









**WHAT IS YOUR FAVOURITE MULTISENSORY EXPERIENCE IN 2026?**





**REALITY  
IS BETTER  
THAN  
VIRTUALITY**









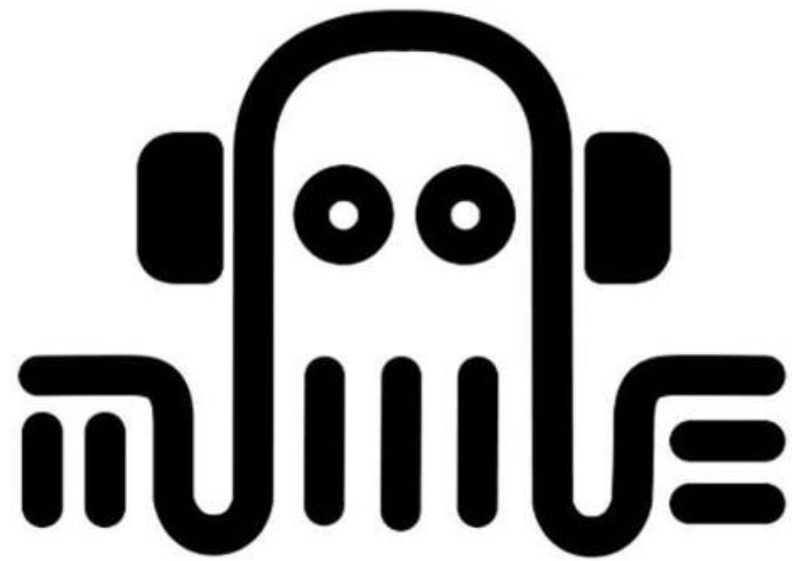




DTU

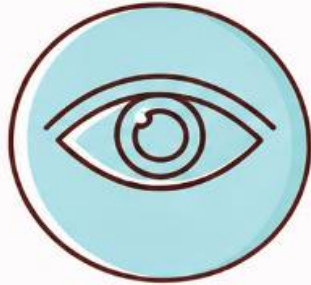


Danmarks Tekniske Universitet  
Technical University of Denmark



M E L A B

**HOW MANY SENSES DO WE HAVE?**



VISUAL



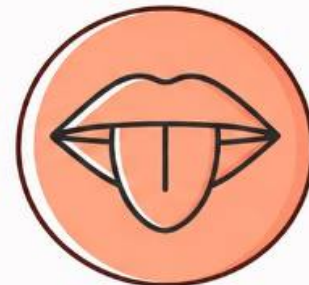
AUDITORY



TACTILE



OLFACTORY

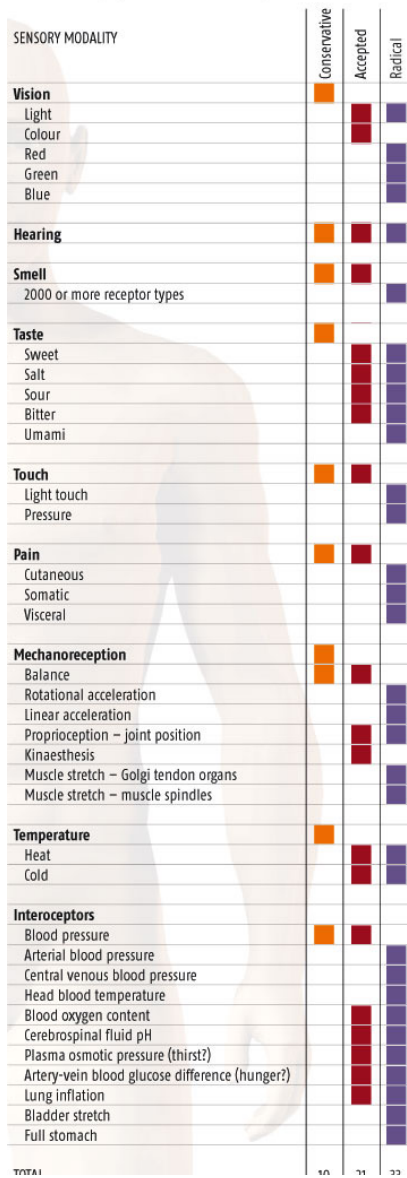


GUSTATORY

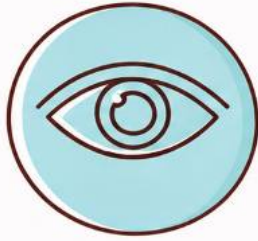


**OPINIONS ARE DIVIDED**

There are many opinions about how many senses we have



[New Scientist article](#), 29 Jan 2005 by Bruce Durie



VISUAL

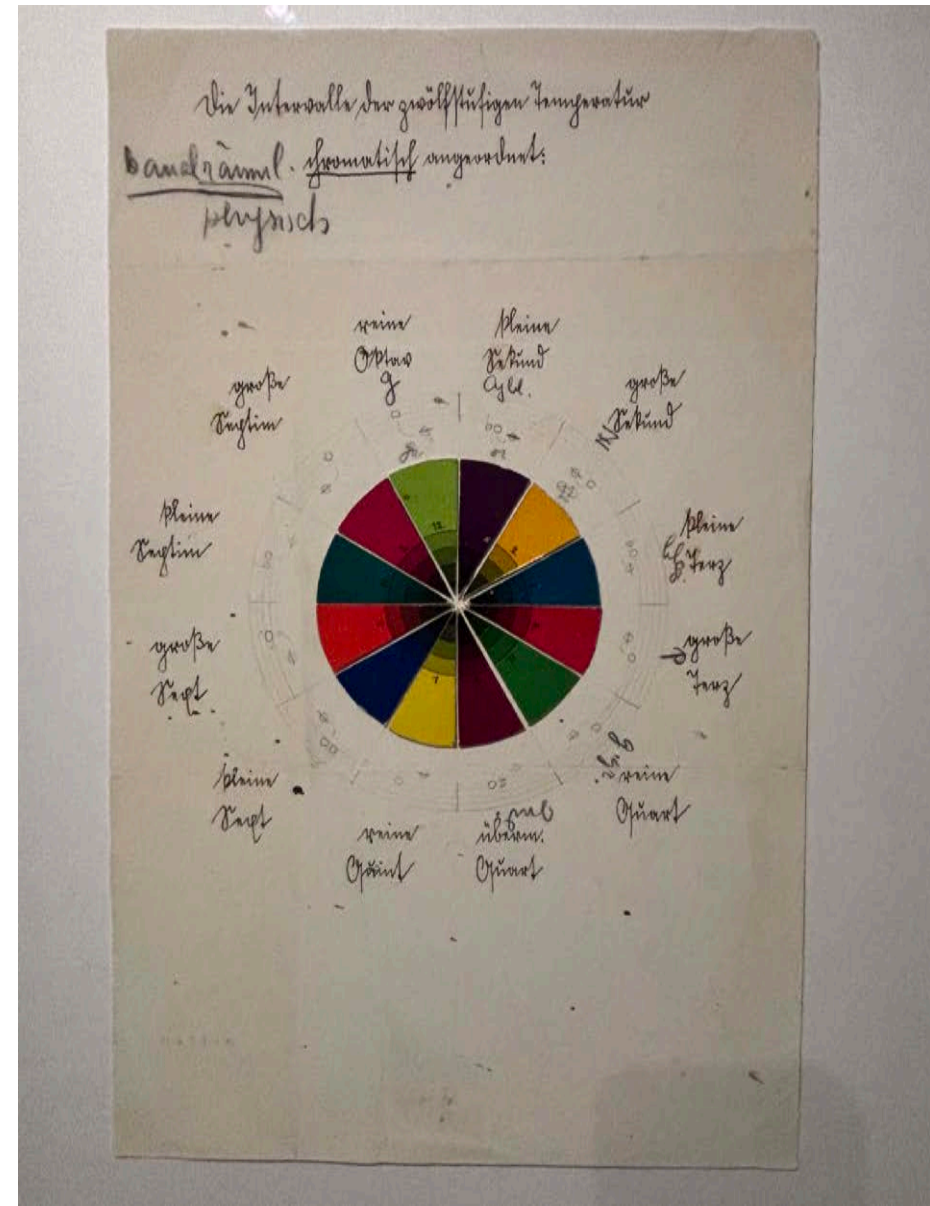


AUDITORY



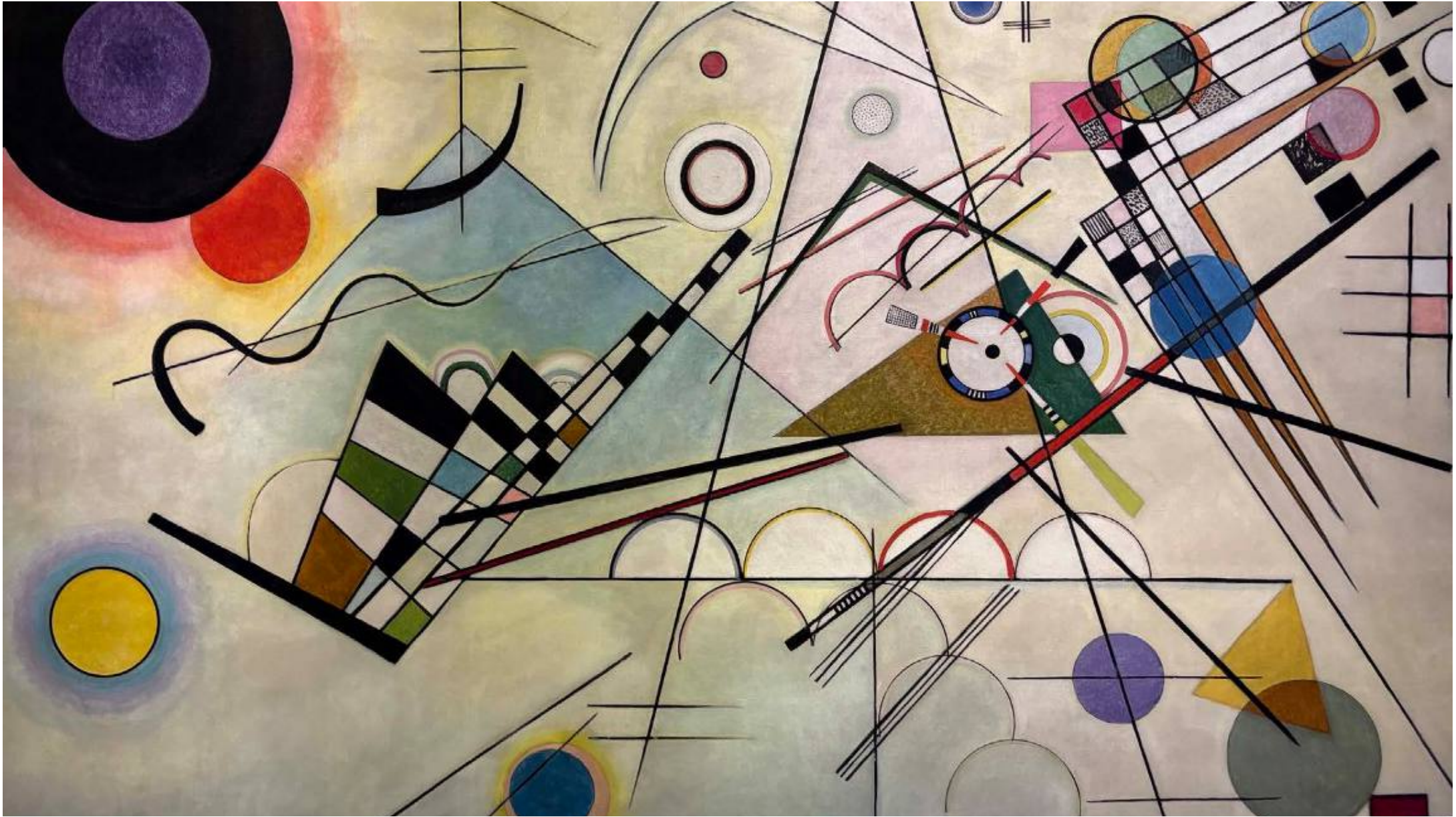
TACTILE

# MULTISENSORY ART









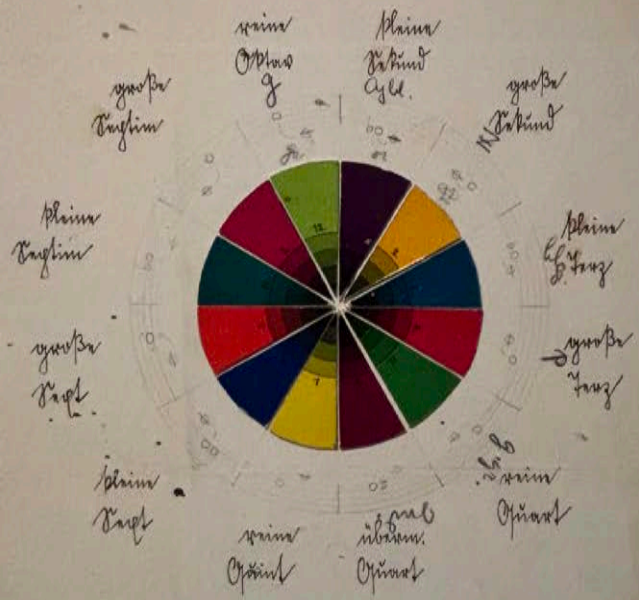
The image shows a handwritten musical score and its analysis. At the top, there are two staves of musical notation. The first staff contains notes and rests, while the second staff contains notes and rests with some markings. Below the musical notation is a graphic notation section consisting of several horizontal lines. The first line has a series of dots and triangles. The second line has a series of dots and triangles. The third line has a series of dots and triangles. The fourth line has a series of dots and triangles. The fifth line has a series of dots and triangles. The sixth line has a series of dots and triangles. The seventh line has a series of dots and triangles. The eighth line has a series of dots and triangles. The ninth line has a series of dots and triangles. The tenth line has a series of dots and triangles. The eleventh line has a series of dots and triangles. The twelfth line has a series of dots and triangles. The thirteenth line has a series of dots and triangles. The fourteenth line has a series of dots and triangles. The fifteenth line has a series of dots and triangles. The sixteenth line has a series of dots and triangles. The seventeenth line has a series of dots and triangles. The eighteenth line has a series of dots and triangles. The nineteenth line has a series of dots and triangles. The twentieth line has a series of dots and triangles.

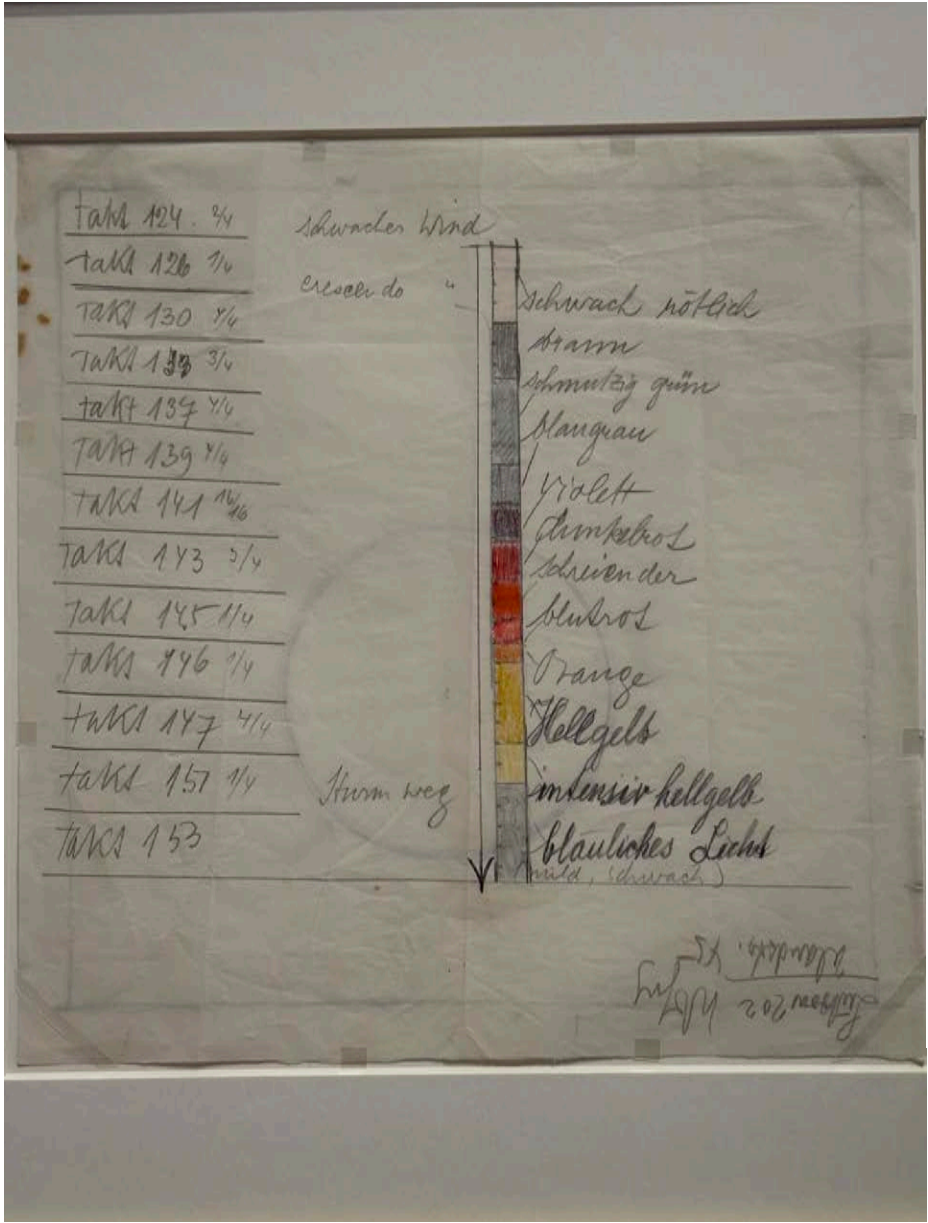
**Farbe / Geräusch**  
 parallel  
 ▲ diametral  
 nur punktlisch parallel  
 ● nur punktlisch diametral  
 oben allein  
 ■ unten allein

analyse eines tonstückes.  
 allegro ritmico / josef haas / op. 69 h.I. takt 7-8.  
 grafisches bild.  
 bearbeitung für einen farbigen geräuschfilm.  
 sichtbarmachung der komponenten.  
 filterung des tons / ergebnis: basis für den vergleich  
 mit optischen kunstwerken.

22.10.38

Die Intervalle der zwölfstimmigen Temperatur  
bandräuml. chromatisch angeordnet:  
 physisch





# COLOR / CHARACTERISTIC / TONE

<b>YELLOW</b>	"warm," "cheeky and exciting," "disturbing for people," "typical earthly color," "compared with the mood of a person it could have the effect of representing madness in color [...] an attack of rage, blind madness, maniacal rage."	loud, sharp trumpets, high fanfares
<b>BLUE</b>	deep, inner, supernatural, peaceful "Sinking towards black, it has the overtone of a mourning that is not human." "typical heavenly color"	light blue: flute darker blue: cello darkest blue of all: organ
<b>GREEN</b>	mixture of yellow and blue, stillness, peace, but with hidden strength, passive "Green is like a fat, very healthy cow lying still and unmoving, only capable of chewing the cud, regarding the world with stupid dull eyes."	quiet, drawn-out, middle position violin
<b>WHITE</b>	"It is not a dead silence, but one pregnant with possibilities." "Harmony of silence",	"Harmony of silence", "pause that breaks temporarily the melody"
<b>BLACK</b>	"Not without possibilities [...] like an eternal silence, without future and hope." Extinguished, immovable.	"final pause, after which any continuation of the melody seems the dawn of another world"
<b>GREY</b>	mixture of white and black	soundlessness.
<b>RED</b>	alive, restless, confidently striving towards a goal, glowing, "manly maturity" Light warm red: strength, energy, joy; Vermilion: glowing passion, sure strength Light cold red: youthful, pure joy, young	"sound of a trumpet, strong, harsh" Fanfare, Tuba deep notes on the cello high, clear violin
<b>BROWN</b>	mixture of red + black dull, hard, inhibited	
<b>ORANGE</b>	mixture of red + yellow radiant, healthy, serious	middle range church bell, alto voice, "an alto violin, singing tone, largo"
<b>VIOLET</b>	mixture of red + blue "morbid, extinguished [...] sad"	english horn, shawm, bassoon

From Wassily Kandinsky, *Concerning the Spiritual in Art*, 1912

17

A musical staff with a treble clef, containing several notes and rests. The notes are mostly quarter and eighth notes.

18

A complex musical diagram featuring overlapping lines, notes, and rests, possibly representing a specific musical structure or analysis. It includes various musical symbols and annotations.

19

A complex musical diagram featuring overlapping lines, notes, and rests, similar to diagram 18. It includes various musical symbols and annotations.

20

20

A complex musical diagram featuring overlapping lines and notes, similar to the previous diagrams. It includes various musical symbols and annotations.

21

A complex musical diagram featuring overlapping lines and notes, similar to the previous diagrams. It includes various musical symbols and annotations.

22

A musical staff with a treble clef, containing several notes and rests. The notes are mostly quarter and eighth notes.

23

A musical staff with a treble clef, containing several notes and rests. The notes are mostly quarter and eighth notes.

24

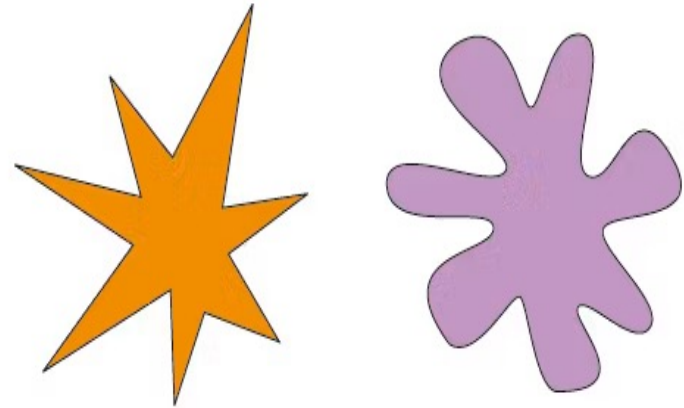
A musical staff with a treble clef, containing several notes and rests. The notes are mostly quarter and eighth notes.

25

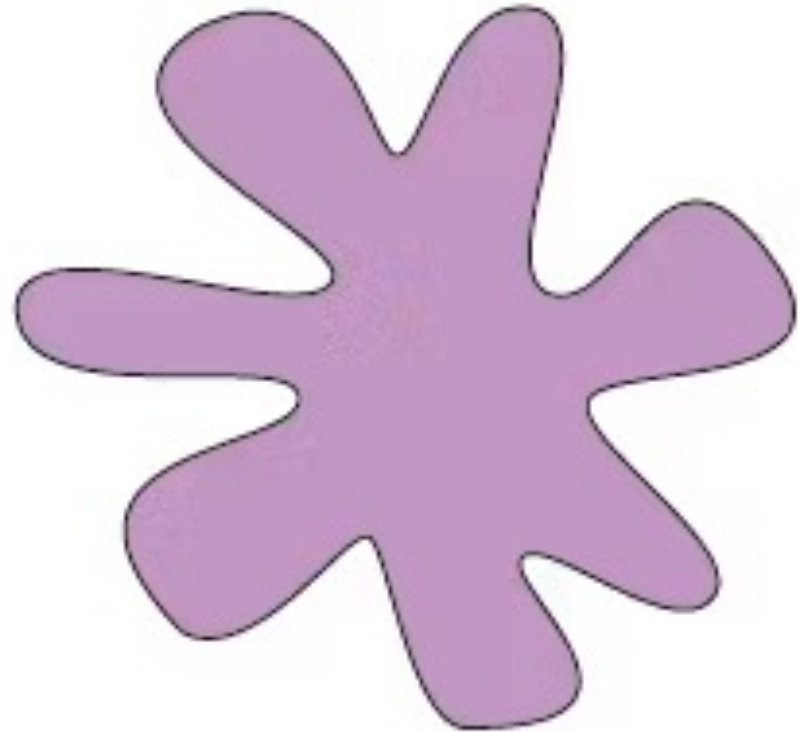
25

A complex musical diagram featuring overlapping lines and notes, similar to the previous diagrams. It includes various musical symbols and annotations.

