



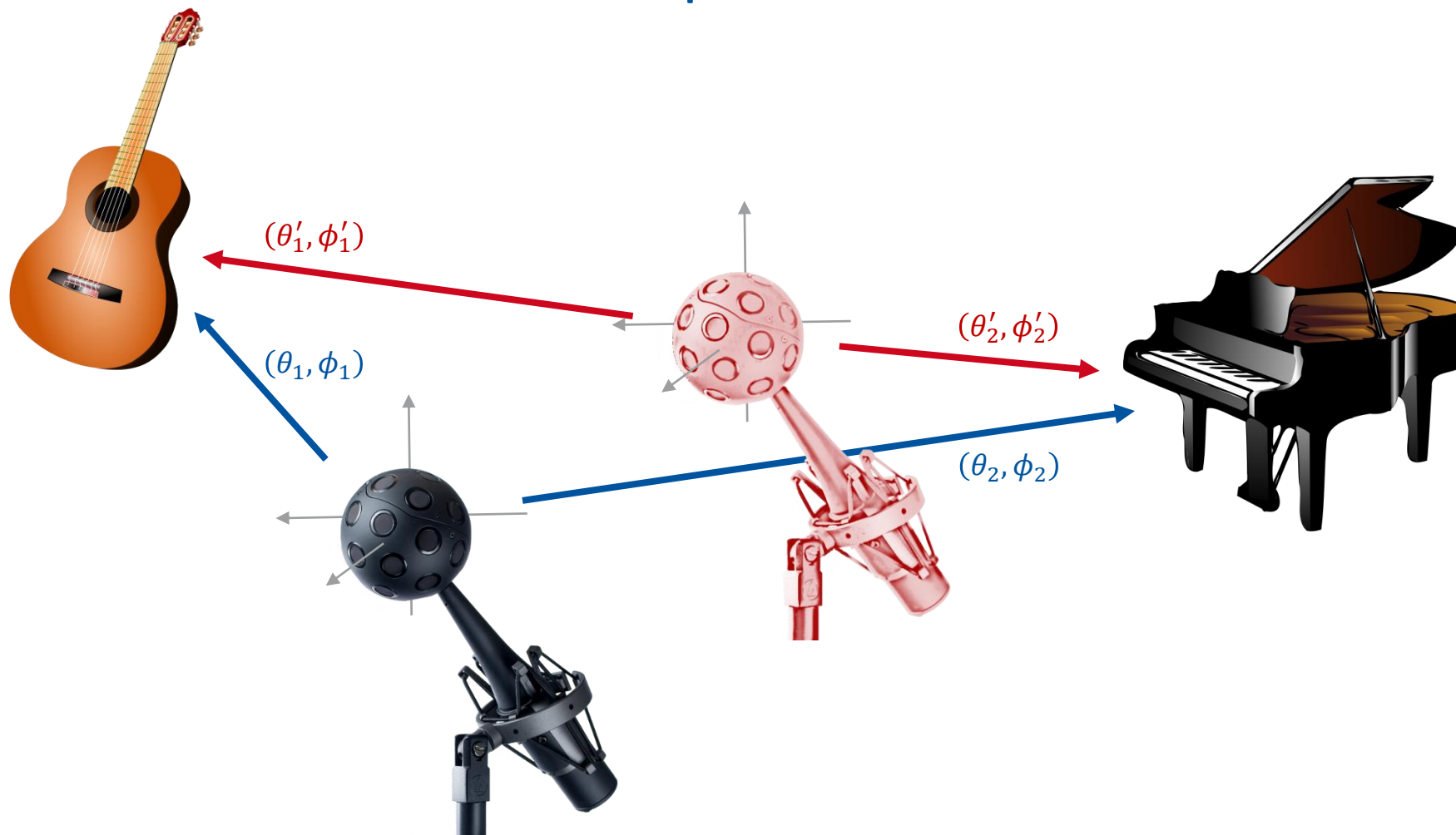
Translation of a Higher Order Ambisonics Sound Scene Based on Parametric Decomposition

Maximilian Kentgens, Andreas Behler, Peter Jax

IEEE ICASSP 2020

From Single Spot Microphone Array Recording...

