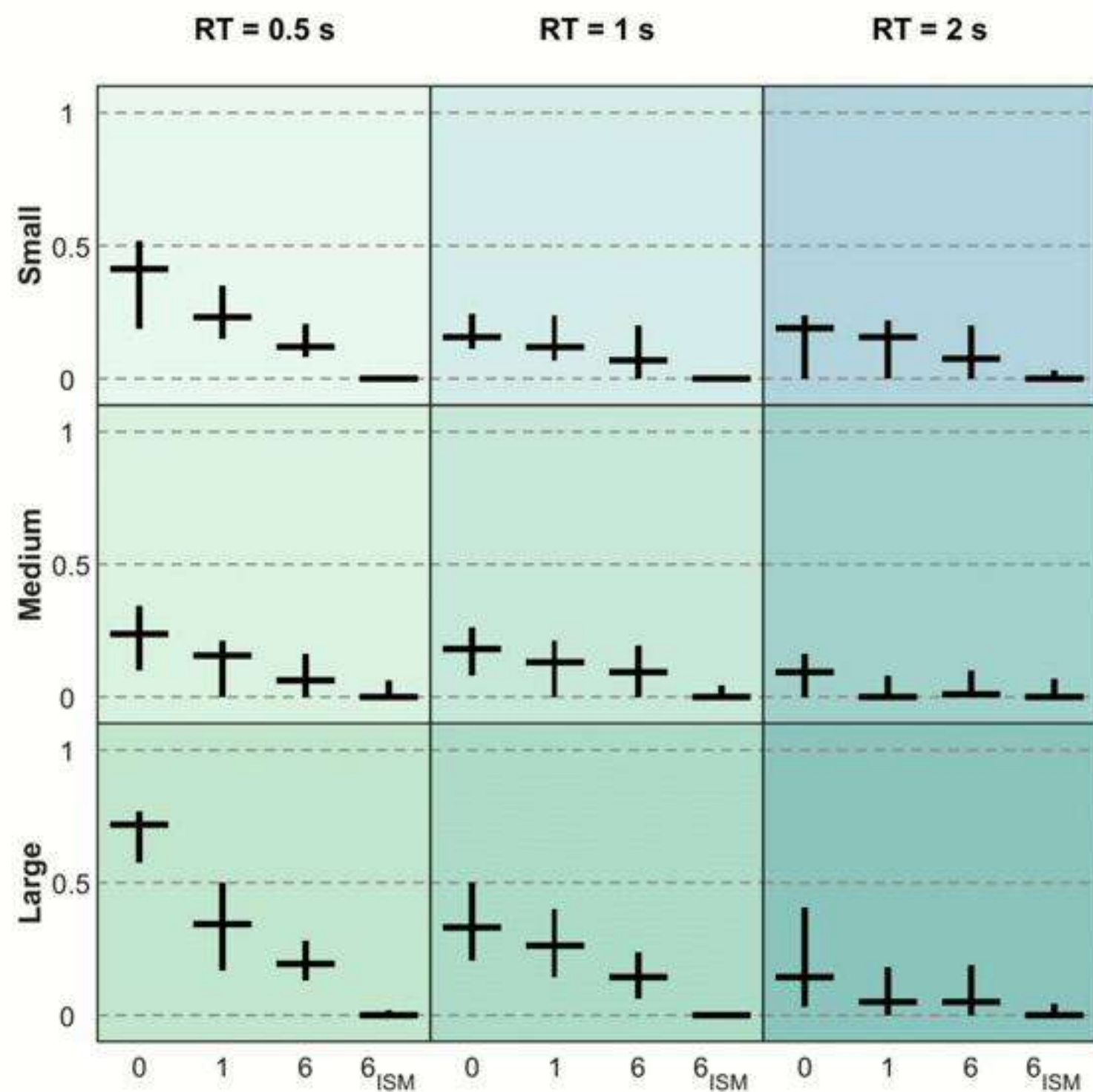


- Differences decrease with **N**

# PERCEPTUAL EVALUATION



