

# Spatial Audio Reproduction Using Primary Ambient Extraction

PhD candidate: **Jianjun HE**

Supervisor: Assoc Prof Gan Woon-Seng

Digital Signal Processing Laboratory  
School of Electrical and Electronic Engineering  
Nanyang Technological University, Singapore



**16<sup>th</sup> Feb, 2016**

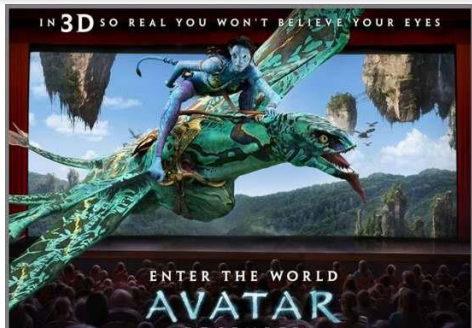


# Outline

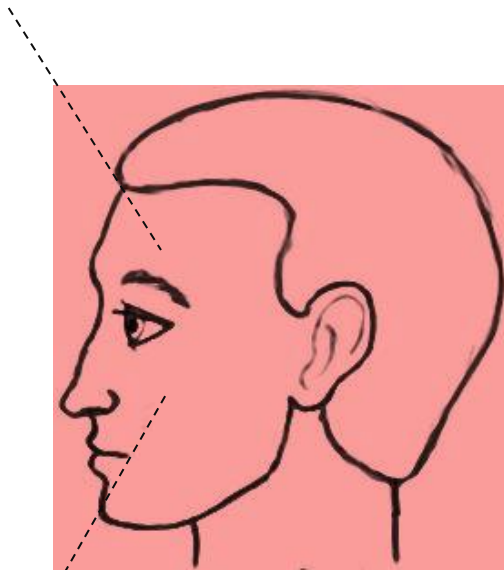
- I. Background of Spatial Audio Reproduction
- II. Overview of Primary Ambient Extraction (PAE)
- III. My Contributions
  1. Linear Estimation based PAE
  2. Ambient Spectrum Estimation based PAE
  3. Time Shifting based PAE
  4. Multi-Source and Multichannel based PAE
  5. Natural Sound Rendering for Headphones
- IV. Conclusions and Future work



# Hearing is Believing



Immersive visual experience



1. What/where is the sound?
2. How is the sound played?
3. How is the sound heard?



Immersive listening experience

Sound content processing

Personal listening cues



# Spatial hearing of sound source

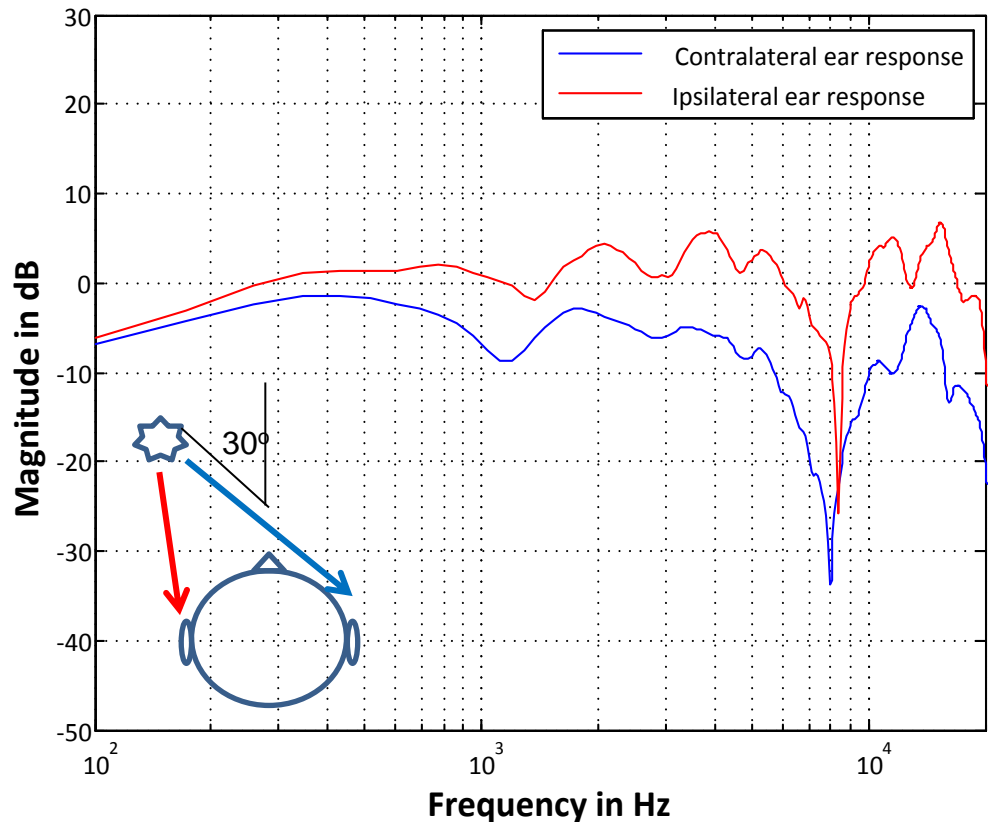
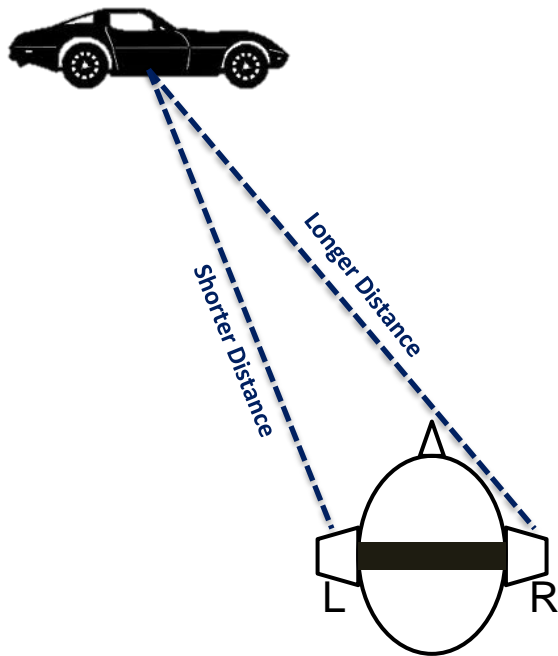
Azimuth, Elevation, Distance



# Head Related Transfer Functions (HRTF)

**ITD:** Interaural Time Difference  
**ILD:** Interaural Level Difference  
**SC:** Spectral Cues

} Head Related Transfer Functions (**HRTF**)



HRTF of KEMAR dummy head for an angle of 30 degree azimuth\*

# Spatial hearing of sound scene



# Source-medium-receiver view of spatial audio reproduction

## Immersive Audio

“being there”