...to User-Movement-Enabled Immersive Reproduction



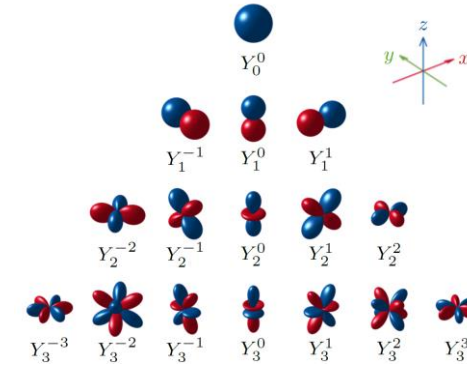
θ inclination angle
 ϕ azimuth angle



- Higher Order Ambisonics (HOA) signal $\mathbf{x}(\lambda, \mu) \in \mathbb{C}^{(N+1)^2}$ with Spherical Harmonics truncation order N :

$$\mathbf{x}(\lambda, \mu) = \mathbf{x}_s(\lambda, \mu) + \mathbf{x}_a(\lambda, \mu)$$

λ : frame index
 μ : frequency bin



- Direct part** $\mathbf{x}_s(\lambda, \mu)$: variable number of $I(\lambda, \mu) \in \{0, 1, 2, \dots, (N+1)^2\}$ plane wave sources:

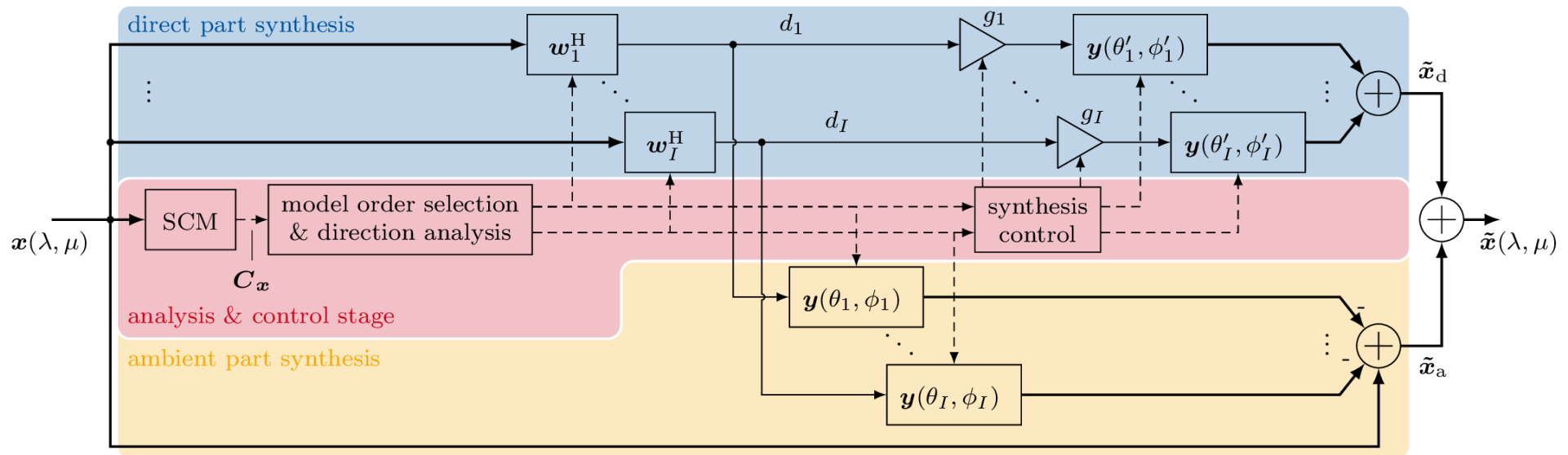
$$\mathbf{x}_s(\lambda, \mu) = \underbrace{\left[\mathbf{y}(\theta_1, \phi_2) \mid \dots \mid \mathbf{y}(\theta_{I(\lambda, \mu)}, \phi_{I(\lambda, \mu)}) \right]}_{= \text{array manifold matrix } \mathbf{Y}_s} \cdot \begin{pmatrix} s_1(\lambda, \mu) \\ \vdots \\ s_{I(\lambda, \mu)}(\lambda, \mu) \end{pmatrix}$$

- Ambient part** $\mathbf{x}_a(\lambda, \mu)$: spatially diffuse

- Spatial covariance matrix (SCM):

