

# A UNIFIED APPROACH TO NUMERICAL AUDITORY SCENE SYNTHESIS USING LOUDSPEAKER ARRAYS

JOSHUA ATKINS, ISMAEL NAWFAL, DANIELE GIACOBELLO

SEPTEMBER 4, 2014



# MOTIVATION

## SPATIAL SOUND SYSTEMS

- **Spatial sound overview**
  - Physical Reconstruction: wave-field synthesis (WFS), near-field compensated higher-order ambisonics (NFC-HOA)
    - **Issues**: not flexible in speaker arrangement, challenging for full-band audio
  - Interpolation: vector base amplitude panning (VBAP)
    - **Issues**: not flexible in speaker arrangement, sources located on surface of array, coloration of sources
  - Numerical Optimization: equivalent source method (ESM), mode-matching, crosstalk cancellation
    - **Issues**: no inclusion of perception, filter design is left as separate problem
- **Proposal in this work (Numerical Auditory Scene Synthesis)**
  - Goal: correct reproduction of perceived auditory scene (not wave field)
  - Convex numerical framework: flexible speaker layouts, listener positions, and error-norms
  - Inherently broadband: time-domain filter generation
  - Spatio-temporal projection: include perception, spatial error distribution

# OPTIMIZATION FRAMEWORK

## PROBLEM STATEMENT

**problem:** design filters to best approximate response at target locations

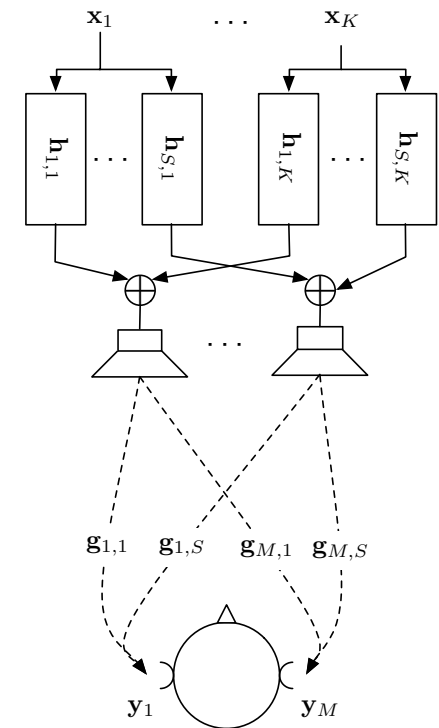
**assumptions:** 1 source, M target points, S speakers

- reproduction:  $\mathbf{y}$  (signal at target points),  $\mathbf{G}$  (acoustic impulse responses, convolution matrices),  $\mathbf{H}$  (unknown filters, convolution matrices),  $\mathbf{x}$  (signal)

$$\begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_M \end{bmatrix} = \begin{bmatrix} \mathbf{G}_{1,1} & \dots & \mathbf{G}_{1,S} \\ \vdots & \ddots & \vdots \\ \mathbf{G}_{M,1} & \dots & \mathbf{G}_{M,S} \end{bmatrix} \begin{bmatrix} \mathbf{H}_1 \\ \vdots \\ \mathbf{H}_S \end{bmatrix} \mathbf{x}$$

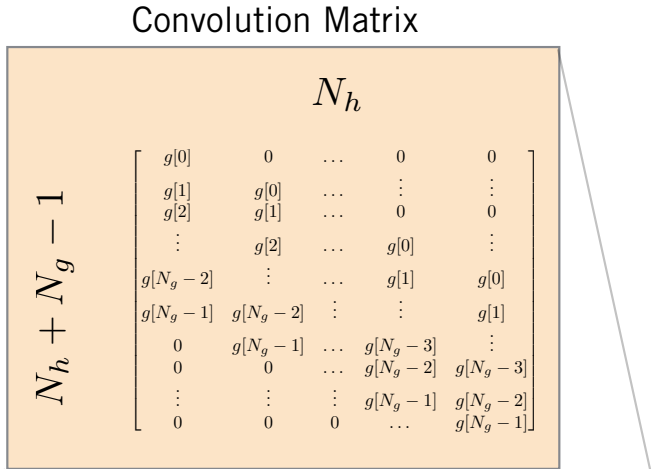
- design:  $\mathbf{t}$  (desired impulse response at target points),  $\mathbf{h}$  (unknown filters)

$$\begin{bmatrix} \mathbf{t}_1 \\ \vdots \\ \mathbf{t}_M \end{bmatrix} = \begin{bmatrix} \mathbf{G}_{1,1} & \dots & \mathbf{G}_{1,S} \\ \vdots & \ddots & \vdots \\ \mathbf{G}_{M,1} & \dots & \mathbf{G}_{M,S} \end{bmatrix} \begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_S \end{bmatrix}$$



