

Parametric Spatial Audio Compression



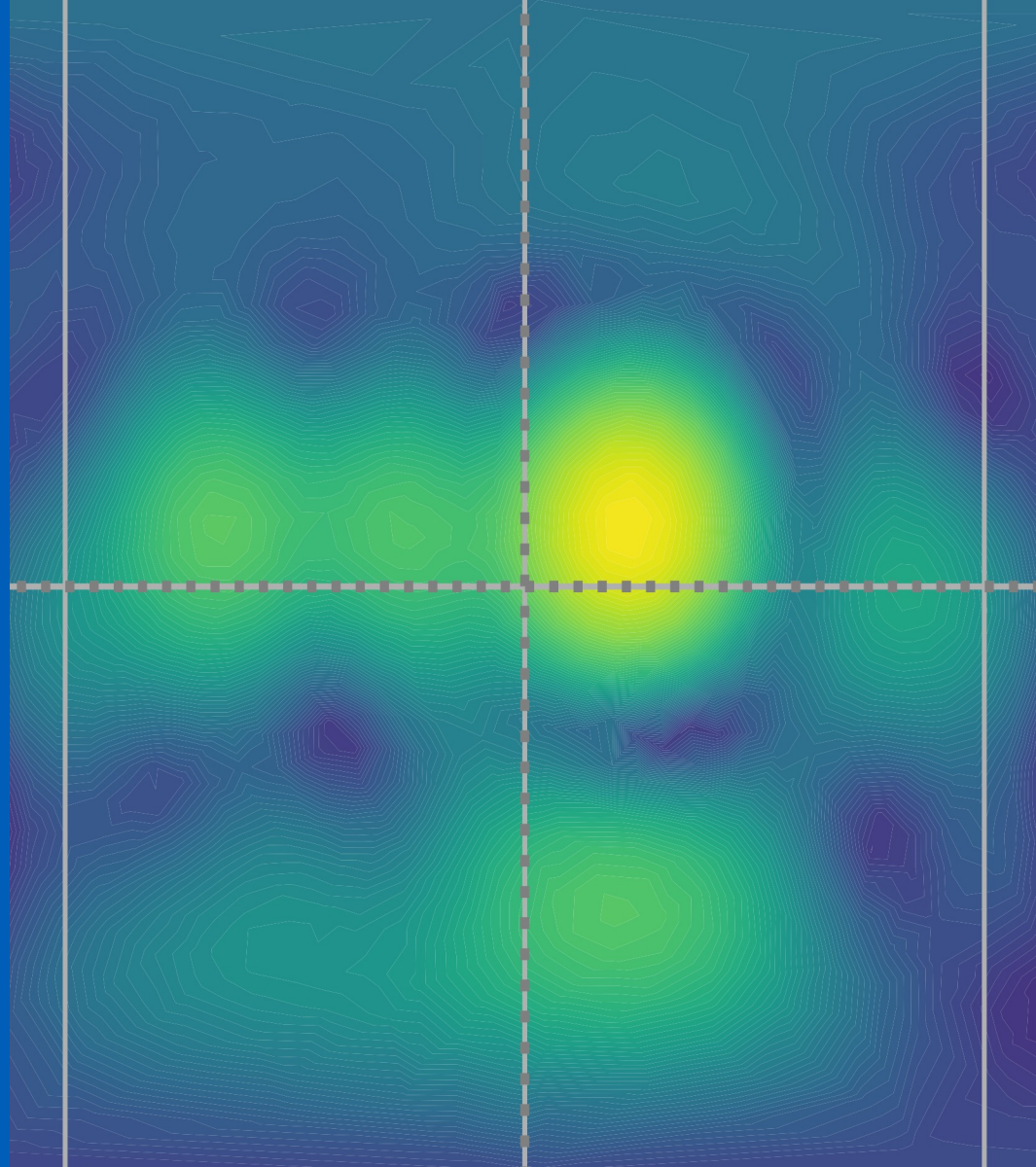
Using Higher-Order- Directional Audio Coding