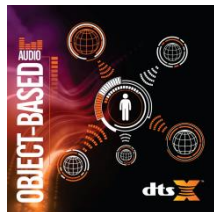
### Source

Audio content

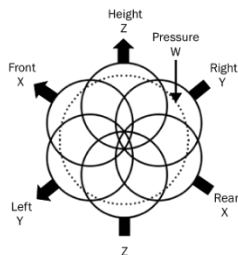
Channel-based



Object-based



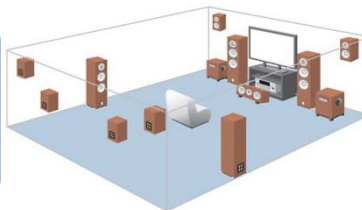
Transform-domain-based



### Medium

Playback system

Loudspeakers



Headphones



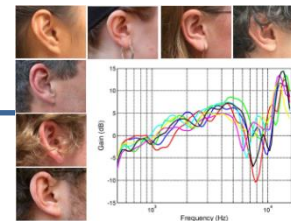
### Receiver

Human ear

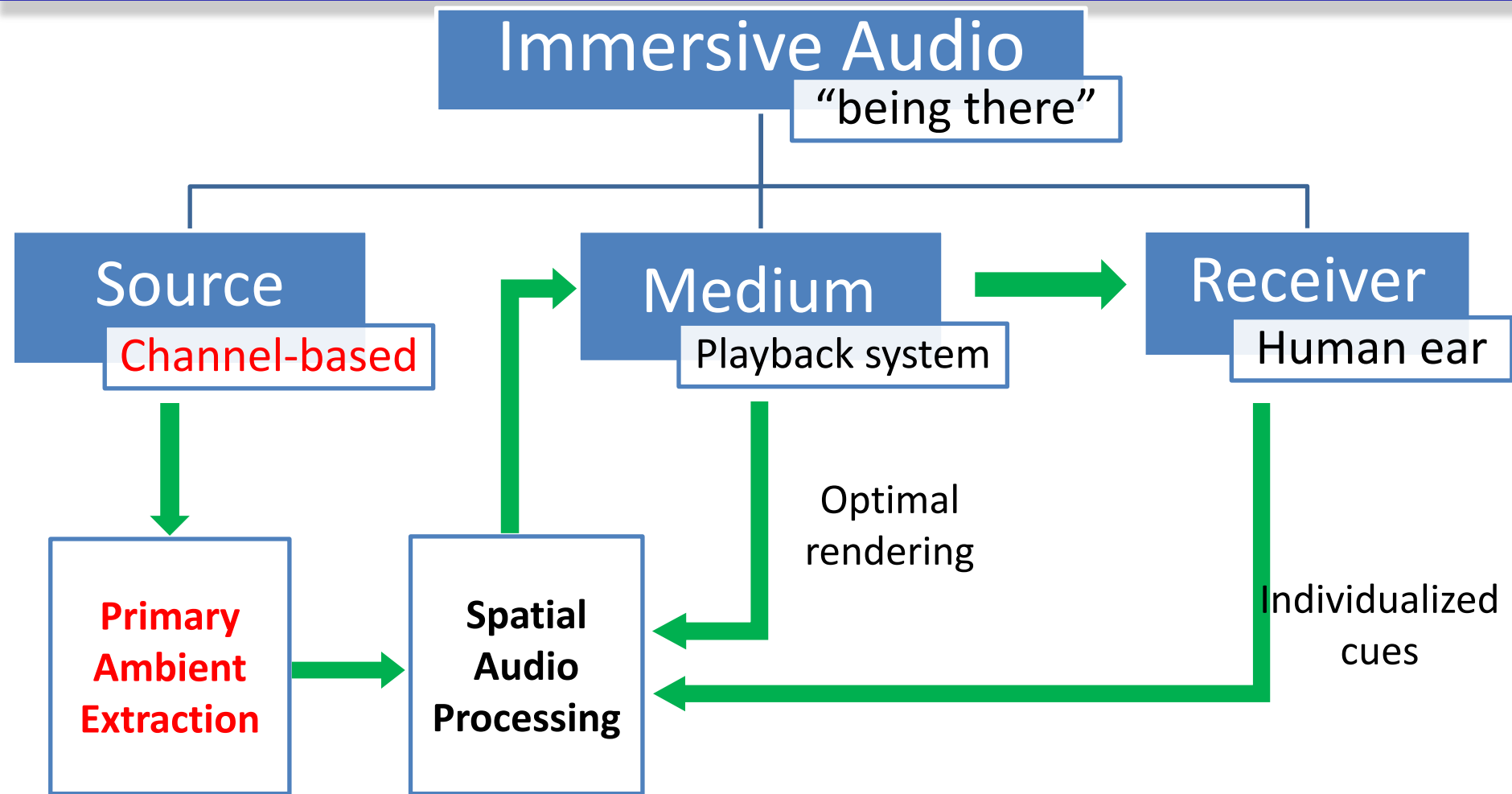
Ear A



Ear B



# Achieving consistency in source and medium in spatial audio



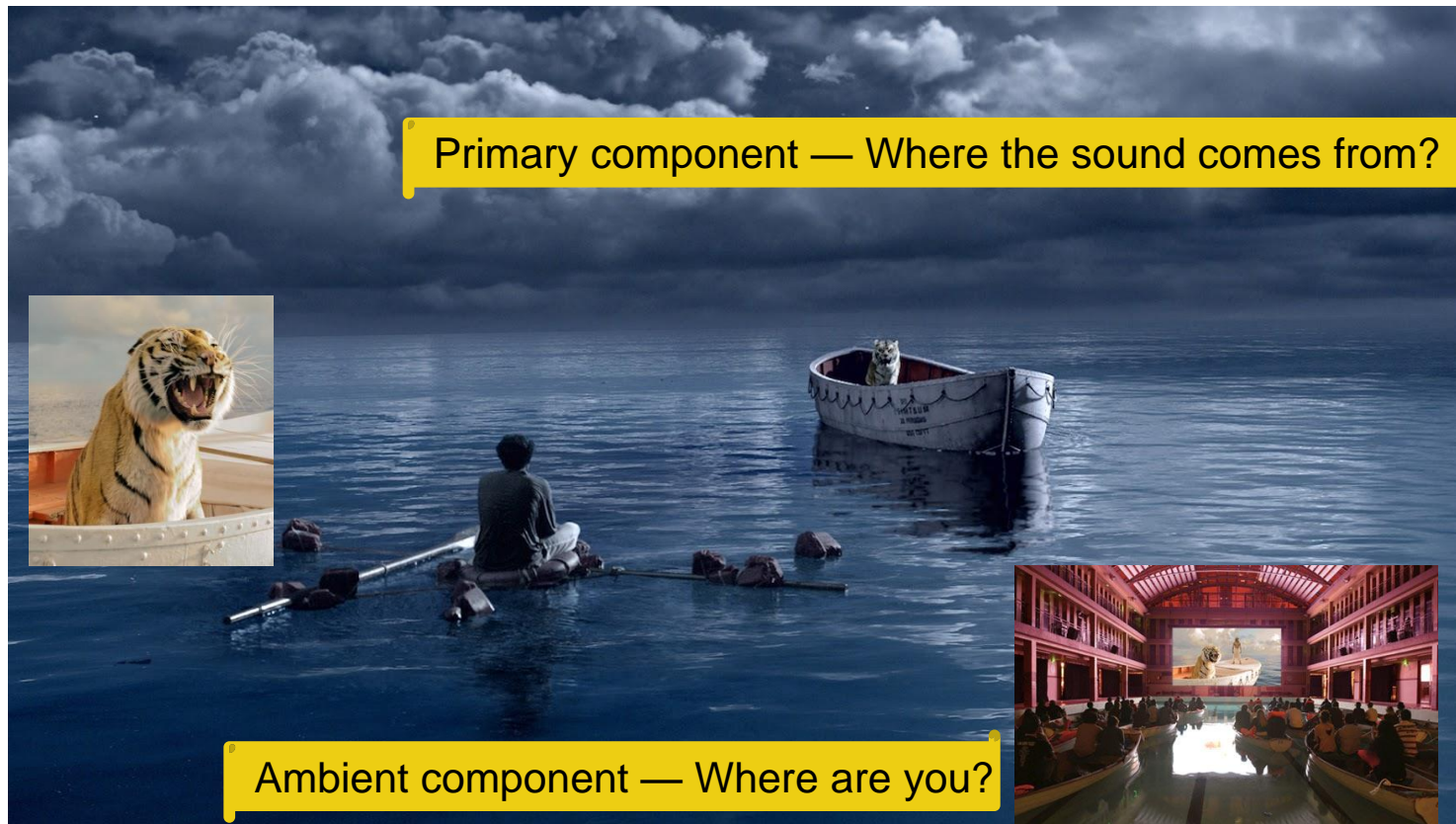
Essentially, PAE serves as a front-end to facilitate **flexible**, **efficient**, and **immersive** spatial audio reproduction.



# What are primary and ambient components?

## Main characteristics

- Primary components are directional
- Ambient components are diffuse



# What does PAE do?

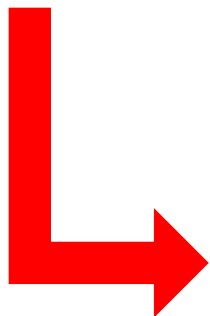


Sound scenes in movies and games often consist of:

**directional sound objects**

and **background ambience**, e.g.,

- Circling aircraft in the rain;
- Flying bee at the waterfall;
- Soaring tiger on the sea;
- Attacking wolves in the forest.



**Primary  
Ambient  
Extraction**



**Directional  
sound**



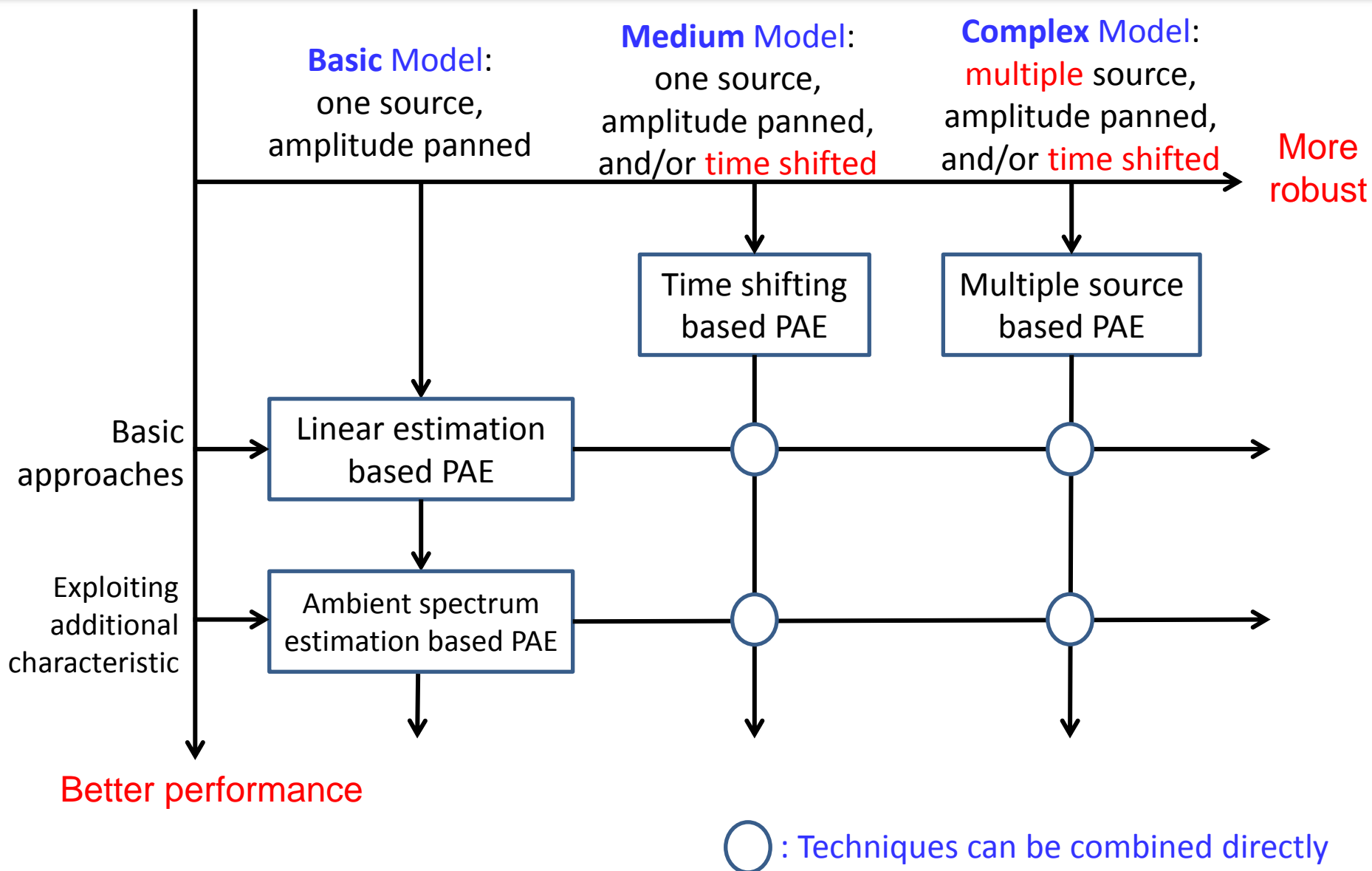
**Background  
ambience**



# A brief literature of PAE

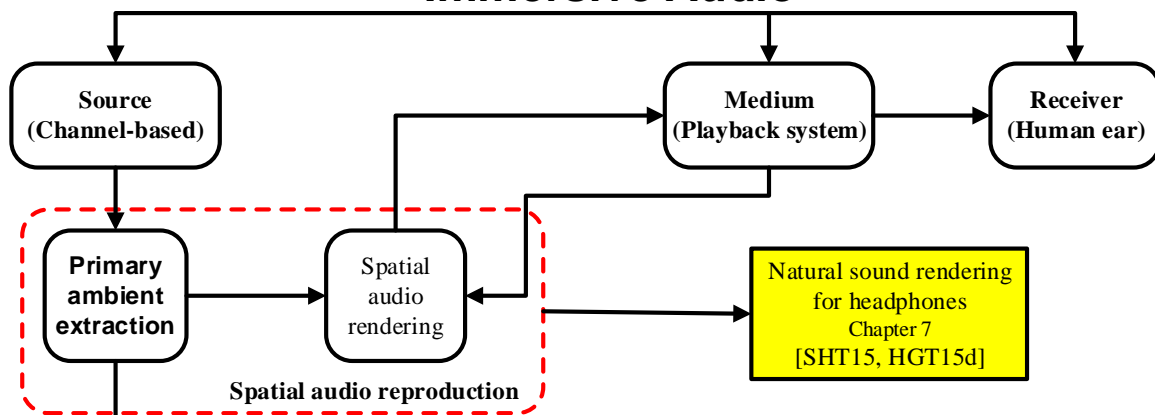
No. of channels	Complexity of audio scenes		
	Basic (amplitude panned single source)	Medium (single source)	Complex (multiple sources)
Stereo (2)	<p><b>Time-frequency masking:</b> [AvJ02], [AvJ04], [MGJ07], [Pul07]</p> <p><b>PCA:</b>[IrA02], [BVM06], [MGJ07], [GoJ07b], [BaS07], [God08], [JHS10], [BJP12], [TaG12], [TGC12], [LBP14]</p> <p><b>Least-squares:</b>[Fal06], [Fal07], [JPL10], [FaB11], [Uhh15]</p> <p><b>Linear estimation:</b> [HTG14]</p> <p><b>Ambient spectrum estimation:</b> [HGT15a], [HGT15b]</p> <p><b>Others:</b> [BrS08], [MeF10], [Har11]</p>	<p><b>LMS:</b> [UsB07]</p> <p><b>Shifted PCA:</b> [HTG13]</p> <p><b>Time-shifting:</b> [HGT15c]</p>	<p><b>PCA:</b> [DHT12], [HGT14], [HeG15]</p>
Multi-channel	<p><b>PCA:</b> [GoJ07b]</p> <p><b>Others:</b> [GoJ07a], [WaF11], [TGC12], [CCK14]</p>	<p><b>ICA &amp; time-frequency masking:</b> [SAM06]</p> <p><b>Pairwise correlations:</b> [TSW12]</p> <p><b>Pairing:</b> [HeG15b]</p> <p><b>Others:</b> [StM15]</p>	<p><b>ICA:</b> [HKO04]</p>
Single	<b>NMF:</b> [UWH07]	<b>Neural network:</b> [Uhp08]	