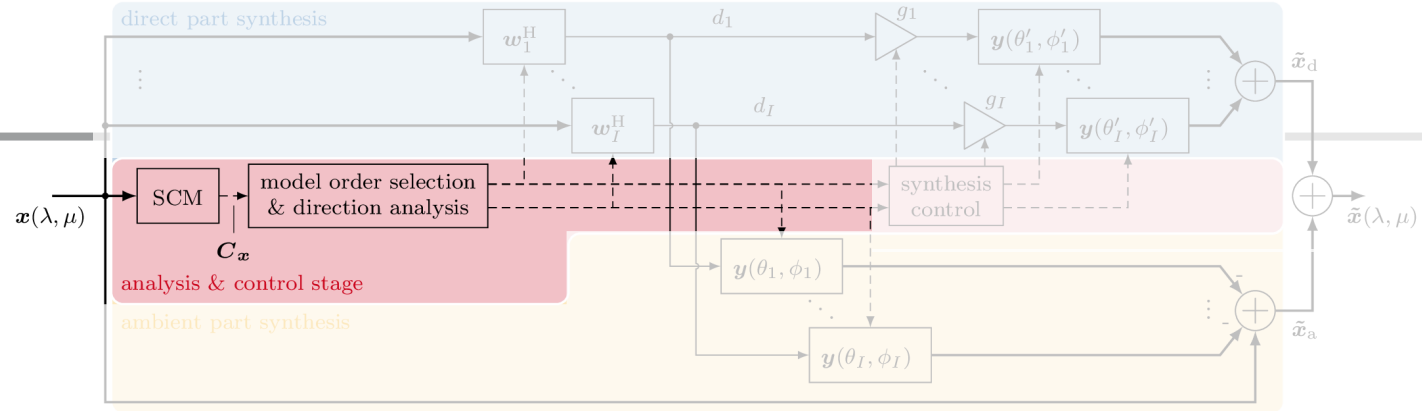
$$\begin{aligned} \mathbf{C}_x(\lambda, \mu) &= \mathbb{E}\{\mathbf{x}(\lambda, \mu)\mathbf{x}^H(\lambda, \mu)\} \\ &= \mathbb{E}\{\mathbf{x}_s(\lambda, \mu)\mathbf{x}_s^H(\lambda, \mu)\} + \text{diag}(\boldsymbol{\sigma}_a^2) \end{aligned}$$

Proposed System Overview



SCM Estimation of Spatial Covariance Matrix
 (θ, ϕ) Angles: (inclination, azimuth)
 λ Frame index
 μ Frequency index
 w_i Beamforming weights
 $y(\theta, \phi)$ Vector of Spherical Harmonics evaluated at (θ, ϕ)

Analysis Stage

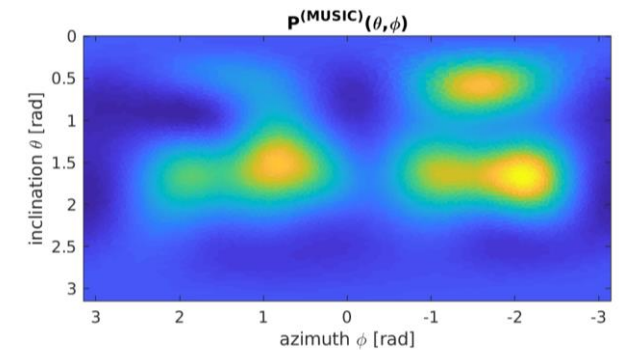
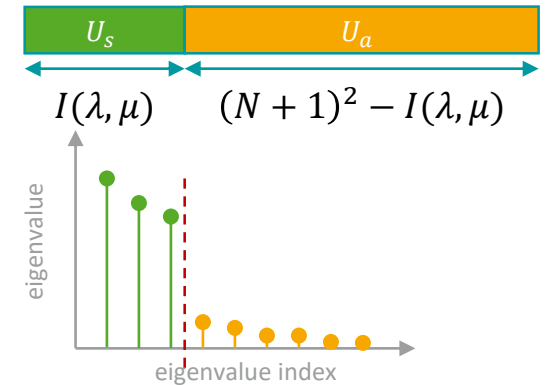