Christoph Hold

25.03.2024



Aalto-yliopisto
Aalto-universitetet
Aalto University



Outline



Introduce Parametric Spatial Audio



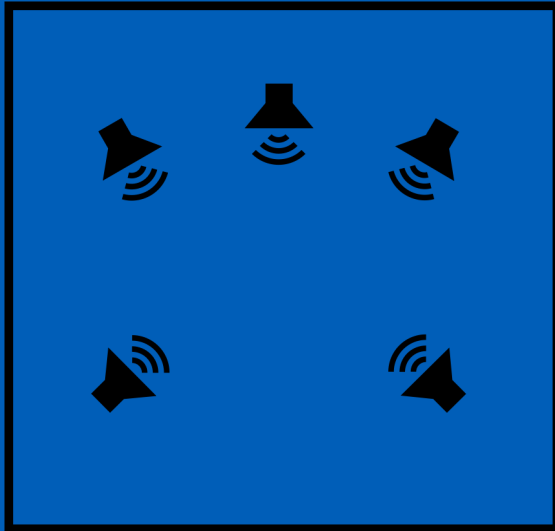
Introduce a Spatial Audio Codec



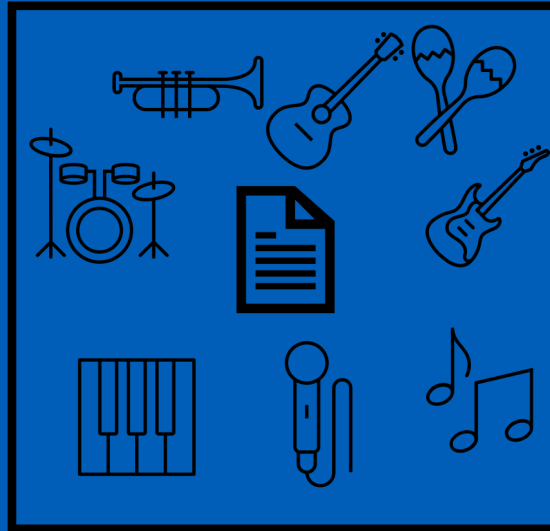
Evaluation

Spatial Audio Formats

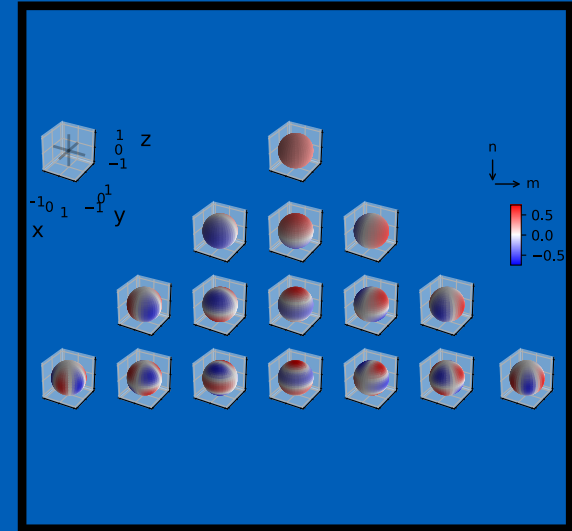
Channel



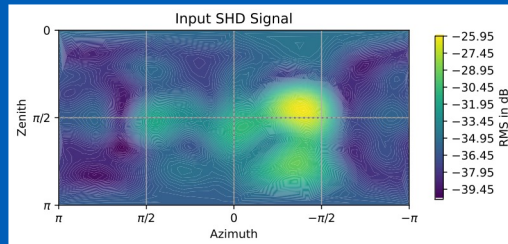
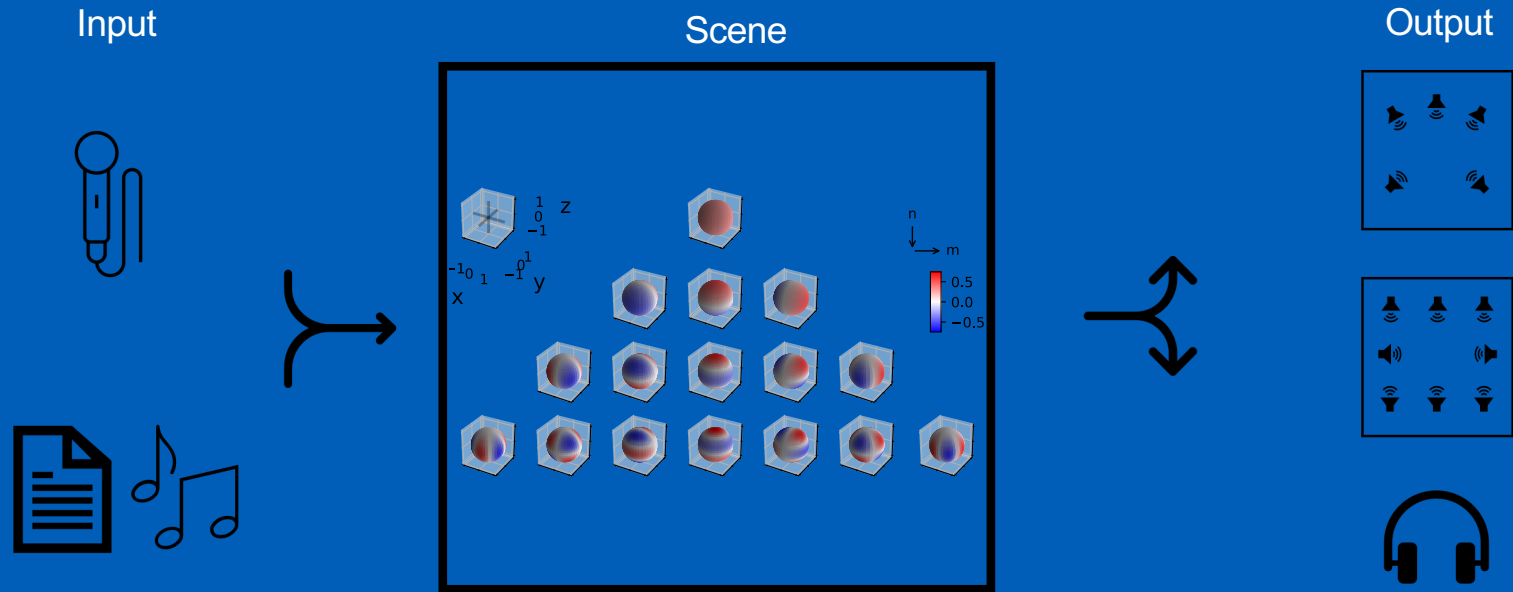
Object



Scene

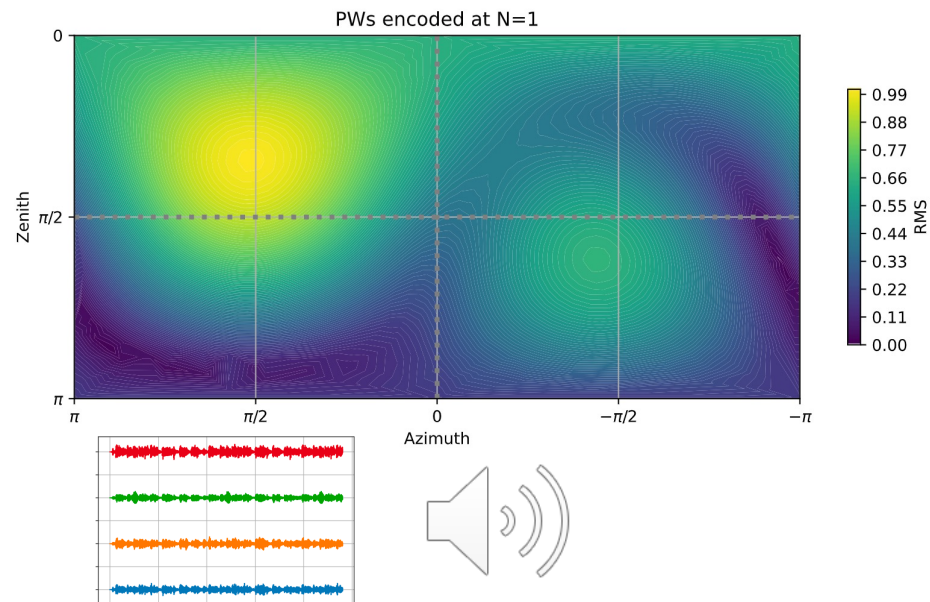


Spatial Audio - Ambisonics

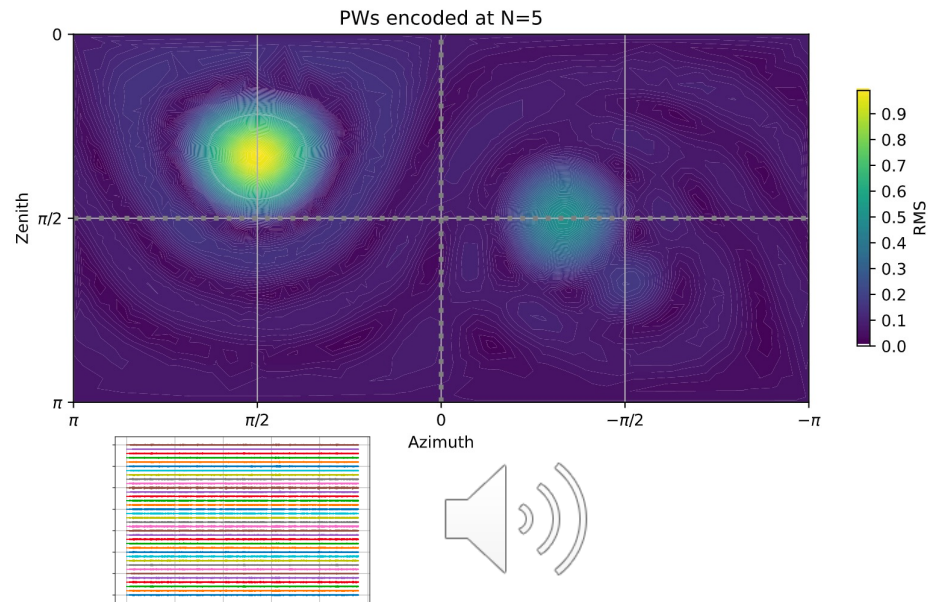


First Order vs. Higher Order

- 4 Audio channels

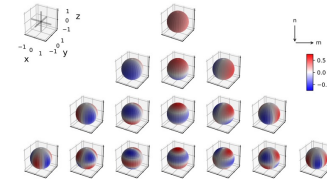
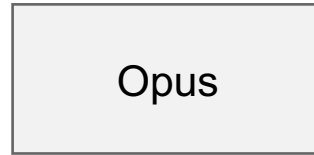
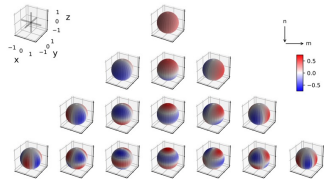