# My contributions

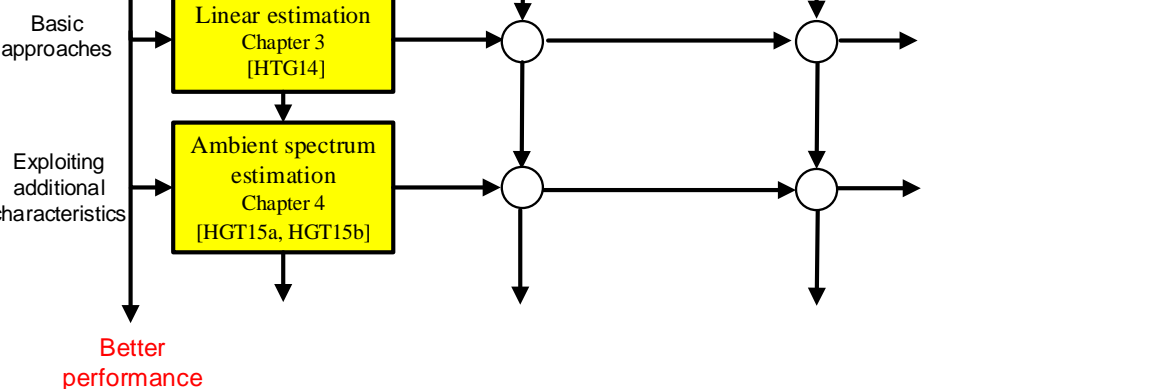


# My contributions

## Immersive Audio



Model: One source, amplitude panned  
 Model: One source, amplitude panned, and/or time shifted  
 Model: Multiple sources, amplitude panned, and/or time shifted  
 More robust in handling complex cases



**I:** Investigation of linear estimation based PAE under the basic signal model.

**II:** Improving PAE performance for strong ambient power cases using ambient spectrum estimation techniques.

**III:** Employing time-shifting techniques for PAE with partially correlated primary components

**IV:** Adaptation of conventional PAE approaches to deal with multiple sources.

**V:** Applying PAE in natural sound rendering headphone systems.

# My contributions in publications

**I:** Investigation of linear estimation based PAE under the basic signal model.

[J1] **J. He**, E. L. Tan, and W. S. Gan, "Linear estimation based primary-ambient extraction for stereo audio signals," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 22, no. 2, pp. 505-517, Feb. 2014.

**II:** Improving PAE performance for strong ambient power cases using ambient spectrum estimation techniques.

[J2] **J. He**, W. S. Gan, and E. L. Tan, "Primary-ambient extraction using ambient phase estimation with a sparsity constraint," *IEEE Signal Process. Letters*, vol. 22, no. 8, pp. 1127-1131, Aug. 2015.

[J3] **J. He**, E. L. Tan, and W. S. Gan, "Primary-ambient extraction using ambient spectrum estimation for immersive spatial audio reproduction," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 9, pp. 1430-1443, Sept. 2015.

**III:** Employing time-shifting techniques for PAE with partially correlated primary components

[J4] **J. He**, W. S. Gan, and E. L. Tan, "Time-shifting based primary-ambient extraction for spatial audio reproduction," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 10, pp. 1576-1588, Oct. 2015.

[C1] **J. He**, E. L. Tan, and W. S. Gan, "Time-shifted principal component analysis based cue extraction for stereo audio signals," in *Proc. ICASSP*, Vancouver, Canada, 2013, pp. 266-270.

**IV:** Adaptation of conventional PAE approaches to deal with multiple sources.

[C2] **J. He**, W. S. Gan, and E. L. Tan, "A study on the frequency-domain primary-ambient extraction for stereo audio signals," in *Proc. ICASSP*, Florence, Italy, 2014, pp. 2892-2896.

[C3] **J. He**, and W. S. Gan, "Multi-shift principal component analysis based primary component extraction for spatial audio reproduction," in *Proc. ICASSP*, Brisbane, Australia, Apr. 2015, pp. 350-354.

**V:** Applying PAE in natural sound rendering headphone systems.

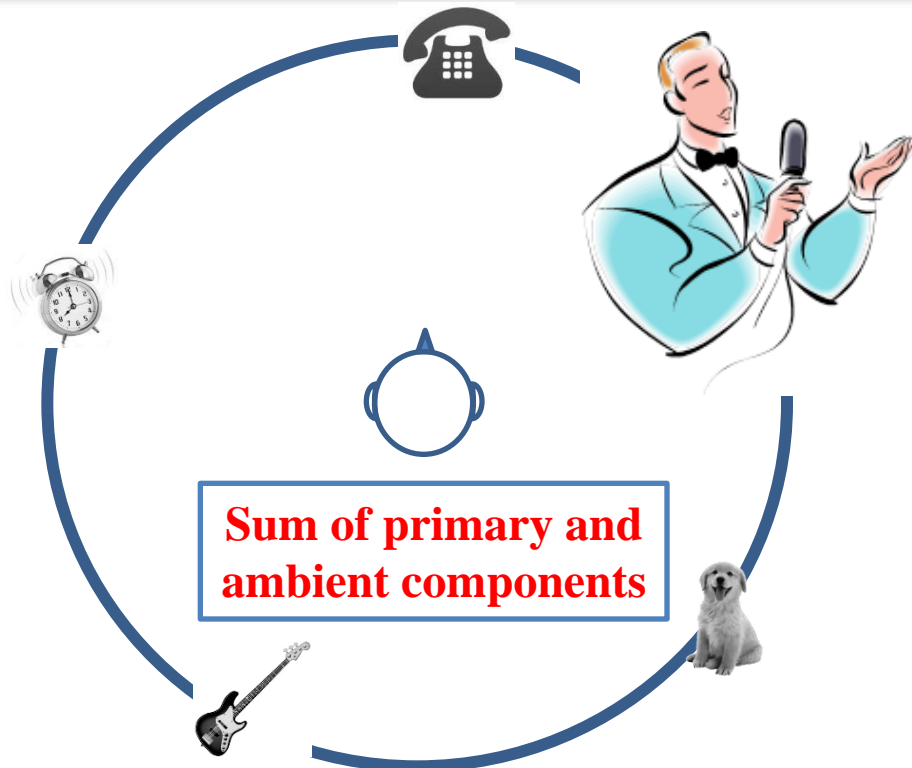
[J5] K. Sunder, **J. He**, E. L. Tan, and W. S. Gan, "Natural sound rendering for headphones: Integration of signal processing techniques," *IEEE Signal Process. Magazine*, vol. 32, no. 2, Mar 2015, pp. 100-113.

# Other related publications

- [C4] **J. He**, W. S. Gan, and E. L. Tan, “On the preprocessing and postprocessing of HRTF individualization based on sparse representation of anthropometry features,” in *Proc. ICASSP*, Brisbane, Australia, Apr. 2015, pp. 639-643.
- [C5] **J. He**, and W. S. Gan, “Applying primary ambient extraction for immersive spatial audio reproduction,” *2015 Asia Pacific Signal and Information Processing Association (APSIPA) Annual Summit and Conference (invited)*, Hong Kong, Dec. 2015.
- [C6] **J. He**, R. Ranjan, and W. S. Gan, “Fast continuous HRTF acquisition with unconstrained movements of human subjects,” in *Proc. ICASSP*, Shanghai, China, Mar. 2016, pp.
- **[Tutorial]** W. S. Gan, and **J. He**, “Assisted listening for headphones and hearing aids: signal processing techniques,” Tutorial at *APSIPA ASC 2015*, Hong Kong, Dec. 2015.
- **[Show & Tell]** D. H. Nguyen, **J. He**, K. K. Phyo, and W. S. Gan, “Real-time audio signal processing platform for natural 3D sound rendering,” Show & Tell at *ICASSP 2016*, Shanghai, China, Mar. 2016.
- [J6] ] **J. He**, R. Ranjan, W. S. Gan, and K. Sunder, “Scalable data reusing normalized LMS for acoustic system identification with short duration signals,” *IEEE Sig. Process. Letters*, under review.

# Objective of PAE

to extract the primary and ambient components from  $M$  mixtures



Mixtures = primary component + ambient component

$$x_m(n) = p_m(n) + a_m(n)$$



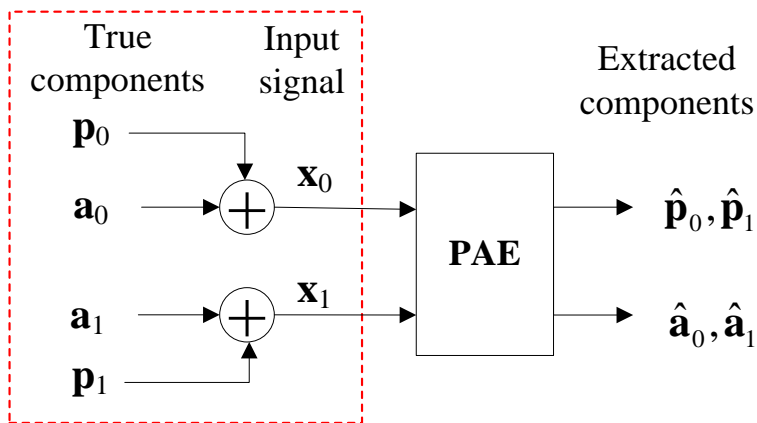
# Definitions with Stereo Signal Model

Signal = Primary + Ambient

$$\mathbf{x}_0 = \mathbf{p}_0 + \mathbf{a}_0$$

$$\mathbf{x}_1 = \mathbf{p}_1 + \mathbf{a}_1$$

Stereo signal model



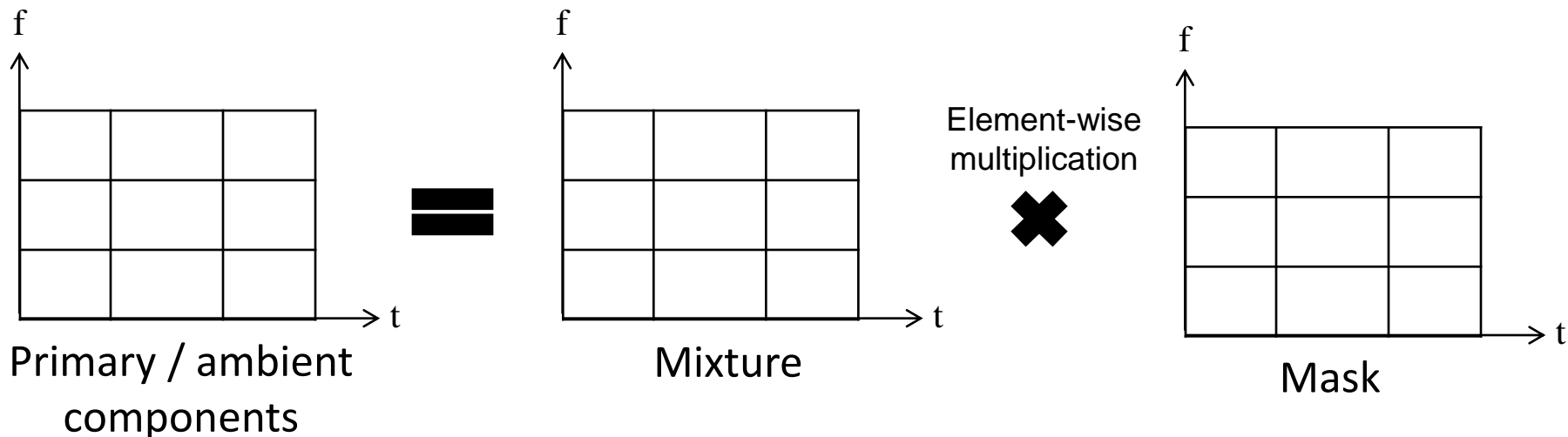
## Assumptions

Primary components highly correlated	$\mathbf{p}_1 = k\mathbf{p}_0$
Ambient components uncorrelated	$\mathbf{a}_0 \perp \mathbf{a}_1$
Primary ambient components uncorrelated	$\mathbf{p}_i \perp \mathbf{a}_j$ $\forall i, j \in \{0, 1\}$
Ambient power balanced	$P_{\mathbf{a}_0} = P_{\mathbf{a}_1}$