# ACOUSTIC MODELS

FLEXIBILITY IN TARGET AND TRANSMISSION MODELS



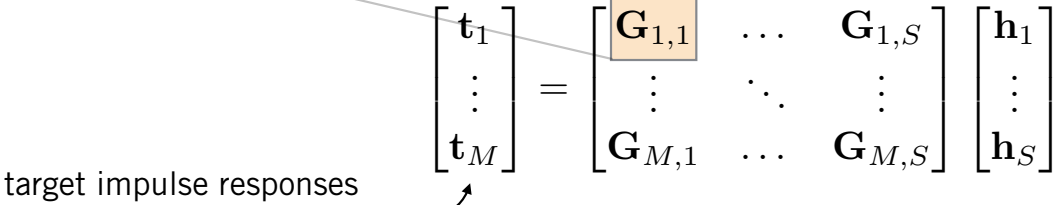
- Plane Wave

$$G(f) = Ae^{i(\mathbf{k} \cdot \mathbf{r} - \omega t)} \xrightarrow{\mathcal{F}^{-1}} g(t) = A\delta\left(\frac{\mathbf{n} \cdot \mathbf{r}}{c} - t\right)$$

- Spherical Wave

$$G(f) = \frac{Ae^{i(kr - \omega t)}}{r} \xrightarrow{\mathcal{F}^{-1}} g(t) = \frac{A}{r}\delta\left(\frac{r}{c} - t\right)$$

- Head-related impulse response (HRIR, BRIR)
- Any other acoustic impulse response



# NUMERICAL AUDITORY SCENE SYNTHESIS PROBLEM

FLEXIBLE CONVEX PROGRAM

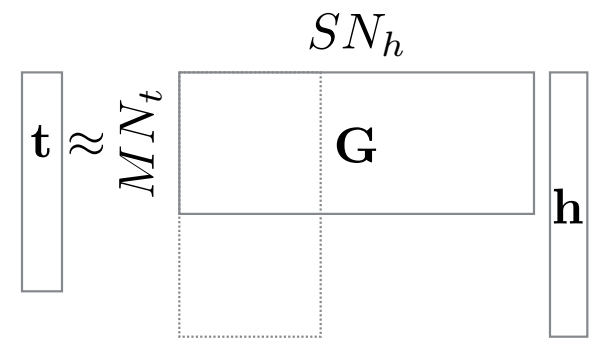
- **Underdetermined** ( $S N_h > M N_t$ , full rank)  
choose one of many exact solutions:

$$\hat{\mathbf{h}} = \arg \min_{\mathbf{h}} \|\mathbf{\Gamma h}\|_q \quad \text{s. t.} \quad \mathbf{Gh} = \mathbf{t}$$

- **Overdetermined** ( $S N_h < M N_t$ , full rank) and/or uncertainty  
approximate solution:

$$\hat{\mathbf{h}} = \arg \min_{\mathbf{h}} \|\mathbf{W}(\mathbf{Gh} - \mathbf{t})\|_p \quad \text{s. t.} \quad \|\mathbf{\Gamma}_i \mathbf{h}\|_{q_i} \leq \gamma_i, \\ \forall i, i = 1 \dots, I$$

- Convex, flexible error/regularizer norm, spatio-temporal projection matrices (incorporate perception)



# SPATIO-TEMPORAL TRANSFORMS

ALTER SOLUTION SPACE AND/OR FILTER SPECIFICATION

$$\hat{\mathbf{h}} = \arg \min_{\mathbf{h}} \|\mathbf{W}(\mathbf{G}\mathbf{h} - \mathbf{t})\|_p \quad \text{s. t.} \|\mathbf{\Gamma}_i \mathbf{h}\|_{q_i} \leq \gamma_i, \\ \forall i, i = 1 \dots, I$$