Segregation into multiple direct time-frequency components & ambient residuum

Variant 1: Subspace approach using SORTE & MUSIC

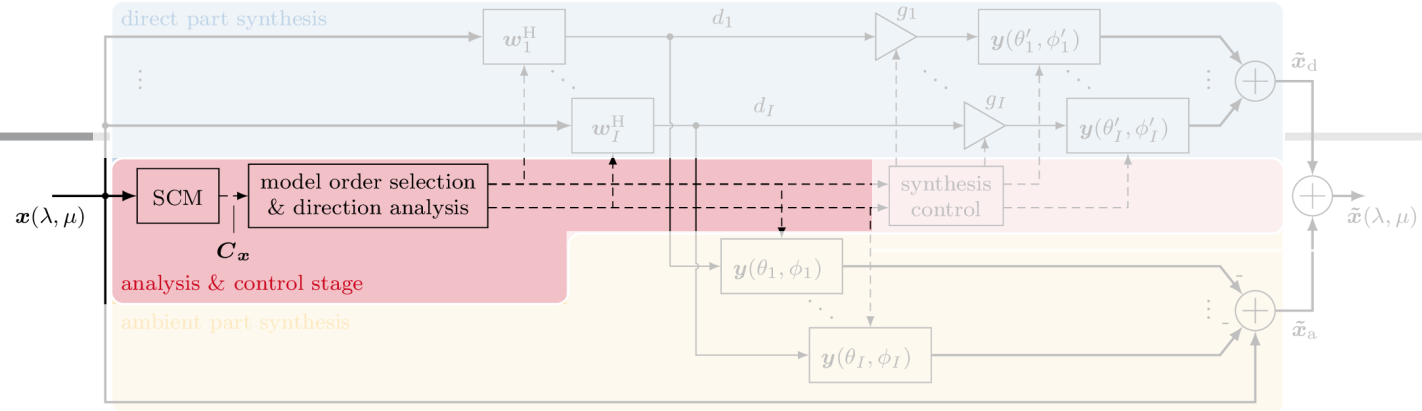
Adopted from [Politis, Tervo & Pulkki, ICASSP'18]

- Use eigenvalue decomposition to segregate $C_x(\lambda, \mu)$ into
 - direct components subspace $U_s(\lambda, \mu)$
 - residual subspace $U_a(\lambda, \mu)$
- Model order selection: determine number $I(\lambda, \mu)$
 - find gap in eigenvalue sequence using *second-order statistic of the eigenvalues* (SORTE)
- Use MUSIC to find $I(\lambda, \mu)$ source directions which are local maxima in



$$P^{(\text{MUSIC})}(\theta, \phi) = \frac{1}{\|y^H(\theta, \phi)U_a\|_2^2}$$

Analysis Stage



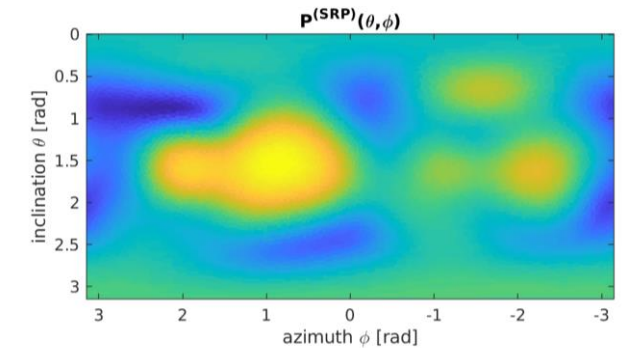
Segregation into multiple direct time-frequency components & ambient residuum