Primary panning factor PPF:  $k = \frac{\mathbf{p}_1}{\mathbf{p}_0}$

Primary power ratio PPR:  $\gamma = \frac{\text{Total primary power}}{\text{Total signal power}}$

# PAE: Time-Frequency Masking



## Mask can be constructed using

- Inter-channel coherence [Avendano and Jot, 2004]
- Pairwise correlation [Thompson et al., 2012]
- Equal level of ambience [Merimaa et al., 2007]
- Diffuseness [Pulkki, 2007]

# Linear Estimation framework in PAE

$$\begin{bmatrix} \hat{\mathbf{p}}_0^T \\ \hat{\mathbf{p}}_1^T \\ \hat{\mathbf{a}}_0^T \\ \hat{\mathbf{a}}_1^T \end{bmatrix} = \begin{bmatrix} w_{P0,0} & w_{P0,1} \\ w_{P1,0} & w_{P1,1} \\ w_{A0,0} & w_{A0,1} \\ w_{A1,0} & w_{A1,1} \end{bmatrix} \begin{bmatrix} \mathbf{x}_0^T \\ \mathbf{x}_1^T \end{bmatrix} = \mathbf{W} \begin{bmatrix} \mathbf{x}_0^T \\ \mathbf{x}_1^T \end{bmatrix}$$

# Performance measures

<p><b>Extraction Accuracy</b></p>	$\mathbf{e}_P = (w_{P0,0} + kw_{P0,1} - 1)\mathbf{p}_0 + 0 + (w_{P0,0}\mathbf{a}_0 + w_{P0,1}\mathbf{a}_1)$ $\mathbf{e}_A = (w_{A0,0} - 1)\mathbf{a}_0 + w_{A0,1}\mathbf{a}_1 + (w_{A0,0}\mathbf{p}_0 + w_{A0,1}\mathbf{p}_1)$ <div style="border: 1px solid black; height: 40px; width: 100%; margin: 10px 0;"></div> $\frac{P_{\text{Error}}}{P_{\text{True signal}}} \approx \frac{P_{\text{Distortion}}}{P_{\text{True signal}}} + \frac{P_{\text{Interference}}}{P_{\text{True signal}}} + \frac{P_{\text{Leakage}}}{P_{\text{True signal}}}$ $\mathbf{ESR} \approx \mathbf{DSR} + \mathbf{ISR} + \mathbf{LSR}$
<p><b>Spatial Accuracy</b></p>	<p><b>ICC, ICLD, ICTD</b>(only for primary)</p>

**ESR**: Error-to-signal ratio

**DSR**: Distortion-to-signal ratio

**ISR**: Interference-to-signal ratio

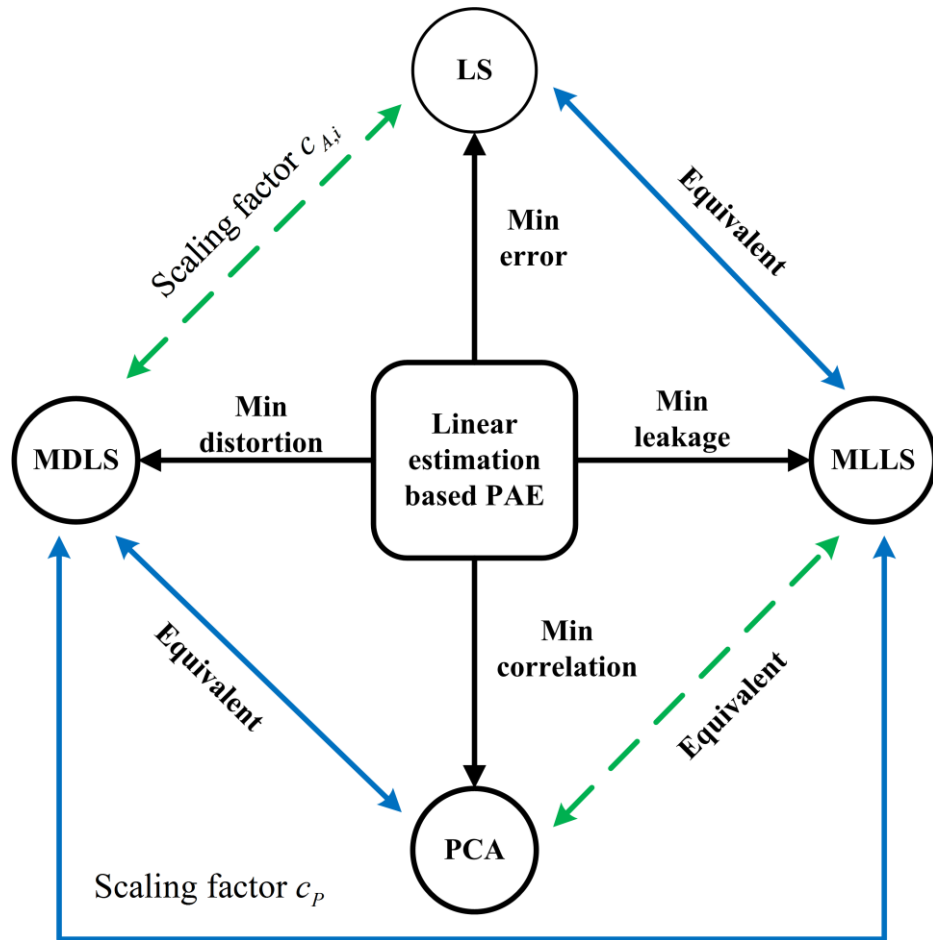
**LSR**: Leakage-to-signal ratio

**ICC** : Inter-channel cross-correlation coefficient

**ICLD**: Inter-channel level difference

**ICTD**: Inter-channel time difference

# PAE: Linear Estimation



$$\begin{bmatrix} \hat{p}_0(n) \\ \hat{p}_1(n) \\ \hat{a}_0(n) \\ \hat{a}_1(n) \end{bmatrix} = \begin{bmatrix} w_{P0,0} & w_{P0,1} \\ w_{P1,0} & w_{P1,1} \\ w_{A0,0} & w_{A0,1} \\ w_{A1,0} & w_{A1,1} \end{bmatrix} \begin{bmatrix} x_0(n) \\ x_1(n) \end{bmatrix}$$

## Objectives and relationships of four linear estimation based PAE approaches.

- **Blue** solid lines represent the relationships in the **primary** component;
- **Green** dotted lines represent the relationships in the **ambient** component.
- **MLLS**: minimum leakage LS
- **MDLS**: minimum distortion LS

# Performance of the four PAE approaches

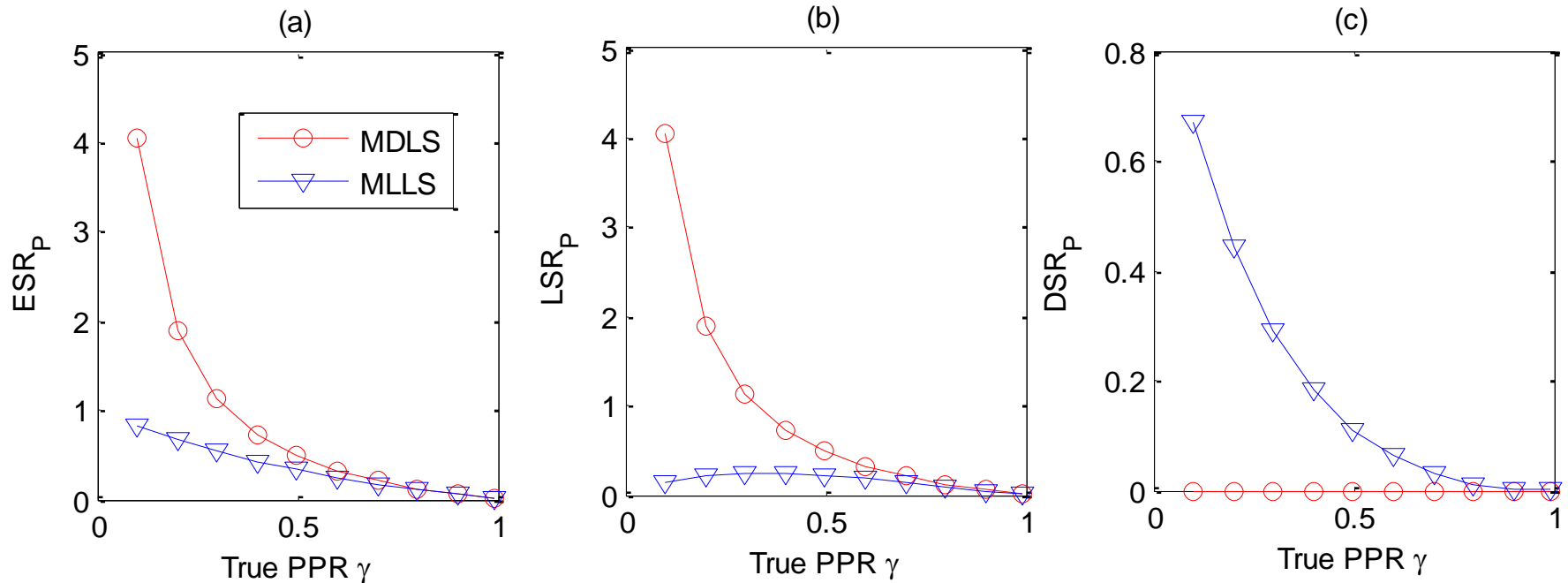
$k$  : PPF

$\gamma$  : PPR

Performance Measure	Primary component		Ambient component		
	MDLS/PCA	MLLS/LS	MLLS/PCA	LS	MDLS
ESR	$\frac{1-\gamma}{2\gamma}$	$\frac{1-\gamma}{1+\gamma}$	$\frac{1}{1+k^2}$	$\frac{1}{1+k^2} \frac{2\gamma}{1+\gamma}$	$\frac{2\gamma}{1+k^2+(k^2-1)\gamma}$
LSR	$\frac{1-\gamma}{2\gamma}$	$\frac{1-\gamma}{2\gamma} \left( \frac{2\gamma}{1+\gamma} \right)^2$	<b>0</b>	$\frac{1}{1+k^2} \frac{2\gamma(1-\gamma)}{(1+\gamma)^2}$	$\frac{(1+k^2)(1-\gamma)2\gamma}{[(1+k^2)(1+\gamma)-2\gamma]^2}$
DSR	<b>0</b>	$\left( \frac{1-\gamma}{1+\gamma} \right)^2$	$\left( \frac{1}{1+k^2} \right)^2$	$\left( \frac{1}{1+k^2} \frac{2\gamma}{1+\gamma} \right)^2$	<b>0</b>
ISR	<b>0</b>		$\left( \frac{k}{1+k^2} \right)^2$	$\left( \frac{k}{1+k^2} \frac{2\gamma}{1+\gamma} \right)^2$	$\left[ \frac{2k\gamma}{(1+k^2)(1+\gamma)-2\gamma} \right]^2$
ICC(ICTD)	1(0)		<b>1</b>	$\frac{2k\gamma}{\sqrt{(1+k^2)^2 - (1-k^2)^2 \gamma^2}}$	
ICLD	$k^2$		$\frac{1}{k^2}$	$\frac{1}{k^2} \frac{1+\gamma+k^2(1-\gamma)}{1+\gamma+\frac{1}{k^2}(1-\gamma)}$	$\frac{1}{k^2} \frac{1-\gamma+k^2(1+\gamma)}{1-\gamma+\frac{1}{k^2}(1+\gamma)}$