- **Time-frequency transform** (DFT, filter banks, and time/frequency weighting, averaging, interpolation)

$$\mathbf{W}_t = \begin{bmatrix} \mathbf{F}_1 & & 0 \\ & \ddots & \\ 0 & & \mathbf{F}_M \end{bmatrix} \quad \mathbf{F} = \begin{bmatrix} f_{1,1} & \dots & f_{1,N_t} \\ \vdots & \ddots & \vdots \\ f_{N_f,N_t} & \dots & f_{N_f,N_t} \end{bmatrix}$$

- **Space-wavenumber transform** (spherical/cylindrical harmonics and spatial weighting, averaging, ...)

$$\mathbf{W}_s = \begin{bmatrix} y_{1,1}\mathbf{I} & \dots & y_{1,M}\mathbf{I} \\ \vdots & \ddots & \vdots \\ y_{C,1}\mathbf{I} & \dots & y_{C,M}\mathbf{I} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} y_{1,1} & \dots & y_{1,M} \\ \vdots & \ddots & \vdots \\ y_{C,1} & \dots & y_{C,M} \end{bmatrix}$$

# NUMERICAL AUDITORY SCENE SYNTHESIS PROBLEM

## OPTIONS FOR SYSTEM DESIGN

Which acoustic model,  $\mathbf{G}$ ?

Which error transform,  $\mathbf{W}$ ?

Which constraint norm,  $q$ ?

Which constraint value,  $\gamma$ ?

- Not studied here:
  - Which error norm,  $p$ ?
  - How many target points?
  - How many speakers?
  - Which speaker locations?

$$\hat{\mathbf{h}} = \arg \min_{\mathbf{h}} \|\mathbf{W}(\mathbf{G}\mathbf{h} - \mathbf{t})\|_p \quad \text{s. t.} \|\mathbf{\Gamma}_i \mathbf{h}\|_{q_i} \leq \gamma_i, \\ \forall i, i = 1 \dots, I$$

Which acoustic target,  $\mathbf{t}$ ?

Which constraint transform,  $\mathbf{\Gamma}$ ?

# CASE STUDY 1: THE EFFECT OF CONSTRAINT NORM

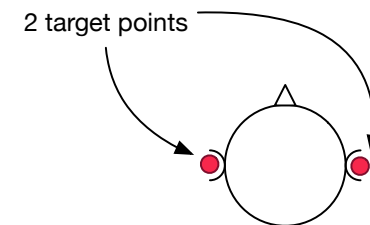
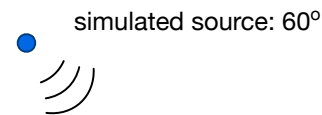
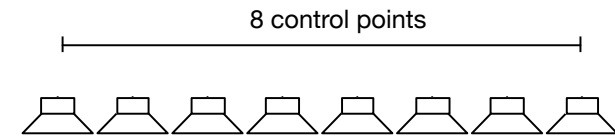
## EXPERIMENT SETUP

$$\hat{\mathbf{h}} = \arg \min_{\mathbf{h}} \|\mathbf{G}\mathbf{h} - \mathbf{t}\|_2 \quad \text{s. t.} \|\mathbf{\Gamma}\mathbf{h}\|_q \leq \gamma$$

Which constraint norm,  $q$ ?

Which constraint value,  $\gamma$ ?

- spherical wave acoustic model ( $\mathbf{G}$ ,  $\mathbf{t}$ )
- $l_2$ -norm error ( $p = 2$ )
- **DFT projection matrix ( $\mathbf{\Gamma}$ )**
- filter length = 1024, modeling delay = 100
- **4 systems:**
  - unconstrained
  - $l_1$ -norm constraint ( $q = 1$ )
  - $l_2$ -norm constraint ( $q = 2$ )
  - $l_\infty$ -norm constraint ( $q = \infty$ )