Variant 2: Steered Response Power (SRP) Map

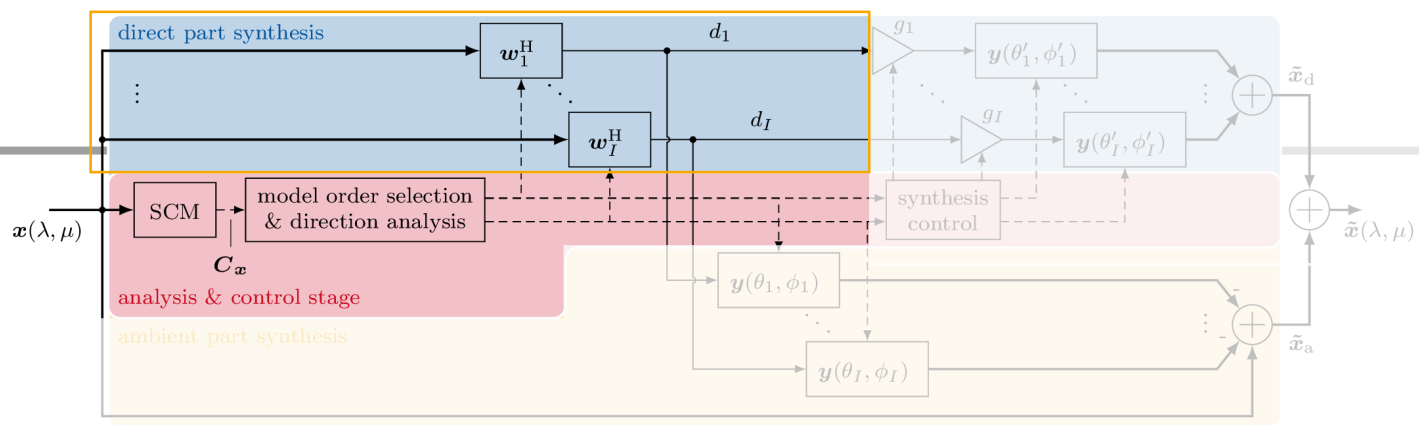
- Source directions: find all local maxima in

$$P^{(SRP)}(\theta, \phi) = \mathbf{y}^H(\theta, \phi) \mathbf{C}_x \mathbf{y}(\theta, \phi)$$

- Model order selection: $I = \#$ minima found



Analysis Stage



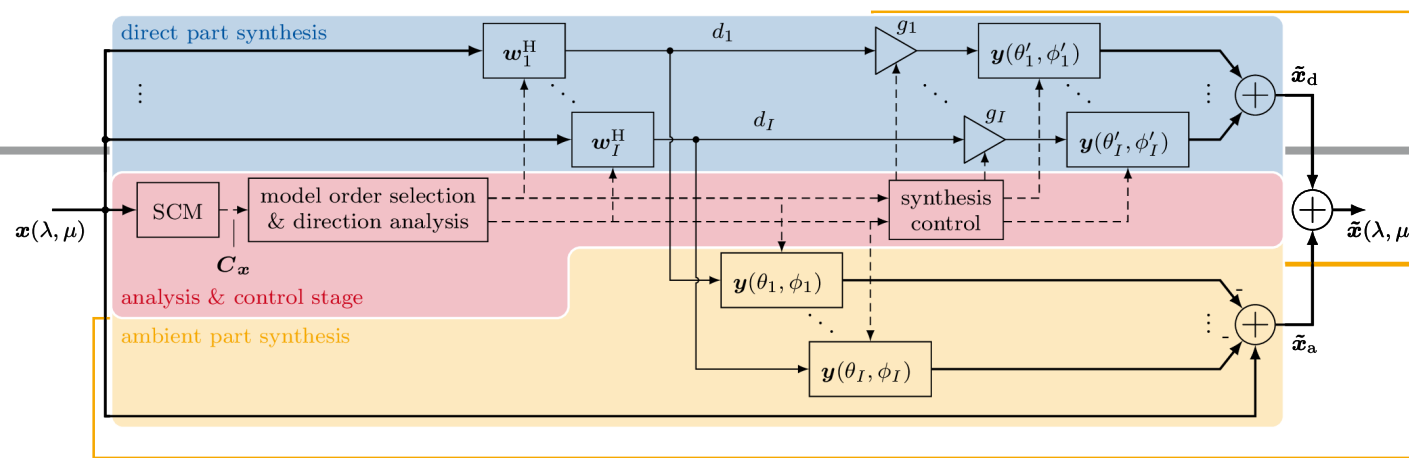
- **Beamforming** to extract direct components $i = 1, 2, \dots, I$

$$d_i(\lambda, \mu) = \mathbf{w}_i^H(\lambda, \mu) \mathbf{x}(\lambda, \mu)$$

weights \mathbf{w}_i : distortionless-response constraint in direction (θ_i, ϕ_i) and null constraints for (θ_j, ϕ_j) , $j \neq i$

- No semantic relation between direct components in time and frequency!