- Differences decrease with **N**
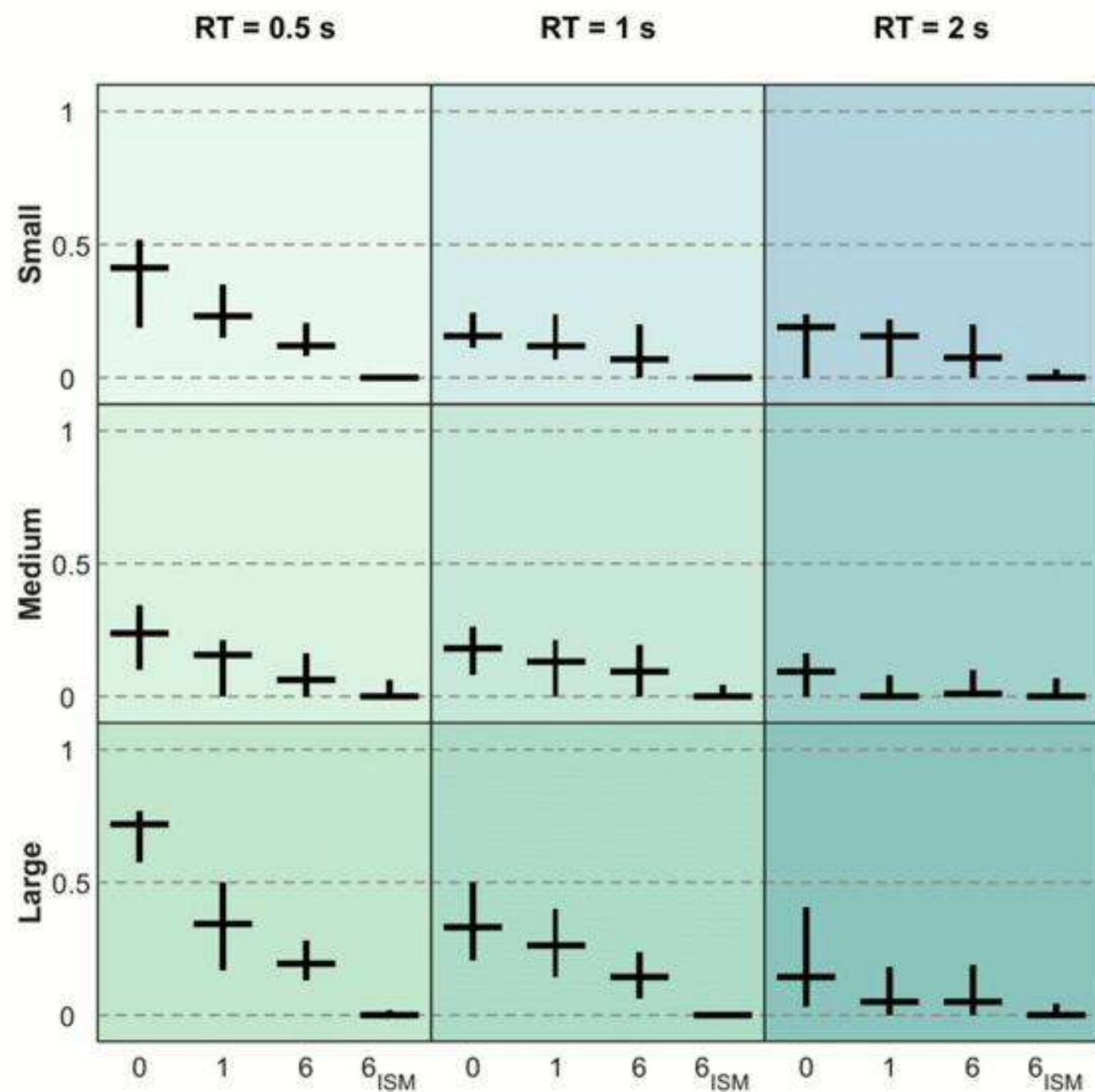
- $6_{ISM}$ is always rated 0