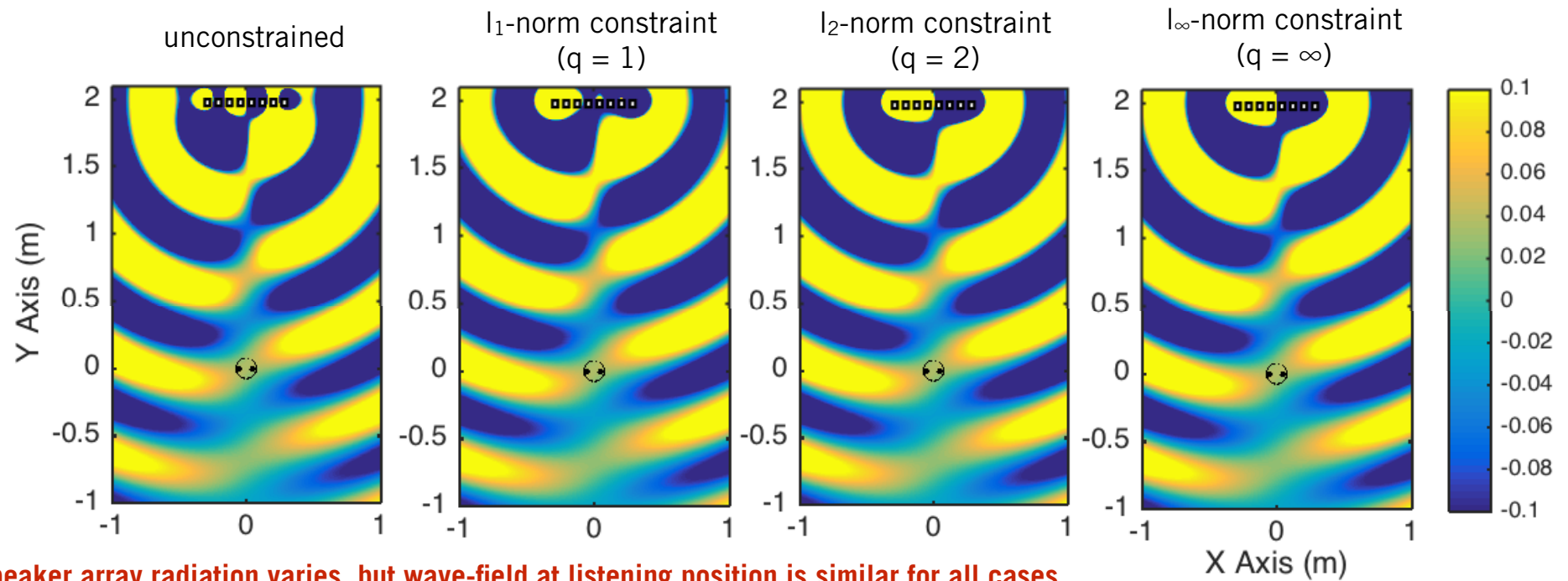




# CASE STUDY 1: THE EFFECT OF CONSTRAINT NORM

WAVEFIELD AT 500 HZ

$$\hat{\mathbf{h}} = \arg \min_{\mathbf{h}} \|\mathbf{G}\mathbf{h} - \mathbf{t}\|_2 \quad \text{s. t.} \|\mathbf{\Gamma}\mathbf{h}\|_q \leq \gamma$$



speaker array radiation varies, but wave-field at listening position is similar for all cases

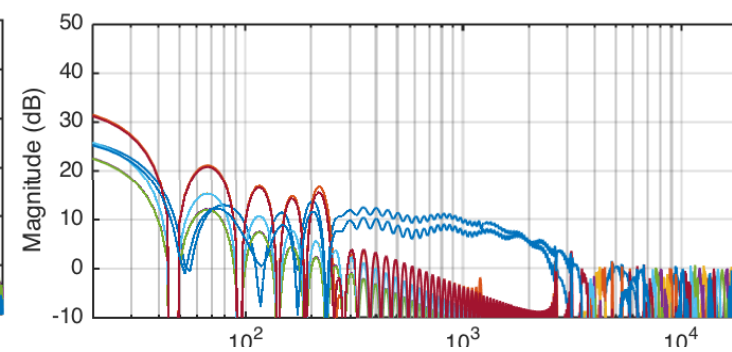
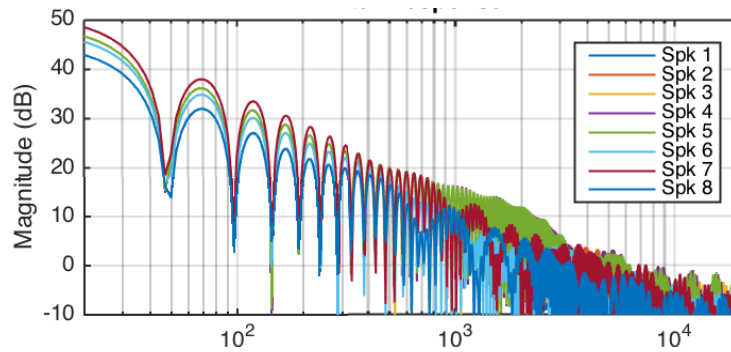
# CASE STUDY 1: THE EFFECT OF CONSTRAINT NORM