- 36 Audio channels

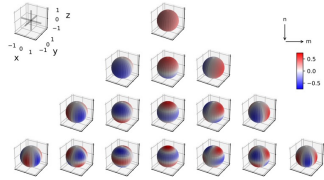


Opus Ambisonics Codec

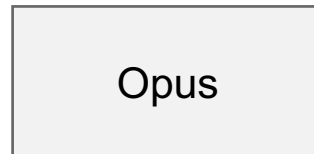
Opus CMF 2



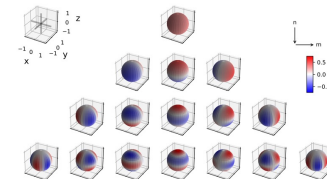
Opus CMF 3



A



B



Transmitting Higher-Order-Ambisonics

Input (MagLS5)




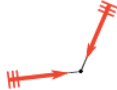

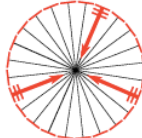

Opus @ 768kbit/s (MagLS5)



Parametric Spatial Audio

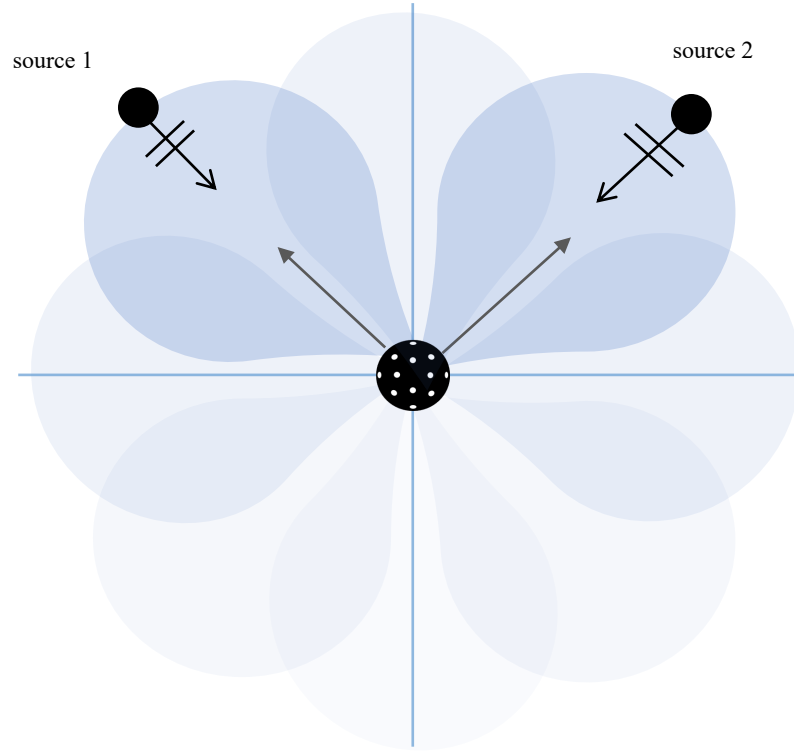
Extract and utilize additional information from the input signals.

Spatial Parameterization Models

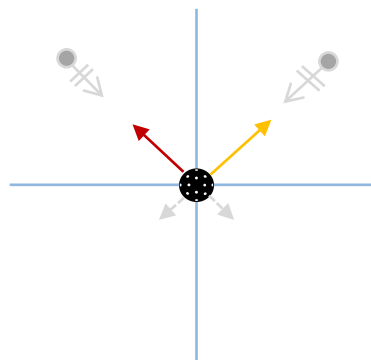
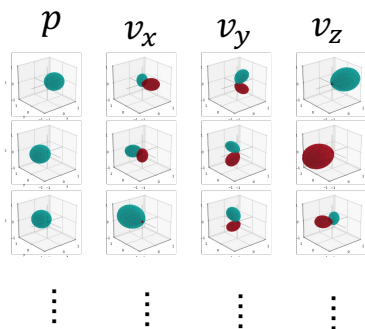
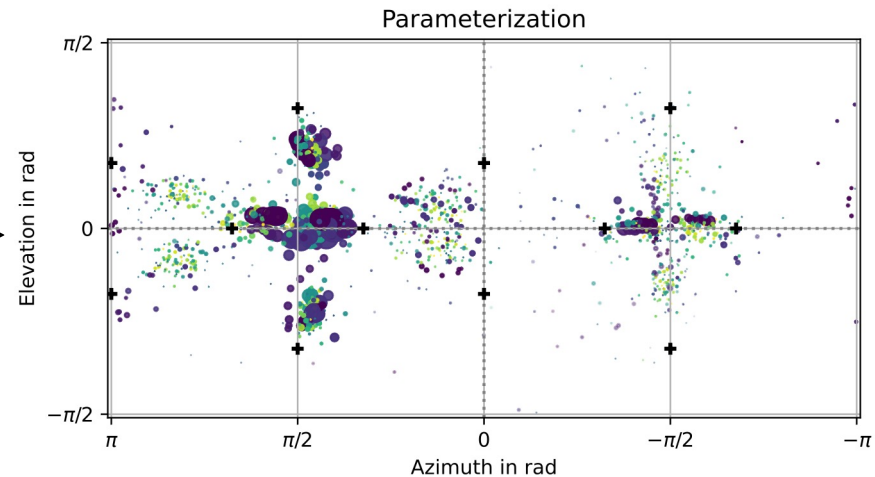
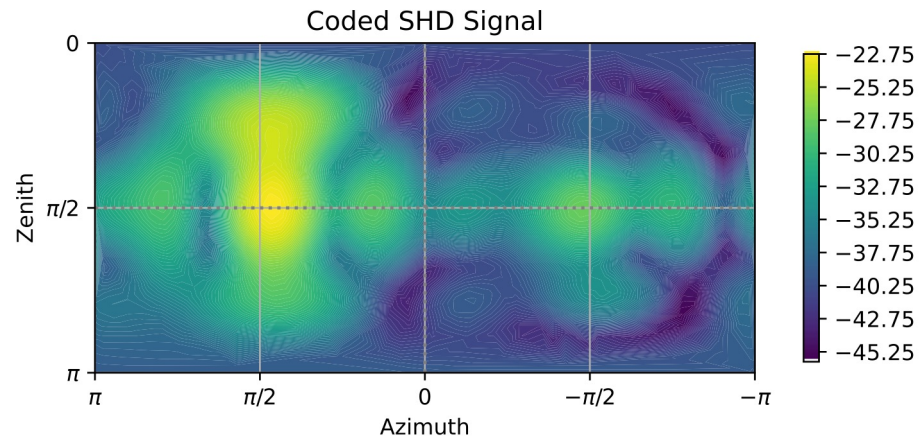
Method	Input	Model	
DirAC (Pulkki, 2006)	FOA (4ch)		1 source component + 1 iso. diffuse component
HARPEX (Berge, 2010)	FOA (4ch)		2 source components
HO-DirAC (Politis et.al., 2015)	HOA (9+ch)		~M sector source + ~M sector diffuse components
Sparse Recovery (Wabnitz, Jin, 2012)	FOA/HOA (4+ch)		$\leq M/2$ source components
COMPASS (Politis et.al., 2018)	FOA/HOA (4+ch)		$\leq M/2$ source components + spatial ambient component