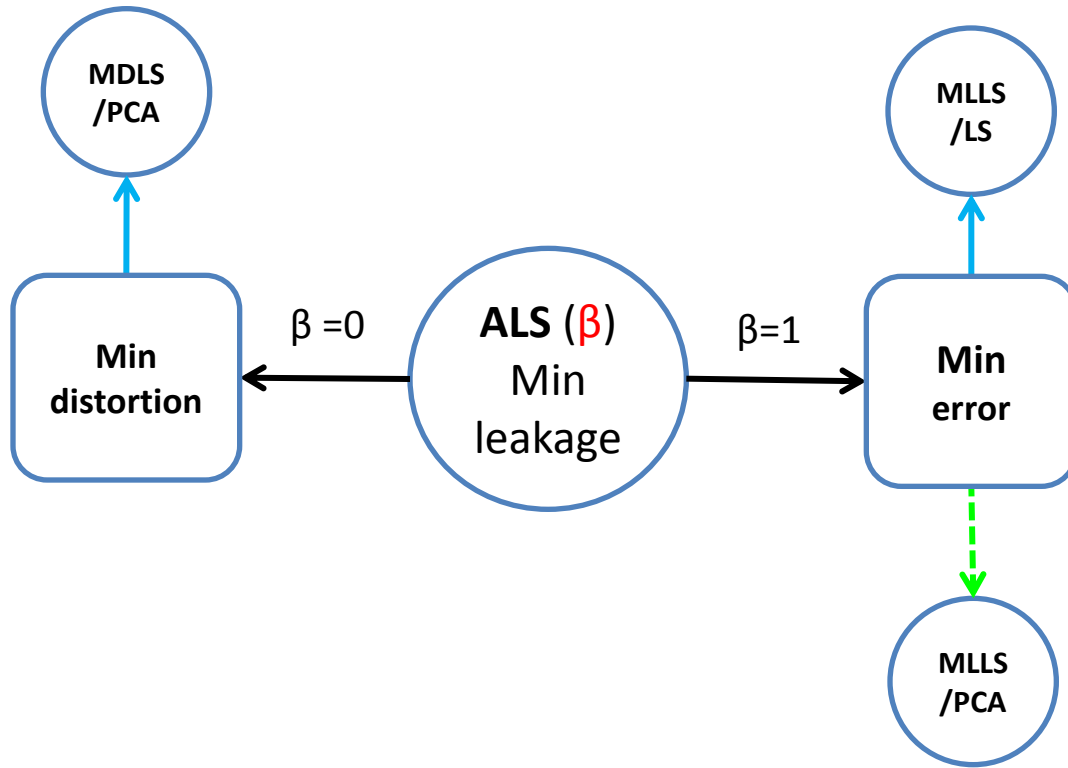
Theoretical results are verified in our experiments!

# Experimental results– Primary extraction



Performance of MDLS (or PCA) and MLLS (or LS) in primary extraction  
(a) ESR, (b) LSR, (c) DSR.

# PAE using Adjustable Least-squares (ALS)



Adjustable factor  $\beta \in [0,1]$

$$\mathbf{W}_{\text{ALS}} = \begin{bmatrix} \frac{1}{1+k^2} \left(1 - \beta \frac{1-\gamma}{1+\gamma}\right) & \frac{k}{1+k^2} \left(1 - \beta \frac{1-\gamma}{1+\gamma}\right) \\ \frac{k}{1+k^2} \left(1 - \beta \frac{1-\gamma}{1+\gamma}\right) & \frac{k^2}{1+k^2} \left(1 - \beta \frac{1-\gamma}{1+\gamma}\right) \\ 1 - \beta \frac{1}{1+k^2} & -\frac{1}{k} \left(1 - \beta \frac{1}{1+k^2}\right) \\ -k \left(1 - \beta \frac{k^2}{1+k^2}\right) & 1 - \beta \frac{k^2}{1+k^2} \end{bmatrix}$$

$$\text{ESR}_P = \frac{1-\gamma}{2\gamma} + \beta(\beta-2) \frac{(1-\gamma)^2}{2\gamma(1+\gamma)}, \quad \text{DSR}_P = \beta^2 \left(\frac{1-\gamma}{1+\gamma}\right)^2, \quad \text{LeSR}_P = \frac{1-\gamma}{1+\gamma},$$

$$\text{ESR}_A = \frac{1}{k^2} + \beta(\beta-2) \frac{1}{k^2(k^2+1)}, \quad \text{DSR}_A = \beta^2 \left(\frac{1}{1+k^2}\right)^2, \quad \text{LeSR}_A = 0.$$

Blue solid lines: primary component; Green dotted lines: ambient component.

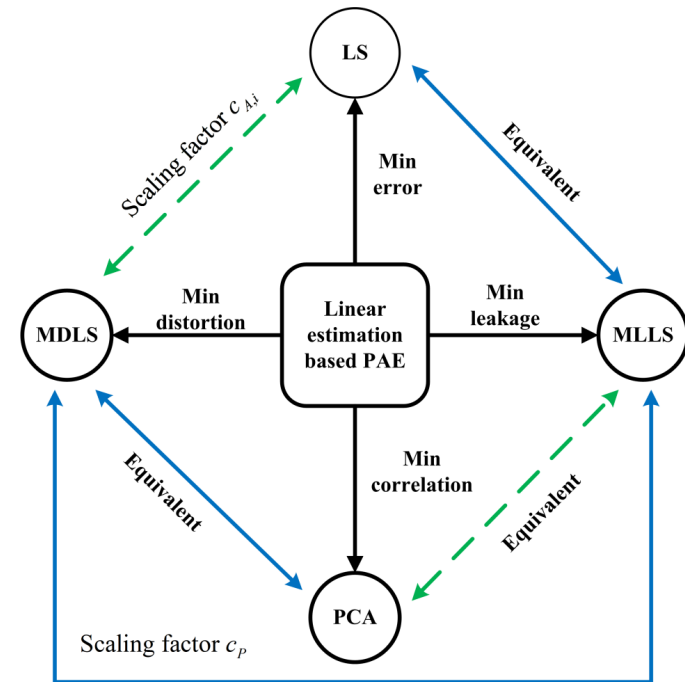


# Recommendations of Linear Estimation based PAE

Approach	Strengths	Weaknesses	Recommendations
PCA	<ul style="list-style-type: none"> <li>No distortion in the extracted primary component;</li> <li>No primary leakage in the extracted ambient component;</li> <li>Primary and ambient components are uncorrelated;</li> </ul>	<ul style="list-style-type: none"> <li>Ambient component severely panned;</li> </ul>	Spatial audio coding and interactive audio in gaming, where the <b>primary component is more important</b> than the ambient component.
LS	<ul style="list-style-type: none"> <li>Minimum MSE in the extracted primary and ambient components;</li> </ul>	<ul style="list-style-type: none"> <li>Severe primary leakage in the extracted ambient component;</li> </ul>	Applications in which both the primary and ambient components are extracted, processed, and <b>finally mixed together</b> .
MLLS	<ul style="list-style-type: none"> <li>Minimum leakage in the extracted primary and ambient components;</li> <li>Primary and ambient components are uncorrelated;</li> </ul>	<ul style="list-style-type: none"> <li>Ambient component severely panned;</li> </ul>	Spatial audio enhancement systems, and applications in which <b>different rendering</b> or playback techniques are employed on the extracted primary and ambient components.
MDLS	<ul style="list-style-type: none"> <li>No distortion in the extracted primary and ambient components;</li> </ul>	<ul style="list-style-type: none"> <li>Severe interference and primary leakage in the extracted ambient component;</li> </ul>	High-fidelity applications in which <b>timbre</b> is of high importance.
ALS	<ul style="list-style-type: none"> <li><b>Performance adjustable;</b></li> </ul>	<ul style="list-style-type: none"> <li>Need to adjust the value of the adjustable factor;</li> </ul>	For applications without explicit requirements.

# Contributions on linear estimation based PAE

1. Formulated the linear estimation framework for PAE.
2. Introduced two groups of performance measures.
  - Extraction accuracy: ESR, DSR, ISR, LSR
  - Spatial accuracy : ICC, ICTD, ICLD
3. Proposed MLLS, MDLS, ALS and compared them with PCA and LS in PAE.
4. Different approaches are recommended in different spatial audio applications.



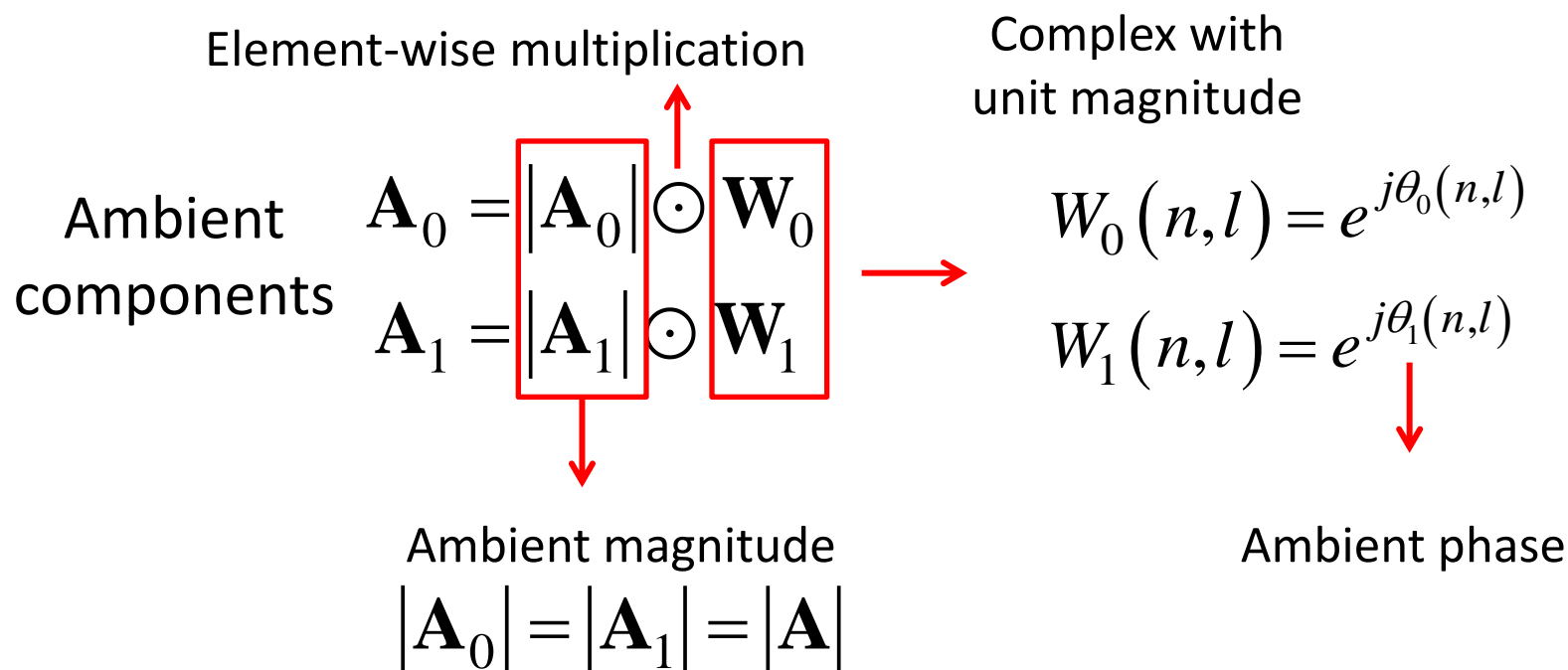
[J1] J. He, E. L. Tan and W. S. Gan, "Linear estimation based primary-ambient extraction for stereo audio signals," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 22, no. 2, pp. 505-517, Feb. 2014.