# PERCEPTUAL EVALUATION



- Differences decrease with **N**
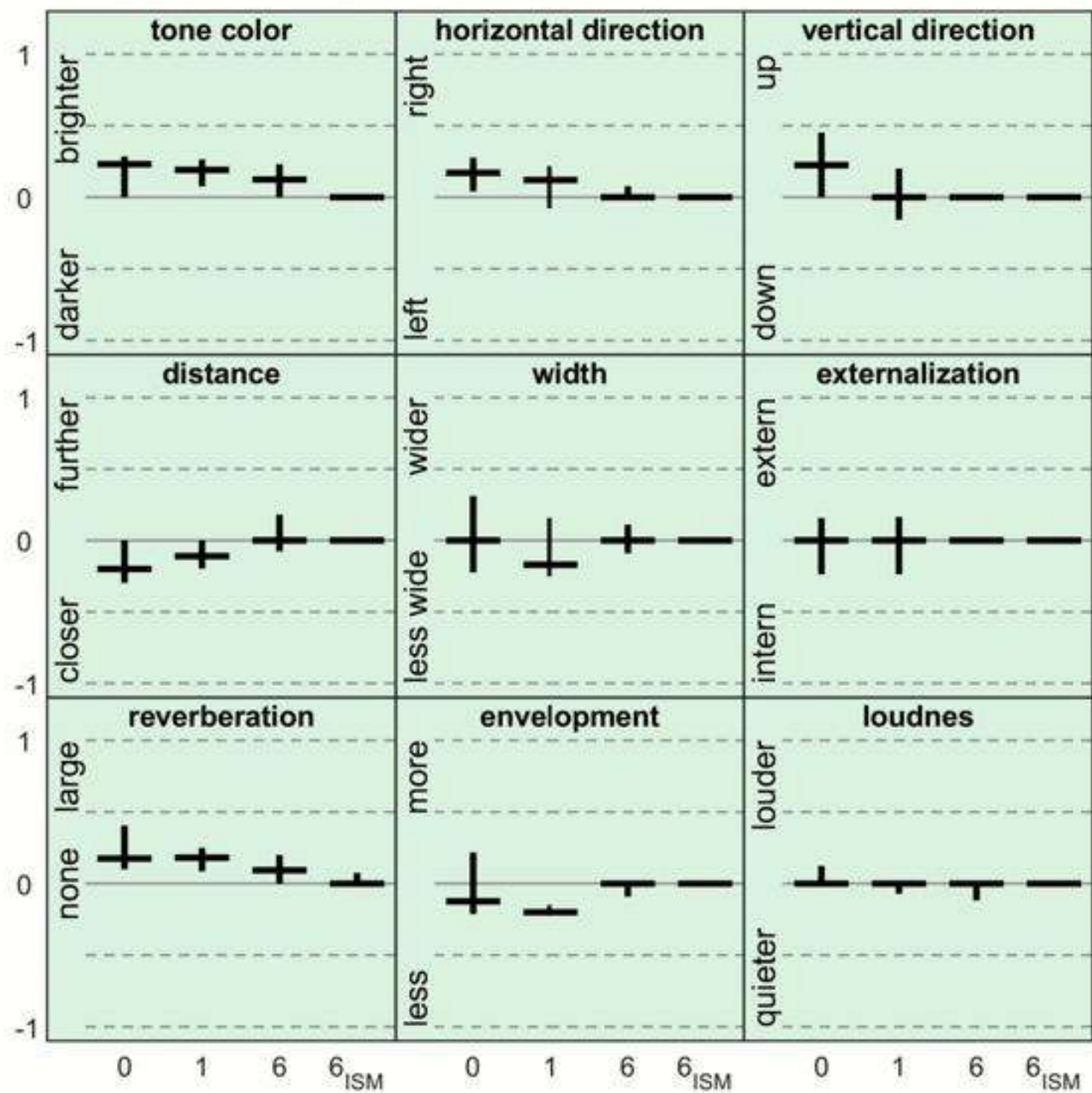
- $6_{ISM}$ is always rated 0