Higher-Order Directional Audio Coding (HODirAC)

Parametric Spatial Audio - HODirAC



Parameter Estimation

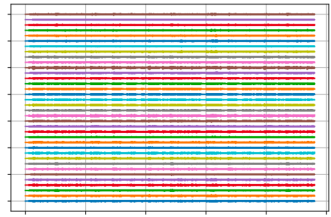


$$i_\xi \propto \Re\{p_\xi^H v_\xi\}$$

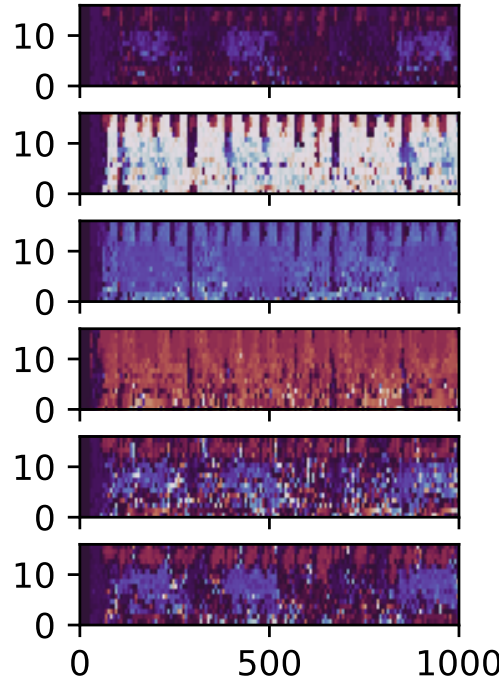
$$\Omega_\xi^{\text{DoA}} = \angle i_\xi$$

$$\psi_\xi = 1 - \frac{\|i_\xi\|}{E_\xi} = 1 - \frac{2\|i_\xi\|}{|p_\xi|^2 + v_\xi^H v_\xi}$$

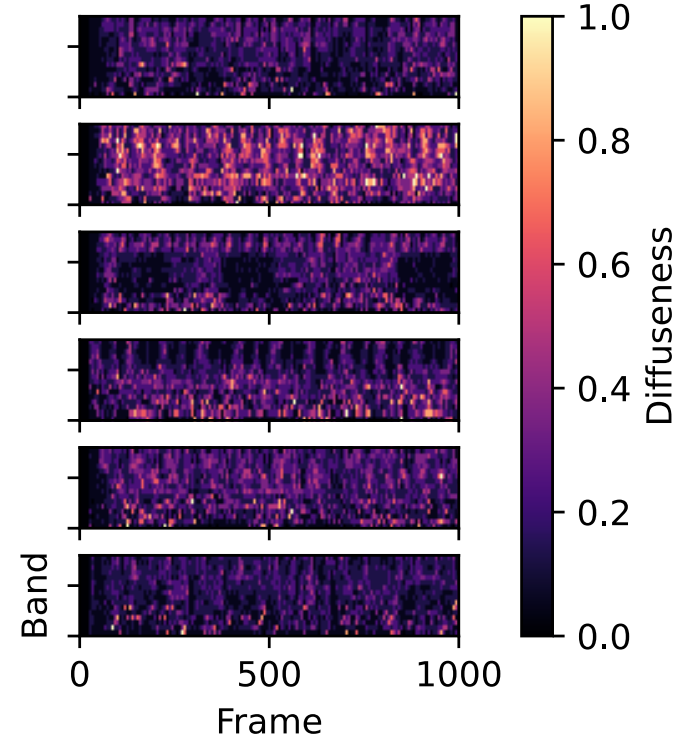
Parameter Estimation - Stream



HOA χ



Azimuth

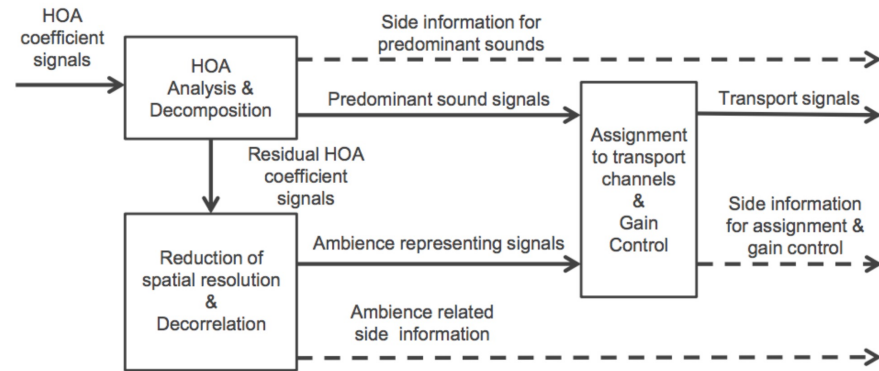


Making a codec

State of the Art: MPEG-H 3D Audio

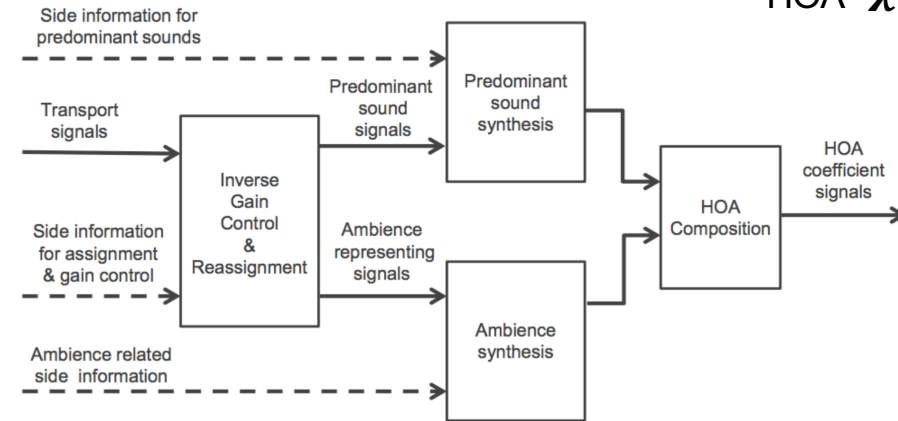
- designed to provide universal coding of channel-based, object-based and Higher Order Ambisonics input