# Observation

Ambient components are decorrelated using:

- A small delay
- All-pass filtering with random phase
- Artificial reverberation, and binaural reverberation

Magnitude of the ambient components are usually kept the same.



# From PAE to APE: ambient phase estimation

X: Mixed signal  
P: Primary components  
k: Primary panning factor

PAE

$$\mathbf{X}_0 = \mathbf{P}_0 + \mathbf{A}_0$$

$$\mathbf{X}_1 = \mathbf{P}_1 + \mathbf{A}_1$$

with

$$\mathbf{P}_1 = k\mathbf{P}_0, \mathbf{A}_0 \perp \mathbf{A}_1$$

$$\mathbf{P}_i \perp \mathbf{A}_j, \forall i, j \in \{0,1\}$$

$$P_{P_1} = k^2 P_{P_0}, P_{A_1} = P_{A_0}$$



APE

$$|\mathbf{A}| = (\mathbf{X}_1 - k\mathbf{X}_0) ./ (\mathbf{W}_1 - k\mathbf{W}_0)$$

and

$$\mathbf{A}_c = (\mathbf{X}_1 - k\mathbf{X}_0) ./ (\mathbf{W}_1 - k\mathbf{W}_0) \odot \mathbf{W}_c,$$

$$\mathbf{P}_c = \mathbf{X}_c - (\mathbf{X}_1 - k\mathbf{X}_0) ./ (\mathbf{W}_1 - k\mathbf{W}_0) \odot \mathbf{W}_c$$

$$\forall \text{channel index } c \in \{0,1\}.$$

Find  $\mathbf{W}_0, \mathbf{W}_1$

Find  $\theta_0, \theta_1$

Find  $\theta_1$

$$\theta_0 = \theta + \arcsin[k^{-1} \sin(\theta - \theta_1)] + \pi$$

$$\text{where } \theta = \angle(\mathbf{X}_1 - k\mathbf{X}_0)$$

# From PAE to APE: ambient phase estimation

PAE



APE

$$\mathbf{X}_0 = \mathbf{P}_0 + \mathbf{A}_0$$

$$\mathbf{X}_1 = \mathbf{P}_1 + \mathbf{A}_1$$

$$|\mathbf{A}| = (\mathbf{X}_1 - k\mathbf{X}_0) ./ (\mathbf{W}_1 - k\mathbf{W}_0)$$



**Is this the only way?**

For each  
time-frequency  
bin

$$\begin{aligned} kX_0 &= P_1 + kA_0 \\ X_1 &= P_1 + A_1 \end{aligned}$$



Complex spectrum  
shown in  
complex plane?

# From PAE to AME: ambient magnitude estimation

Put the complex spectrum in  
2 dimensional vector form

$$k \vec{X}_0 = \vec{P}_1 + k \vec{A}_0$$

$$\vec{X}_1 = \vec{P}_1 + \vec{A}_1$$



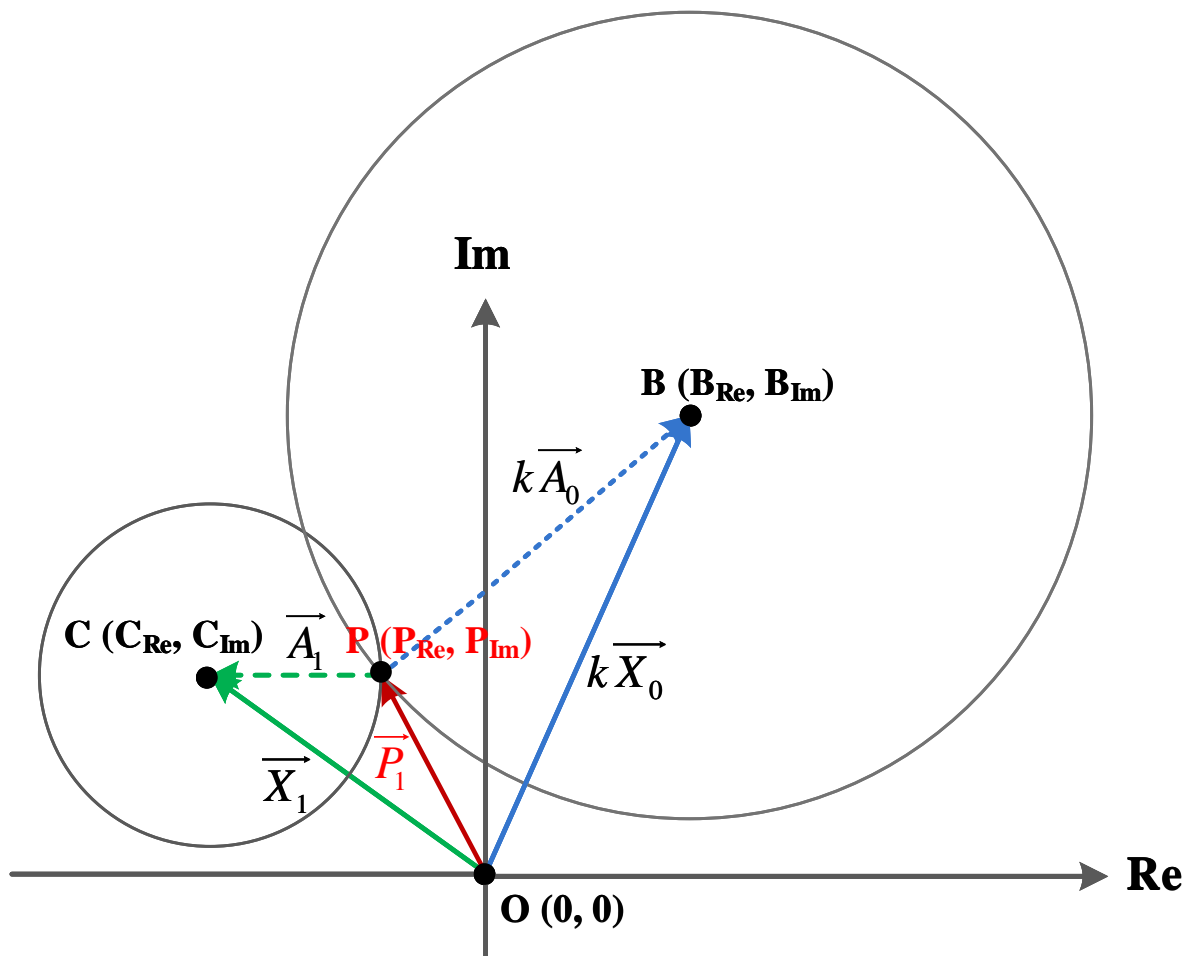
$$k \vec{X}_0 = \vec{OB} = (B_{\text{Re}}, B_{\text{Im}})$$

$$\vec{X}_1 = \vec{OC} = (C_{\text{Re}}, C_{\text{Im}})$$

$$\vec{P}_1 = \vec{OP} = (P_{\text{Re}}, P_{\text{Im}})$$

$$k \vec{A}_0 = \vec{PB}$$

$$\vec{A}_1 = \vec{PC}$$



# From PAE to AME: ambient magnitude estimation

$$r = |\vec{A}_0| = |\vec{A}_1| \longrightarrow |\vec{PC}| = r, \quad |\vec{PB}| = kr.$$

$$(P_{\text{Re}} - B_{\text{Re}})^2 + (P_{\text{Im}} - B_{\text{Im}})^2 = k^2 \hat{r}^2,$$

$$(P_{\text{Re}} - C_{\text{Re}})^2 + (P_{\text{Im}} - C_{\text{Im}})^2 = \hat{r}^2.$$

$$\hat{P}_{\text{Re}} = \frac{B_{\text{Re}} + C_{\text{Re}}}{2} + \frac{(C_{\text{Re}} - B_{\text{Re}})(k^2 - 1)\hat{r}^2 \pm (B_{\text{Im}} - C_{\text{Im}})\beta}{2|\overline{BC}|^2},$$

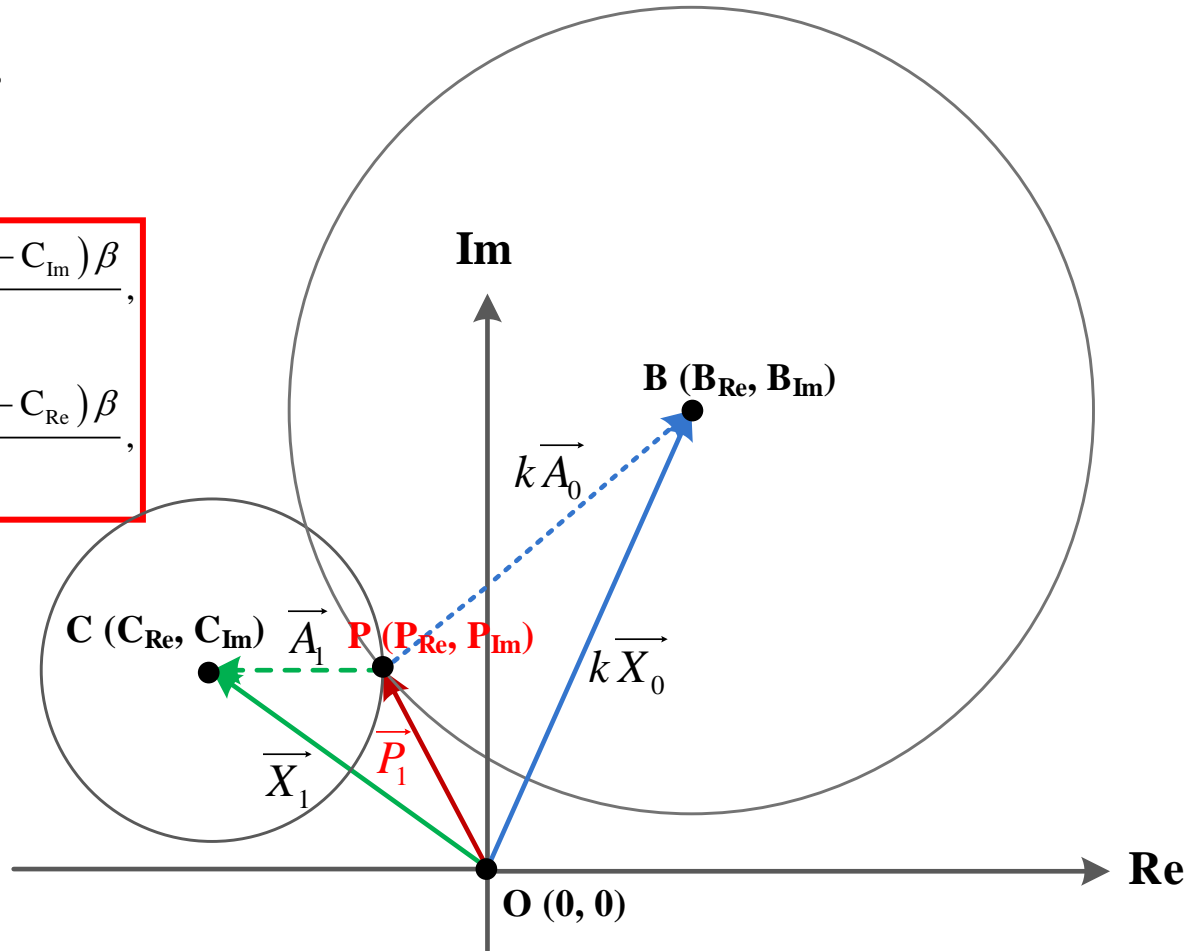
$$\hat{P}_{\text{Im}} = \frac{B_{\text{Im}} + C_{\text{Im}}}{2} + \frac{(C_{\text{Im}} - B_{\text{Im}})(k^2 - 1)\hat{r}^2 \mp (B_{\text{Re}} - C_{\text{Re}})\beta}{2|\overline{BC}|^2},$$