- 6 Cis overlap with 0

# PERCEPTUAL EVALUATION



- Differences decrease with **N**

- $6_{ISM}$ is always rated 0

- 6 Cis overlap with 0

- Small differences in every tested Quality

# PERCEPTUAL EVALUATION



- Differences decrease with **N**

- $6_{ISM}$ is always rated 0

- 6 Cis overlap with 0

- Small differences in every tested Quality

- Loudness not problematic

INTRODUCTION

METHOD

RESULTS

DISCUSSION

# SUMMARY

- End-to-end parametric audio system including early reflections

# SUMMARY

- End-to-end parametric audio system including early reflections

- Evaluated on two different models

# SUMMARY

- End-to-end parametric audio system including early reflections

- Evaluated on two different models

- Detected reflections agree across models

# SUMMARY

- End-to-end parametric audio system including early reflections

- Evaluated on two different models

- Detected reflections agree across models

- Early reflections important in large and dry rooms

# SUMMARY

- End-to-end parametric audio system including early reflections

- Evaluated on two different models

- Detected reflections agree across models

- Early reflections important in large and dry rooms

- Differences decrease with increasing reflections

# SUMMARY

- End-to-end parametric audio system including early reflections

- Evaluated on two different models

- Detected reflections agree across models

- Early reflections important in large and dry rooms

- Differences decrease with increasing reflections

- 6 early reflections sufficient in most cases

# SUMMARY

- End-to-end parametric audio system including early reflections

- Evaluated on two different models

- Detected reflections agree across models

- Early reflections important in large and dry rooms

- Differences decrease with increasing reflections

- 6 early reflections sufficient in most cases

- Floor reflections seems important

# CONTRIBUTIONS

- Detecting and selecting early reflections