BROADBAND FILTER RESPONSE, FIXED  $L_2$ -NORM=39DB FOR ALL CASES

Large gain required at low frequencies

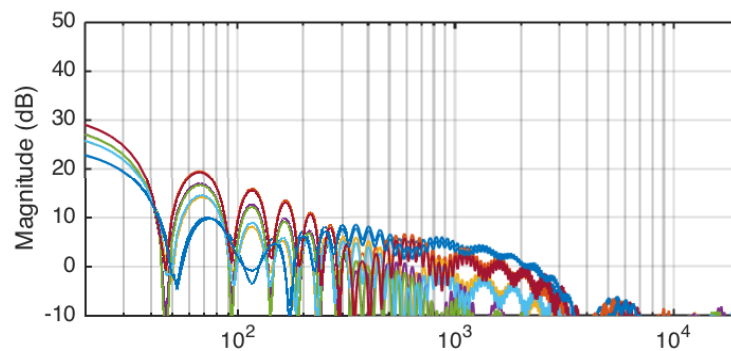
Few active speakers per frequency band

unconstrained

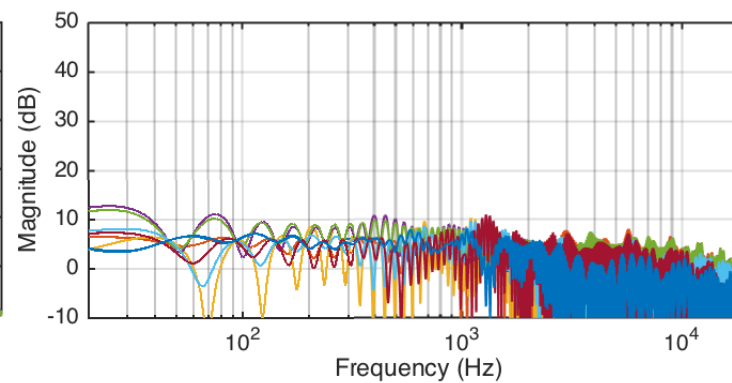


$l_1$ -norm

$l_2$ -norm



Limited required power



$l_\infty$ -norm

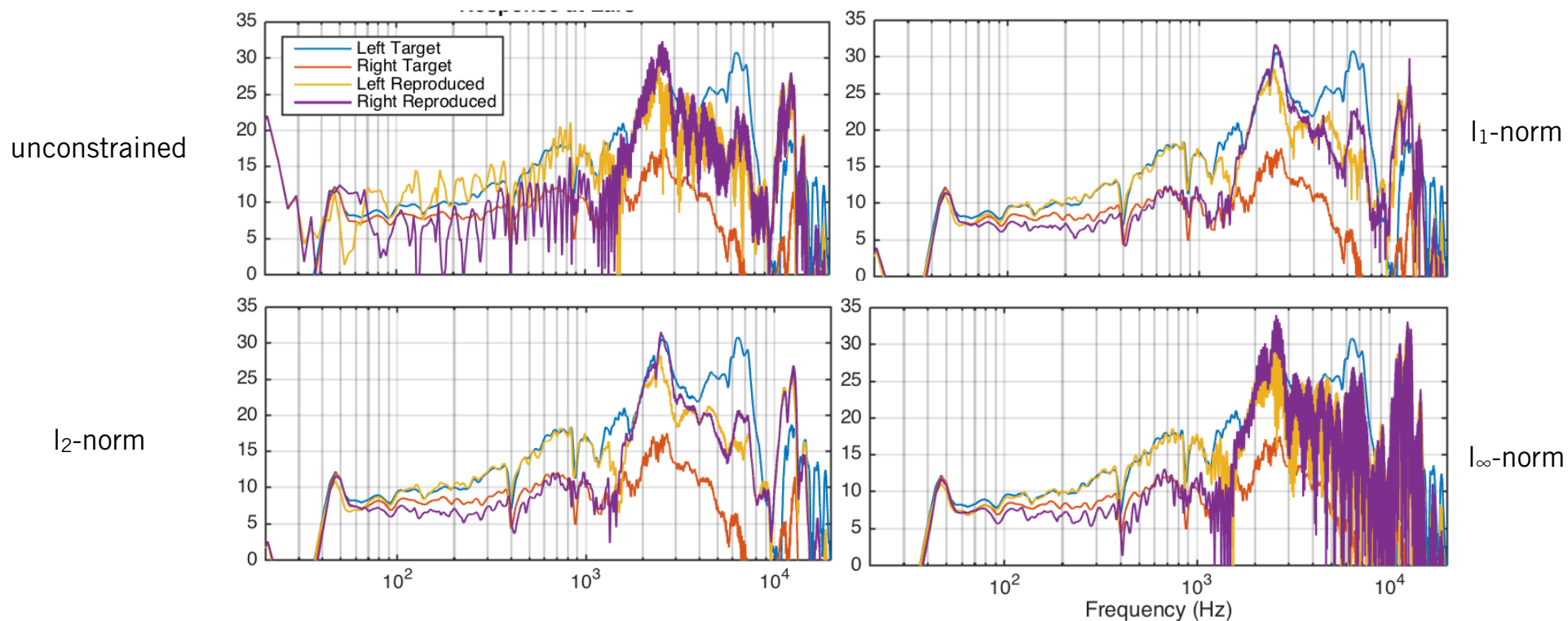
Limited maximum magnitude



# CASE STUDY 1: THE EFFECT OF CONSTRAINT NORM

RESPONSE SIMULATED AT EAR DRUM REFERENCE (TARGET POINTS)

All fail above 1-2 kHz, high frequency coloration