$$\hat{P}_1 = \hat{P}_{\text{Re}} + j\hat{P}_{\text{Im}}, \quad \hat{P}_0 = k^{-1}(\hat{P}_{\text{Re}} + j\hat{P}_{\text{Im}}),$$

$$\hat{A}_1 = X_1 - (\hat{P}_{\text{Re}} + j\hat{P}_{\text{Im}}),$$

$$\hat{A}_0 = X_0 - k^{-1}(\hat{P}_{\text{Re}} + j\hat{P}_{\text{Im}}).$$

**Find  $r$**



$$|\overline{BC}| = \sqrt{(C_{\text{Re}} - B_{\text{Re}})^2 + (C_{\text{Im}} - B_{\text{Im}})^2},$$

$$\beta = \sqrt{\left[ (k+1)^2 \hat{r}^2 - |\overline{BC}|^2 \right] \left[ (k-1)^2 \hat{r}^2 - |\overline{BC}|^2 \right]}.$$

# APE $\leftrightarrow$ AME

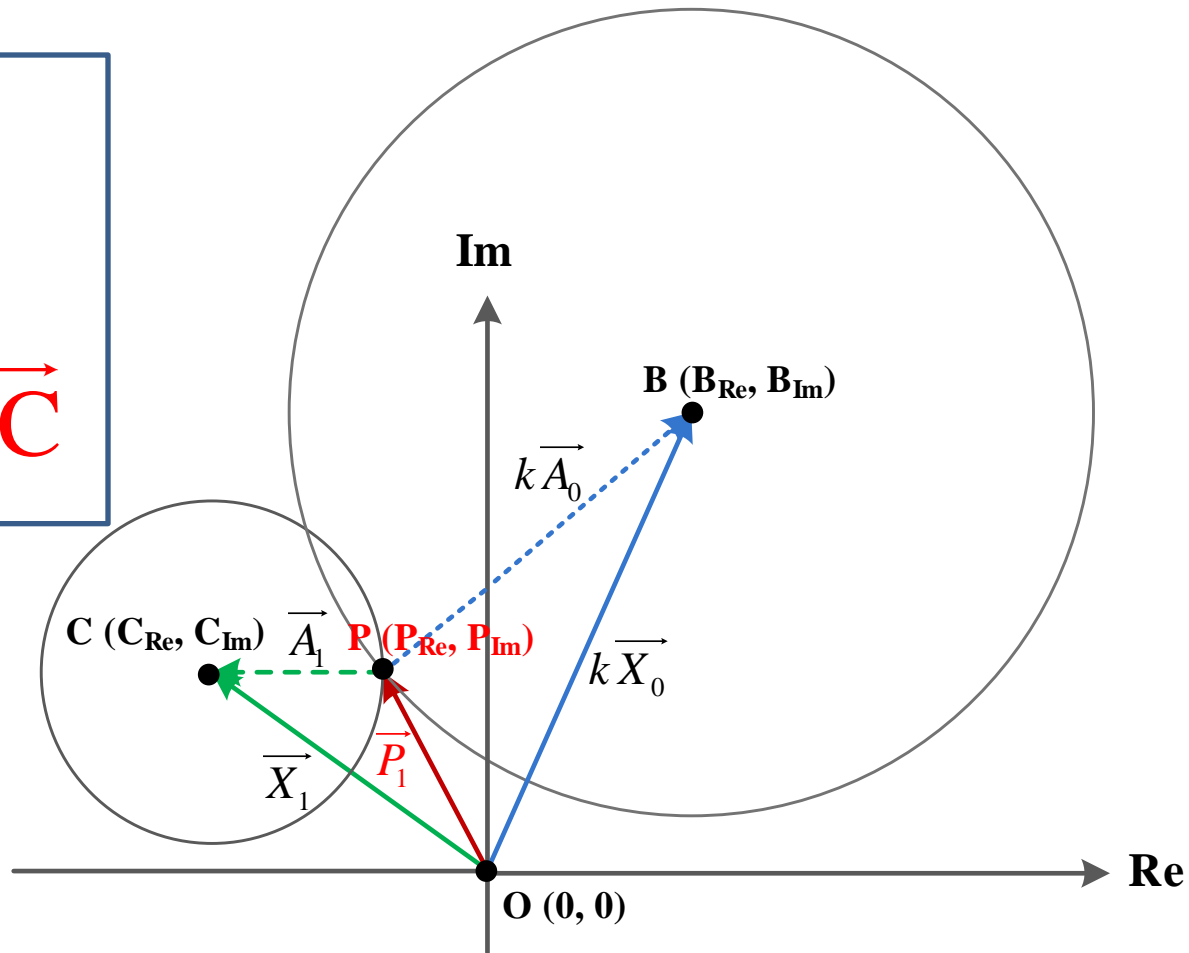
AME: Find  $r = |\overrightarrow{PC}|$



APE: Find  $\theta_1 = \angle \overrightarrow{PC}$

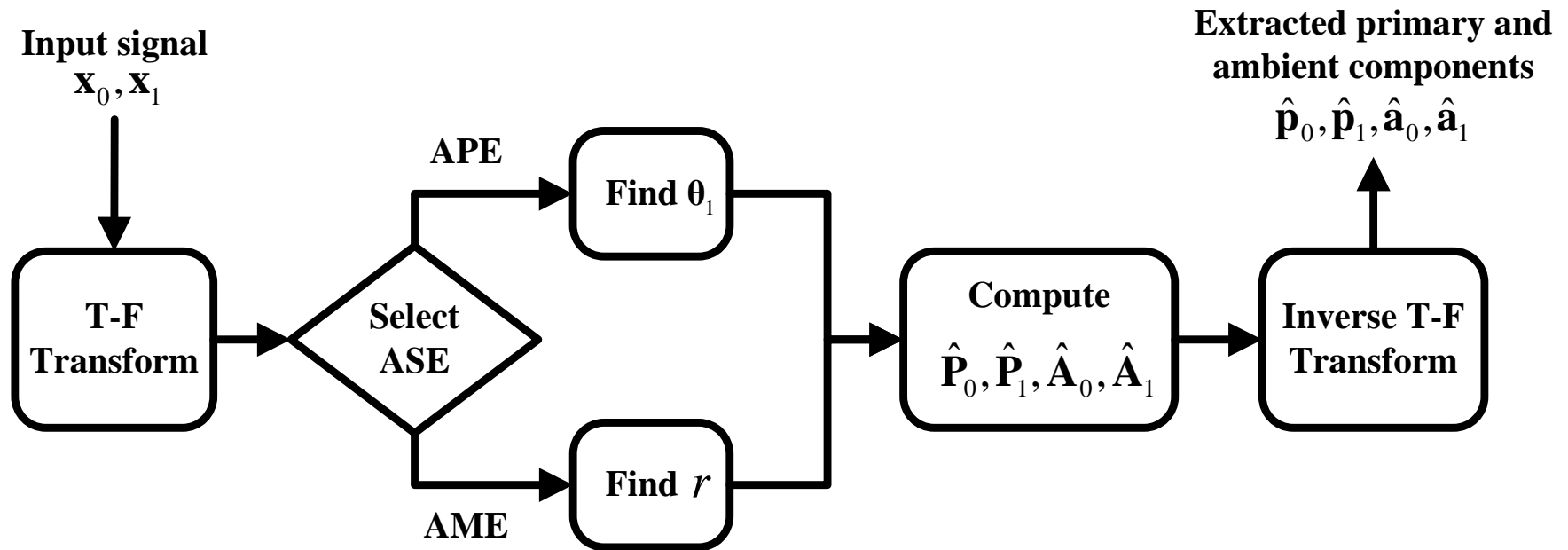


**A**mbient  
**S**pectrum  
**E**stimation





# Ambient Spectrum Estimation (ASE)



# ASE: how to estimate? Using a Sparsity Constraint

Objective

Find  $\theta_1$

Find  $r$

How

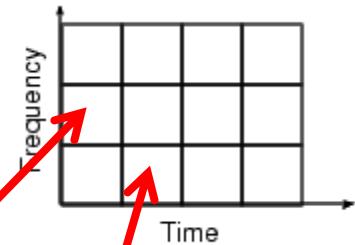
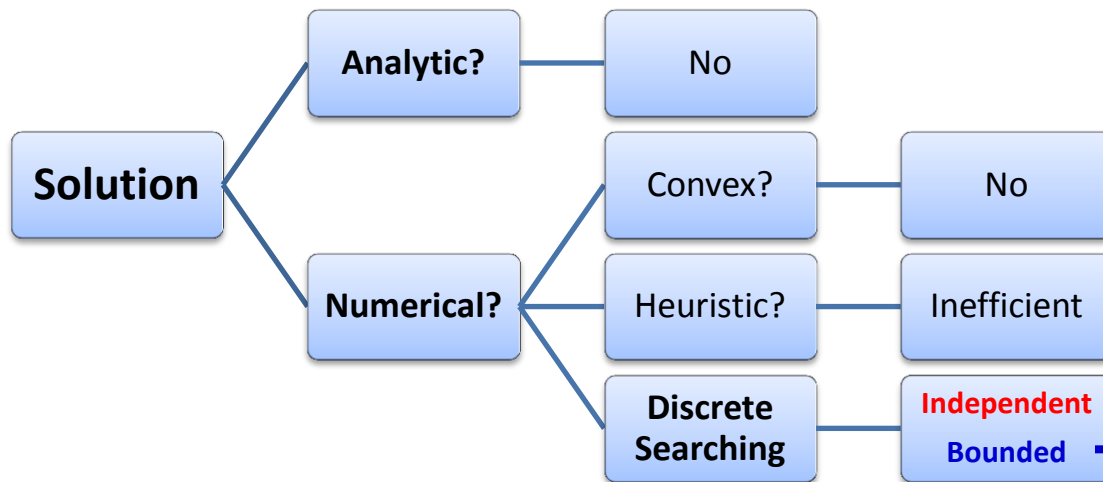
$$\hat{\theta}_1^* = \arg \min_{\hat{\theta}_1} \|\hat{\mathbf{P}}_1\|_1$$

$$\hat{\mathbf{r}}^* = \arg \min_{\hat{\mathbf{r}}} \|\hat{\mathbf{P}}_1\|_1$$

Method

APES

AMES



$$|\hat{P}_1(n_i, l_i)| \perp A_1(n_j, l_j) \quad \forall n_i \neq n_j \text{ or } l_i \neq l_j$$

$$\hat{\theta}_1 \in [-\pi, \pi]$$

$$\hat{\mathbf{r}} \in [r_{lb}, r_{ub}]$$

# Approximate solution APEX and computational cost

Approximate close-form solution APEX :

$$\hat{\boldsymbol{\theta}}_1^* = \begin{cases} \angle \mathbf{X}_1 & , \forall k > 1 \\ \angle (\mathbf{X}_1 - \mathbf{X}_0), \forall k = 1 \end{cases}$$