HOA χ



HOA - Encoder

HOA $\tilde{\chi}$



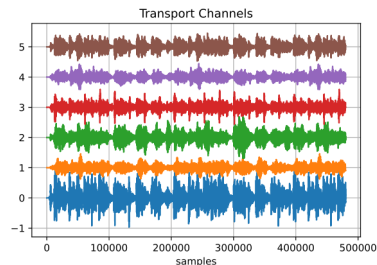
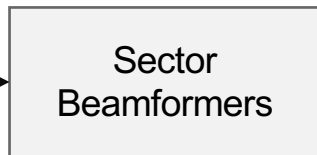
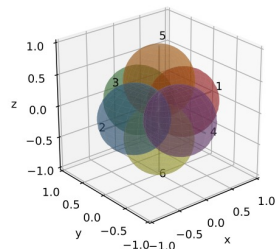
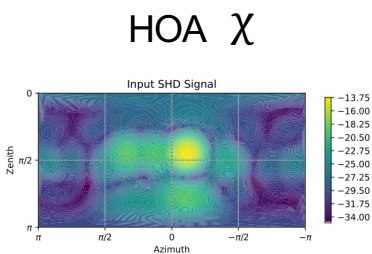
Decoder

Sen, D., Peters, N., Kim, M. Y., & Morrell, M. (2016). Efficient compression and transportation of scene based audio for television broadcast. *Proceedings of the AES International Conference, 2016-July*.

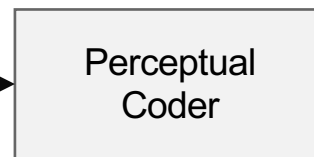
Herre, J., Hilpert, J., Kuntz, A., & Plogsties, J. (2015). MPEG-H 3D Audio—The New Standard for Coding of Immersive Spatial Audio. *IEEE Journal of Selected Topics in Signal Processing*.

Encoder

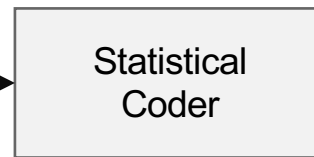
Encoder



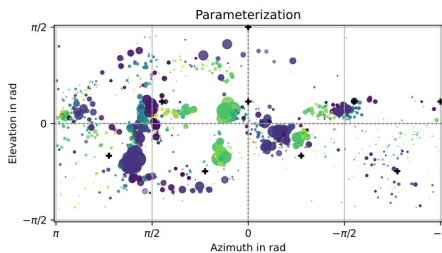
Audio Signals



Parameters

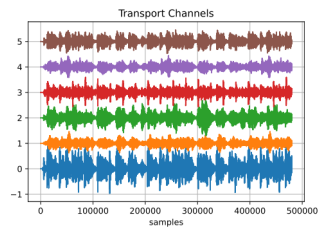


*.hoac



Decoder

Decoder



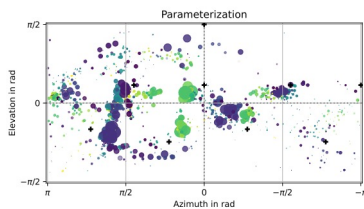
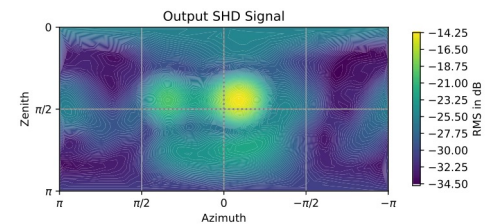
$\left(\begin{array}{l} \text{*} \\ \text{.hoac} \end{array} \right)$

Audio \mathbf{X}

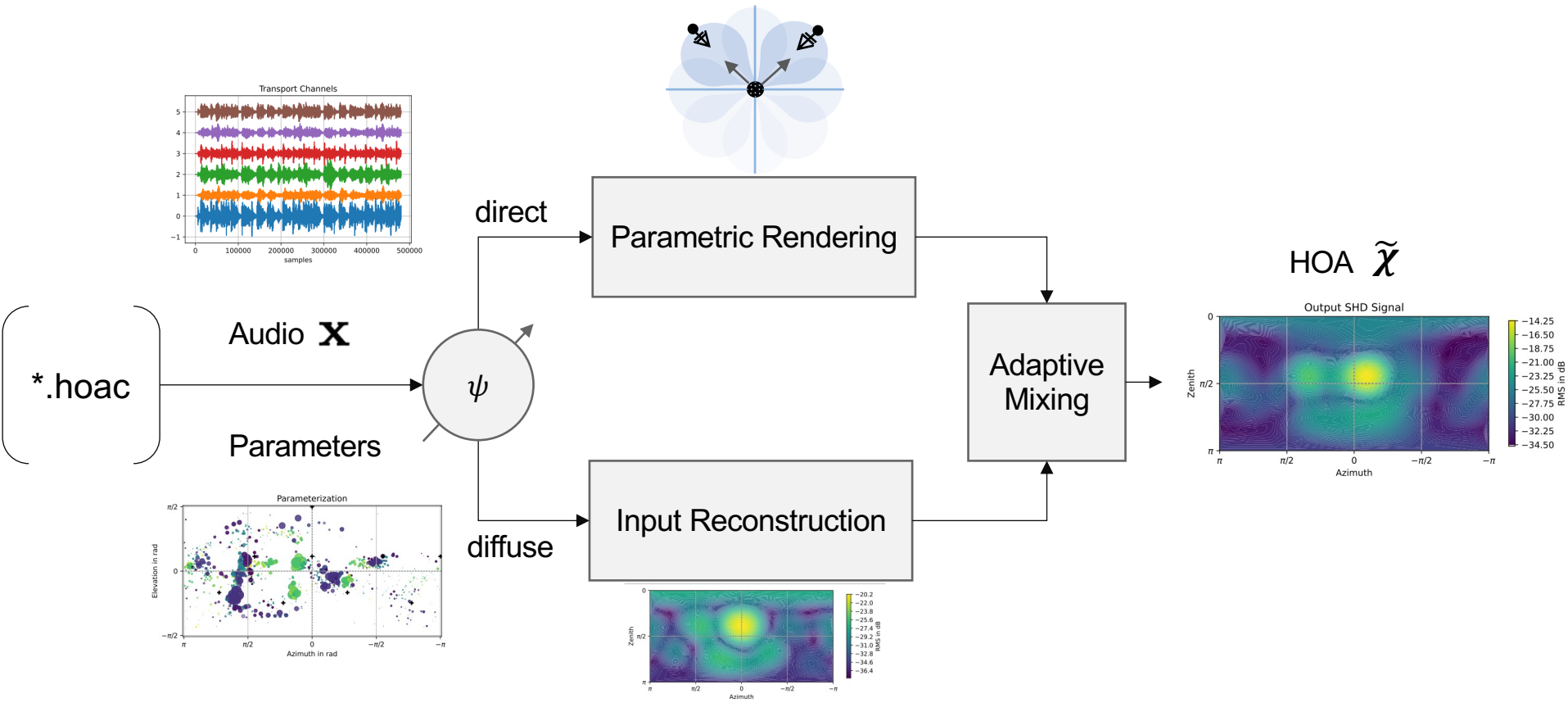
Parameters

Spatial
Reconstruction

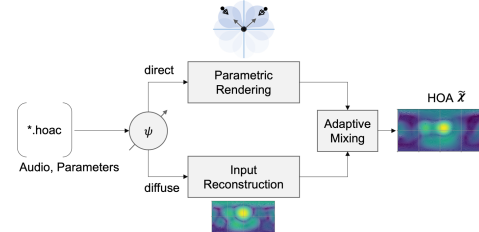
HOA $\tilde{\mathbf{X}}$



Decoder - Spatial Reconstruction



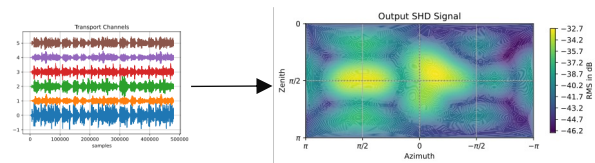
Recover and Resynthesize



- Reconstruction (low orders) :