# **CROSSMODAL CORRESPONDANCES**



# BOUBA KIKI OR TAKETE MALUMA?



1929, Wolfgang Köhler

# Matching sounds to shapes: Evidence of the Bouba-Kiki effect in naïve baby chicks

<sup>1</sup>\*Maria Loconsole, <sup>2</sup>Silvia Benavides-Varela, <sup>1</sup>Lucia Regolin

<sup>1</sup>Department of General Psychology, University of Padova, Padova, Italy

<sup>2</sup>Department of Developmental Psychology and Socialisation, University of Padova, Padova, Italy

\* Corresponding author: Maria Loconsole

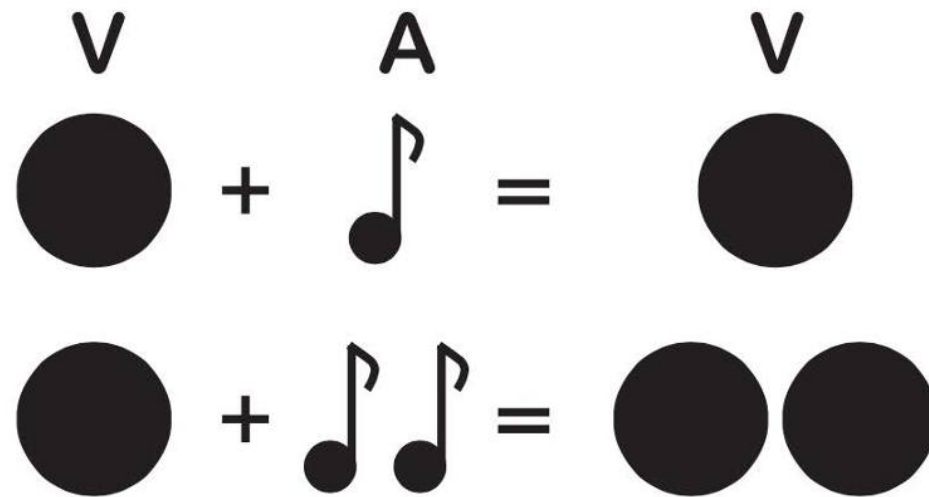
Email: [maria.loconsole@unipd.it](mailto:maria.loconsole@unipd.it)

**Keywords:** Sound-symbolism; Crossmodal associations; Bouba-Kiki effect; Domestic chicken; Predispositions

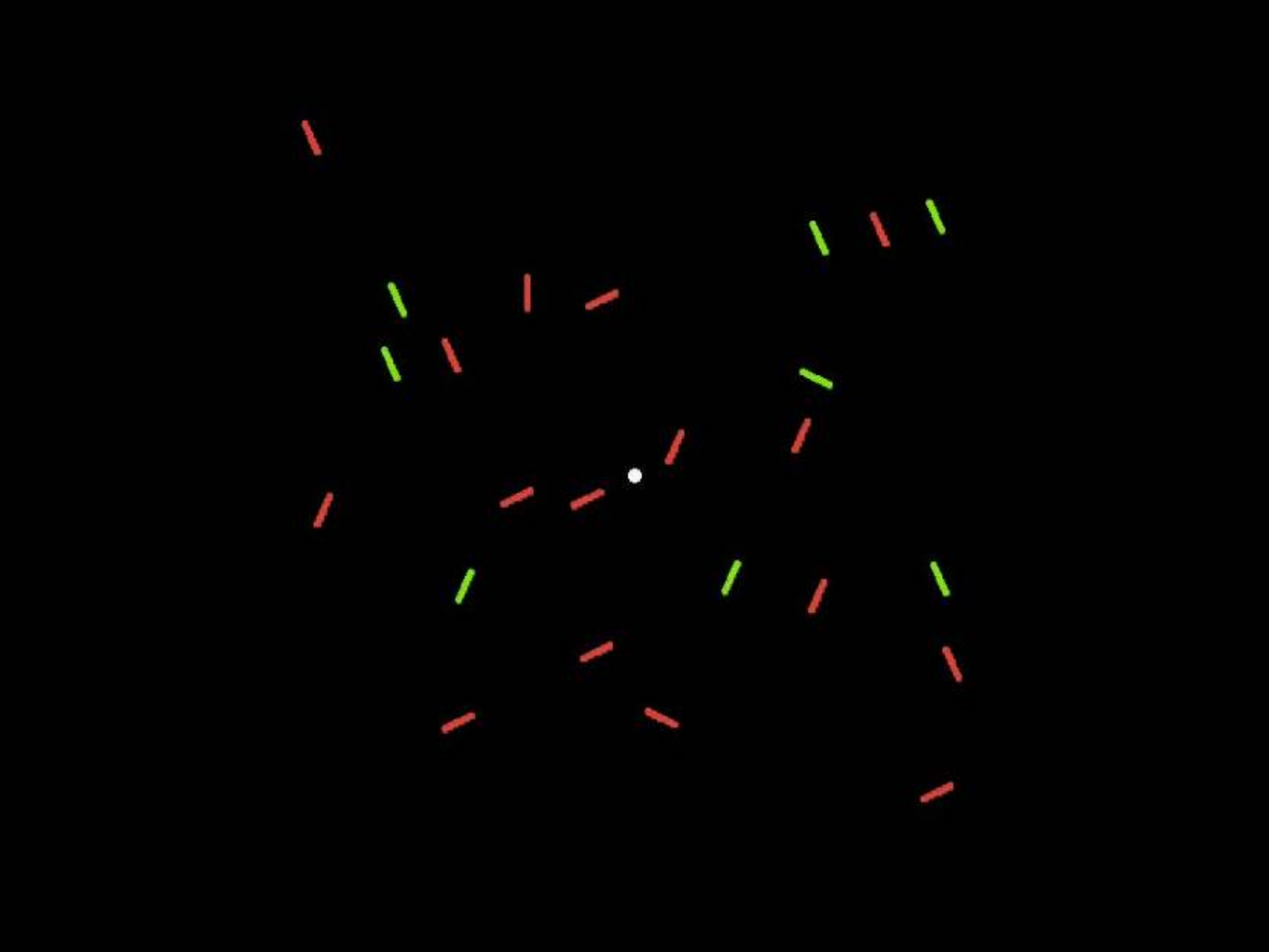
You are about to be FLASHED with black DISKS

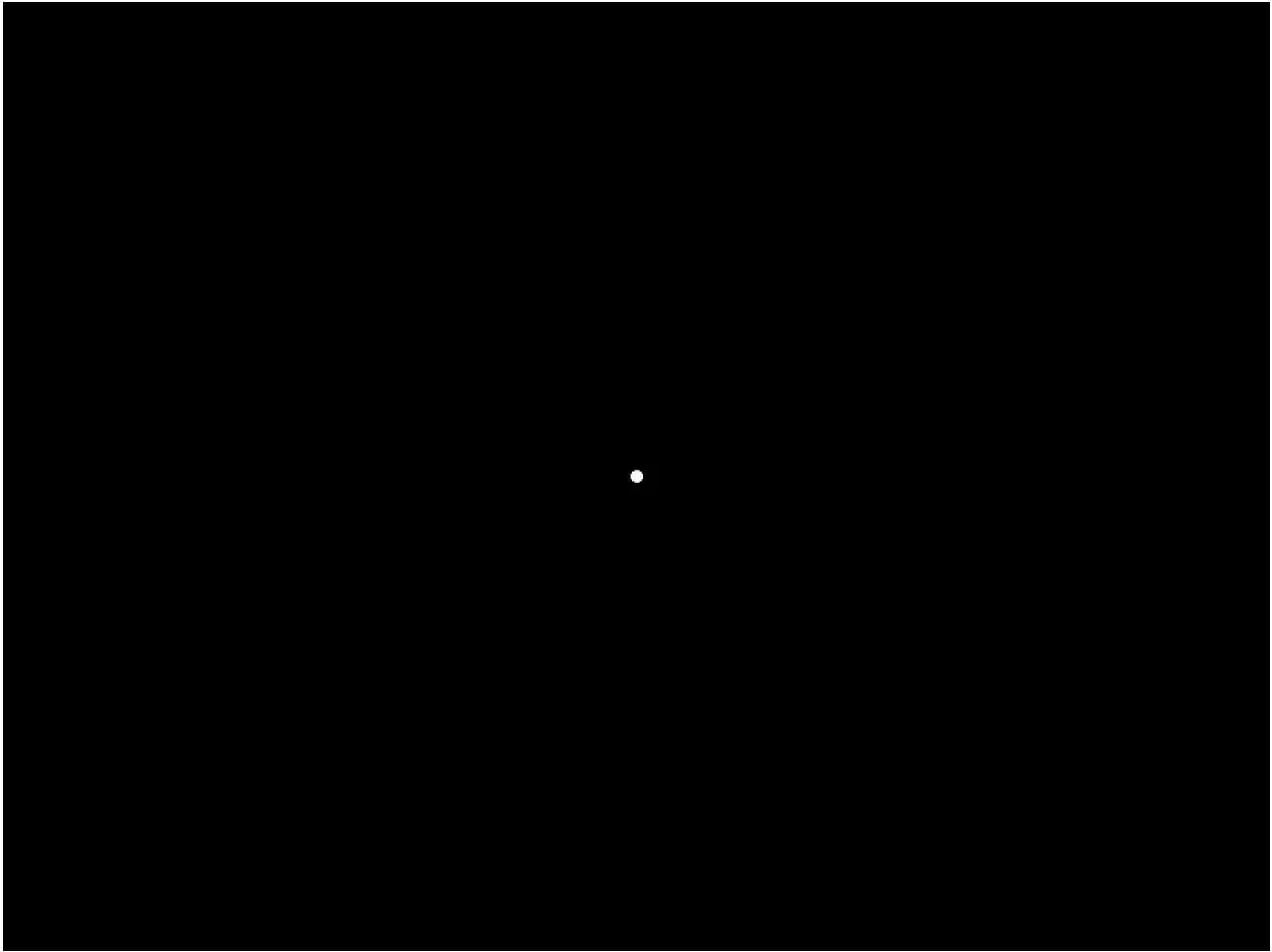
How many times does each disk flash?

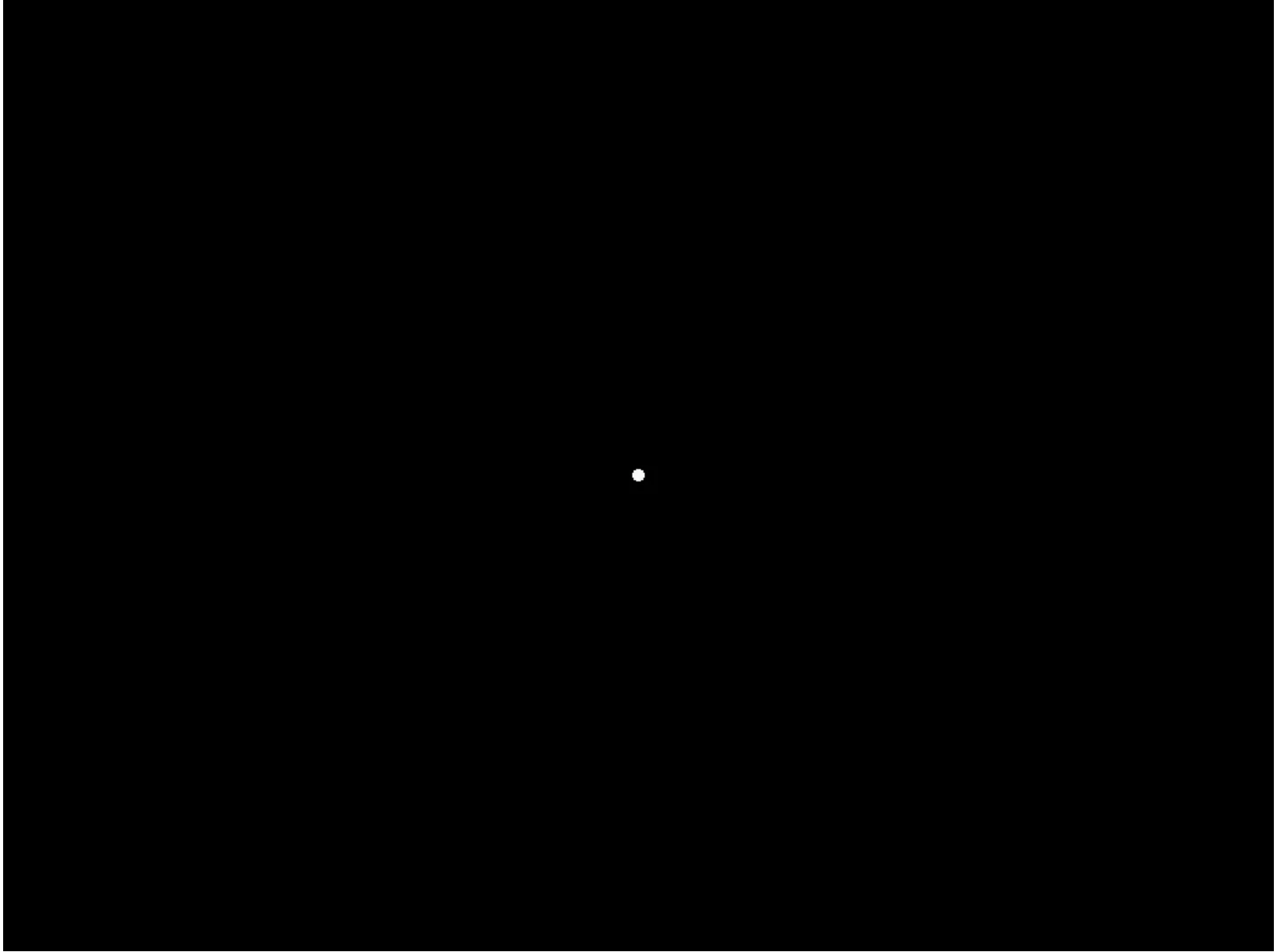
# ILLUSIONARY FLASHES

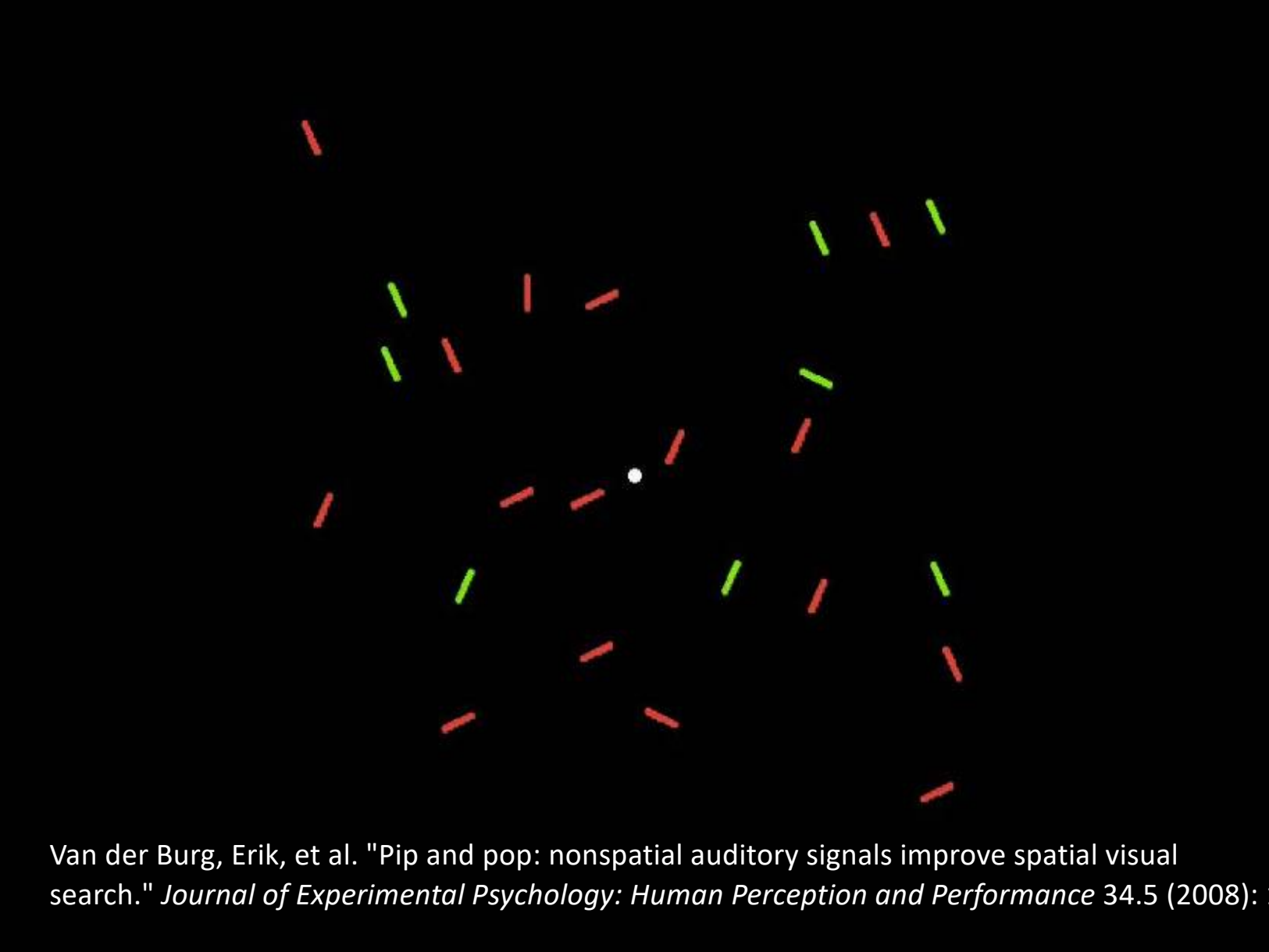


Shams, Ladan, Yukiyasu Kamitani, and Shinsuke Shimojo. "What you see is what you hear." *Nature* 408.6814 (2000): 788-788.









Van der Burg, Erik, et al. "Pip and pop: nonspatial auditory signals improve spatial visual search." *Journal of Experimental Psychology: Human Perception and Performance* 34.5 (2008): 1

# Missing The Point: An Exploration of How to Guide Users' Attention During Cinematic Virtual Reality

Lasse T. Nielsen, Matias B. Møller, Sune D. Hartmeyer, Troels C. M. Ljung\*, Niels C. Nilsson, Rolf Nordahl, and Stefania Serafin<sup>†</sup>  
Aalborg University Copenhagen

## Abstract

Recent technological advances have brought Virtual Reality (VR) into the homes of consumers, and there is a growing interest in bringing cinematic experiences from the screen and into VR. However, cinematic VR limits filmmakers' ability to effectively guide the audience's attention. In this paper we present a taxonomy of approaches to guiding users' attention, and present a study comparing two such approaches with a control condition devoid of guidance. One approach guides users by controlling their body's orientation, and the other implicitly directs their attention by encouraging them to follow a firefly with their gaze. The results revealed interesting, albeit statistically insignificant, indications that assuming control of the user's action may negatively influence presence, whereas the firefly was perceived as significantly more helpful.

**Keywords:** virtual reality; film, attention; presence; cinematic VR

**Concepts:** •Computing methodologies → Virtual reality;  
•Human-centered computing → Empirical studies in HCI;

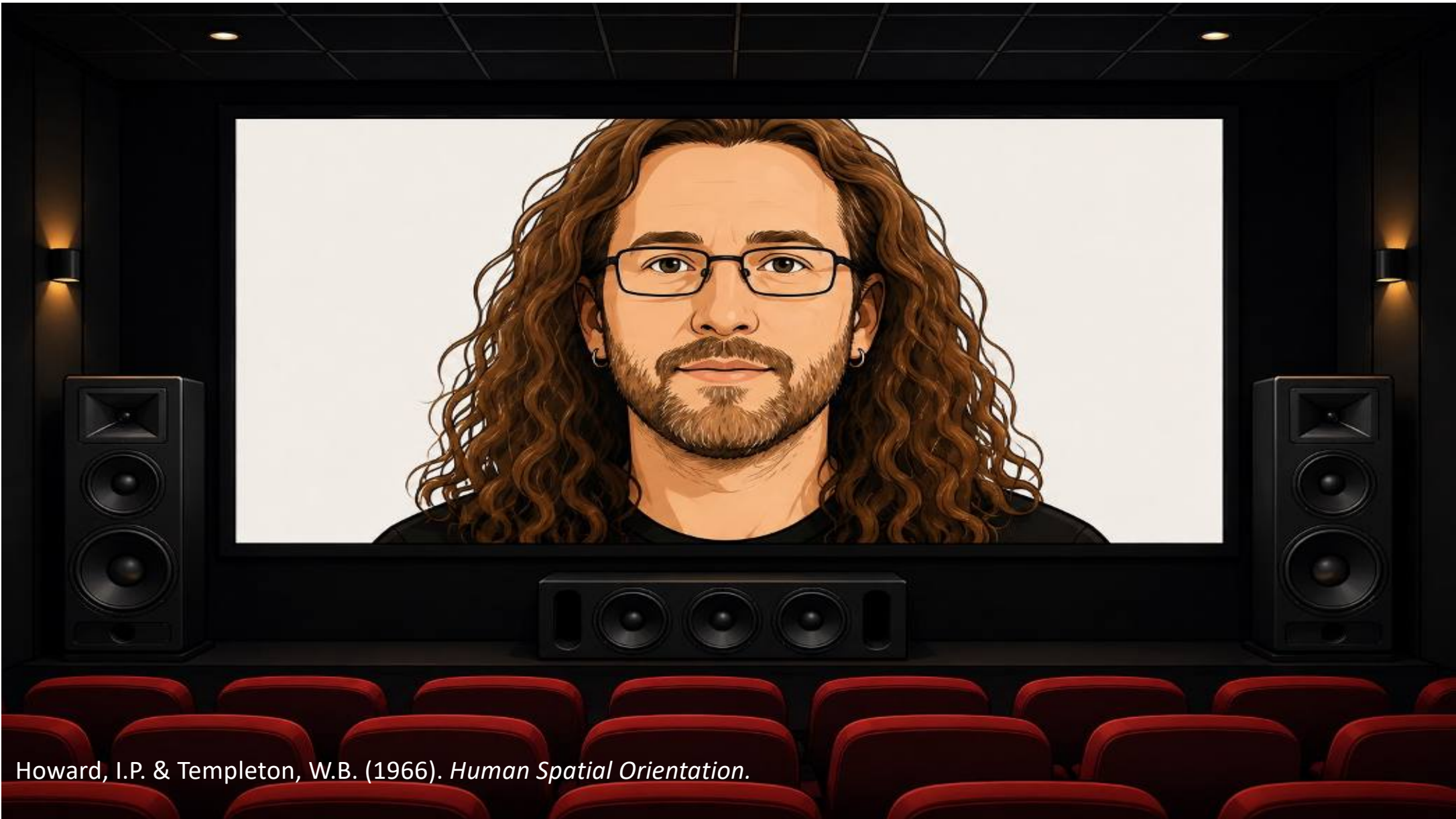
## 1 Introduction

Most contemporary filmmakers adhere to conventions that have

simply by turning their heads, or even change the vantage point by moving around the scene, and it remains to be seen if editing can be used effectively in relation to cinematic VR without disorienting the audience. This inability to control the audience's attention naturally poses a problem when one aspires to create a coherent narrative since coherence is contingent upon "careful selection and presentation of actions whose causal and temporal relationships highlight an underlying plot" [Young 2000]. In relation to cinematic VR, the filmmaker can no longer rely on cinematography to show the audience the building blocks of the plot since the camera is identified with the user and no longer is under authorial control [Aylett and Louchart 2003]. In relation to pre-authored narratives this leaves the filmmaker with at least three different, albeit not mutually exclusive, options: 1) progression of the story is halted until the user's head or gaze direction makes it reasonable to assume that important events and objects have been observed; 2) the system dynamically presents events and objects within the user's field of view; and 3) the filmmaker uses cues to steer the user's attention towards relevant events and objects (e.g., using *mise-en-scène* and sound). The work documented in the current paper is focussed on the third option. Section 2 outlines a novel taxonomy of different categories of cues which can be used to guide users' attention during cinematic VR, and section 3 presents a user study aimed at exploring how two different types of cues influence the sensation of presence in the virtual environment (VE) and recollection of the scene. Finally,



McGurk, Harry, and John MacDonald. "Hearing lips and seeing voices." *Nature* 264.5588 (1976): 746-748.



Howard, I.P. & Templeton, W.B. (1966). *Human Spatial Orientation*.

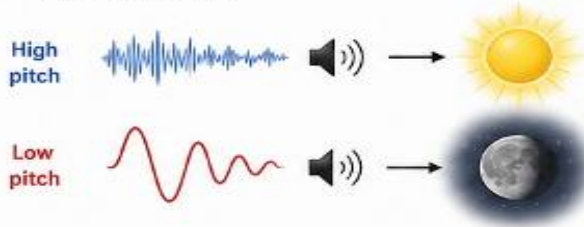
### 1. PITCH ↔ SIZE

High-pitched sounds are associated with small objects; low-pitched sounds with large objects.