# CONTRIBUTIONS

- Detecting and selecting early reflections

- Double sloped parametric late reverberation

# CONTRIBUTIONS

- Detecting and selecting early reflections

- Double sloped parametric late reverberation

- It doesn't take to much (early reflections) to trick the brain

# CONTRIBUTIONS

- Detecting and selecting early reflections

- Double sloped parametric late reverberation

- It doesn't take to much (early reflections) to trick the brain

- Inclusion into Triton work-flow possible

# CONTRIBUTIONS

- Detecting and selecting early reflections

- Double sloped parametric late reverberation

- It doesn't take to much (early reflections) to trick the brain

- Inclusion into Triton work-flow possible

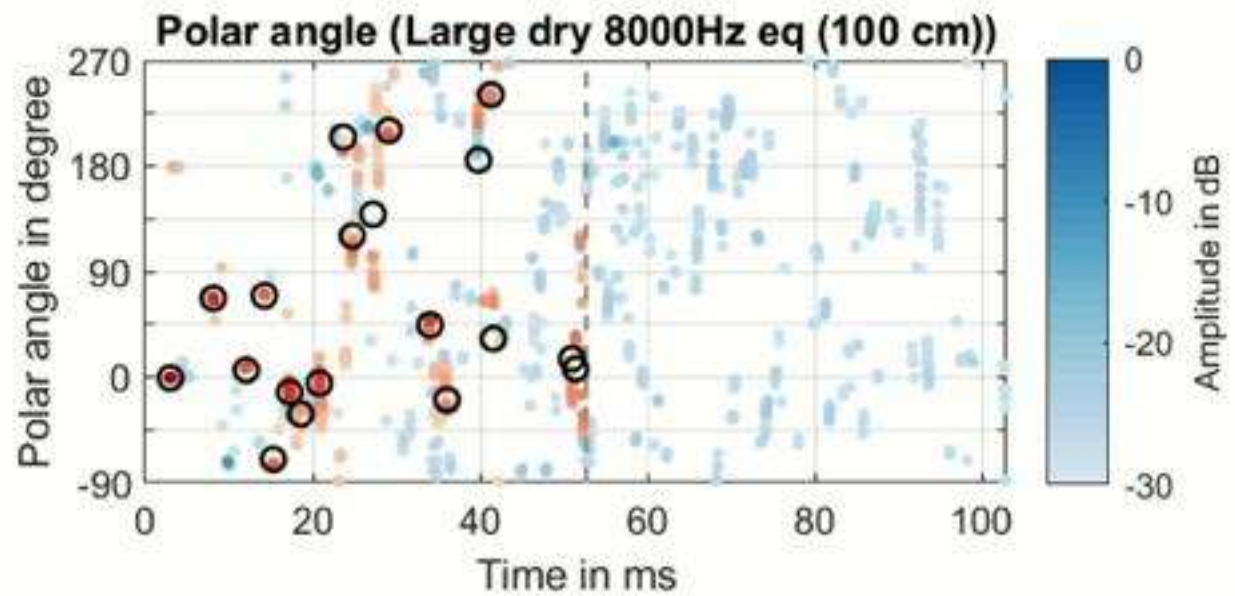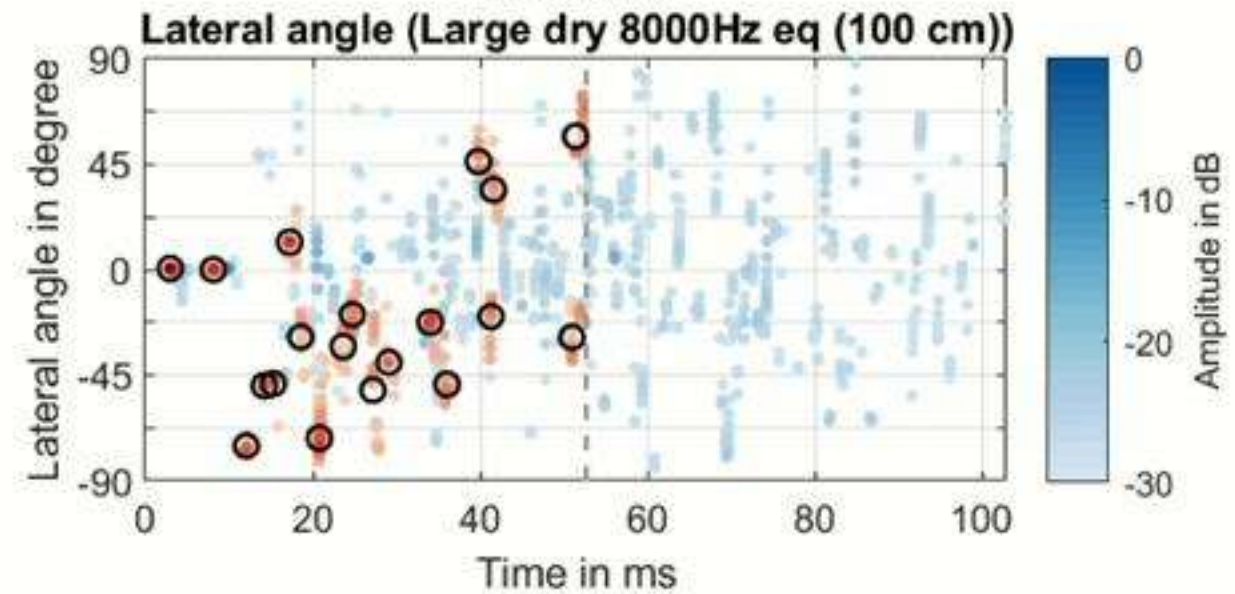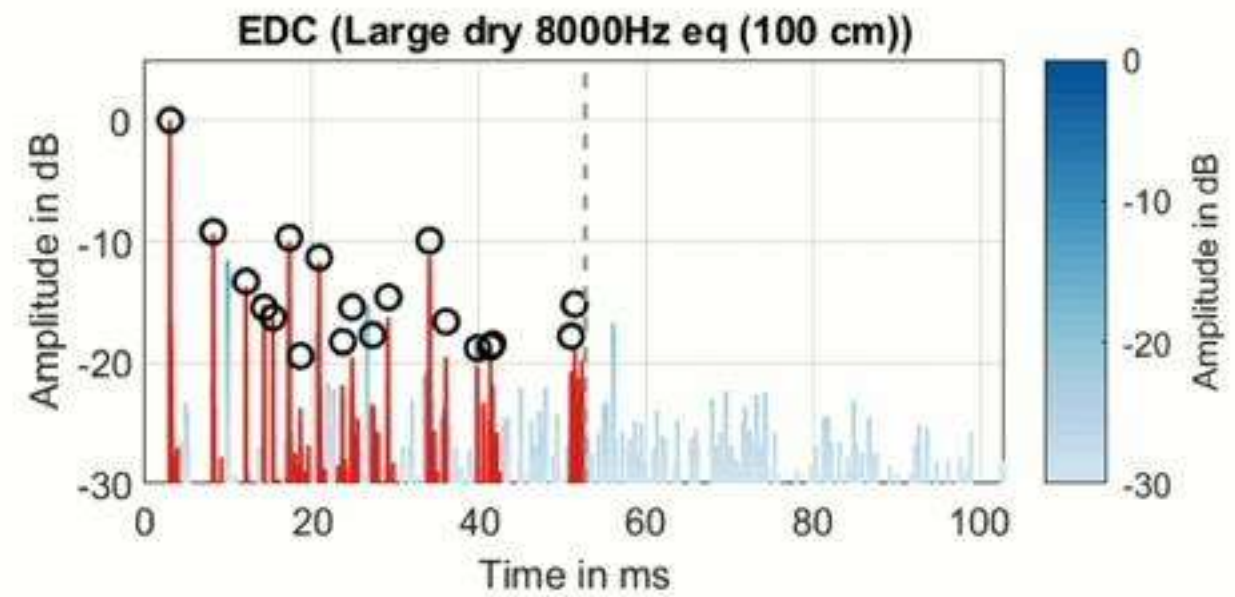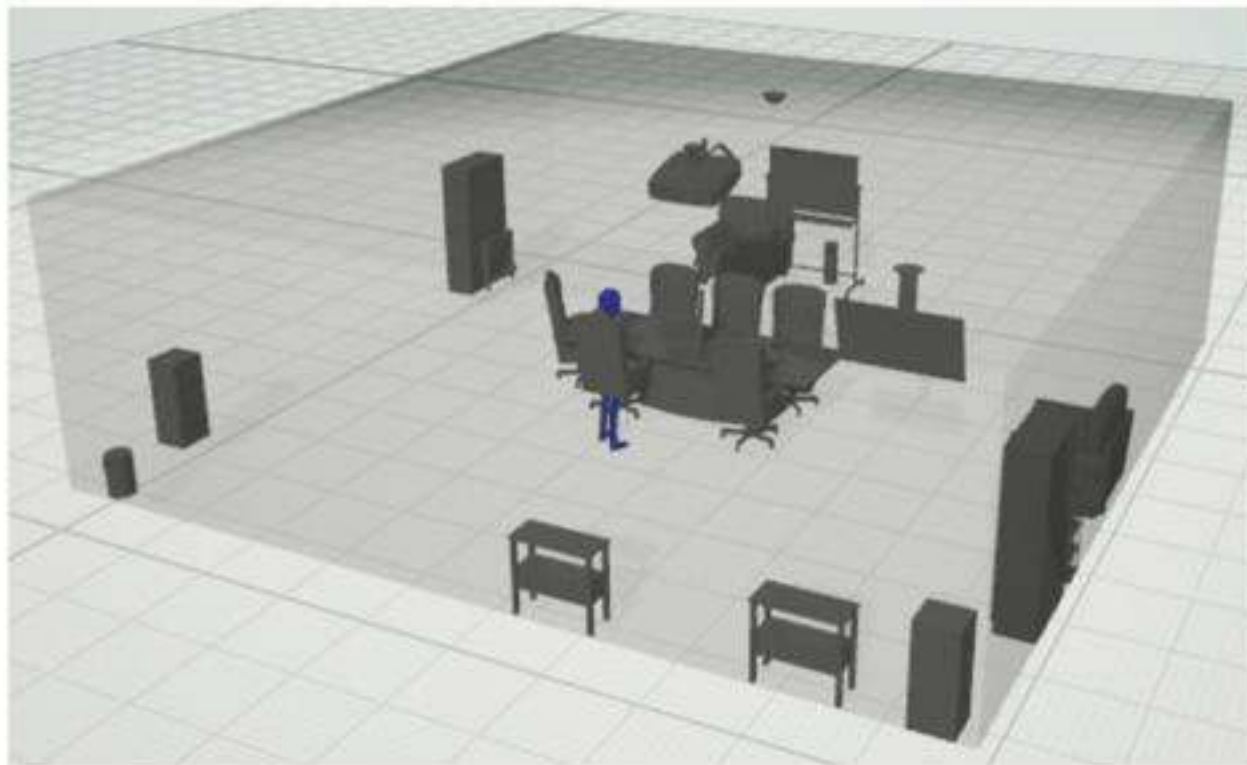- Currently FIR for late reverberation, new approach allows IIR