$$\text{SFB } \mathbf{x} = \mathbf{A}\boldsymbol{\chi}, \text{ and } \tilde{\boldsymbol{\chi}} = \mathbf{B}\mathbf{x} \rightarrow$$

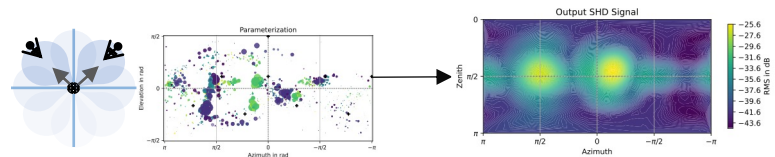
$$\mathbf{B} = [\mathbf{A}]^\dagger$$



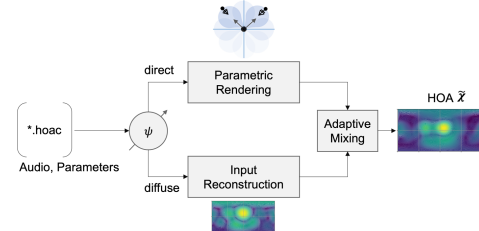
- Resynthesis (high orders) :

$$\hat{\boldsymbol{\chi}} = \underbrace{\mathbf{B} \text{diag}(\boldsymbol{\psi}_s)}_{\text{diffuse}} \mathbf{x} + \underbrace{\beta_A \mathbf{Y}(\boldsymbol{\Omega}_s) \text{diag}(1 - \boldsymbol{\psi}_s)}_{\text{directional}} \mathbf{x}$$

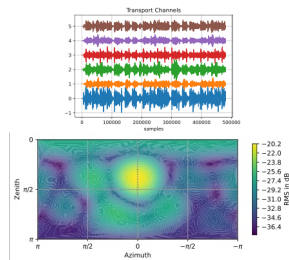
$$\hat{\boldsymbol{\chi}} = \mathbf{Q}\mathbf{x}$$



Mix and Match



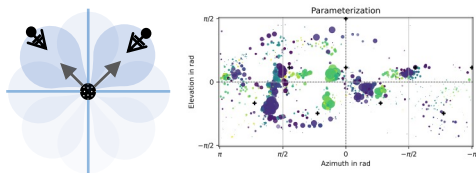
Measurements



$$\mathbf{C}_x = \mathcal{E}\{\mathbf{x}\mathbf{x}^H\}$$

$$\tilde{\mathbf{C}}_\chi = \mathcal{E}\{\tilde{\chi}\tilde{\chi}^H\} = \mathbf{B}\mathbf{C}_x\mathbf{B}^H$$

Target + Model



$$\mathbf{C}_\chi = \mathbf{C}_{\text{dir}} + \mathbf{C}_{\text{dif}}$$

$$\mathcal{E}\{\mathbf{Q}\mathbf{x}(\mathbf{Q}\mathbf{x})^H\} = \mathbf{Q}\mathbf{C}_x\mathbf{Q}^H$$

Solution

$$\tilde{\chi} = \mathbf{M}\mathbf{x} + \mathbf{r}$$

Evaluation

Input (MagLS5)



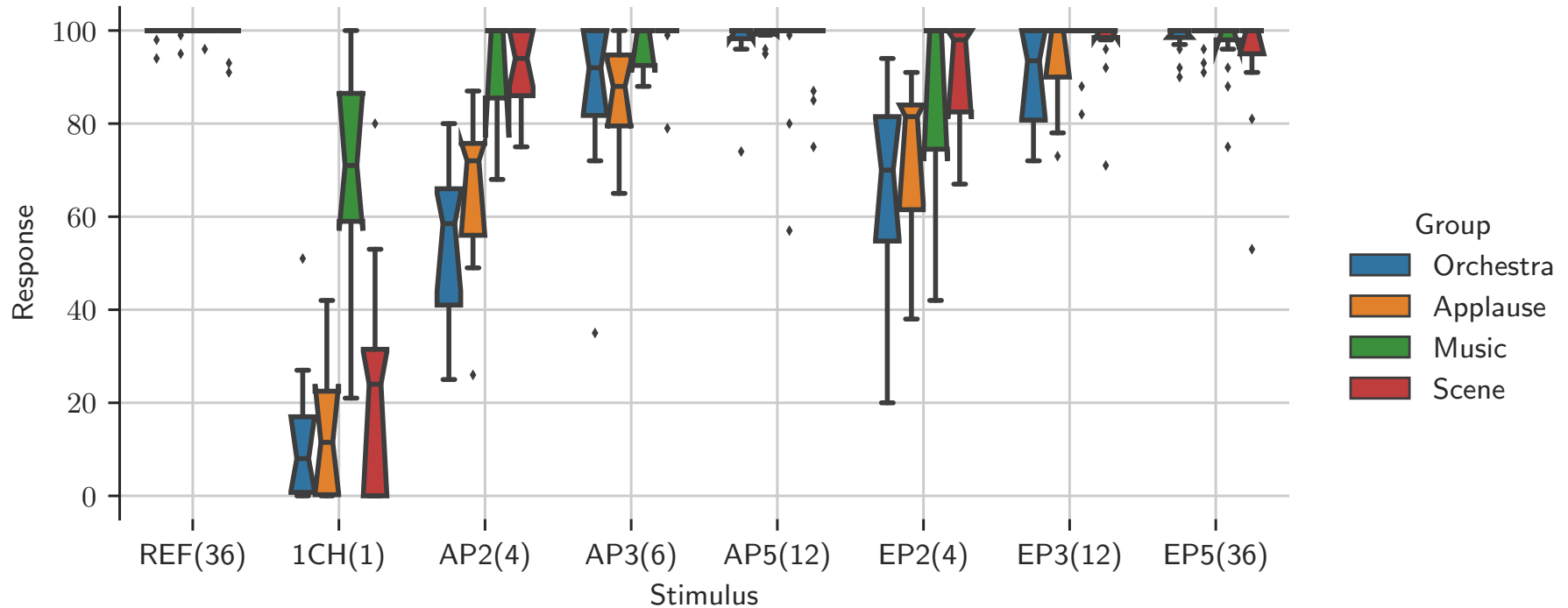
Opus @ 768kbit/s (MagLS5)