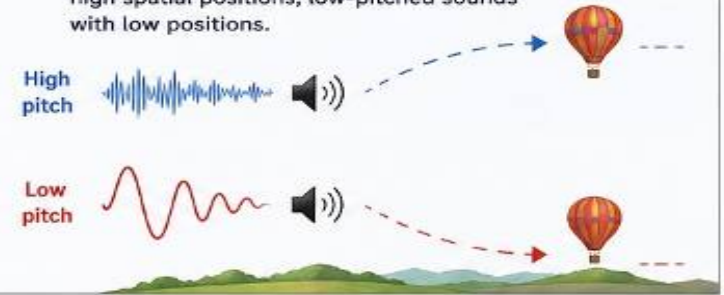
### 2. PITCH ↔ BRIGHTNESS (LIGHTNESS)

High-pitched sounds are matched with bright or light stimuli; low-pitched sounds with dark stimuli.



### 3. PITCH ↔ SPATIAL ELEVATION

High-pitched sounds are associated with high spatial positions; low-pitched sounds with low positions.



### 4. PITCH ↔ SHAPE

High-pitched sounds are associated with angular shapes; low-pitched sounds with rounded shapes.



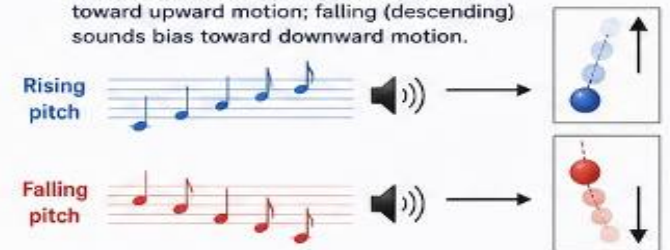
### 5. SOUND SYMBOLISM: BOUBA / KIKI EFFECT

People consistently match invented words with shapes.



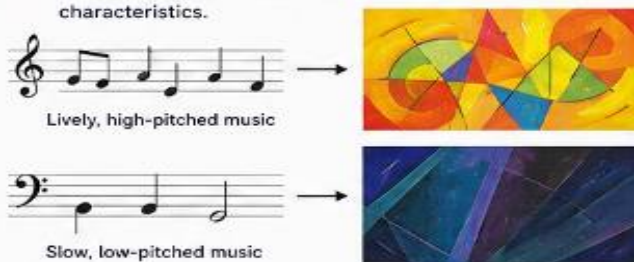
### 6. PITCH ↔ MOTION DIRECTION

Rising (ascending) sounds bias perception toward upward motion; falling (descending) sounds bias toward downward motion.



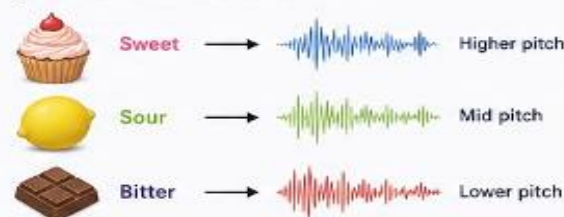
### 7. MUSIC ↔ VISUAL ART / COLOR

People reliably match music with paintings or colors and visual forms with musical characteristics.



### 8. TASTE ↔ SOUND

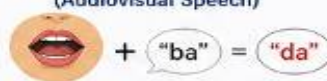
Tastes and flavors are associated with particular sound characteristics.



Flavors can also be matched with musical characteristics (e.g., brighter vs. darker timbre).


### 9. CLASSIC MULTISENSORY ILLUSIONS & EFFECTS

**McGurk Effect (Audiovisual Speech)**




What we hear is influenced by what we see.

**Ventriloquism Effect**



Sounds are perceived as coming from the visual location.

**Pip-and-Pop Effect**

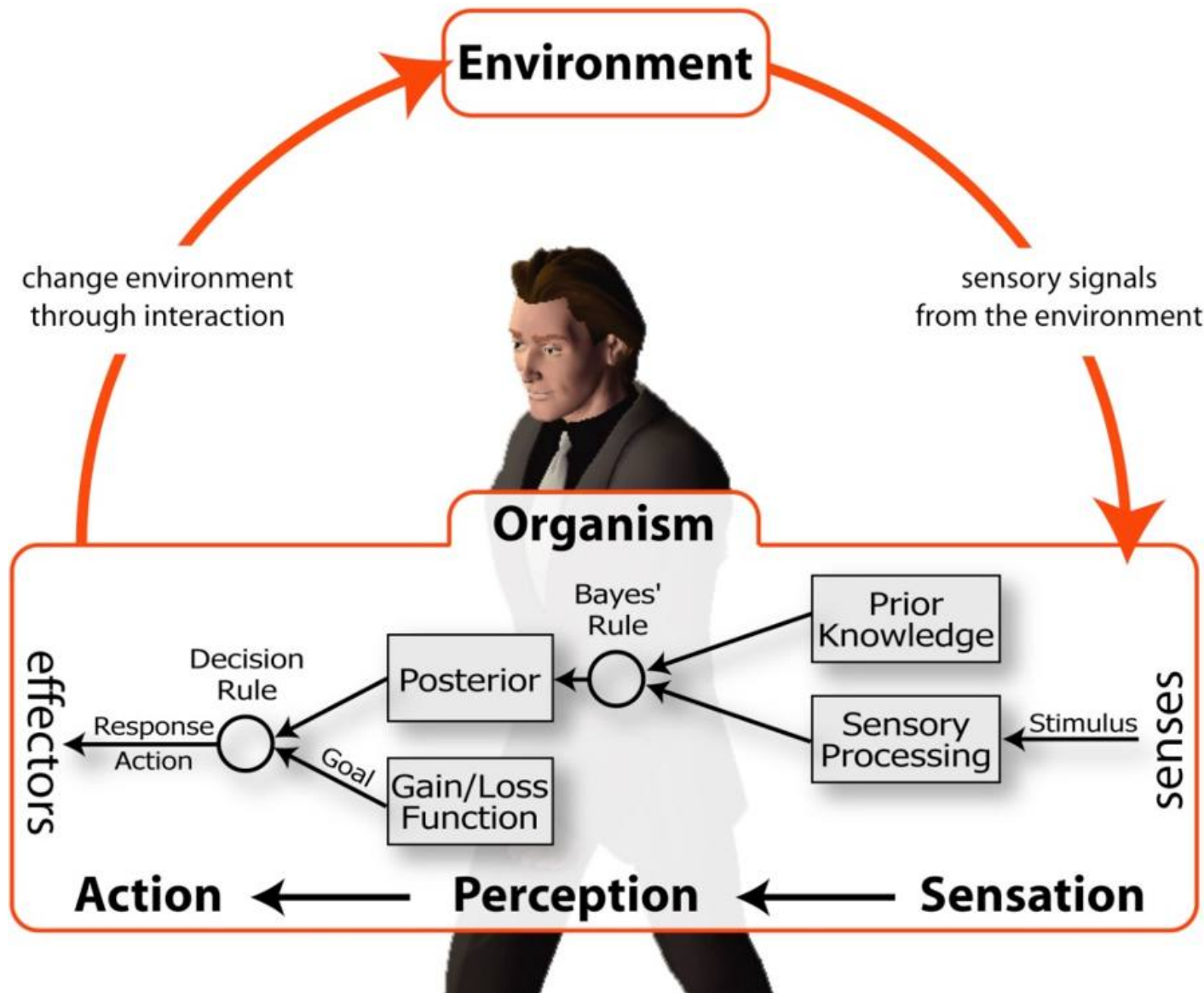


A sound ("pop") makes a brief visual cue ("pip") easier to find.

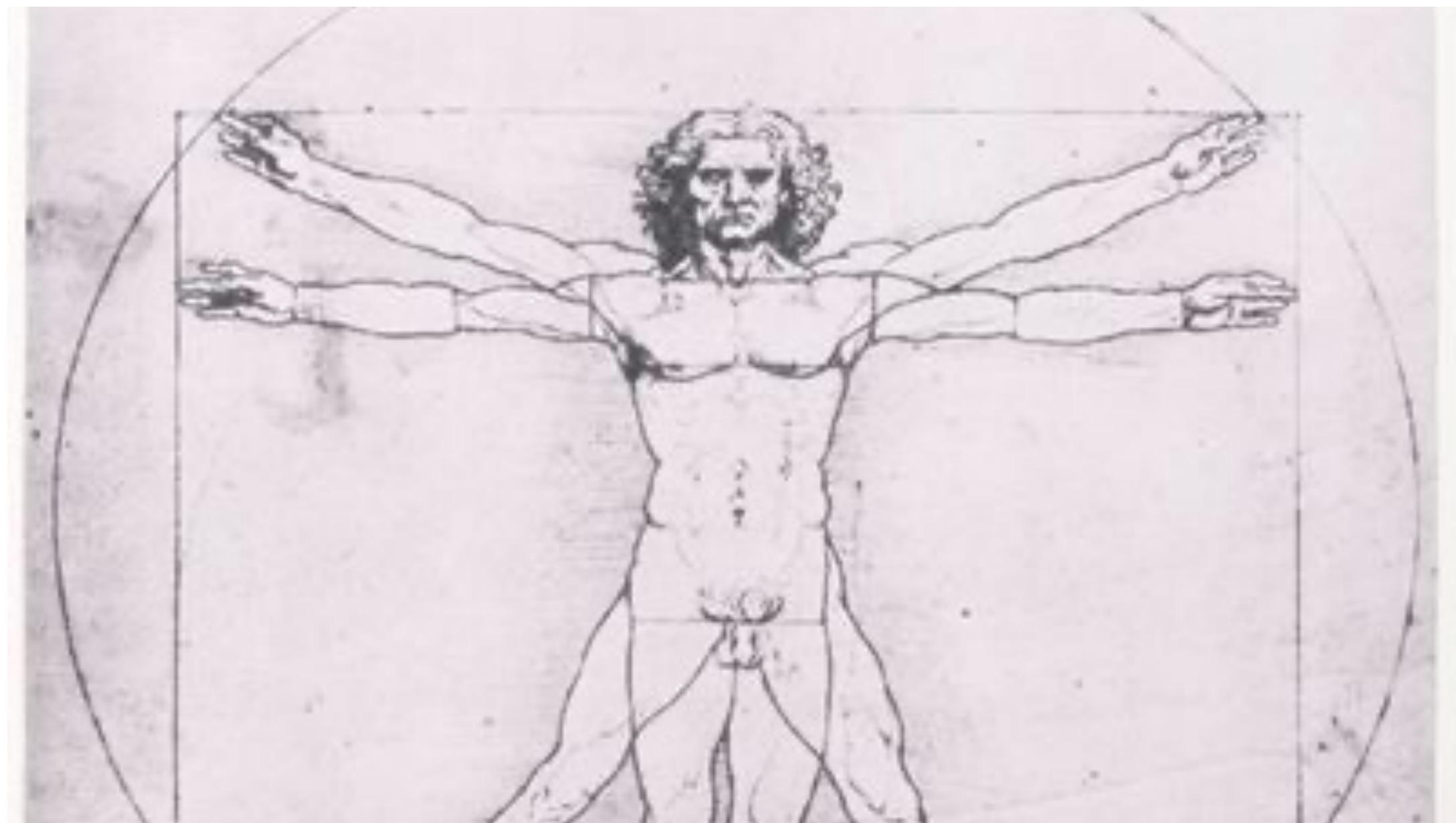
**Cross-modal Plasticity**

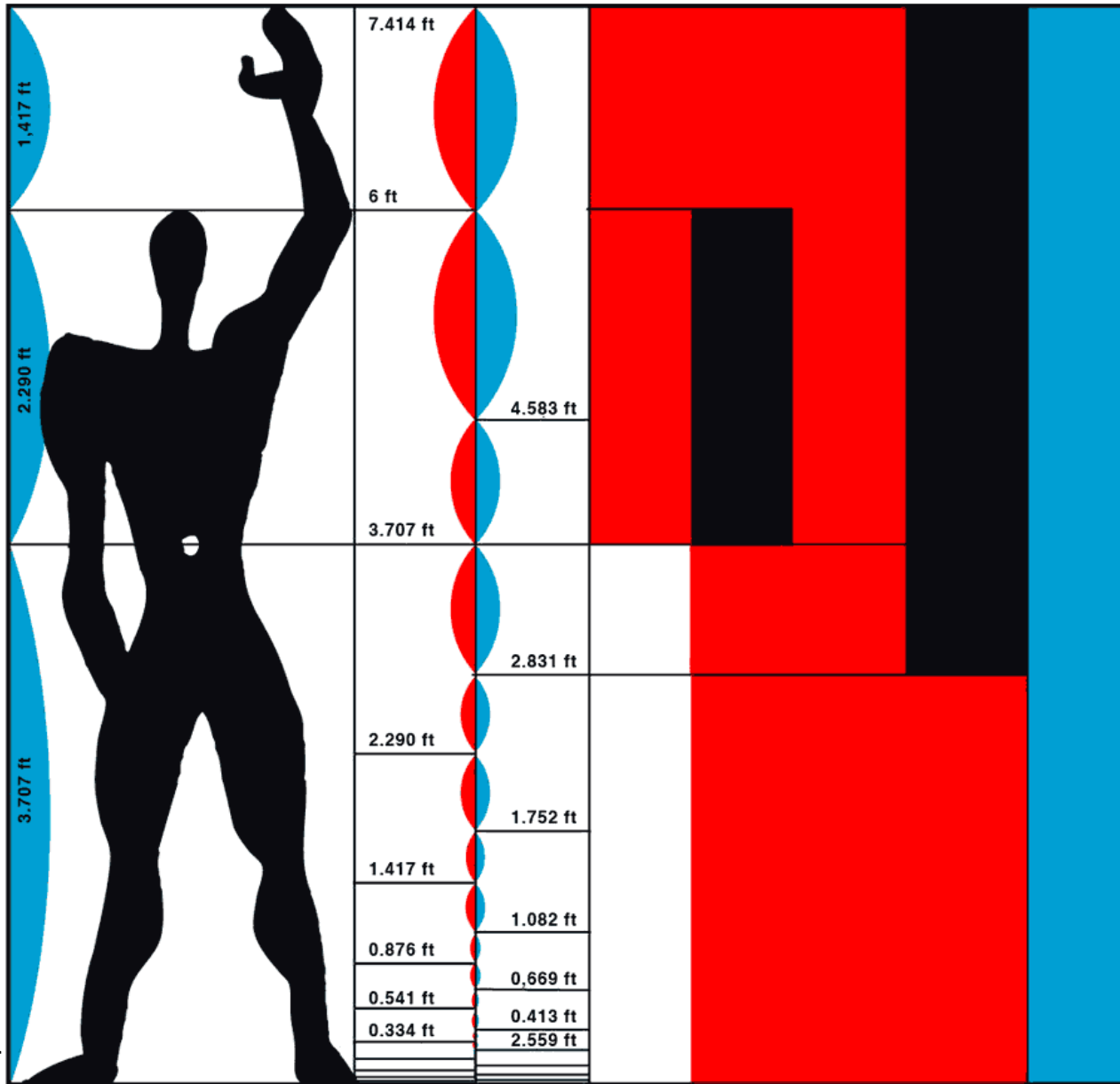


In deaf individuals, visual information can recruit auditory cortex.



Ernst, Marc O. "A Bayesian view on multimodal cue integration." *Human body perception from the inside out* 131 (2006): 105-131.





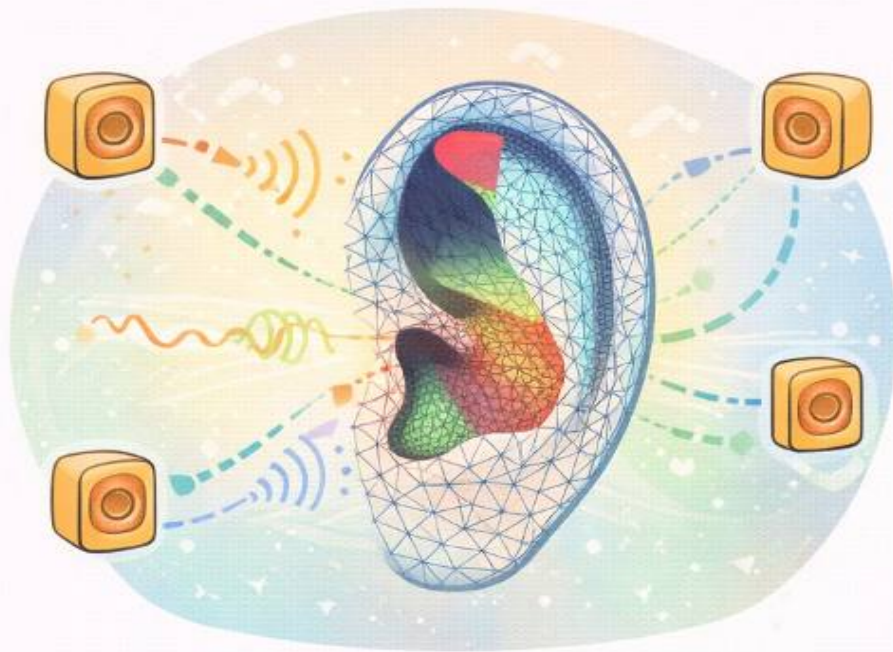
Modulor man LeCorbusier







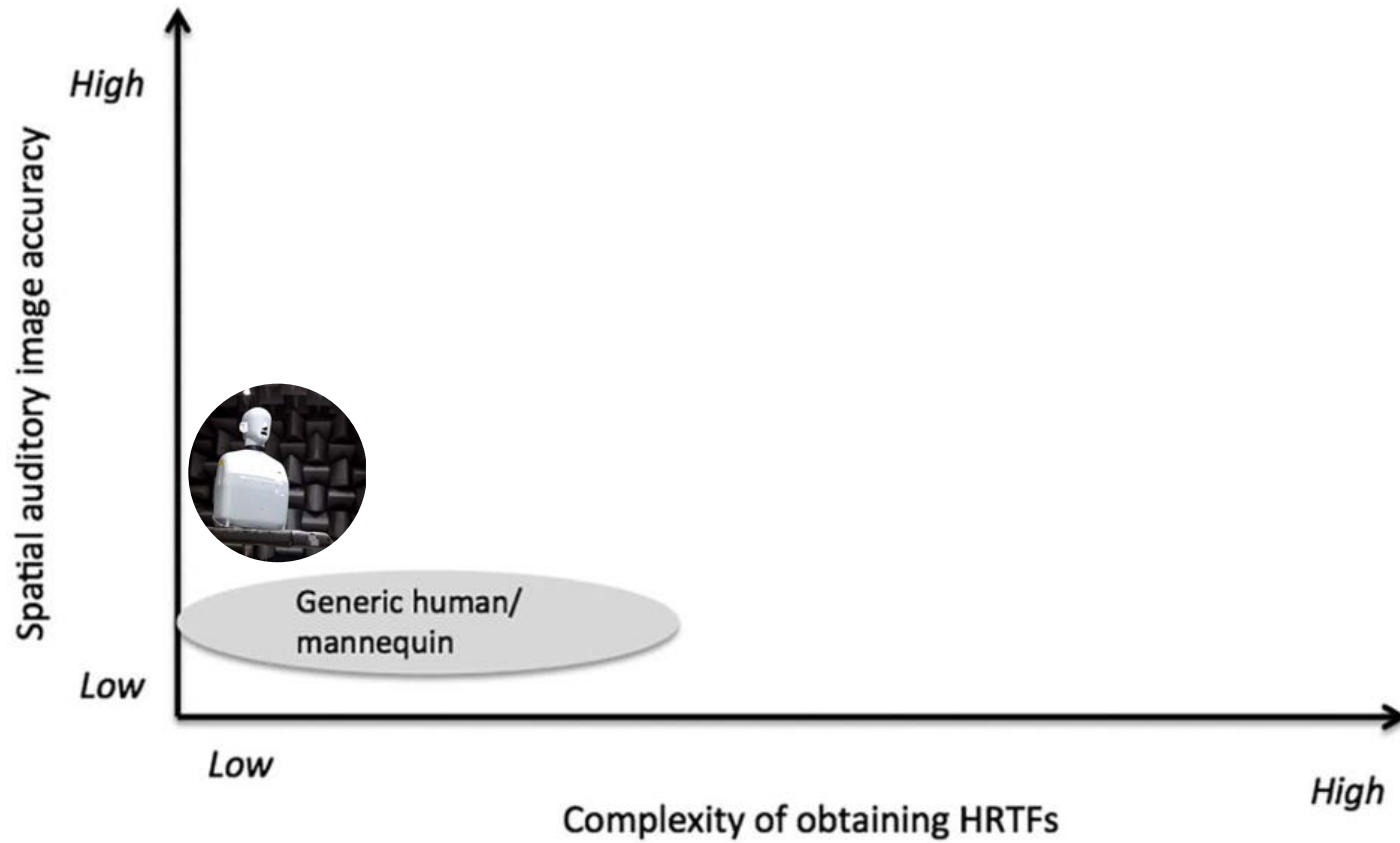
## Your Ear is an Acoustic Fingerprint



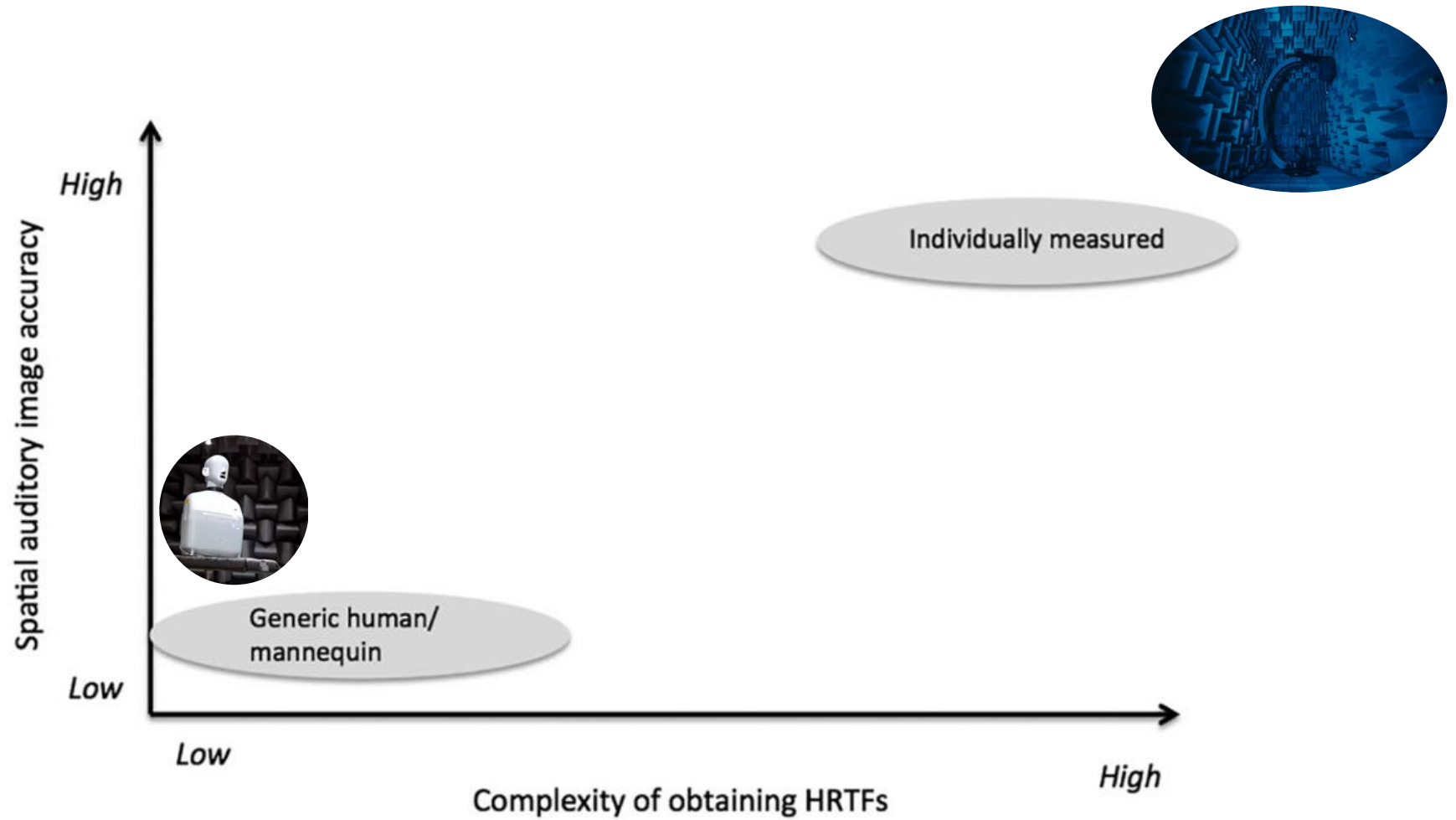
The ear acts as a directional acoustic filter that encodes spatial information.