# CASE STUDY 2: PERCEPTUAL ERROR TRANSFORM

## EXPERIMENTAL SETUP

$$\hat{\mathbf{h}} = \arg \min_{\mathbf{h}} \|\mathbf{W}(\mathbf{G}\mathbf{h} - \mathbf{t})\|_2 \quad \text{s. t.} \quad \|\mathbf{\Gamma}\mathbf{h}\|_q \leq \gamma$$

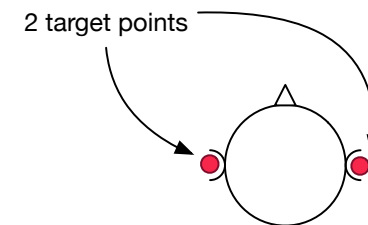
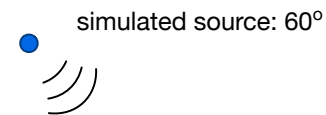
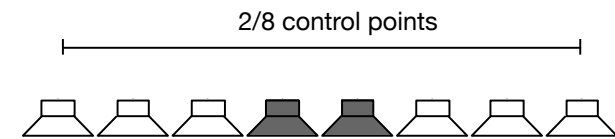
Which acoustic model,  $\mathbf{G}$ ?

Which error transform,  $\mathbf{W}$ ?

Which acoustic target,  $\mathbf{t}$ ?

Which constraint transform,  $\mathbf{\Gamma}$ ?