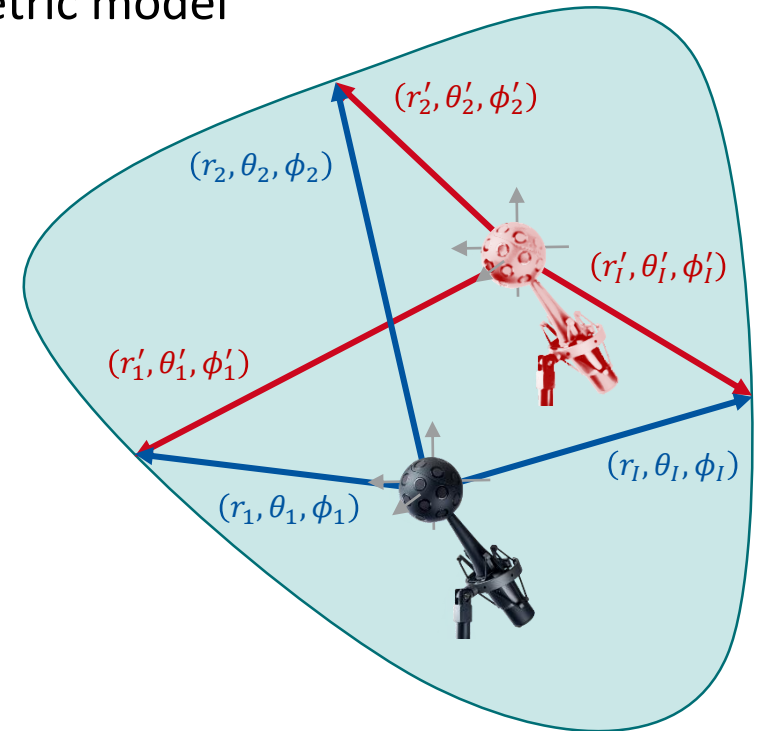
Synthesis Stage



- **Direct components: repanning & gain adjustment** according to geometric model

$$\tilde{\mathbf{x}}_d = [\mathbf{y}(\theta'_1, \phi'_1) \mid \dots \mid \mathbf{y}(\theta'_I, \phi'_I)] \cdot \begin{pmatrix} \frac{r_1}{r'_1} & & 0 \\ & \ddots & \\ 0 & & \frac{r_I}{r'_I} \end{pmatrix} \cdot \begin{pmatrix} d_i(\lambda, \mu) \\ \vdots \\ d_I(\lambda, \mu) \end{pmatrix}$$

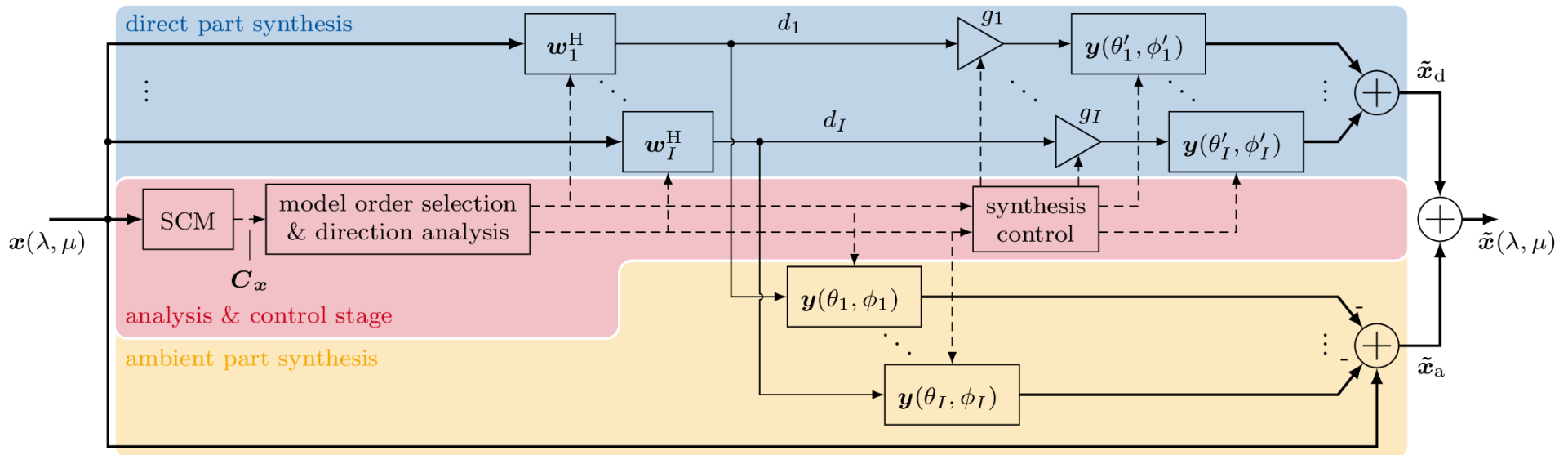
a priori information on source distances r_1, r_2, \dots, r_I needed