## Generic HRTF vs Personalized HRTF



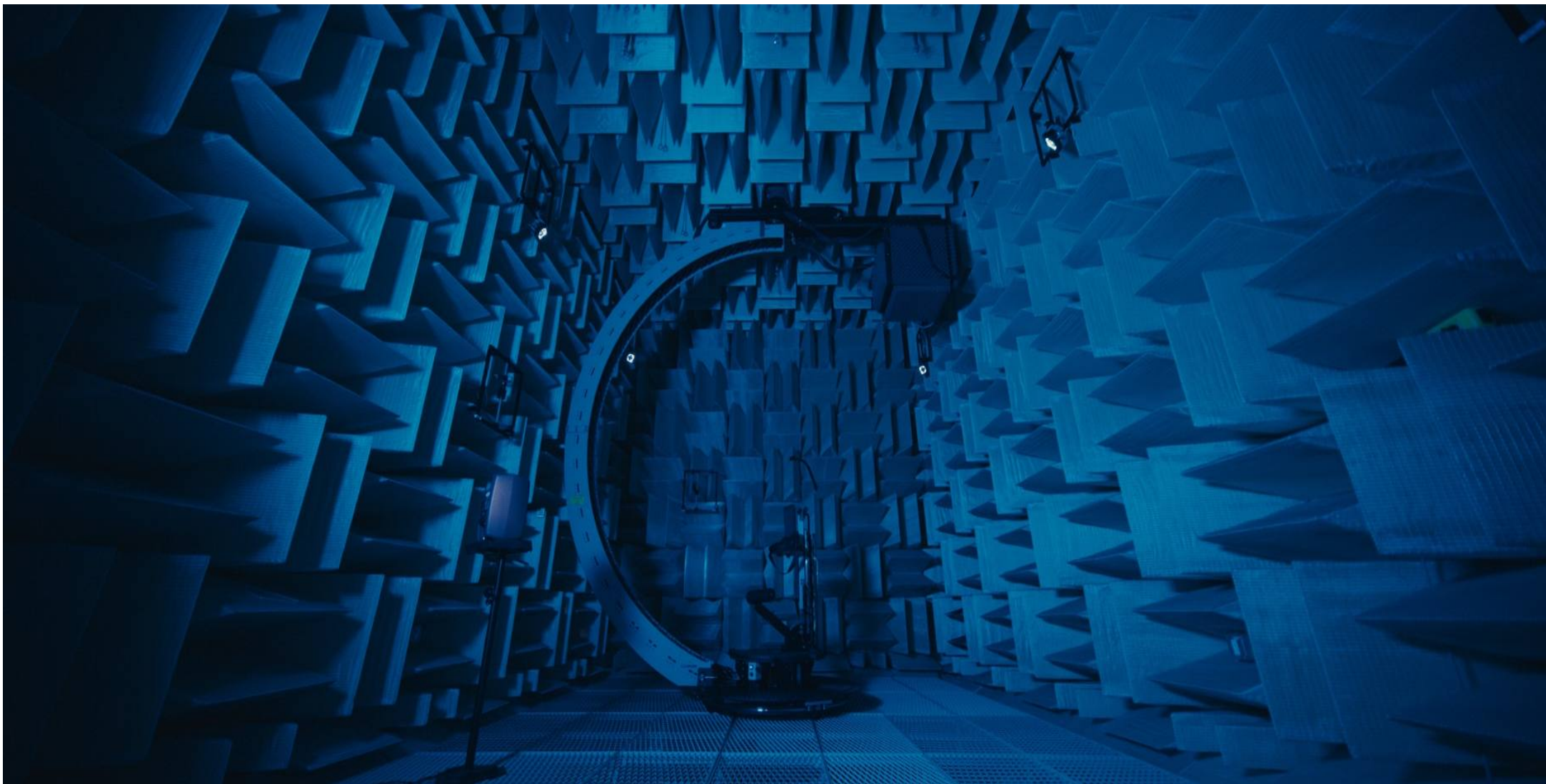


Adapted from: Roginska, A., & Geluso, P. (2017). *Immersive Sound*. Focal Press.



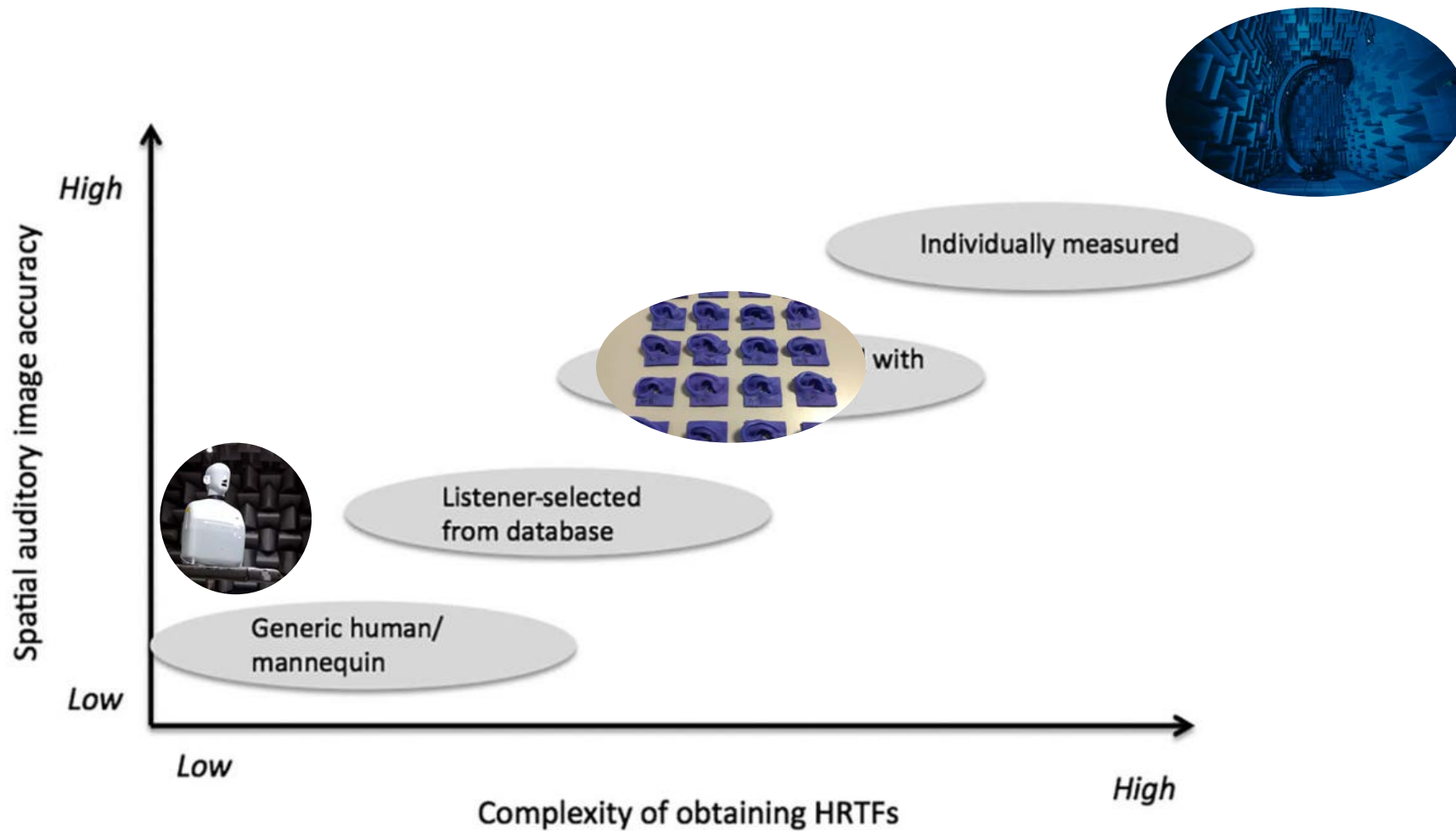
Adapted from: Roginska, A., & Geluso, P. (2017). *Immersive Sound*. Focal Press.

- [https://facebookresearch.github.io/SS2\\_HRTF/](https://facebookresearch.github.io/SS2_HRTF/)



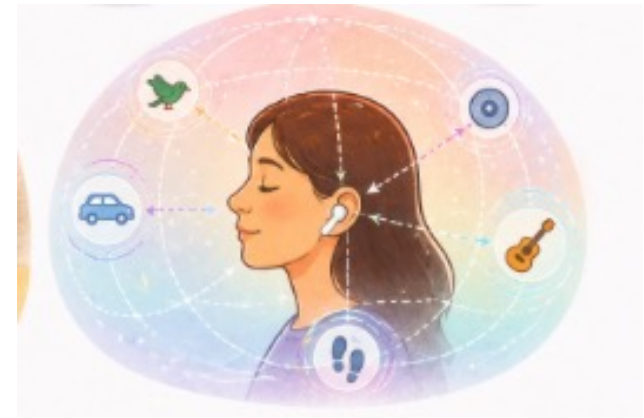


Simone Spagnol, 2021



Adapted from: Roginska, A., & Geluso, P. (2017). *Immersive Sound*. Focal Press.

**Which HRTF FOR WHOM?**



## Generic HRTFs May be Good Enough in Virtual Reality. Improving Source Localization through Cross-Modal Plasticity

 Christopher C. Berger<sup>1,2†</sup>  Mar Gonzalez-Franco<sup>1,†</sup>

 Ana Tajadura-Jiménez<sup>3,4</sup>  Dinei Florencio<sup>1</sup>  Zhengyou Zhang<sup>1,5</sup>

<sup>1</sup> Microsoft Research, Redmond, WA, United States

<sup>2</sup> Division of Biology and Biological Engineering, California Institute of Technology, Pasadena, CA, United States

<sup>3</sup> UCL Interaction Centre, University College London, London, United Kingdom

<sup>4</sup> Interactive Systems DEI-Lab, Universidad Carlos III de Madrid, Madrid, Spain

<sup>5</sup> Department Electrical Engineering, University of Washington, Seattle, WA, United States

SCIENTIFIC REPORTS

**OPEN** Auditory Accommodation to Poorly Matched Non-Individual Spectral Localization Cues Through Active Learning

received: 8 May 2018  
accepted: 17 December 2018  
published online: 31 January 2019

Peter Stitt<sup>1</sup>, Lorenzo Picinali<sup>2</sup> & Brian F. G. Katz<sup>3</sup>



Brain plasticity



Ventriloquism effect

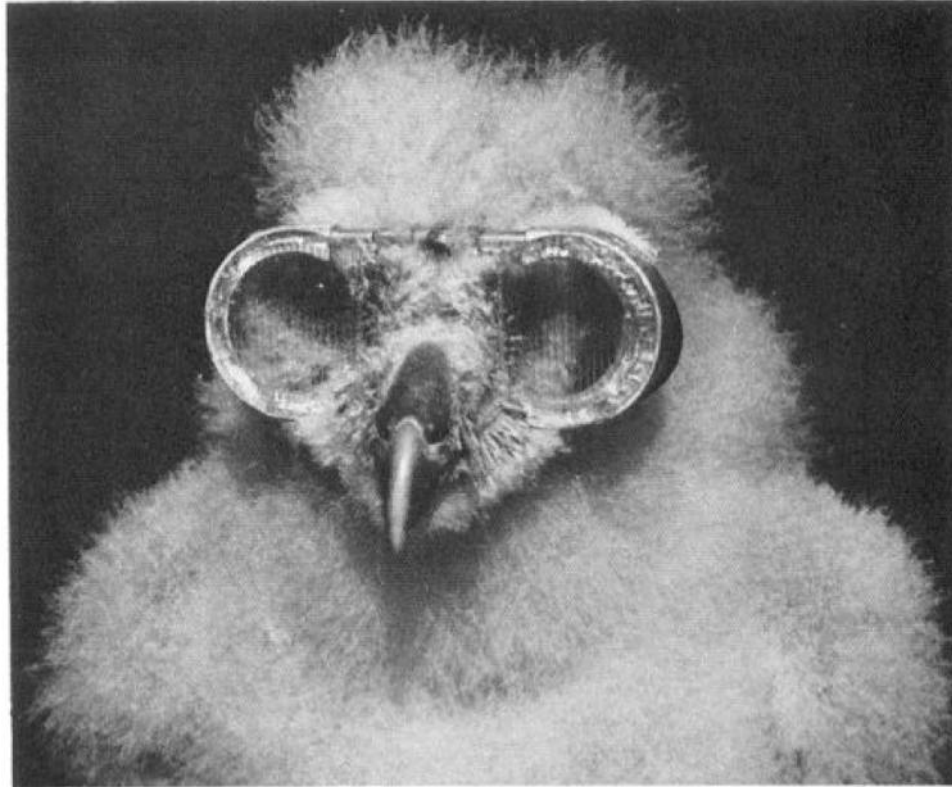
## A new step to optimize sound localization adaptation through the use of vision

Tristan-Gael Bara<sup>1,2</sup>, Alma Guilbert<sup>2</sup>, and Tifanie Bouchara<sup>1</sup>

<sup>1</sup> CEDRIC (EA4626), CNAM, HeSam Université, 75003 Paris, France

<sup>2</sup> VAC Laboratory (EA 7326), Université de Paris, 92774 Boulogne-Billancourt, France

It depends



*Figure 1.* A baby barn owl, 28 d old, wearing binocular Fresnel prisms that displaced the visual field  $34^\circ$  to the right. Photograph indicates the optical quality of the prisms and the large visual field afforded by the spectacles. The vertical lines on the prisms are the edges of the individual prism elements that allow large optical displacements to be achieved with thin prisms.

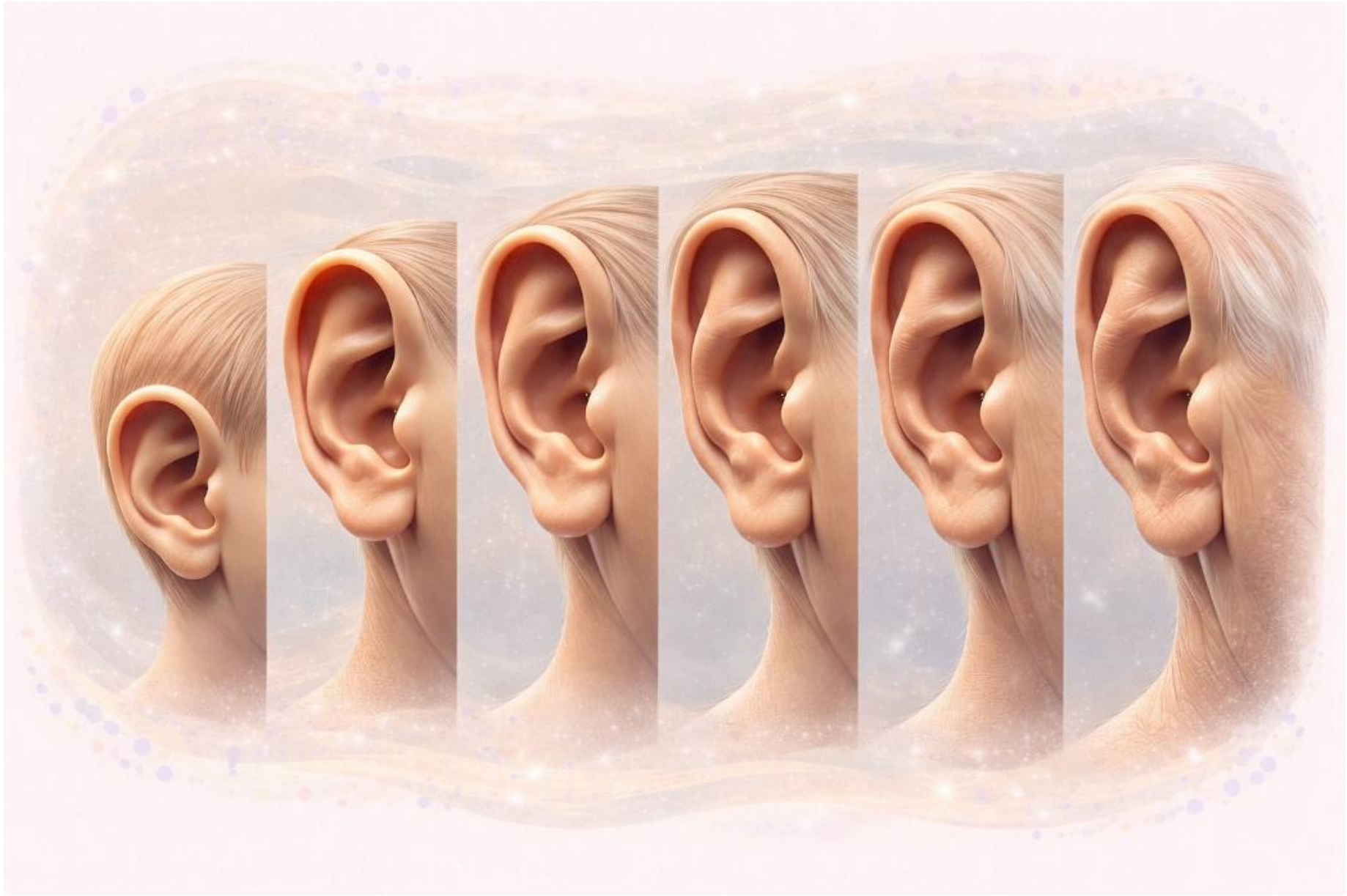
Eric I. Knudsen & Phyllis F. Knudsen (1985), *Vision guides the adjustment of auditory localization in young barn owls.*

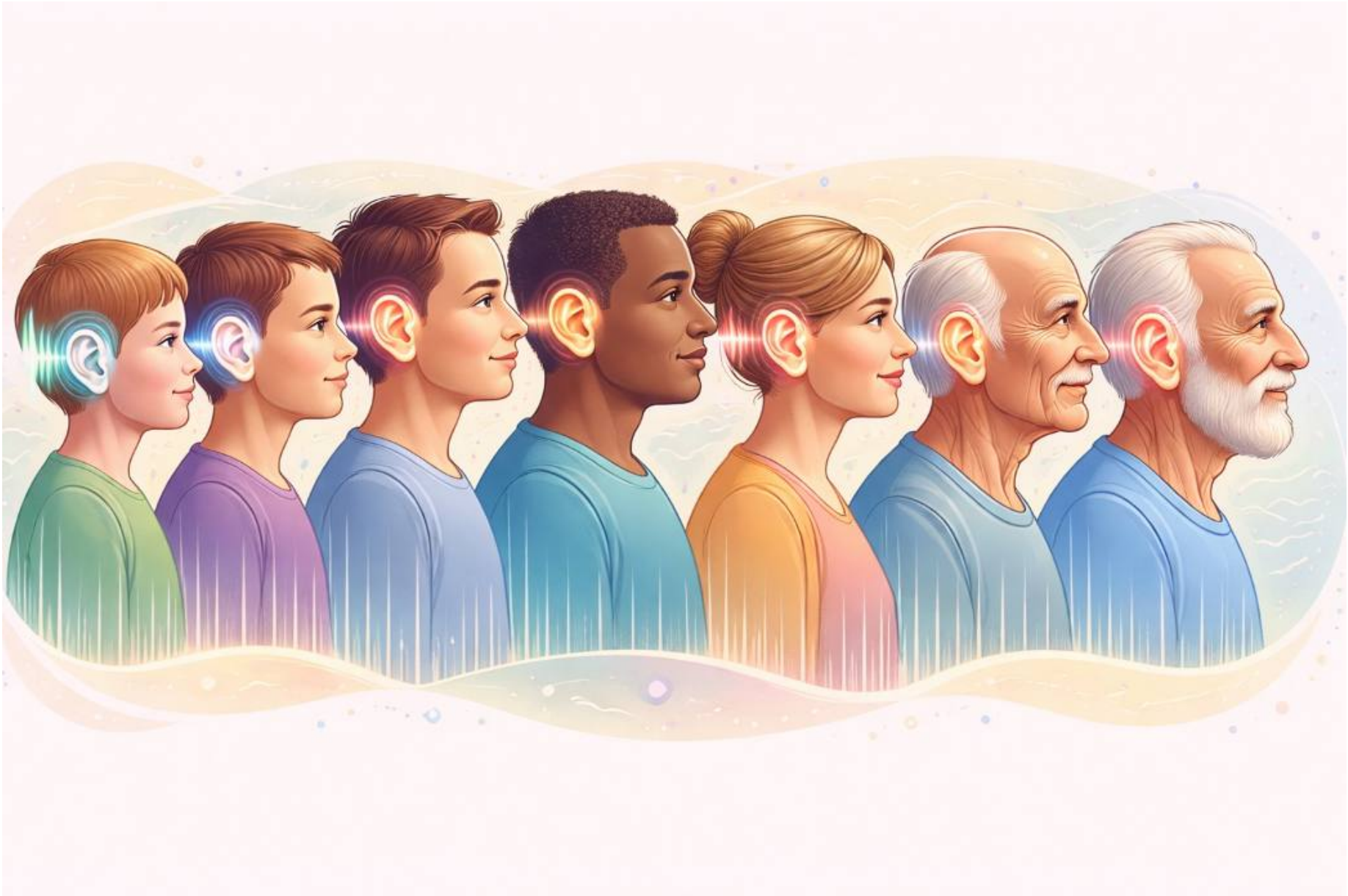


**DANMARKS FRIE  
FORSKNINGSFOND**  
INDEPENDENT RESEARCH  
FUND DENMARK



What about blind or low vision individuals?