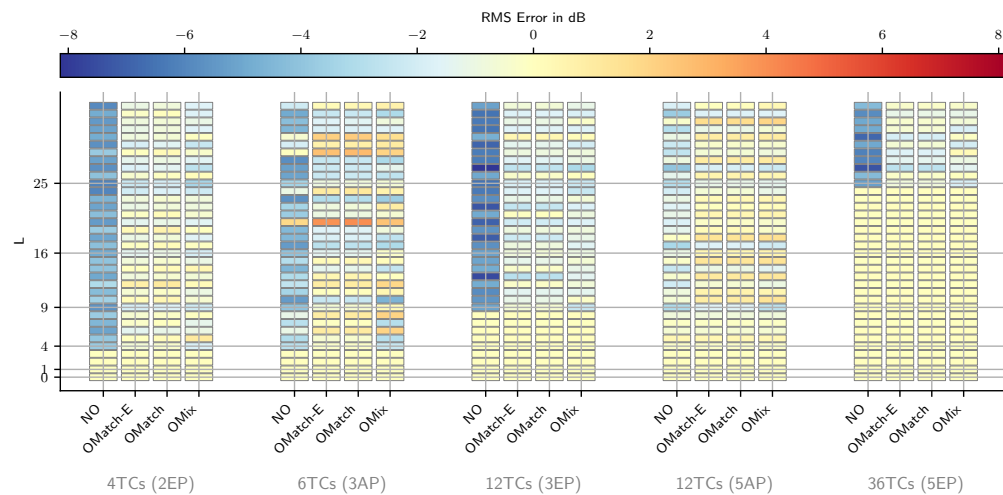
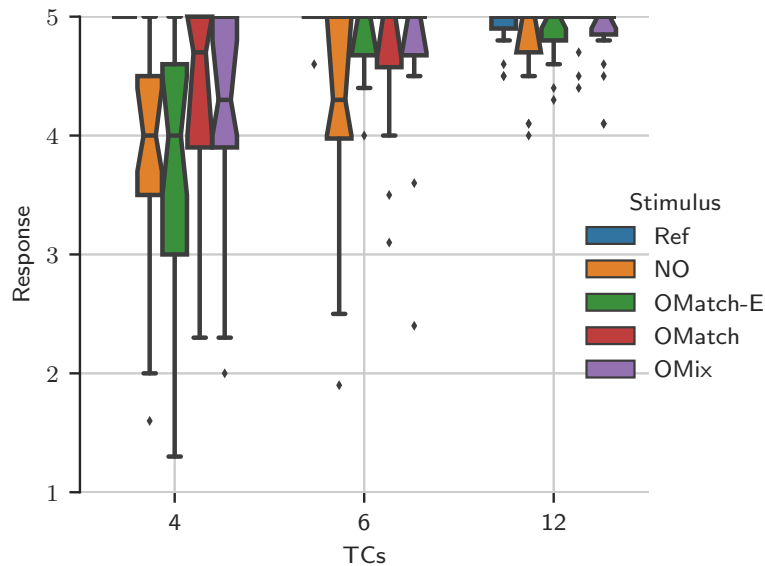
Hoac @ “Opusbitrate” (MagLS5)