- **Ambient part:** cancel out direct components from input signal $\mathbf{x}(\lambda, \mu)$

- In total:

$$\tilde{\mathbf{x}}(\lambda, \mu) = \tilde{\mathbf{x}}_d(\lambda, \mu) + \tilde{\mathbf{x}}_a(\lambda, \mu)$$

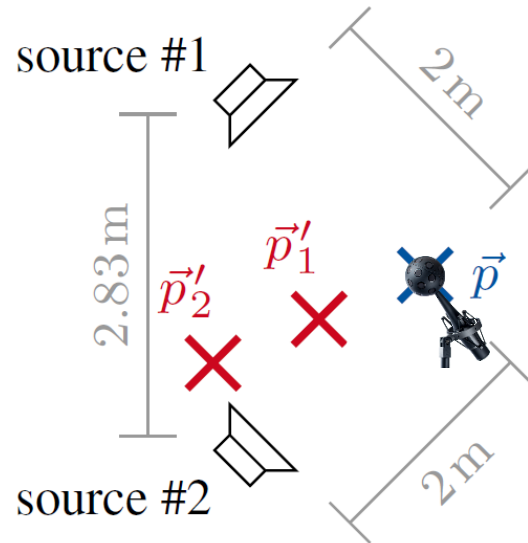


■ Early reflections not considered in signal model

- Likely to be modified as if they were direct sound
- **Precedence effect:** early reflections do not contribute to perceived source direction



MUSHRA Listening Test for Localization & Overall Audio Quality | 18 Participants



- HOA scene (order $N=4$) with two sound sources
- Translation to \vec{p}'_1 and \vec{p}'_2
- Two variants of the proposed algorithm:
 - MUSIC/SORTE
 - SRP-based
- For reference:
 - ORACLE: direction estimates replaced by actual source directions
 - Low anchor (AL): omni-directional sound, 4 kHz low-pass filtered
 - Mid anchor (AM): non-translated input signal $x(\lambda, \mu)$
 - Hidden reference (HR): ground truth at \vec{p}'_1 and \vec{p}'_2 , respectively

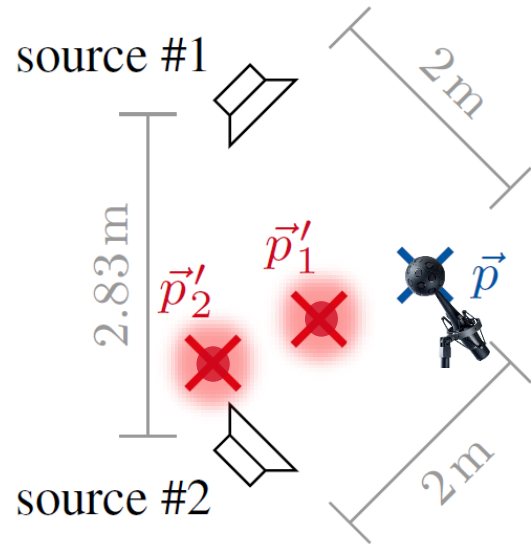
T_{60}	0.4s
source #1	pop music
source #2	male speaker
noise	diffuse (-30dB)
sample duration	8s

sampling rate	16kHz
SH order N	4
frame length	512 samples
frame overlap	50%
sub-band grouping	critical bands
rec. smoothing	0.1 s
decay time	
# spatial sampling points Q	900
oracle distance r_i	2m