# AURAL DIVERSITY



# CoolHear Worskhop

live performance  
multisensory workshop  
interactive demos

19<sup>th</sup> April - 3.00 to 5.00 PM

ROYAL DANISH ACADEMY OF MUSIC  
Rosenørns Alle 22, 1970 Frederiksberg

## HoloBand

Augmented reality for perceptual music training  
Bakka Iványi, Christian Tsilidis, Scott Taylor, Trus Bendik Tjernstrand



## Music training in VR

Virtual Reality application for training listening to music  
Sine Marie Kromann Kristensen, Emi Sønderkov Hansen



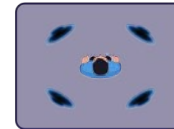
## Uncharted Chants

Mobile game for musical perception training  
Erik Frej Knudsen, Jonas Sim Andersen, Håvard Nuijen



## Lyt Igen

Empowers young people with hearing loss to reach their full potential using VR  
ME-Lab, CHBH and Decibel



## Name the instrument

Listening experiment where the system helps to identify what instruments are playing in a music piece

Doga Buse Cavdik, Francesco Garris



## VAM

High fidelity vibrotactile actuator  
Razvan Paşa

## Effects of spatialization

Listening experiment to examine how extreme spatialization influences musical appreciation and instruments segregation in a string quartet  
Jesper Andersen



## Multisensory live music

Bass, drums and voice duo can be felt in the entire body through custom designed furniture

Antonia Baršić, Peter Williams, Razvan Paşa, Francesco Garris



## Tickle Tuner

Haptic telephone cover for musical training  
Francesco Garris

## For more information and booking

[melcph.create.aau.dk/coolhear-workshop](http://melcph.create.aau.dk/coolhear-workshop)

## Contact us

Stefania Serafin  
sts@create.aau.dk

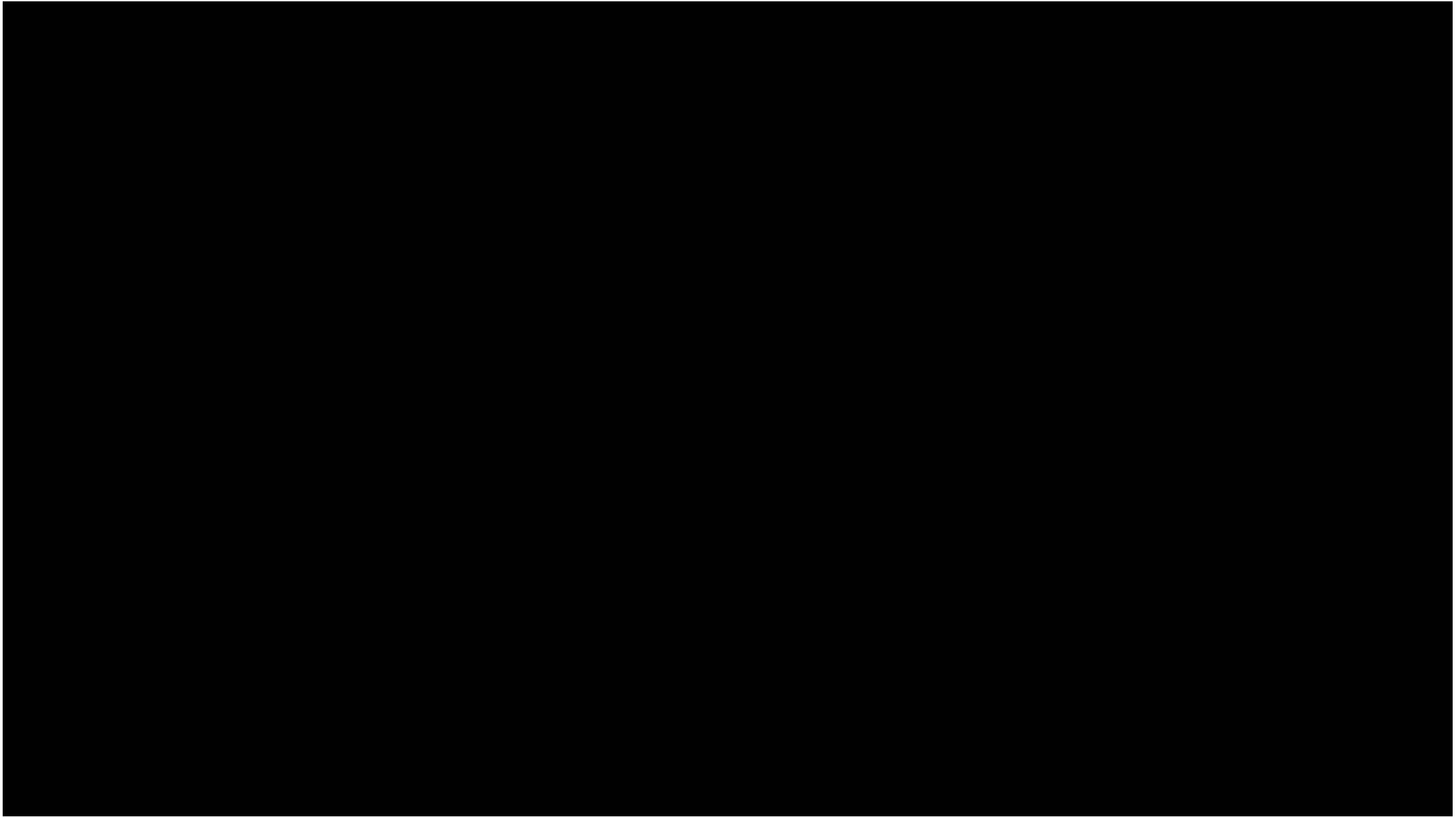
Lone Marianne Percy-Smith  
lone.percy-smith@regionh.dk

## Event organized by



## In collaboration with





# MUSIC VISUALIZATION



# HOLOBAND



Ivanyi, B., Tsalidis, C., Naylor, S., Tjemslund, T. B., Adjorlu, A., Kepp, N. E., & Serafin, S. (2022). HoloBand: An Augmented Reality Experience to Train Music Perception for the hard-of-hearing. In *2022 AES International Conference on Audio for Virtual and Augmented Reality, AVAR 2022* (pp. 1-10). Audio Engineering Society.

### Guitar



Visuals  Solo



### Bass



Visuals  Solo



### Vocals



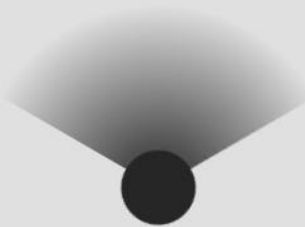
Guitar



Synth

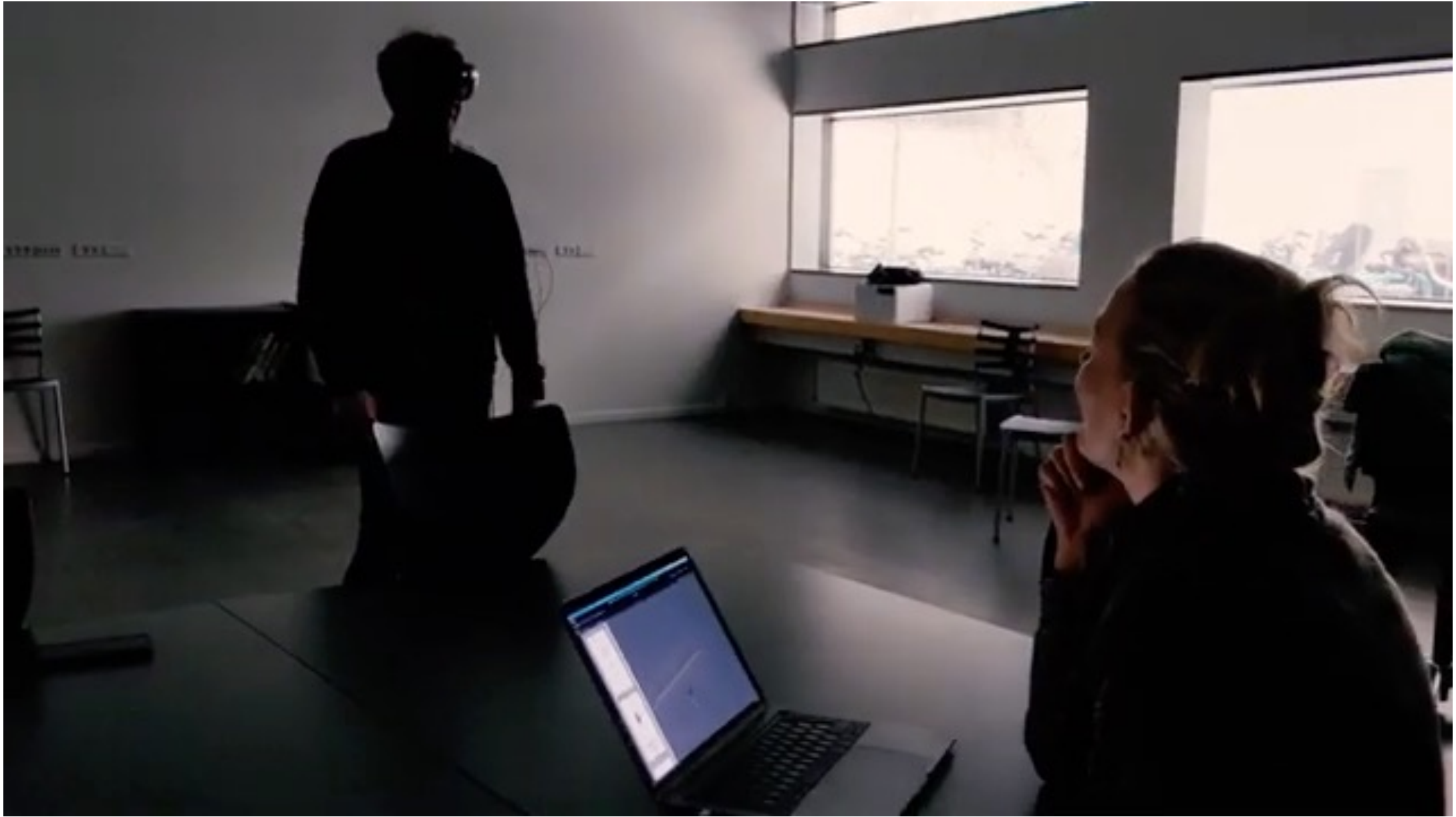


Vocals



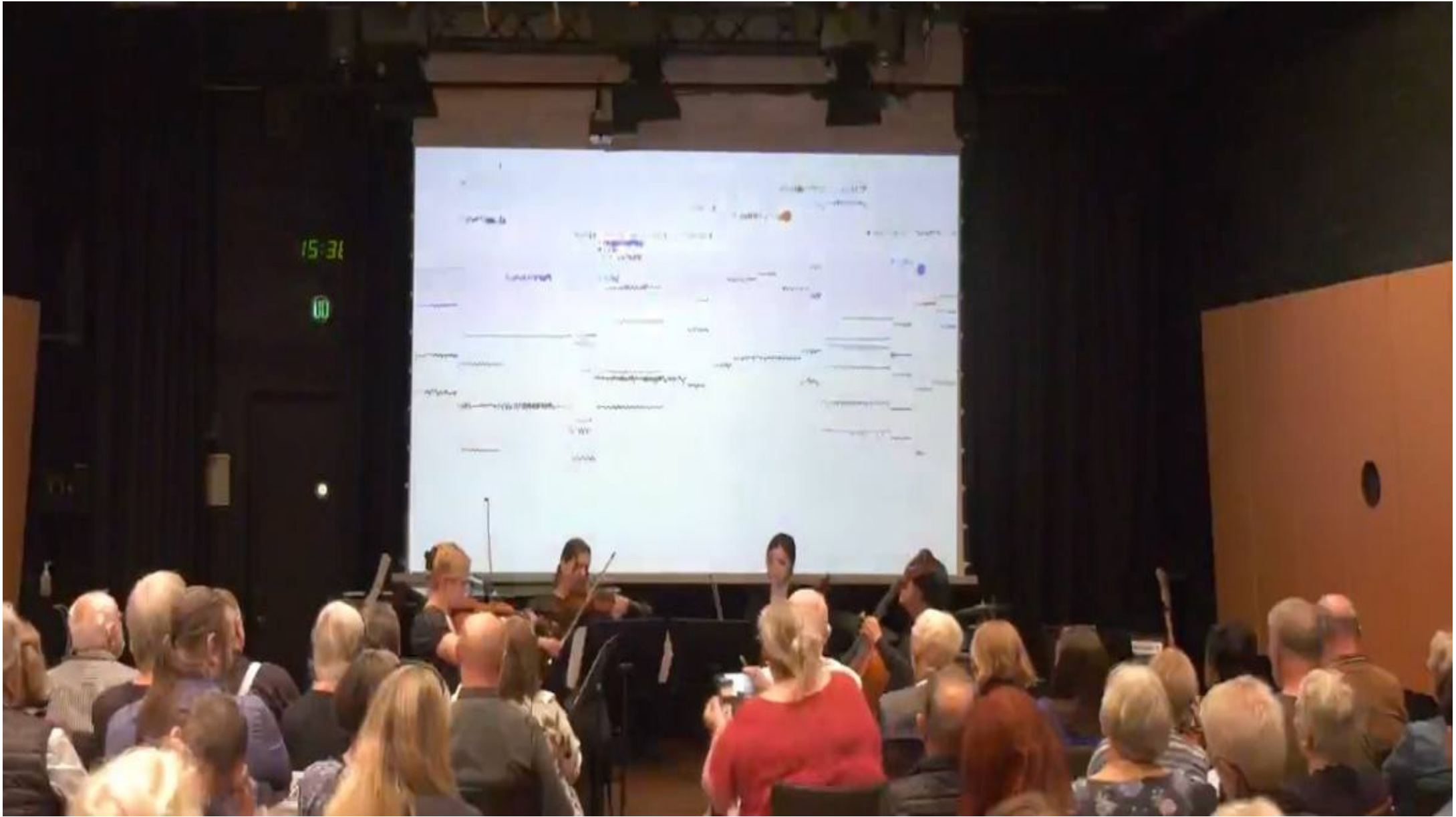
BASS

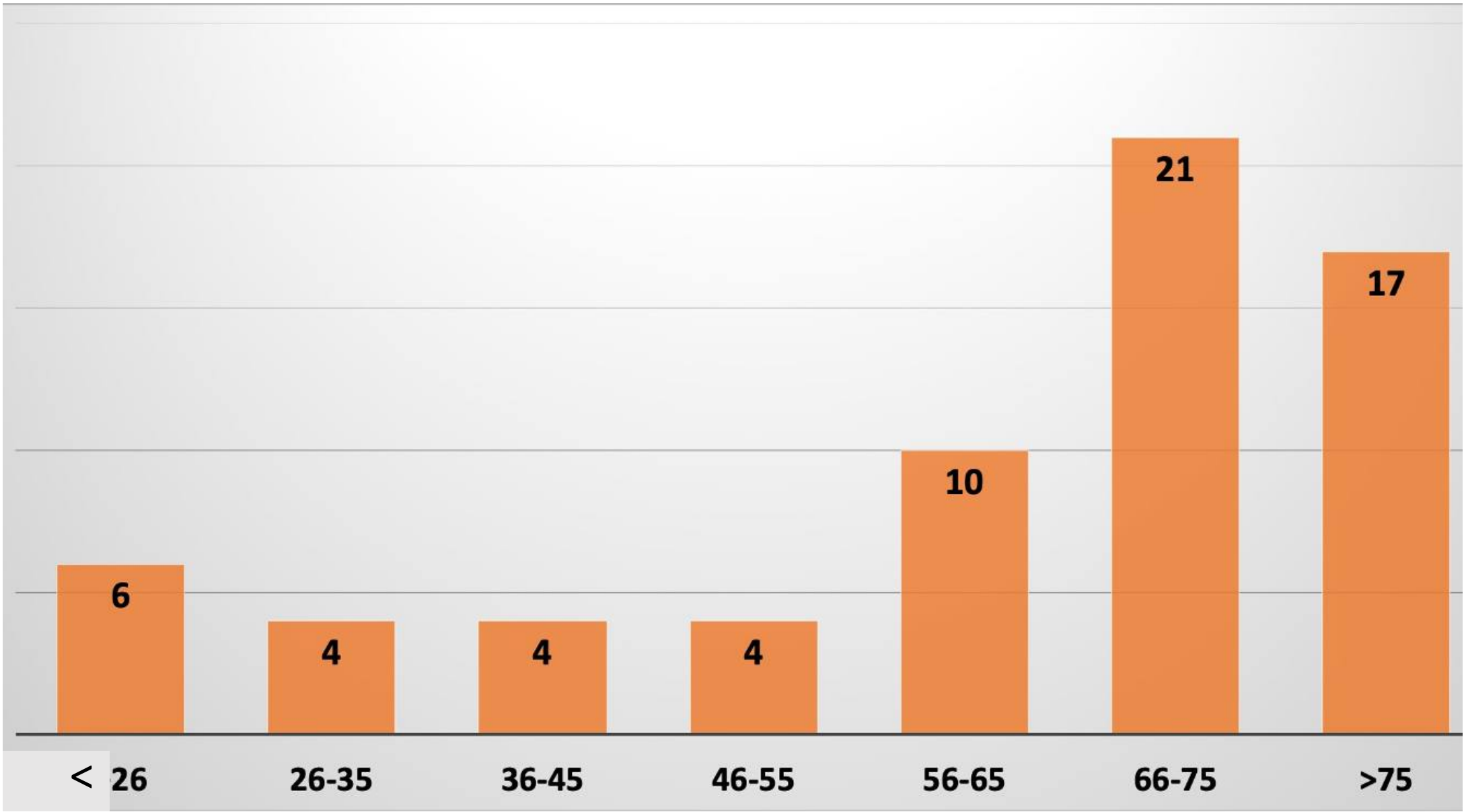


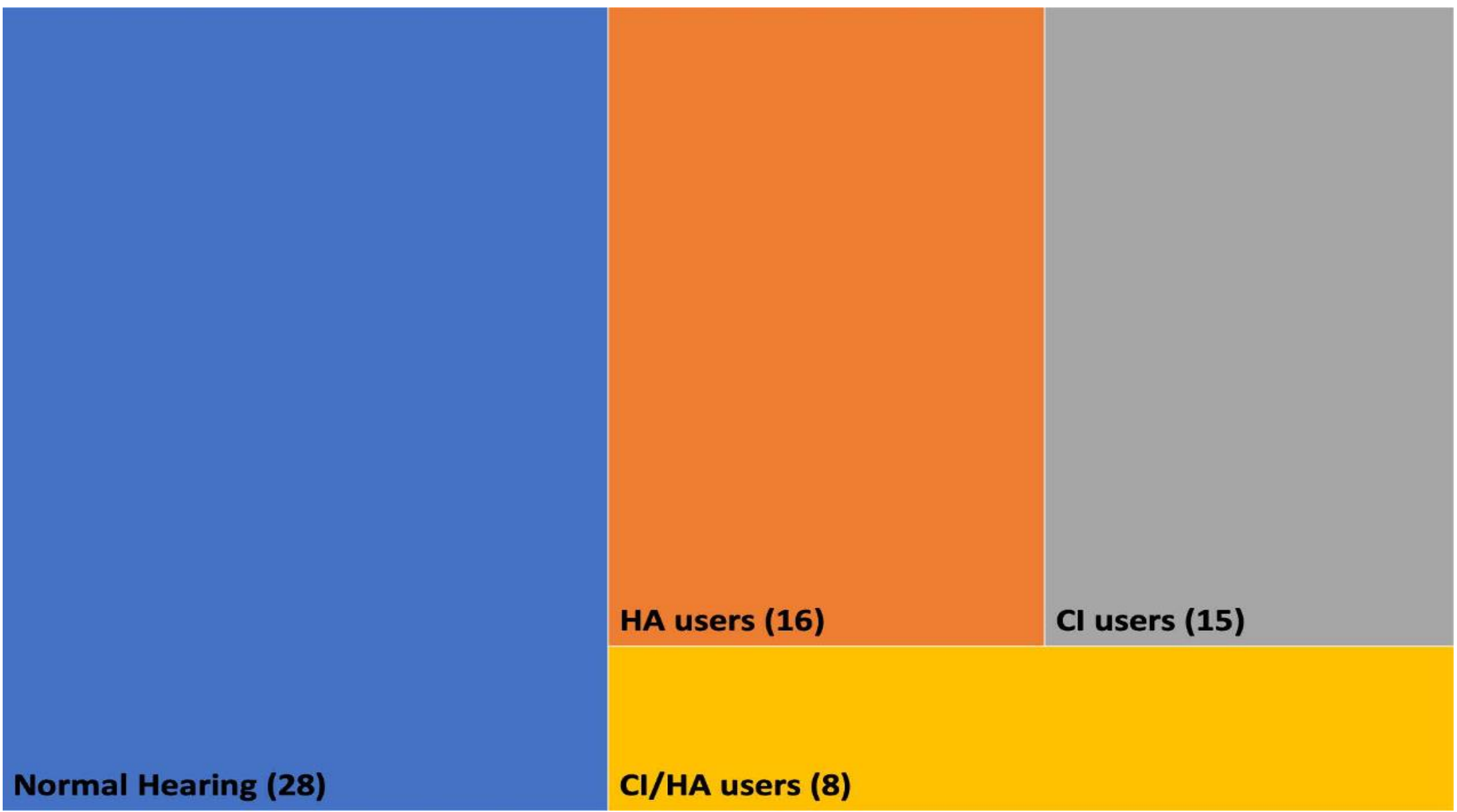


# SOCIAL EXPERIENCES









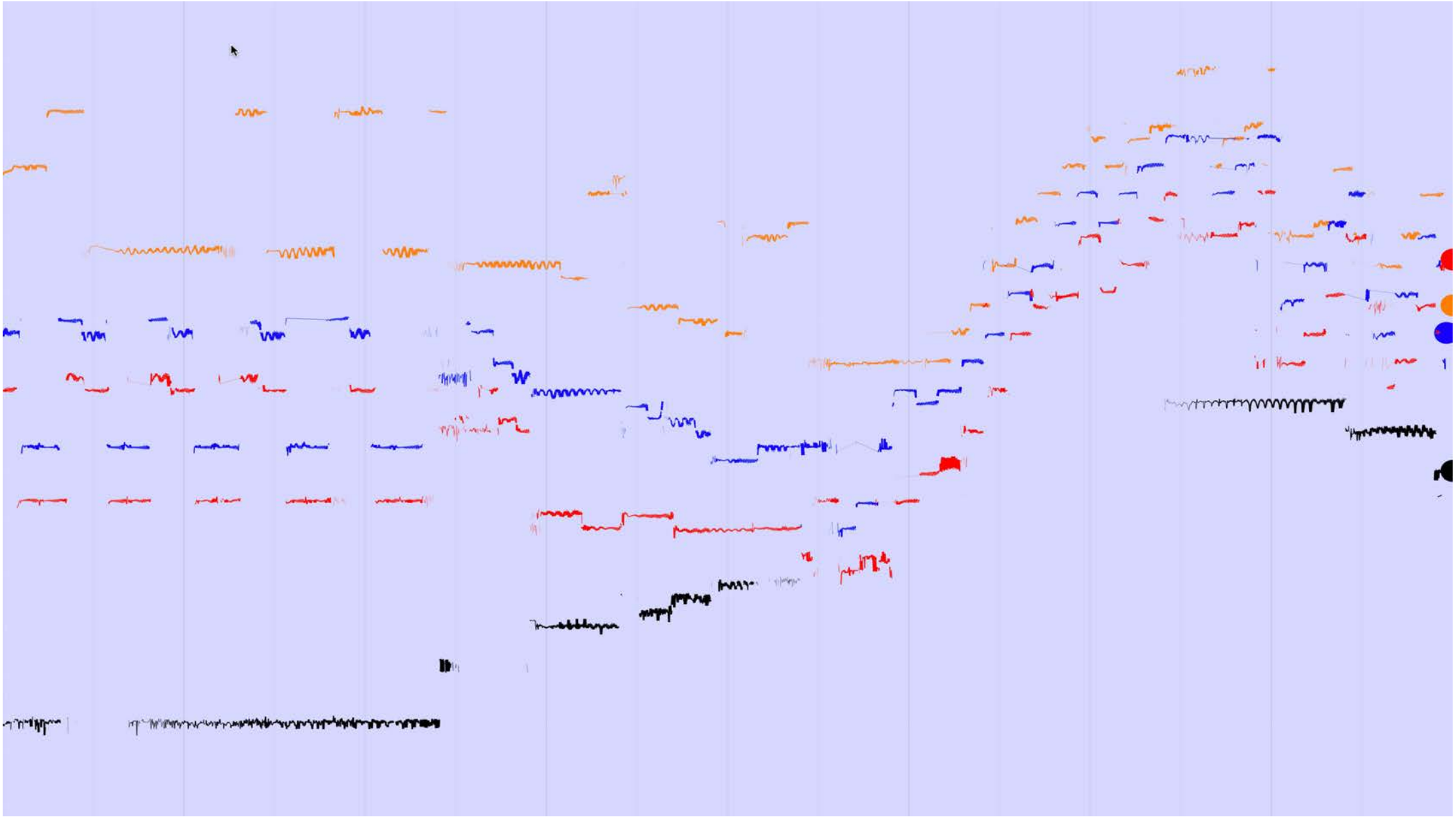
**Normal Hearing (28)**

**HA users (16)**

**CI users (15)**

**CI/HA users (8)**





# RESULTS

Participants filled a questionnaire and 16 participated to interviews.  
Here we focus on interviews

The visualization was useful

The visualization was distracting. I closed my eyes.

The visualization took focus away from the music.

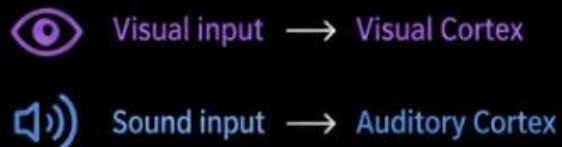
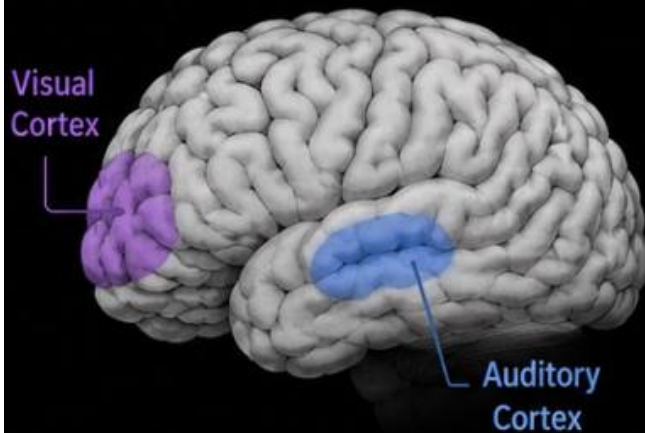
There should have been an introduction.