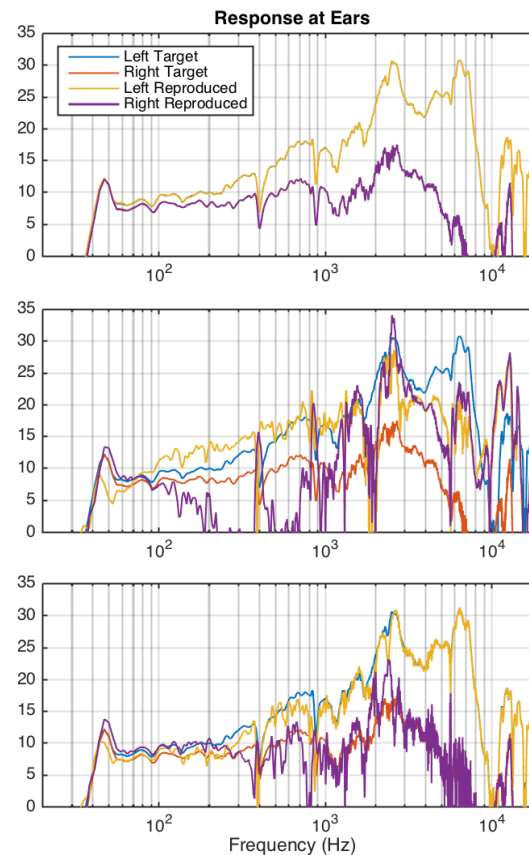
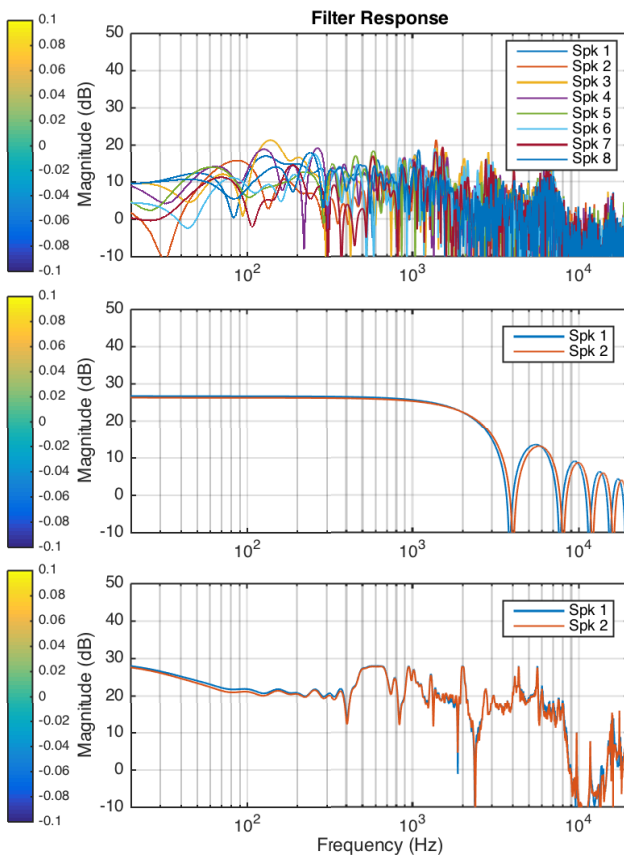
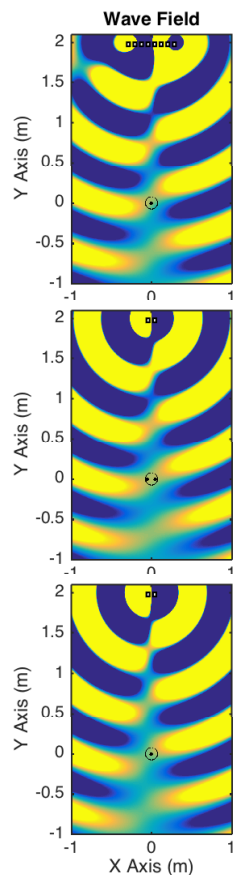
- $l_2$ -norm error ( $p = 2$ )
- filter length = 1024, modeling delay = 100
- **3 systems:**
  - 8 speakers, HRIR, unconstrained
  - 2 speakers, spherical wave, unconstrained
  - 2 speakers, HRIR,  $l_\infty$ -norm constraint, ERB-spaced DFT ( $\mathbf{W}$ ,  $\mathbf{\Gamma}$ )