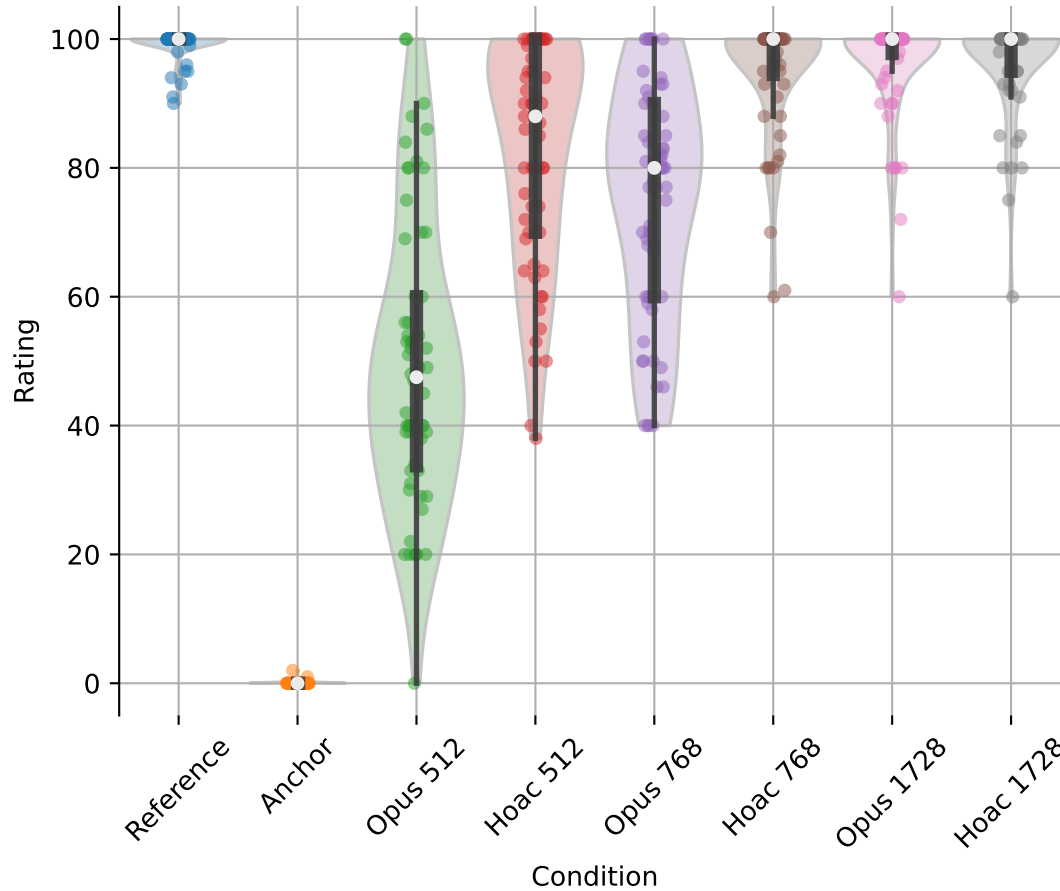
Influence of Transport Channels



Influence of Adaptive Mixing

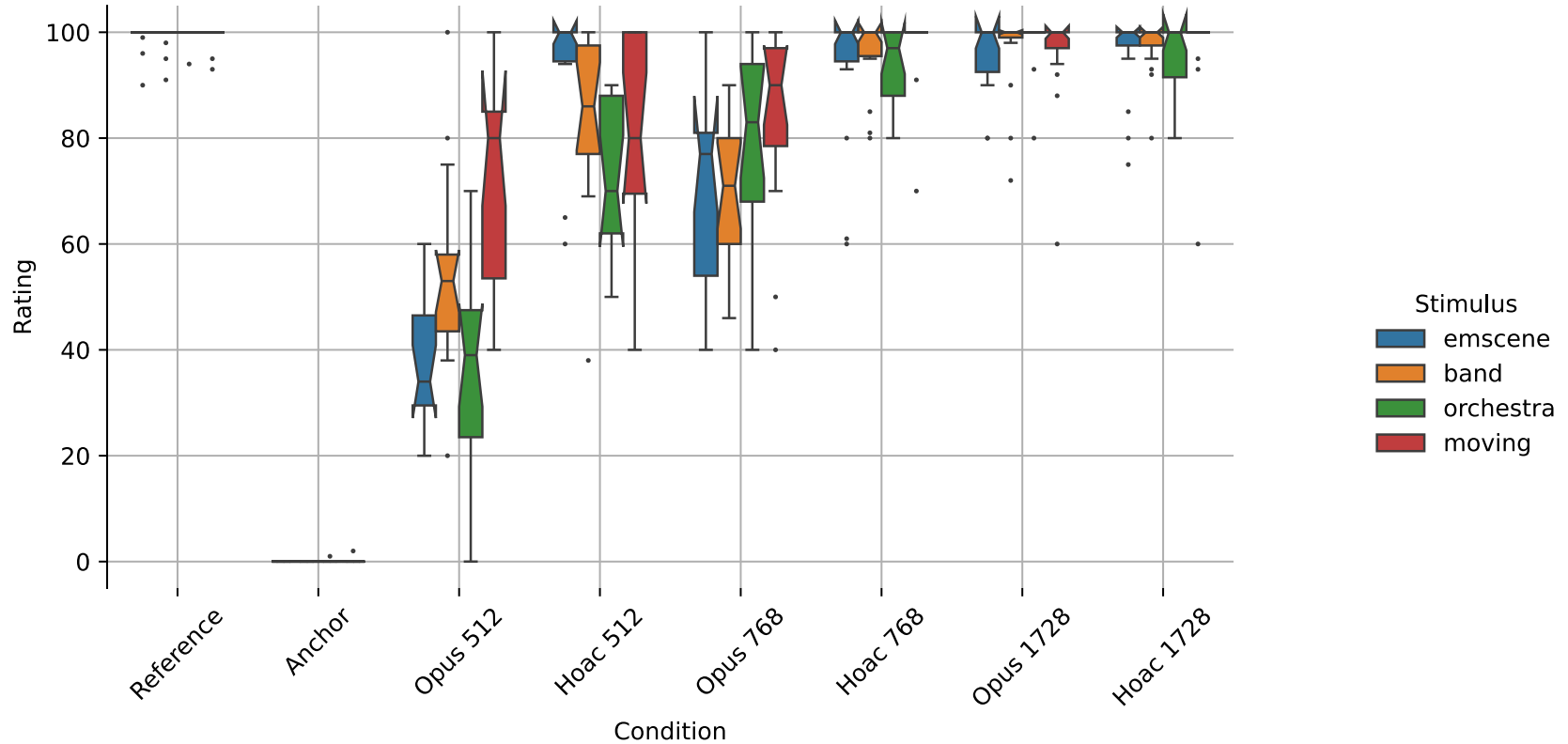


Influence of Coders



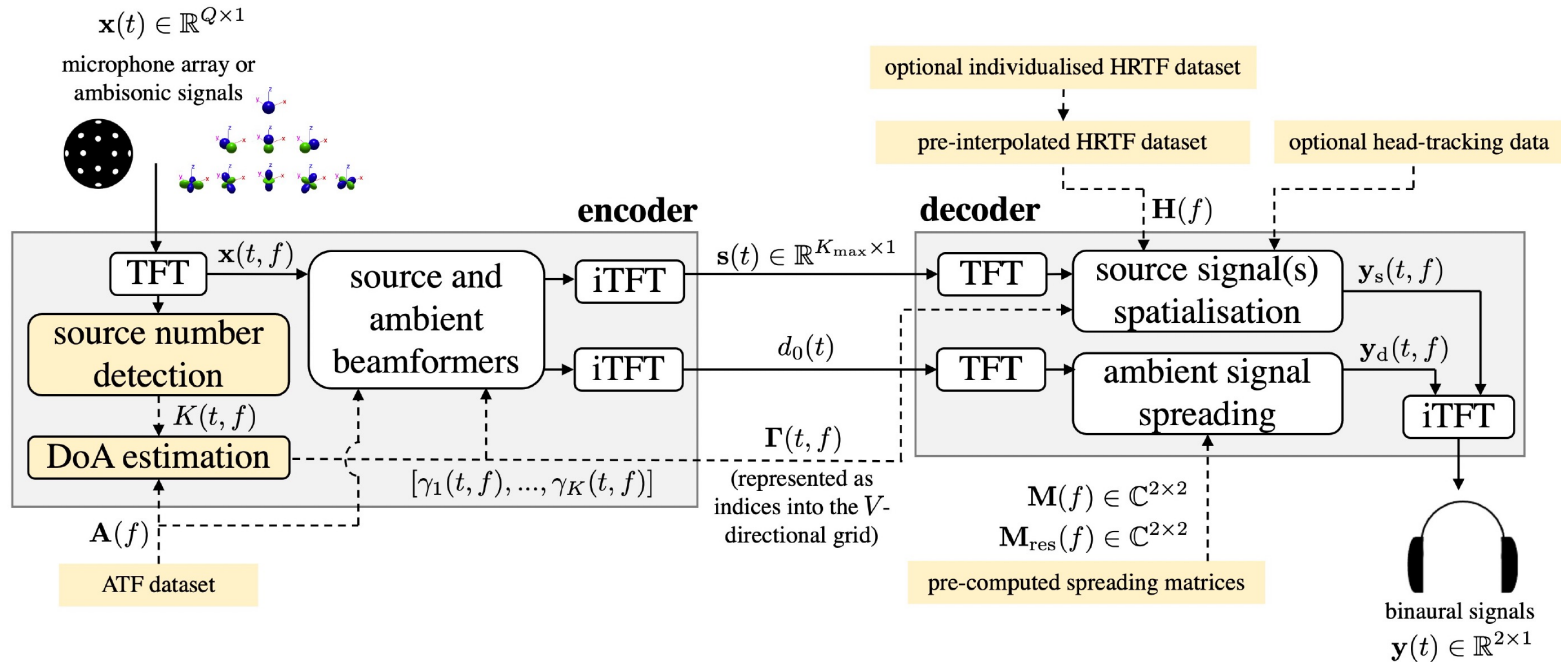
Perceptually-Motivated Spatial Audio Codec for Higher-Order Ambisonics Compression,
Christoph Hold, Leo McCormack, Archontis Politis, Ville Pulkki, IEEE ICASSP, 2024

Influence of Coders



Which model is the best?

It depends!



Implications and Outlook

- Parametric spatial audio can benefit compression
- Allows balancing bitrate between audio core-coders and metadata
 - Fewer audio TCs at higher quality, and quantized metadata
- Input parameterization can improve decoding

- Improvements suggested for Opus (Ambi Map 3)
- Implemented low resource Ambisonics layer in 3GPP SA4 IVAS

More:



<http://research.spa.aalto.fi/publications/papers/hoac/>



aalto.fi



Aalto-yliopisto
Aalto-universitetet
Aalto University