# Visual Stimuli Can Activate Auditory Cortex After Deafness

*Cross-modal plasticity and outcomes with cochlear implantation*

## 1. NORMAL HEARING

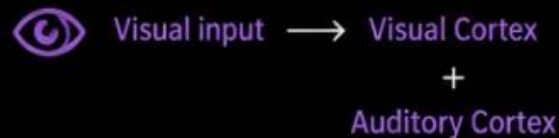
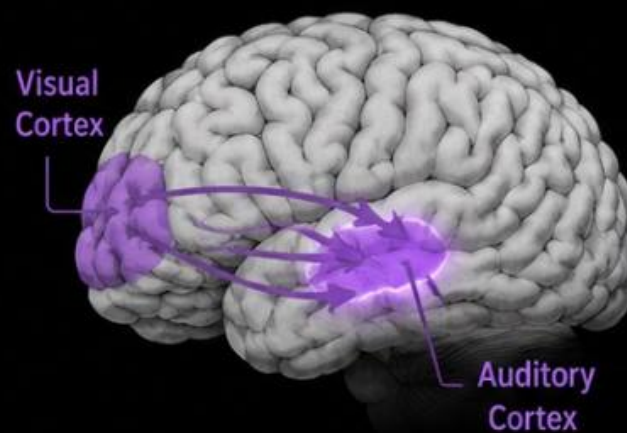
Vision activates visual cortex.  
Sound activates auditory cortex.



Auditory cortex is dedicated

## 2. DEAFNESS

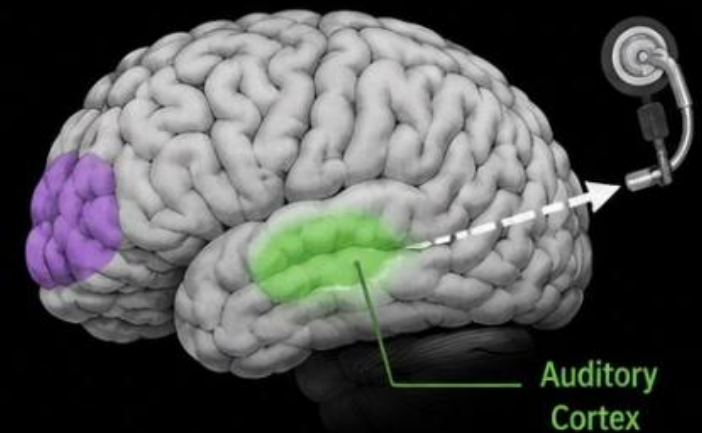
In the absence of sound, visual input  
recruits auditory cortex.



Auditory cortex is partially

## 3. COCHLEAR IMPLANT OUTCOMES

Restoring sound can re-engage auditory cortex.  
Outcome depends on extent of reorganization.



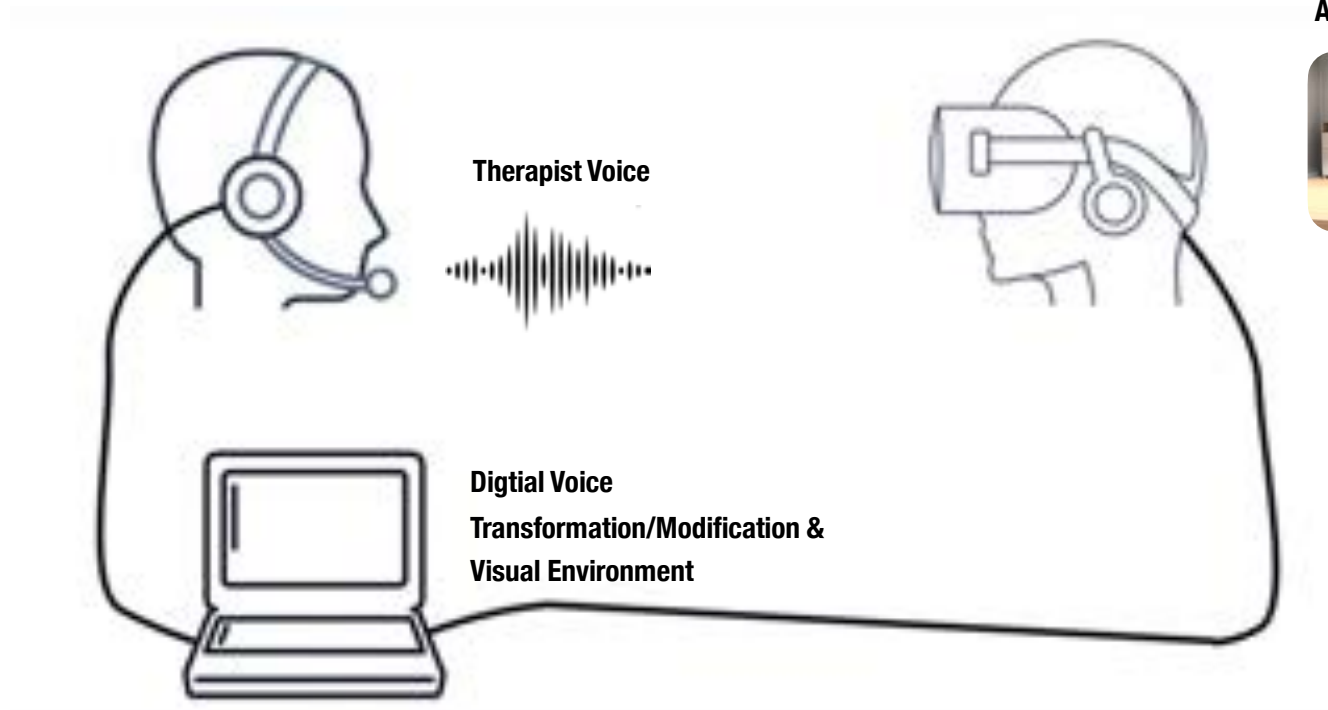
✓ Less visual takeover → Better CI outcomes  
(auditory cortex more available for sound)

⚠ More visual takeover → Poorer CI outcomes  
(auditory cortex less available for sound)

# HEARING VOICES



# AVATAR THERAPY



Avatar/Malevolent Voice

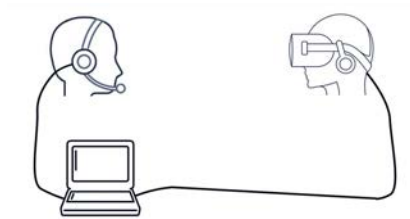


# AVATAR THERAPY

0 1

## Session setup

*The patient is positioned in virtual reality (VR)*



0 2

## Avatar Creation

*Avatar is created matching the perceived voice*



0 3

## Live Interaction

*The therapist speaks and guides the patient*



0 4

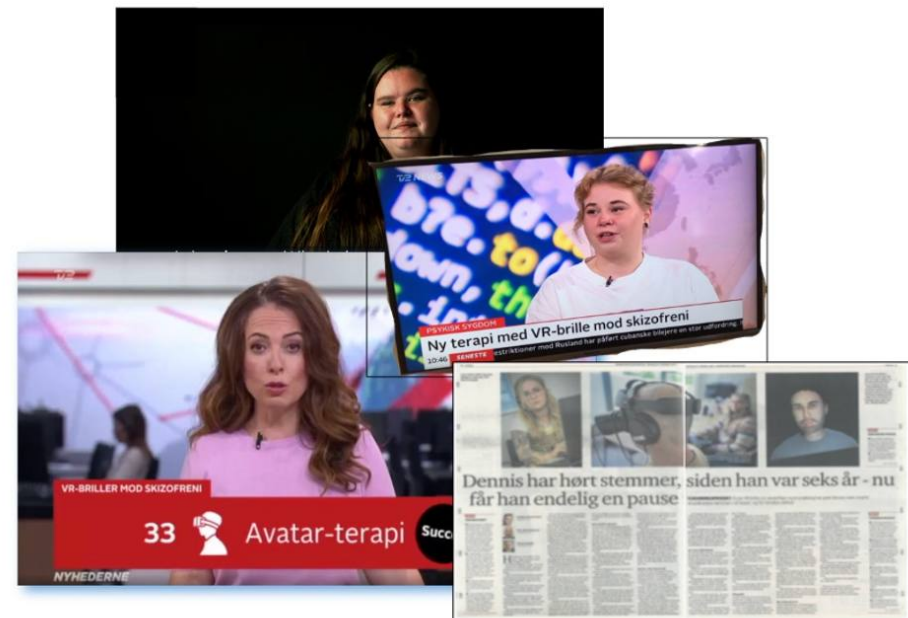
## Observation

*The therapist observes and iterates (domain knowledge)*



# AVATAR THERAPY

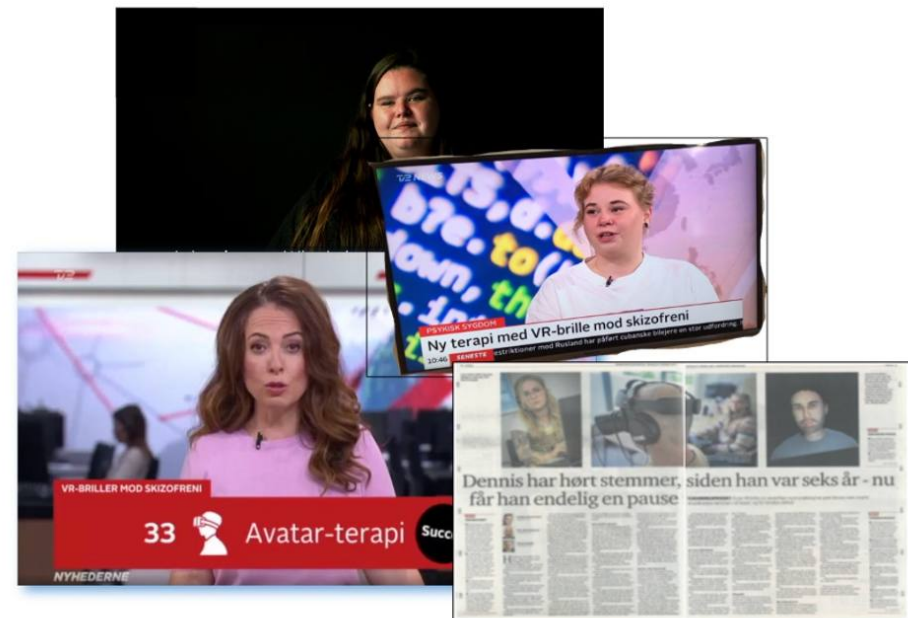
- Reduced severity and frequency of auditory hallucinations after 12 weeks.
- 11 out of 271 treatment resistant patients no longer heard voices by the 24 week follow up.
- Greater satisfaction with treatment.
- Not for everyone, with possibilities of improving general realism, immersion and user-experience.



[1] Glenthøj, et al., "Impact of avatar features and presence on treatment outcomes in virtual reality-assisted therapy for auditory hallucinations", Schizophrenia research, Elsevier, 2025

# AVATAR THERAPY

- Reduced severity and frequency of auditory hallucinations after 12 weeks.
  - 11 out of 271 treatment resistant patients no longer heard voices by the 24 week follow up.
  - Greater satisfaction with treatment.
- Not for everyone, with possibilities of improving general realism, immersion and user-experience.



[ L. Glenthøj, et al., "Impact of avatar features and presence on treatment outcomes in virtual reality-assisted therapy for auditory hallucinations", Schizophrenia research, Elsevier, 2025

# AVATAR THERAPY

“

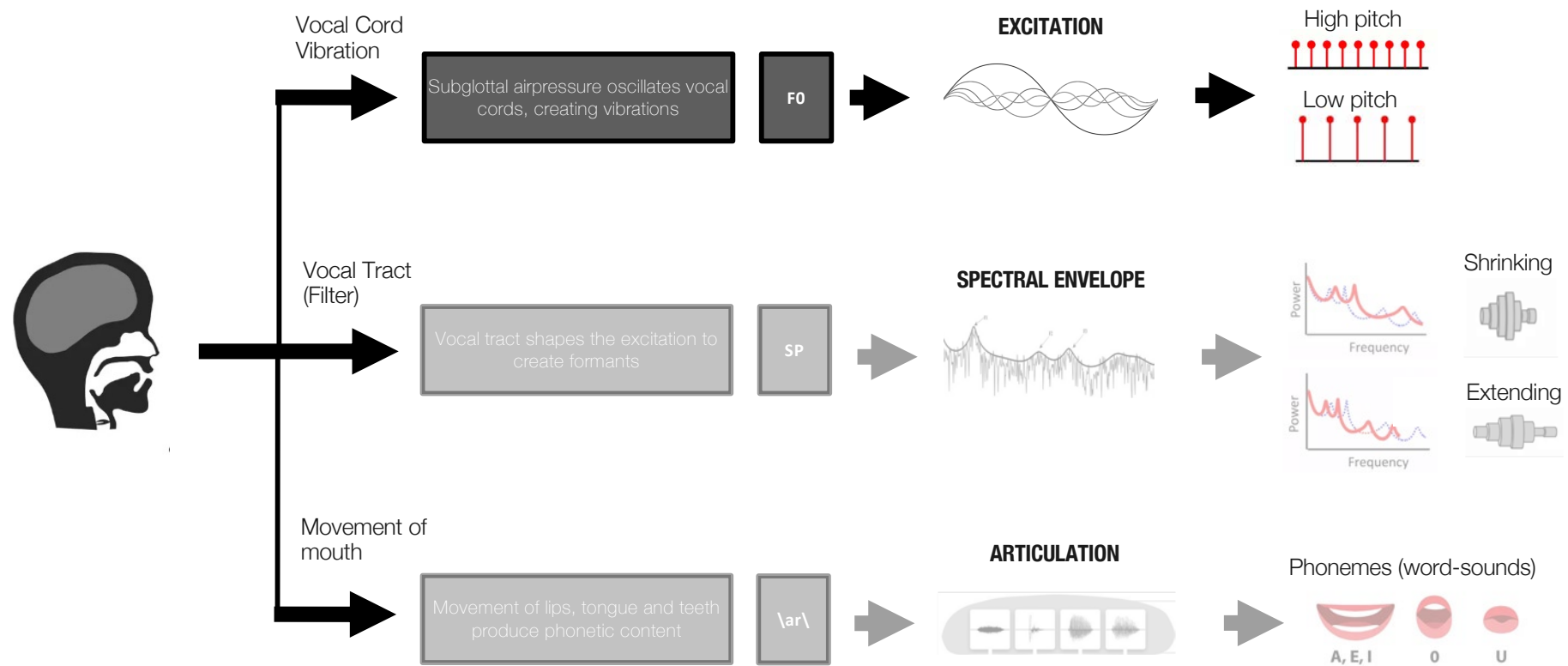
*Realism of the avatar, particularly  
in its auditory dimension, may  
facilitate engagement and  
therapeutic processing.*

”

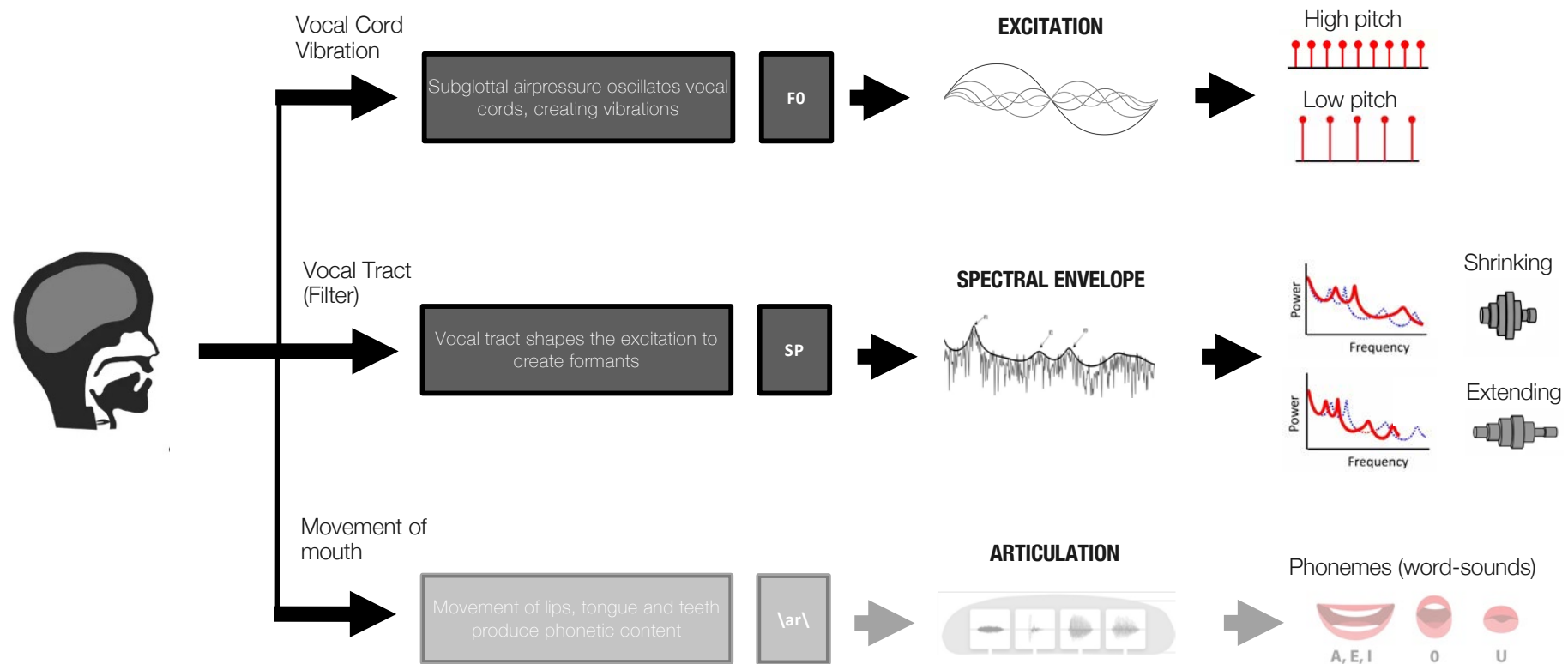
[1] L. Glenthøj, et al., "Impact of avatar features and presence on treatment outcomes in virtual reality-assisted therapy for auditory hallucinations", Schizophrenia research, Elsevier, 2025

---

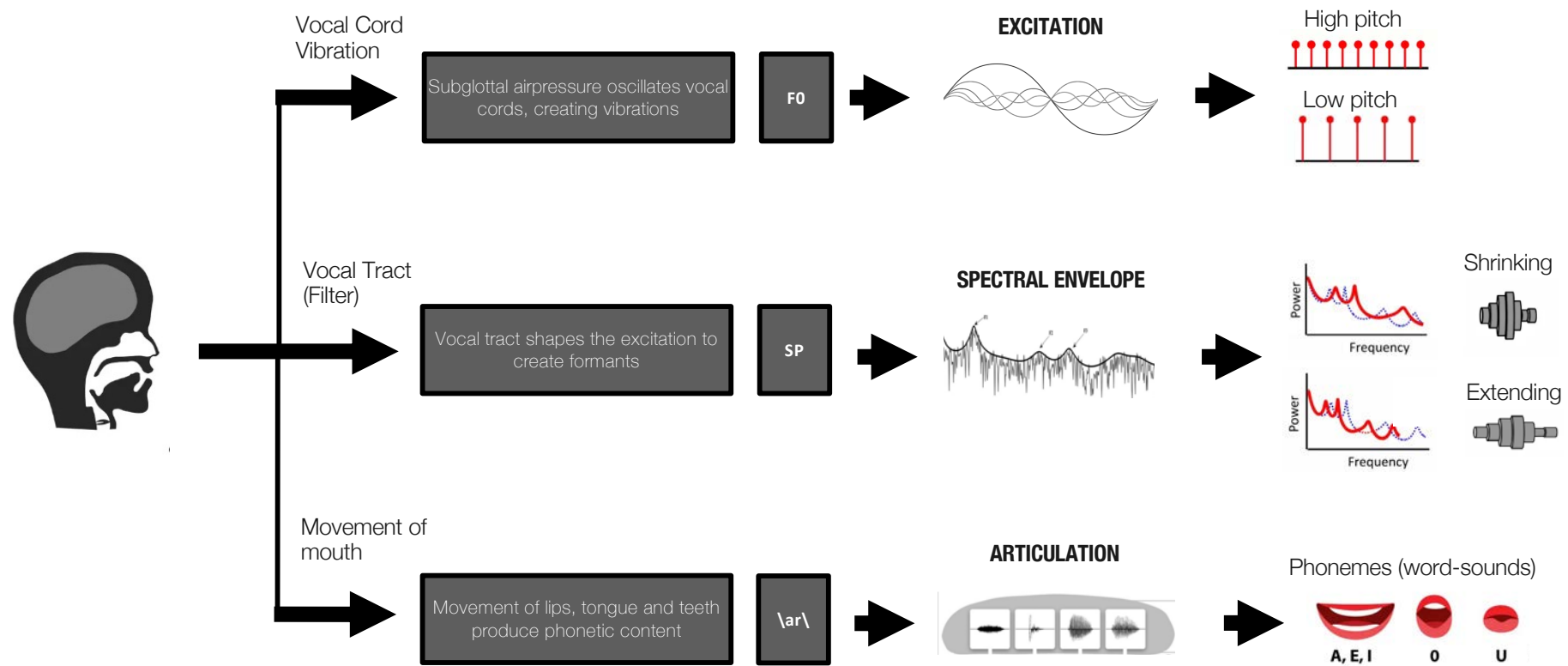
# SPEECH PRODUCTION



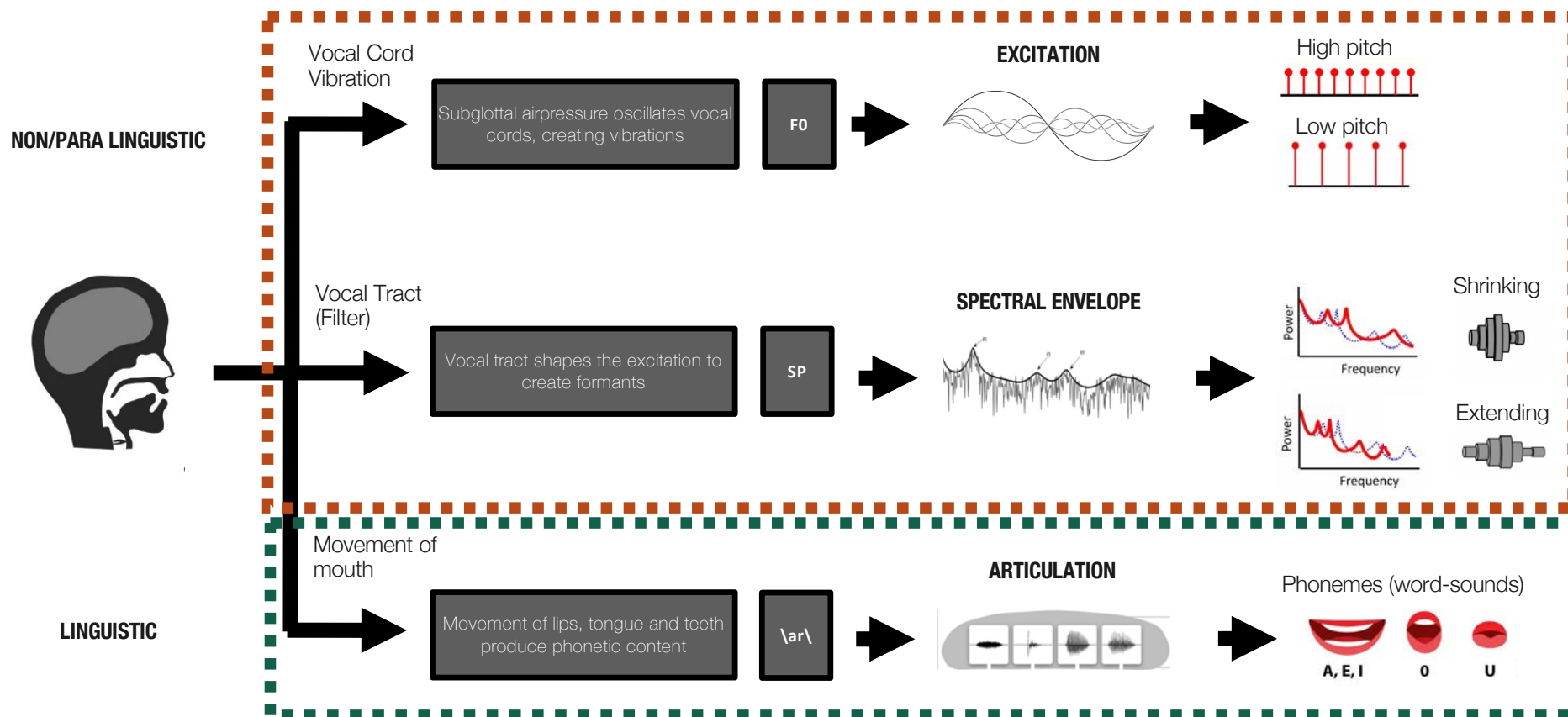
# SPEECH PRODUCTION



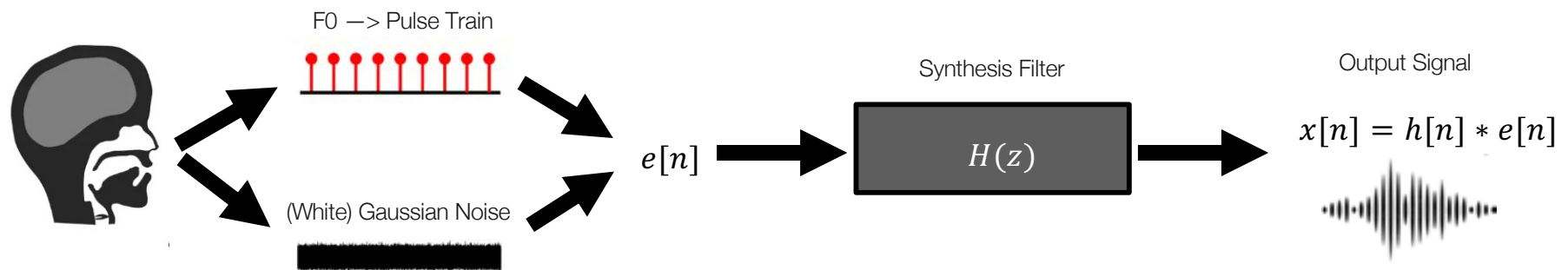
# SPEECH PRODUCTION



# SPEECH PRODUCTION



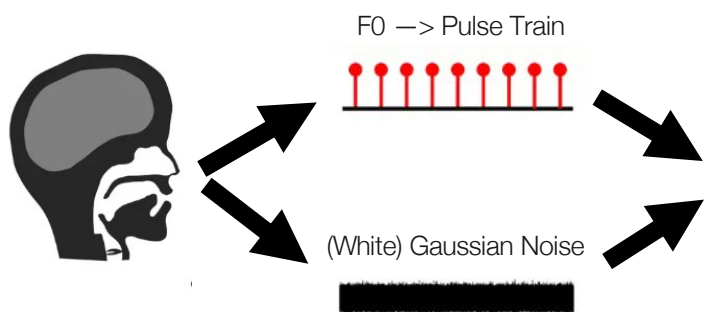
# VOICE CONVERSION (NAIVE APPROACH)