# CONTRIBUTIONS

- Detecting and selecting early reflections

- Double sloped parametric late reverberation

- It doesn't take to much (early reflections) to trick the brain

- Inclusion into Triton work-flow possible

- Currently FIR for late reverberation, new approach allows IIR
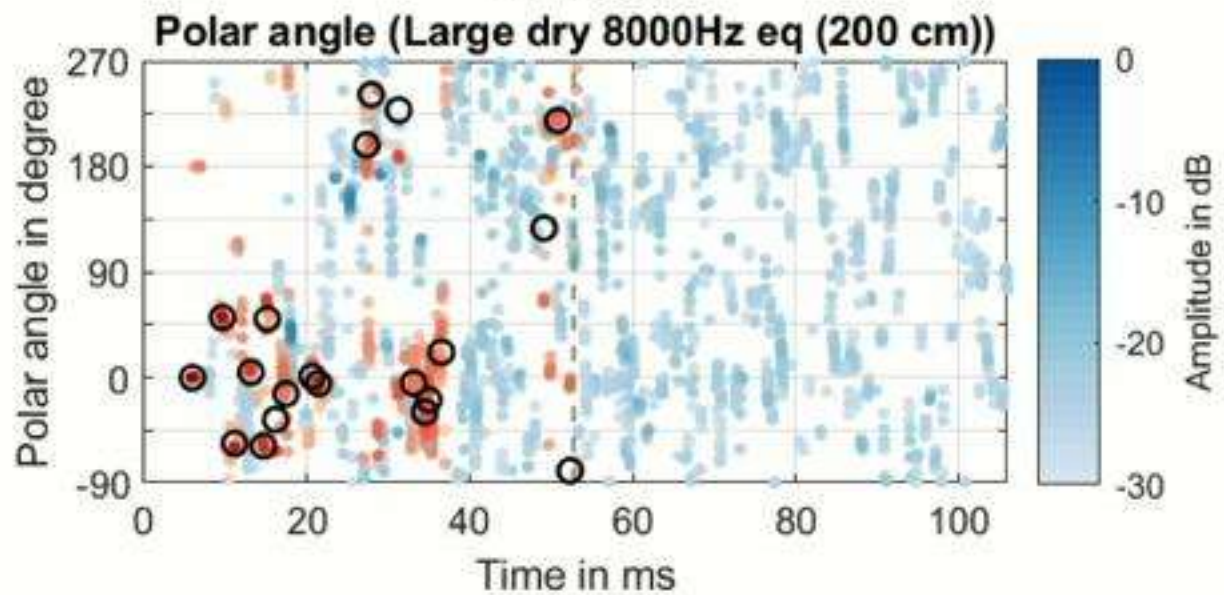
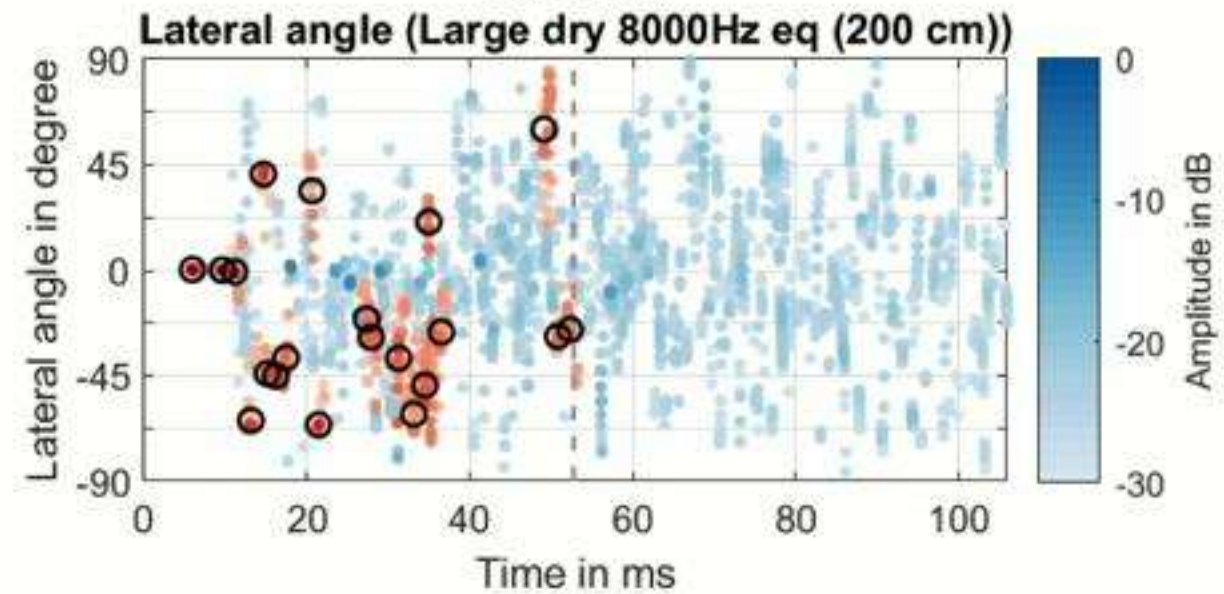➤ More efficient, and higher perceptual quality

# OUTLOOK

- Non-empty rooms

# OUTLOOK

- Non-empty rooms

- Spatial smoothness

# OUTLOOK

- Non-empty rooms

- Spatial smoothness

- Outdoor Environments

# OUTLOOK

- Non-empty rooms

- Spatial smoothness

- Outdoor Environments

- Parametric rendering
  - Discover first order reflections

# OUTLOOK

- Non-empty rooms

- Spatial smoothness

- Outdoor Environments

- Parametric rendering
    - Discover first order reflections
    - Detect reflections with short memory

# OUTLOOK

- Non-empty rooms

- Spatial smoothness

- Outdoor Environments

- Parametric rendering
  - Discover first order reflections
  - Detect reflections with short memory
  - Room dependent search time

# OUTLOOK

- Non-empty rooms

- Spatial smoothness

- Outdoor Environments

- Parametric rendering
  - Discover first order reflections
  - Detect reflections with short memory
  - Room dependent search time
  - Directional variance

# OUTLOOK

- Non-empty rooms

- Spatial smoothness

- Outdoor Environments

- Parametric rendering
  - Discover first order reflections
  - Detect reflections with short memory
  - Room dependent search time
  - Directional variance
  - Frequency dependent rendering

# OUTLOOK

- Non-empty rooms

- Spatial smoothness

- Outdoor Environments

- Parametric rendering
  - Discover first order reflections
  - Detect reflections with short memory
  - Room dependent search time
  - Directional variance
  - Frequency dependent rendering
  - Directional reverberation