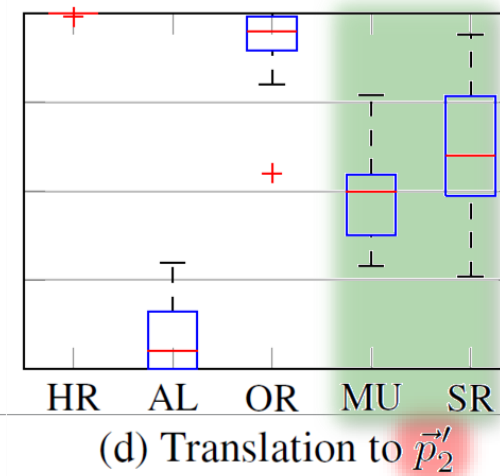
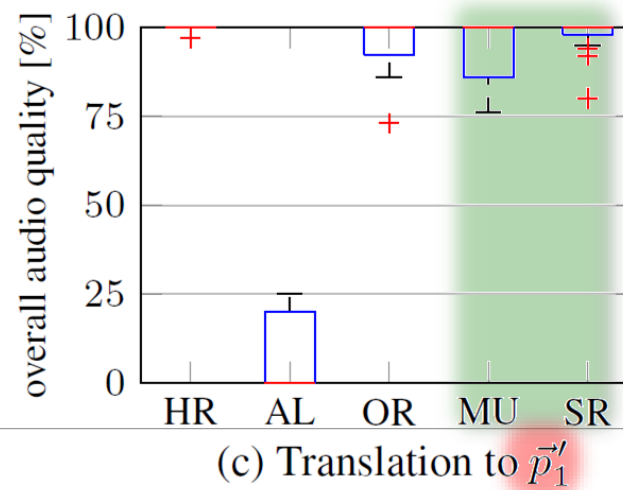
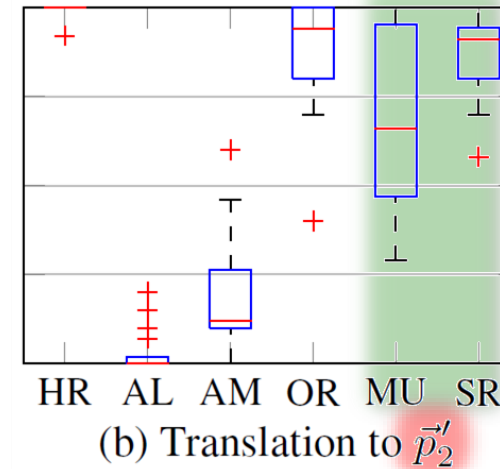
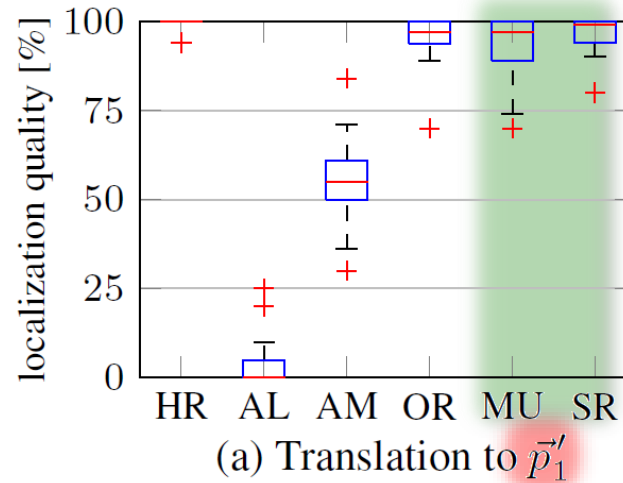
Subjective Evaluation



MUSHRA Listening Test for Localization & Overall Audio Quality | 18 Participants



T_{60}	0.4s
source #1	pop music
source #2	male speaker
noise	diffuse (-30 dB)
sample duration	8s



- HR hidden reference
- AL low anchor
- AM mid anchor (input signal)
- OR proposed method with oracle information on source directions
- MU **proposed method with MUSIC**
- SR **proposed method with SRP**

Conclusions



- Novel 3DoF+ system for translation in scene-based HOA content
 - Only single capture device needed
- Excellent subjective performance...
 - ... despite model error due to early reflections not treated differently than direct sound (→ psychoacoustic precedence effect)
 - Median subject ratings in localization between 91 % and 99 % (SRP method)
- Applications: VR/AR systems, immersive teleconferencing & telepresence