# VOICE CONVERSION (NAIVE APPROACH)

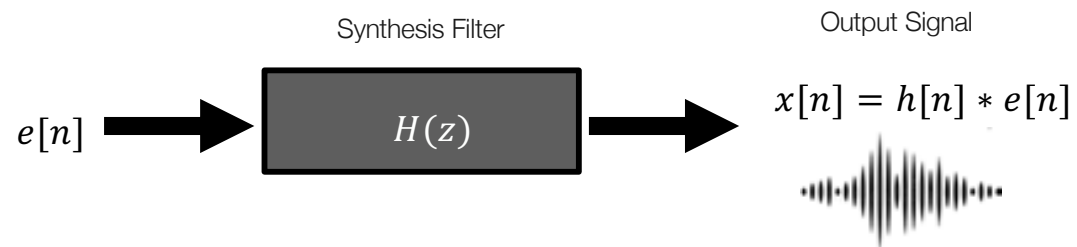
## ANALYSIS

Extract features directly from input



## SPECTRAL PREDICTION

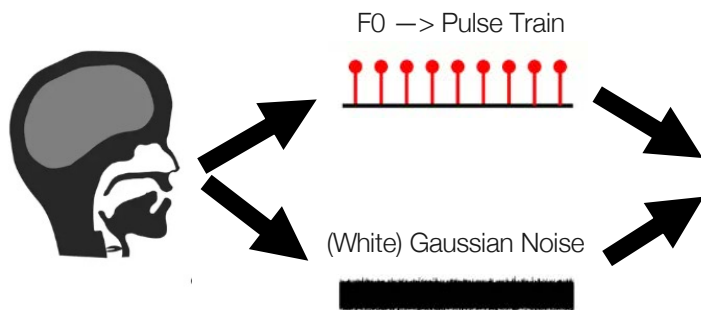
Learn target specific coefficients for a filter



# VOICE CONVERSION (NAIVE APPROACH)

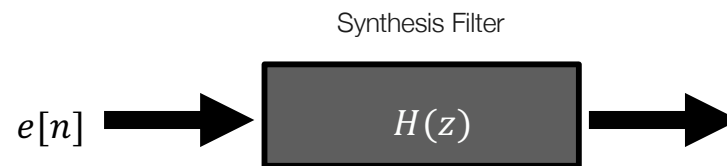
## ANALYSIS

Extract features directly from input



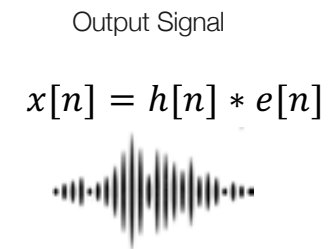
## SPECTRAL PREDICTION

Learn target specific coefficients for a filter

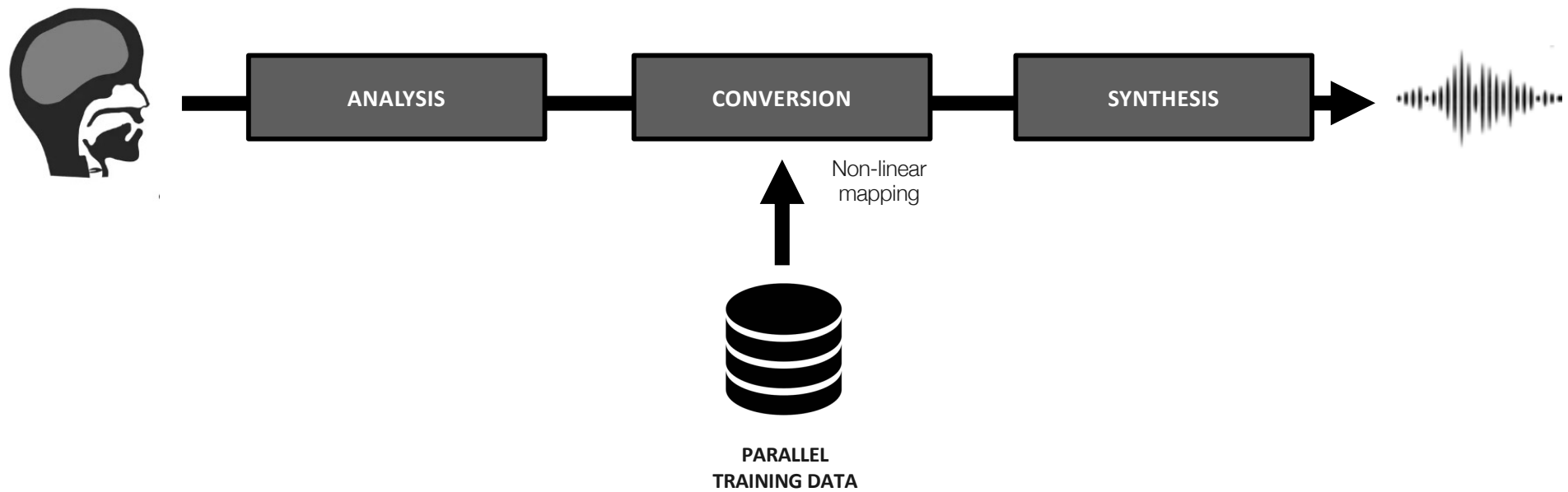


## RULE-BASED SYNTHESIS

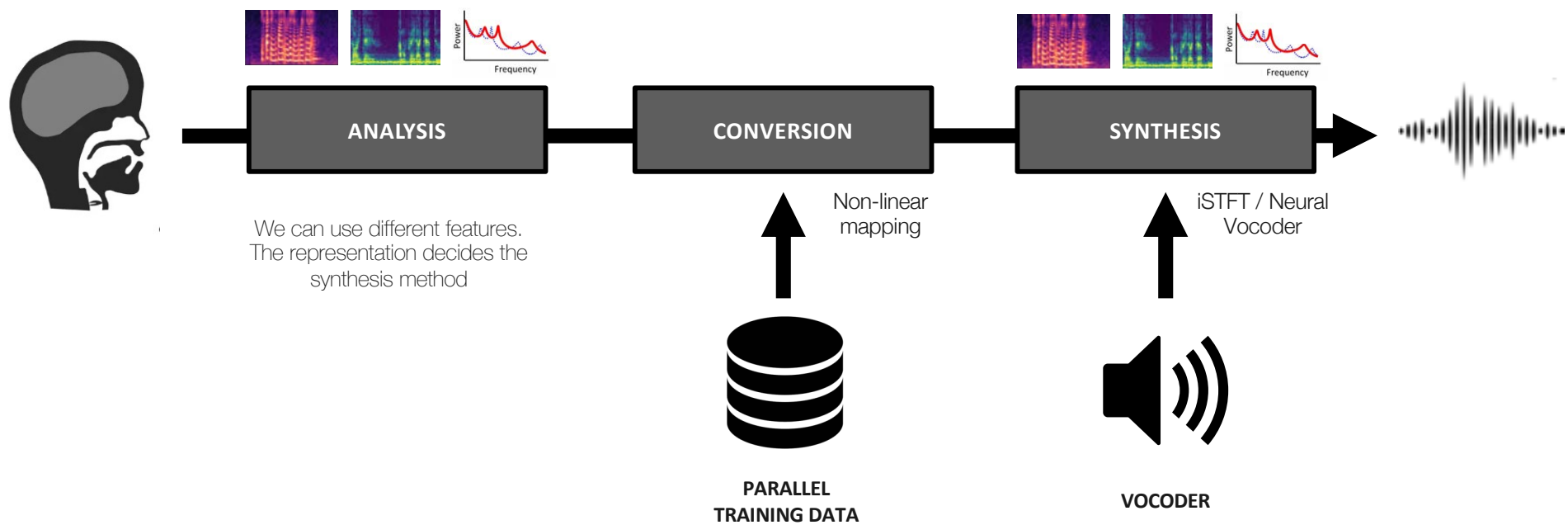
Heuristic modification, frequency domain information (phase-vocoder)



# VOICE CONVERSION



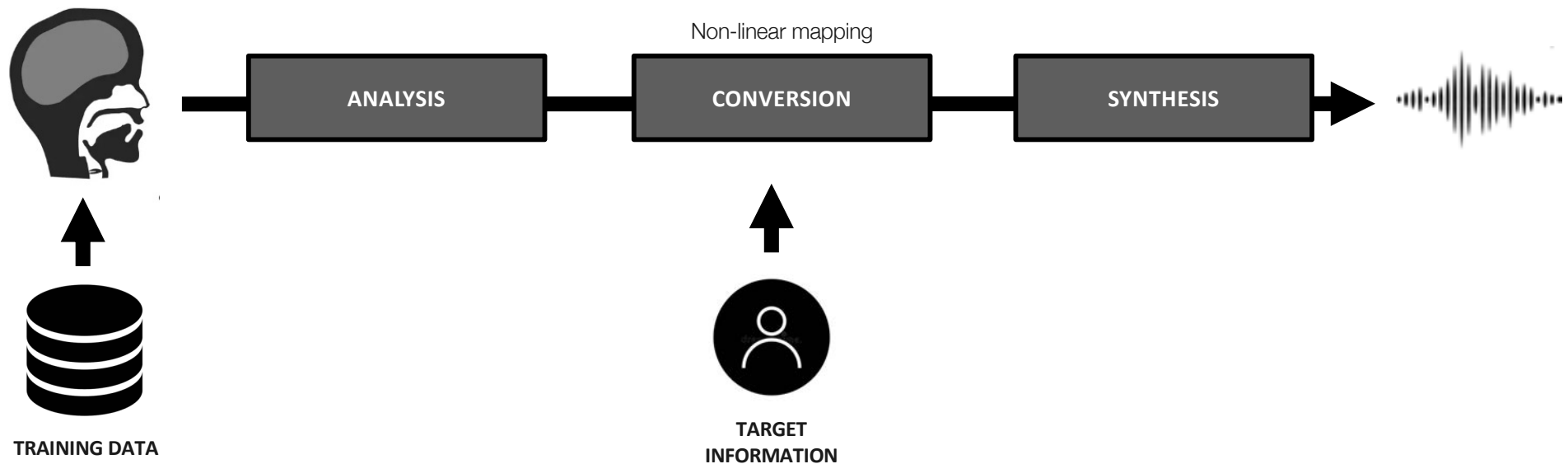
# VOICE CONVERSION



# VOICE CONVERSION (NON-PARALLEL)

**Recognition:** Extract desired information such as linguistics and eliminate unwanted source information.

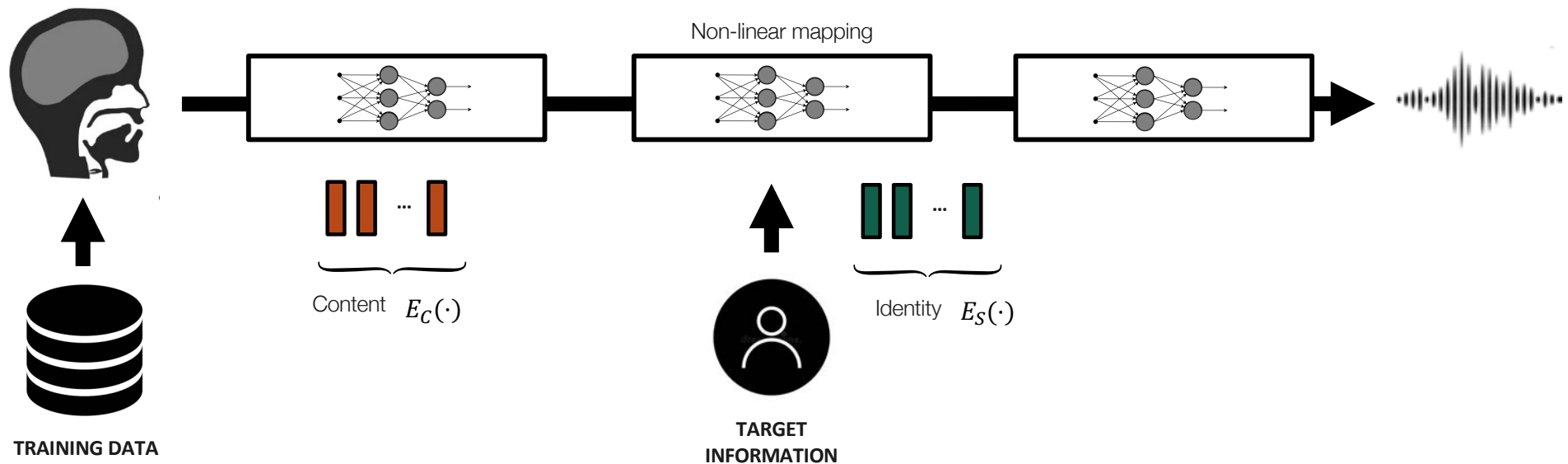
**Synthesis:** Inject and condition the generation on information from the target, such as identity.



# VOICE CONVERSION (NON-PARALLEL)

**Recognition:** Extract desired information such as linguistics and eliminate unwanted source information.

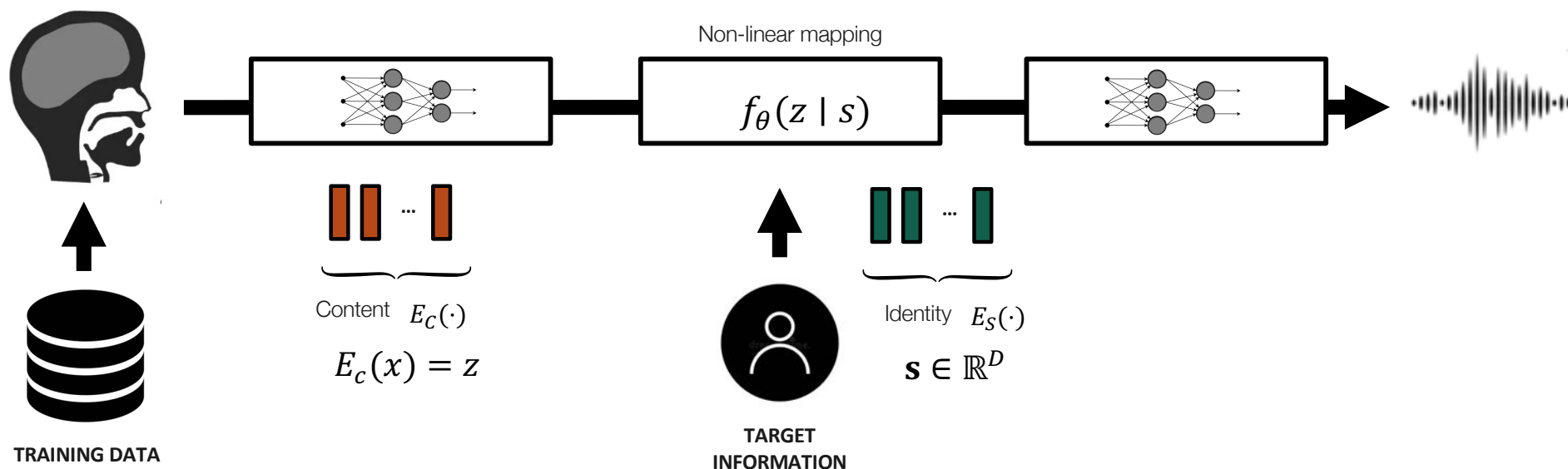
**Synthesis:** Inject and condition the generation on information from the target, such as identity.



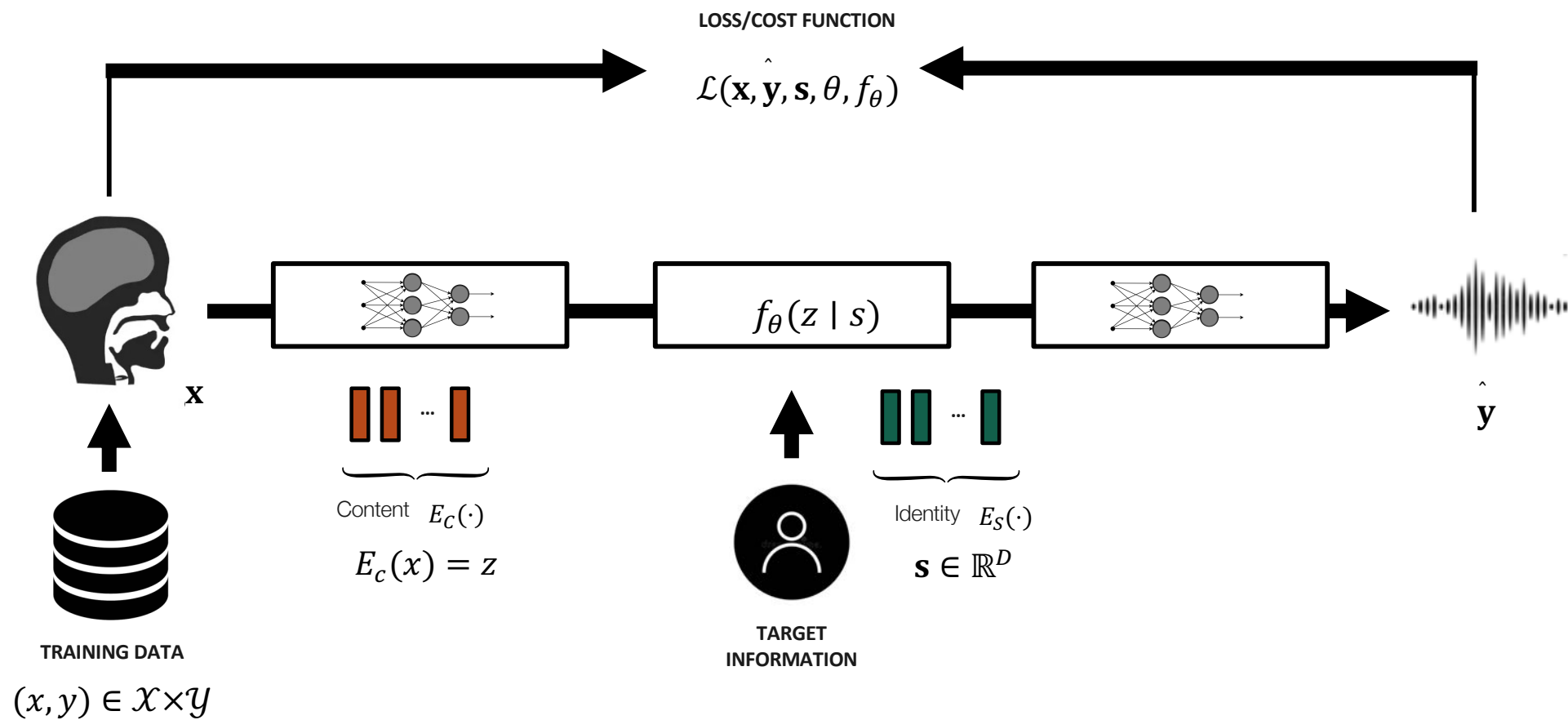
# VOICE CONVERSION (NON-PARALLEL)

**Recognition:** Extract desired information such as linguistics and eliminate unwanted source information.

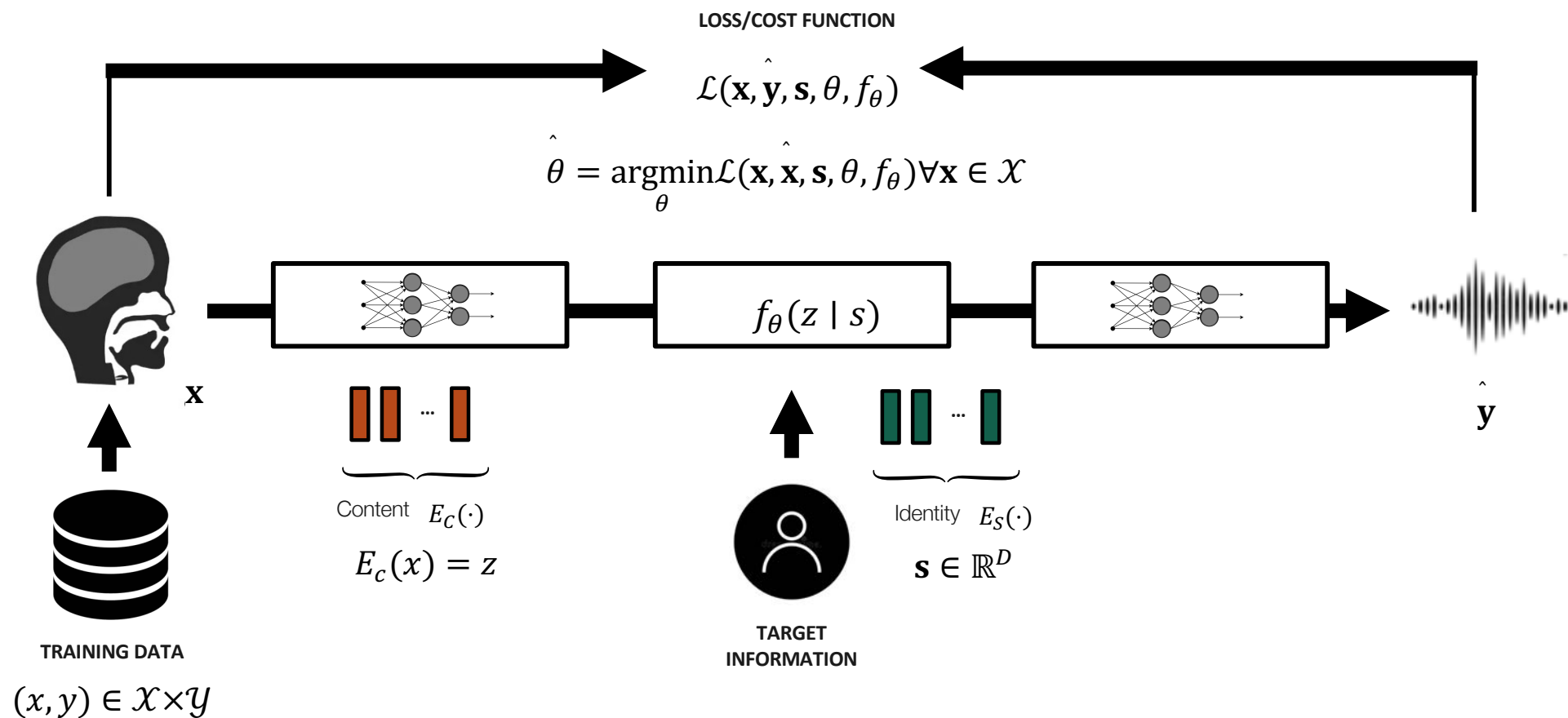
**Synthesis:** Inject and condition the generation on information from the target, such as identity.