# THANKS!

Ivan Tashev for having me

Hannes Gamper for mentoring, critical feedback and help

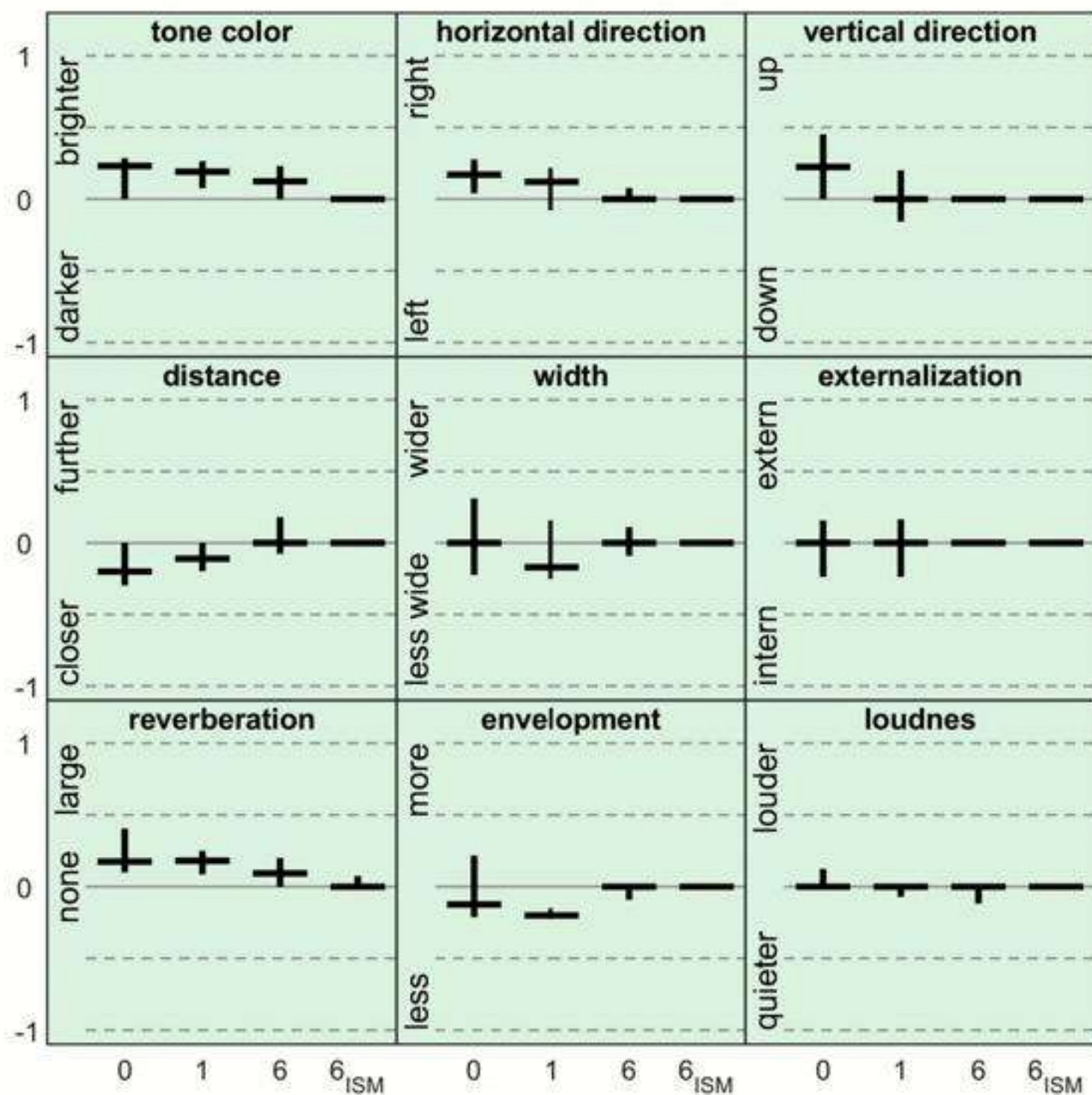Nikunj Raghuvanshi for critical feedback and help

Acoustics Team for help with Triton

Audio and Acoustics Group for the warm welcome

All subjects that participated in the listening test

# PERCEPTUAL EVALUATION



- Differences decrease with **N**

- $6_{ISM}$ is always rated 0

- 6 Cis overlap with 0

- Small differences in every tested Quality

- Loudness not problematic

# THANKS!

Ivan Tashev for having me

Hannes Gamper for mentoring, critical feedback and help

Nikunj Raghuvanshi for critical feedback and help

Acoustics Team for help with Triton

Audio and Acoustics Group for the warm welcome

All subjects that participated in the listening test