# CASE STUDY 2: PERCEPTUAL ERROR TRANSFORM

WAVEFIELD (500 HZ), FILTER RESPONSE, AND RESPONSE AT EAR DRUM

8 speakers  
HRIR  
unconstrained



exact solution

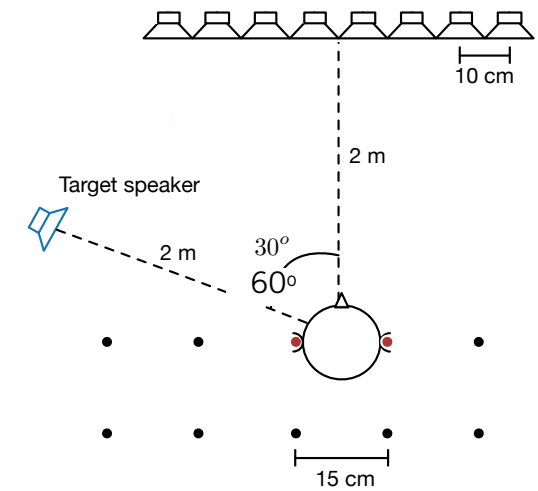
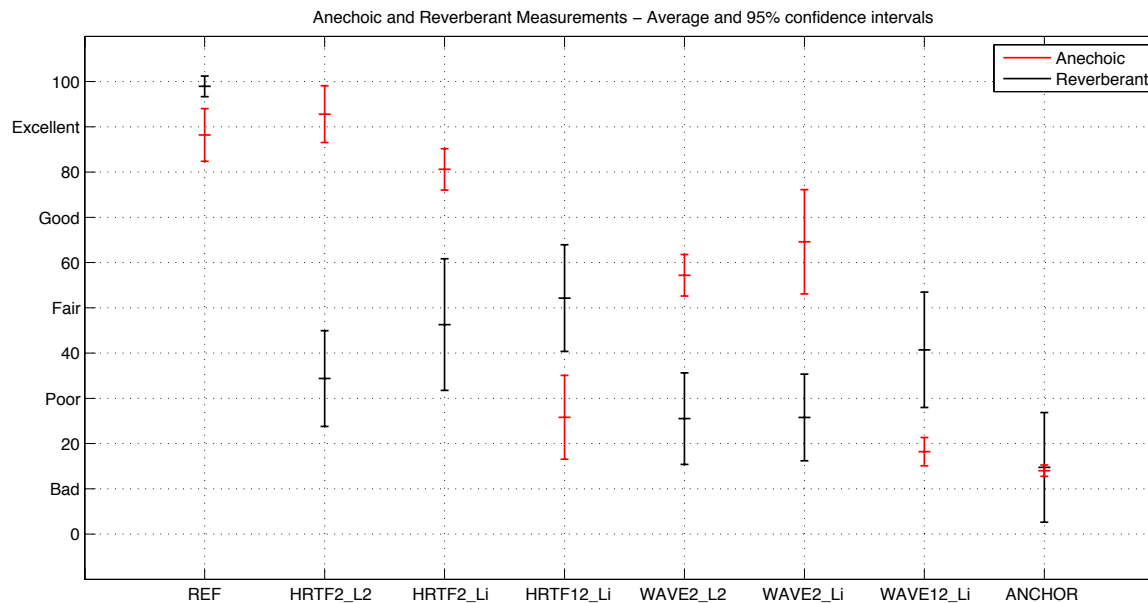
large timbral errors

very close response



# IN PRACTICE: SOME EARLY PERCEPTUAL RESULTS

FROM: AES 55TH CONFERENCE ON SPATIAL AUDIO



**HRIR target is preferred, multipoint (overdetermined) does better in reverberant scenario**

[1] Ismael Nawfal, Joshua Atkins, Daniele Giacobello, Stephen Nimick. "Perceptual Evaluation of Numerical Auditory Scene Synthesis Using Loudspeaker Arrays." Proceedings of the 55th Convention of the Audio Engineering Society. August 2014.

# CONCLUSION

...

- Numerical Auditory Scene Synthesis
  - Flexible spatial rendering method for generating time-domain **broadband filters**
  - Can be used with **arbitrary loudspeaker arrays**
  - **Convex** program guarantees achievable solution
  - Spatio-temporal transform matrices allow for simple inclusion of **perceptual constraints**
- Analysis
  - Showed effect of filter constraint norm on resulting system
    - easily prefer sparse loudspeaker activations or limit maximum gain applied to loudspeaker array
  - Simple perceptual constraints: ERB-spaced transform, HRIR target/acoustic model
    - outperforms spherical wave assumption in objective & subjective tests



**THANKS!**