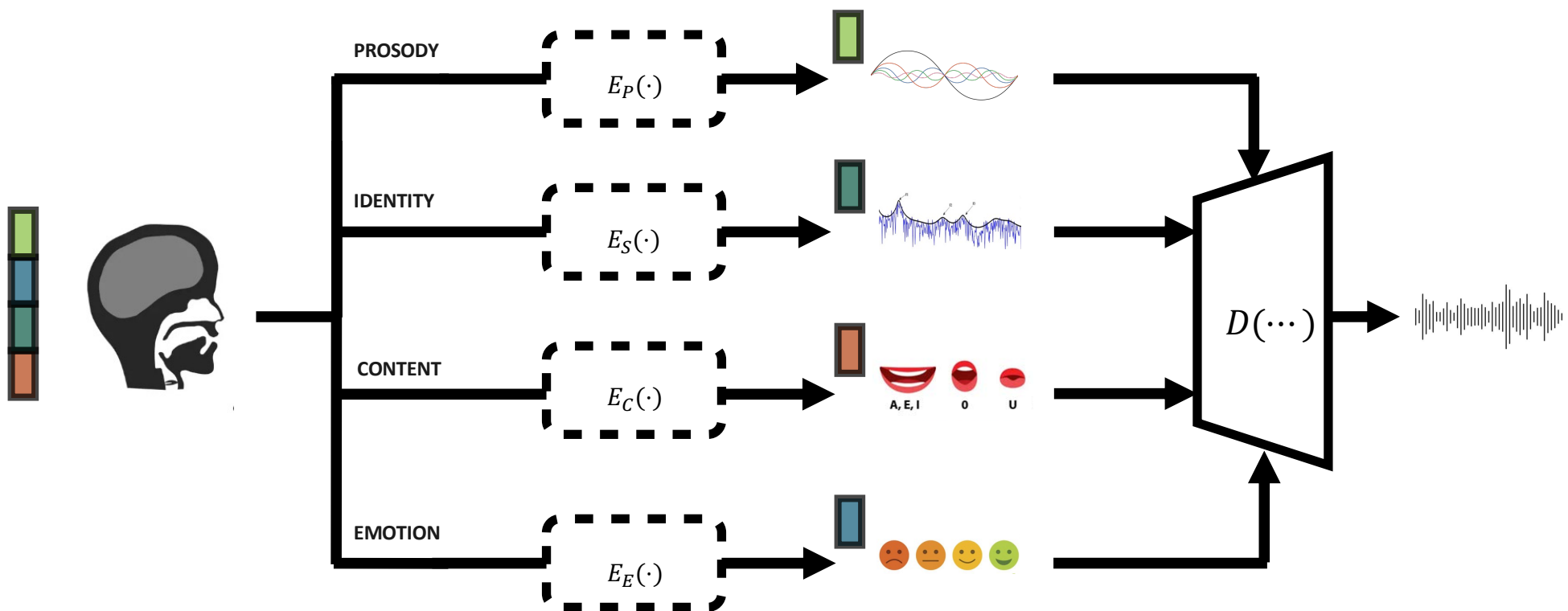
# VOICE CONVERSION (NON-PARALLEL)



# VOICE CONVERSION (NON-PARALLEL)



# VOICE CONVERSION (NON-PARALLEL)



# The CHALLENGE Project: Fighting Auditory Hallucinations by using Virtual Reality

Luis Viera, Peter Fisher and Simon Lajboschitz \*

Khora Virtual Reality

Stefania Serafin †

Multisensory Experience Lab, Aalborg University Copenhagen

Merete Nordentoft ‡

Psykiatrisk Center København

The screenshot shows a virtual reality interface for customizing a character's voice. On the left is a 3D model of a woman's face. On the right is a control panel with a navigation bar at the top containing tabs: Main, Basic, Face, Features, Therapist, Voice (selected), and Session. Below the navigation bar is the title "Voice characteristics". There are six circular sliders for adjusting voice parameters: Gender (Feminine to Masculine), Age (Young to Old), Pitch (Low to High), Size (Small to Large), Presence (Min to Max), and Tremble (Min to Max). The Tremble slider is currently set to a high value. At the bottom of the control panel are two buttons: "Play" and "Show key", flanked by left and right arrow navigation symbols.

# AVATAR THERAPY

0 1

## Voice Manipulation

AI-driven voice synthesis to make the avatar voice sound more natural

TTS / Voice Cloning

Naturalness & Prosody

0 2

## Controllability

The voice should be controllable, natural and match human descriptors.

Controllability

Personification

0 3

## Technical Specs

Latency < 100 ms, with the ability to run the system on a consumer CPU (Unity).

Latency

Deployment

[C. J. Edwards, et al., "Thevoicecharacterisation checklist: psychometric properties of a brief clinical assessment of voices as socialagents," Frontiers in Psychiatry, vol. Volume 14, 2023.]

Main

Basic

Face

Features

**Therapist**

Voice

Session

## Therapist voice

Please select your gender below and record your voice.

Therapist gender ▾

Text language ▾

● Record

Play

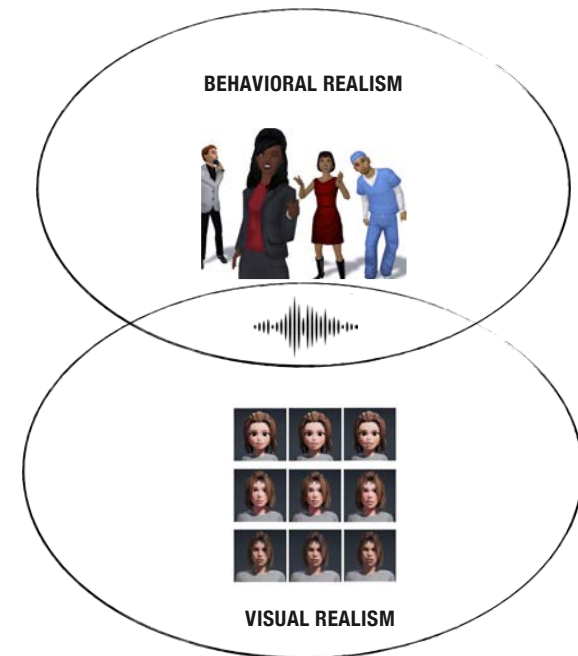


Show key



# MOTIVATION

- A cartoonish avatar can maintain a **strong social presence** if it exhibits human-like behaviour and voice.
- **Synthetically created** voices are perceived as **less appealing** when paired with photorealistic avatars.
- Virtual characters with human voices are generally considered more **understandable** and **expressive** than their **synthetic** counterparts.



Zegaran et. Al, "Does avatar design in educational games promote a positive emotional experience among learners?" *SAGE E-Learning and Digital Media*, 2021.

Higgins et. Al, "Sympathy for the digital: Influence of synthetic voice on affinity, social presence and empathy for photorealistic virtual humans", *Computers & Graphics*, 2022.

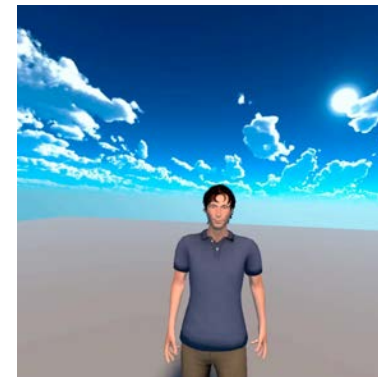
[Cabral et. Al, "The Influence of Synthetic Voice on the Evaluation of a Virtual Character", *ISCA Interspeech*, 2017.

# EXPERIMENT

**1.** Realism, **2.** Social Presence, **3.** Emotional, **4.** Likeability



21 Participants (within subjects)

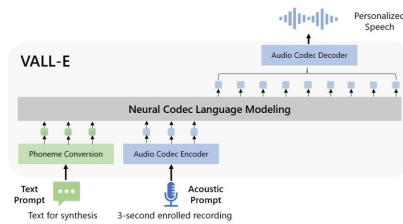


# EXPERIMENT

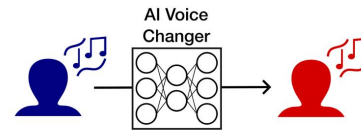
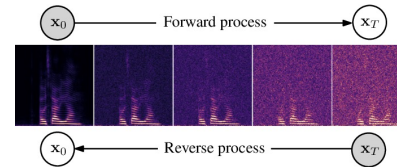
a) Coqui (TTS)



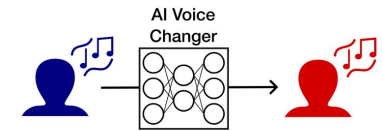
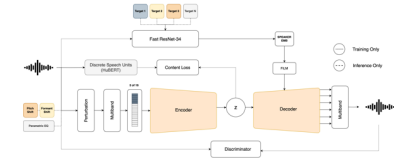
b) VALL-E (TTS)



c) Diff-VC (STS)



d) S-RAVE (STS)

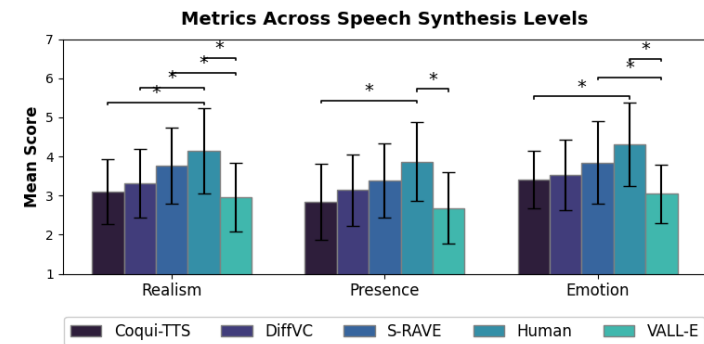
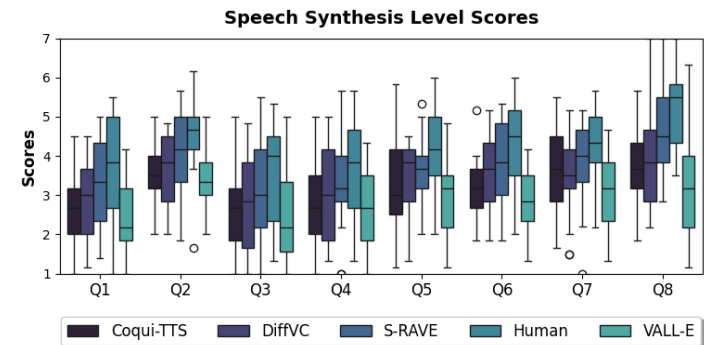


# EXPERIMENT

Category	Item	Question	ID
Realism	Overall realism	<i>"I found the person realistic overall"</i>	Q1
	Appearance	<i>"The appearance of the person matched the voice"</i>	Q2
Social Presence	General Presence	<i>"It feels like there was a presence of another person in the room with me"</i>	Q3
	Believableability	<i>"The person did not seem to be alive"</i>	Q4
Empathy	Emotion	<i>"I did not feel that the person's emotions seemed genuine"</i>	Q5
	Alignment	<i>"The emotions of the person matched the voice"</i>	Q6
	Expressivity	<i>"I felt the person was expressive"</i>	Q7
Voice Trait	Likeability	<i>"I disliked the voice of the person"</i>	Q8

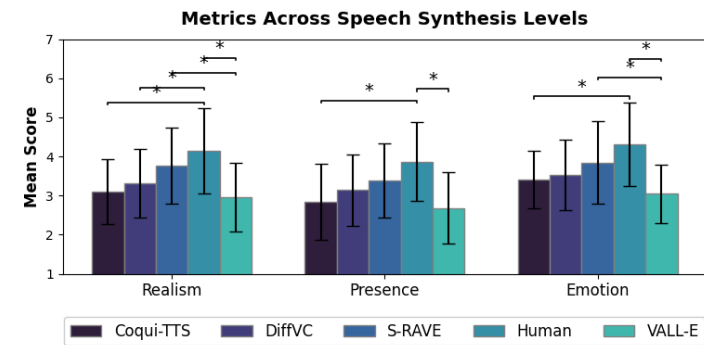
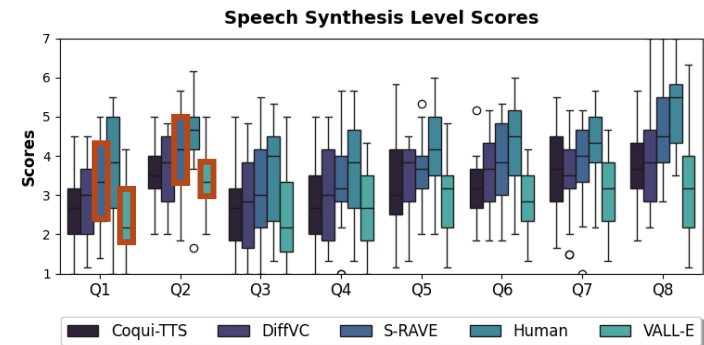
# RESULTS

- None of the avatars were deemed realistic or present.
- Avatars using **STS-based** speech approached the human baseline on all metrics.
- STS methods outperform TTS systems in conveying **emotional subtleties**.
- **TTS systems have emotional limitations:** lack of intonation, pitch and natural phonetic coding within text-based-tokens.



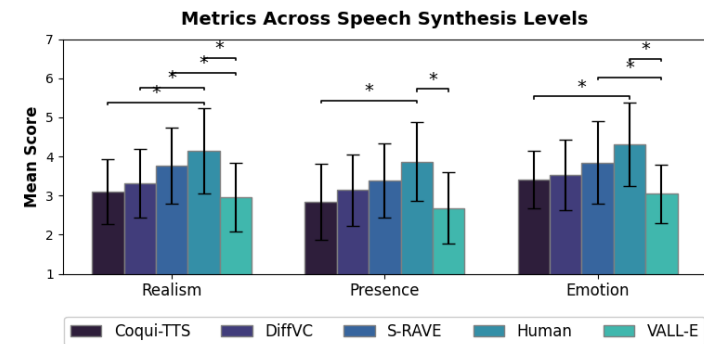
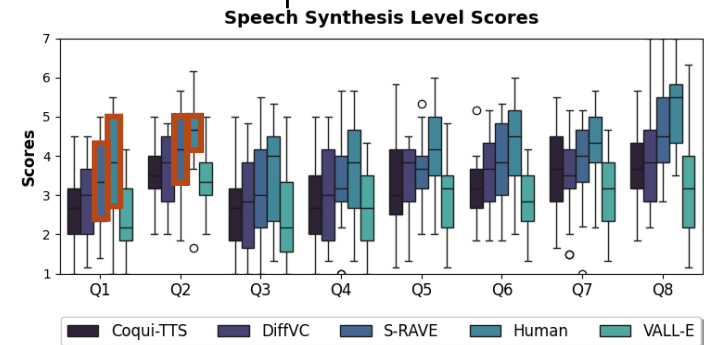
# RESULTS

- None of the avatars were deemed realistic or present.
- Avatars using **STS-based** speech approached the human baseline on all metrics.
- STS methods outperform TTS systems in conveying **emotional subtleties**.
- **TTS systems have emotional limitations:** lack of intonation, pitch and natural phonetic coding within text-based-tokens.



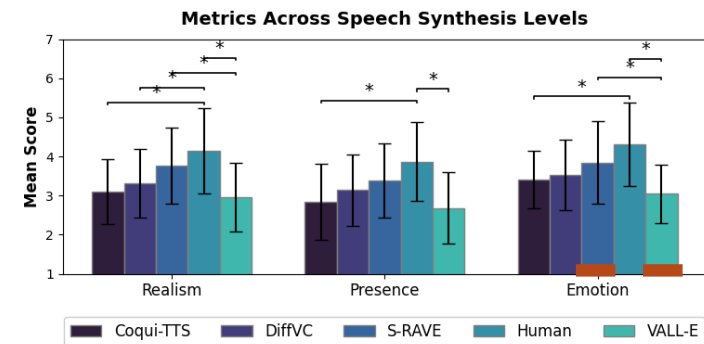
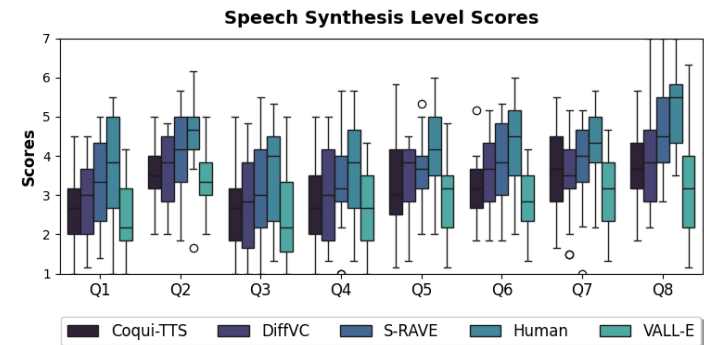
# RESULTS

- None of the avatars were deemed realistic or present.
- Avatars using **STS-based** speech approached the human baseline on all metrics.
- STS methods outperform TTS systems in conveying **emotional subtleties**.
- **TTS systems have emotional limitations:** lack of intonation, pitch and natural phonetic coding within text-based-tokens.



# RESULTS

- None of the avatars were deemed realistic or present.
- Avatars using **STS-based** speech approached the human baseline on all metrics.
- STS methods outperform TTS systems in conveying **emotional subtleties**.
- **TTS systems have emotional limitations:** lack of intonation, pitch and natural phonetic coding within text-based-tokens.



# UNIVERSAL DESIGN



- Ronald Mace “The design of products and environments to be usable by all people, to the greatest extent possible, without the need for adaptation or specialized design.”

# PRINCIPLES OF UNIVERSAL DESIGN



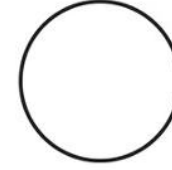
**1. Equitable Use**



**2. Flexibility in Use**



**3. Simple and Intuitive Use**



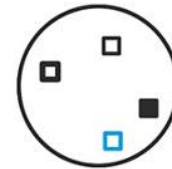
**4. Perceptible Information**



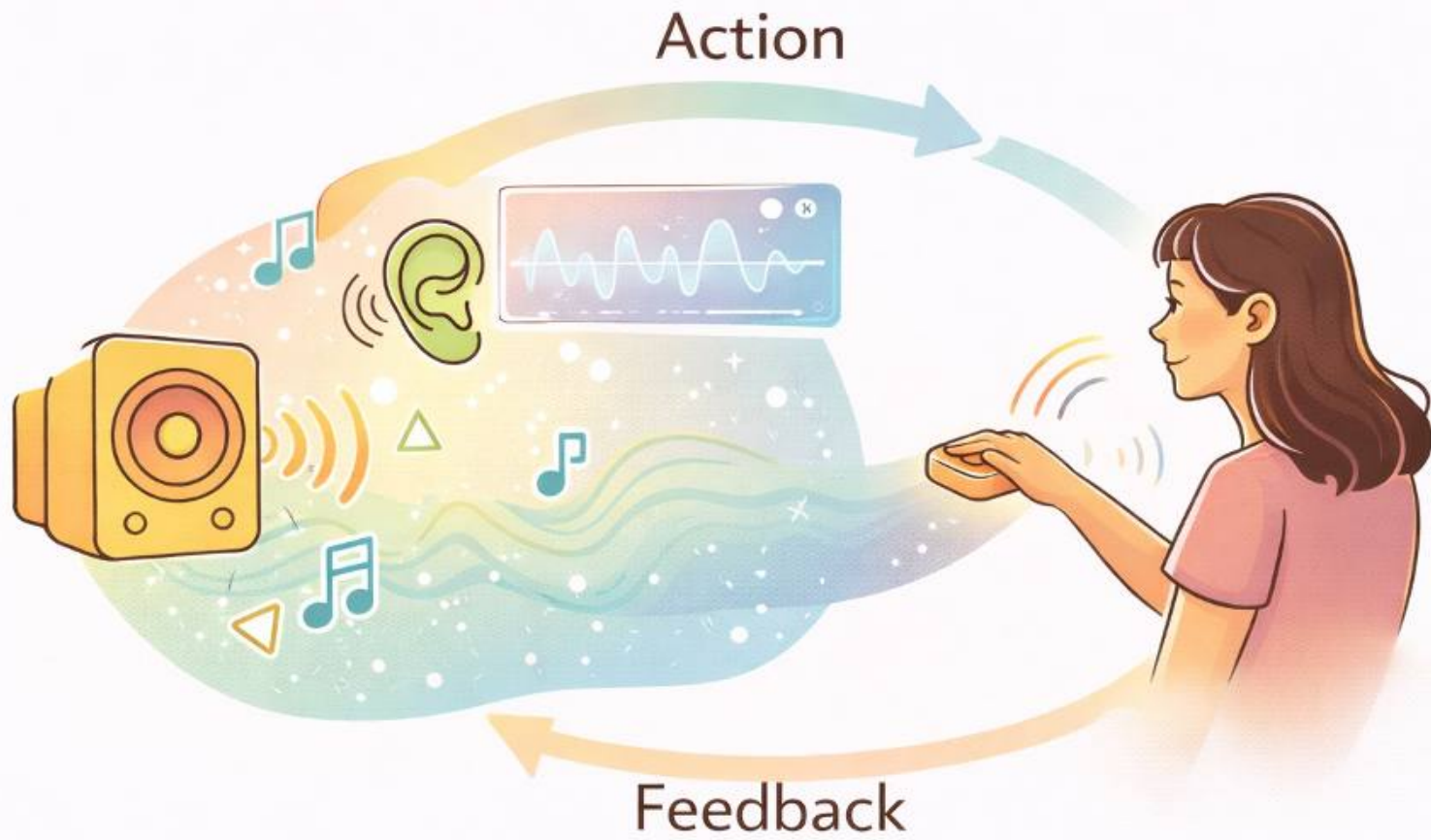
**5. Tolerance for Error**



**6. Low Physical Effort**



**7. Size and Space for Approach and Use**



Sonic Interaction Design

from Sonic Interaction Design



to Sonic Universal Design

# 7 PRINCIPLES OF SONIC UNIVERSAL DESIGN

*Designing sound experiences that are usable, comfortable, and meaningful for all, across diverse hearing abilities, contexts, and technologies.*



## 1. Equitable Use

Sound experiences are usable and meaningful for people with a wide range of hearing abilities and technologies.

**Key idea:**

No one is excluded from the sonic experience.



## 2. Flexibility in Use

Sound can be experienced, customized, and interacted with in multiple ways to suit individual preferences and needs.

**Key idea:**

Provide choice, control, and adaptability.



## 3. Simple and Intuitive Use

Sonic interactions are easy to understand, learn, and use, regardless of experience, language, or hearing ability.

**Key idea:**

Make sound interactions clear and predictable.



## 4. Perceptible Information

Important information is conveyed through sound in ways that are detectable and understandable across diverse listening conditions.

**Key idea:**

Communicate clearly through multiple sonic cues.



## 5. Tolerance for Error

Sonic interactions are forgiving and support recovery from mistakes, misunderstandings, or missed cues.

**Key idea:**

Reduce stress and support confident interaction.

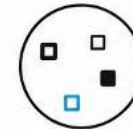


## 6. Low Physical Effort

Sound interactions require minimal physical, cognitive, and listening effort to use and perceive.

**Key idea:**

Make listening and interacting effortless and comfortable.



## 7. Size and Space for Approach and Use

Sonic interfaces scale across devices, environments, and situations, supporting orientation and access for all users.

**Key idea:**

Design for diverse bodies, spaces, and situations.

# SID versus SUD



## Past / Present

Sonic Interaction Design

Sound as interaction modality

Designing interactions with sound

Adding sound to interfaces

Generic users

## Future

Sonic Universal Design

Sound as inclusive infrastructure

Designing for perceptual diversity

Rethinking interfaces through sound

Diverse users

# THANKS!

## Stefania Serafin

[stefse@dtu.dk](mailto:stefse@dtu.dk)

<https://www.linkedin.com/in/stefaniaserafin/>

Copenhagen Hearing  
and Balance Centre