## Computation cost per time-frequency bin

Operations	Square root	Addition	Multiplication	Division	Comparison	Trigonometric operation
APES	D	15D+18	15D+13	4D+6	D-1	<b>7D+6</b>
AMES	2D+2	25D+35	24D+24	9D+13	D-1	<b>0</b>
APEX	0	13	7	4	1	<b>7</b>

D: number of phase or magnitude estimates in discrete searching

# Evaluation of PAE- Extraction accuracy

Performance measure: **Error-to-Signal Ratio (ESR)**

$$\text{ESR}_P = 10 \log_{10} \left\{ \frac{1}{2} \sum_{c=0}^1 \frac{\|\hat{\mathbf{p}}_c - \mathbf{p}_c\|_2^2}{\|\mathbf{p}_0\|_2^2} \right\}, \quad \text{ESR}_A = 10 \log_{10} \left\{ \frac{1}{2} \sum_{c=0}^1 \frac{\|\hat{\mathbf{a}}_c - \mathbf{a}_c\|_2^2}{\|\mathbf{a}_c\|_2^2} \right\}.$$

**Error** = **Distortion** + **Interference** + **Leakage**

**ESR**  $\approx$  **DSR** + **ISR** + **LSR**

When there is no analytic solution, how to compute these measures?

**We propose an optimization technique to compute these measures.**

# Optimization method for PAE Extraction accuracy

$$\begin{aligned}\hat{\mathbf{p}}_c &= \mathbf{p}_c + Leak_{\mathbf{p}_c} + Dist_{\mathbf{p}_c} \\ &= \mathbf{p}_c + (w_{Pc,0}\mathbf{a}_0 + w_{Pc,1}\mathbf{a}_1) + Dist_{\mathbf{p}_c},\end{aligned}$$

$$\begin{aligned}(w_{Pc,0}^*, w_{Pc,1}^*) &= \\ \arg \min_{(w_{Pc,0}, w_{Pc,1})} &\left\| \hat{\mathbf{p}}_c - \mathbf{p}_c - (w_{Pc,0}\mathbf{a}_0 + w_{Pc,1}\mathbf{a}_1) \right\|_2^2,\end{aligned}$$

$$LSR_P = 10 \log_{10} \left\{ \frac{1}{2} \sum_{c=0}^1 \frac{\|w_{Pc,0}^*\mathbf{a}_0 + w_{Pc,1}^*\mathbf{a}_1\|_2^2}{\|\mathbf{p}_c\|_2^2} \right\},$$

$$DSR_P = 10 \log_{10} \left\{ \frac{1}{2} \sum_{c=0}^1 \frac{\|\hat{\mathbf{p}}_c - \mathbf{p}_c - (w_{Pc,0}^*\mathbf{a}_0 + w_{Pc,1}^*\mathbf{a}_1)\|_2^2}{\|\mathbf{p}_c\|_2^2} \right\}.$$

$$\begin{aligned}\hat{\mathbf{a}}_c &= \mathbf{a}_c + Leak_{\mathbf{a}_c} + Intf_{\mathbf{a}_c} + Dist_{\mathbf{a}_c} \\ &= \mathbf{a}_c + w_{Ac,c}\mathbf{p}_c + w_{Ac,1-c}\mathbf{a}_{1-c} + Dist_{\mathbf{a}_c},\end{aligned}$$

$$\begin{aligned}(w_{Ac,c}^*, w_{Ac,1-c}^*) &= \\ \arg \min_{(w_{Ac,c}, w_{Ac,1-c})} &\left\| \hat{\mathbf{a}}_c - \mathbf{a}_c - (w_{Ac,c}\mathbf{p}_c + w_{Ac,1-c}\mathbf{a}_{1-c}) \right\|_2^2,\end{aligned}$$

$$LSR_A = 10 \log_{10} \left\{ \frac{1}{2} \sum_{c=0}^1 \frac{\|w_{Ac,c}^*\mathbf{p}_c\|_2^2}{\|\mathbf{a}_c\|_2^2} \right\},$$

$$ISR_A = 10 \log_{10} \left\{ \frac{1}{2} \sum_{c=0}^1 \frac{\|w_{Ac,1-c}^*\mathbf{a}_{1-c}\|_2^2}{\|\mathbf{a}_c\|_2^2} \right\},$$

$$DSR_A = 10 \log_{10} \left\{ \frac{1}{2} \sum_{c=0}^1 \frac{\|\hat{\mathbf{a}}_c - \mathbf{a}_c - (w_{Ac,c}\mathbf{p}_c + w_{Ac,1-c}\mathbf{a}_{1-c})\|_2^2}{\|\mathbf{a}_c\|_2^2} \right\}.$$

# Objective evaluation

## Stimuli

- Primary component:
  - Speech,  $k = 2$
- Ambient component:
  - Wave lapping sound
- Primary power ratio (PPR):
  - (0, 1) at an interval of 0.1
- FFT size: 4096

## Performance evaluated

1. **Extraction accuracy:** ESR
2. **Spatial accuracy:** ICC, ICLD

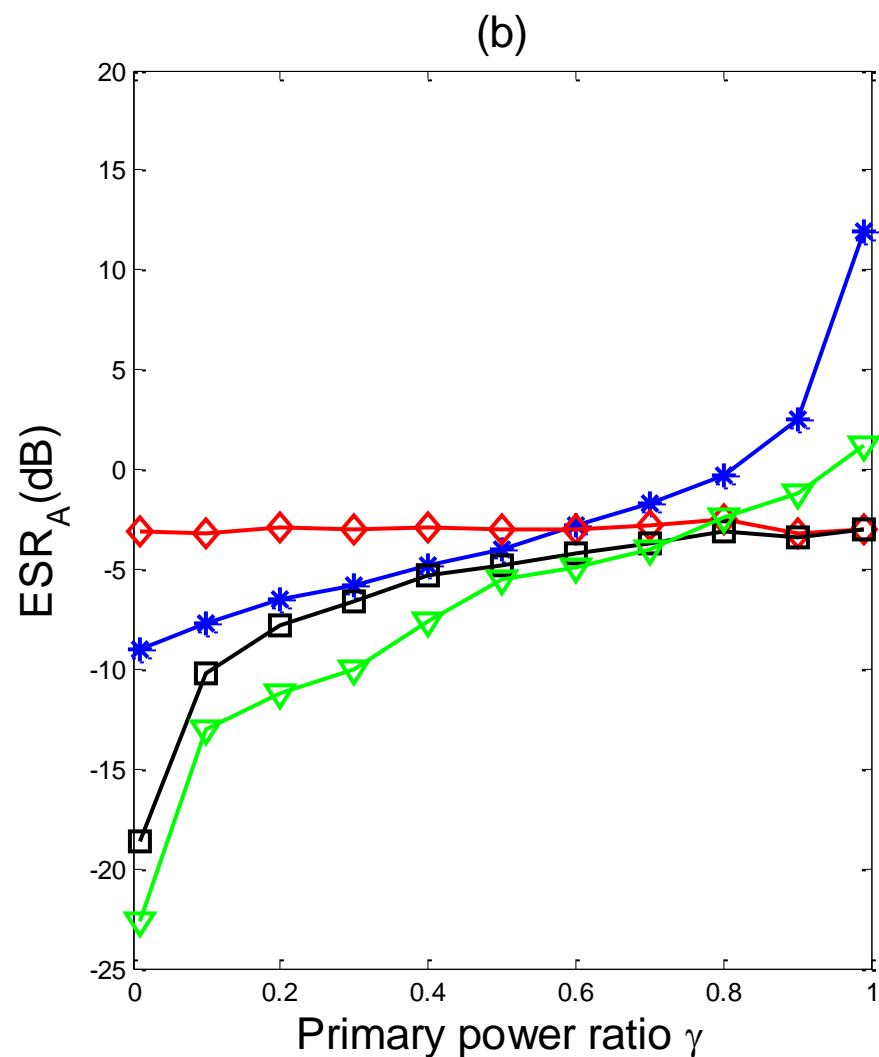
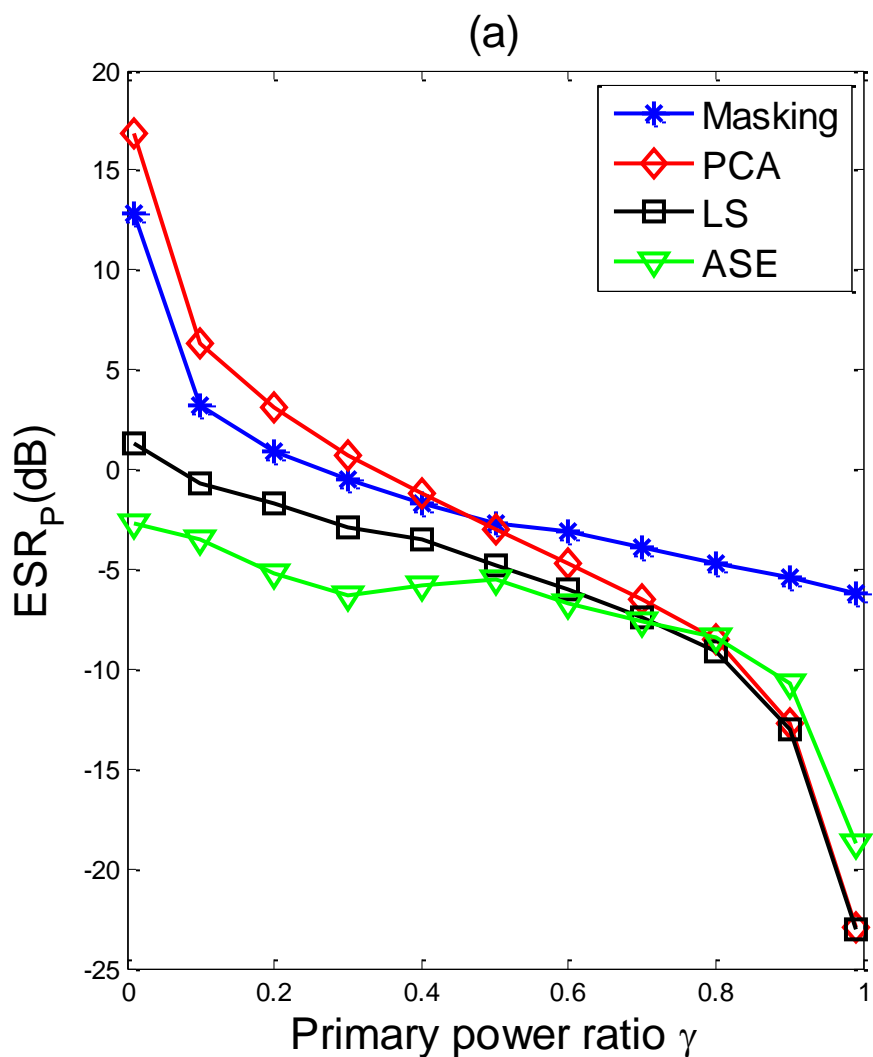
## Approaches compared

- Masking
- PCA
- LS
- ASE (APEX)

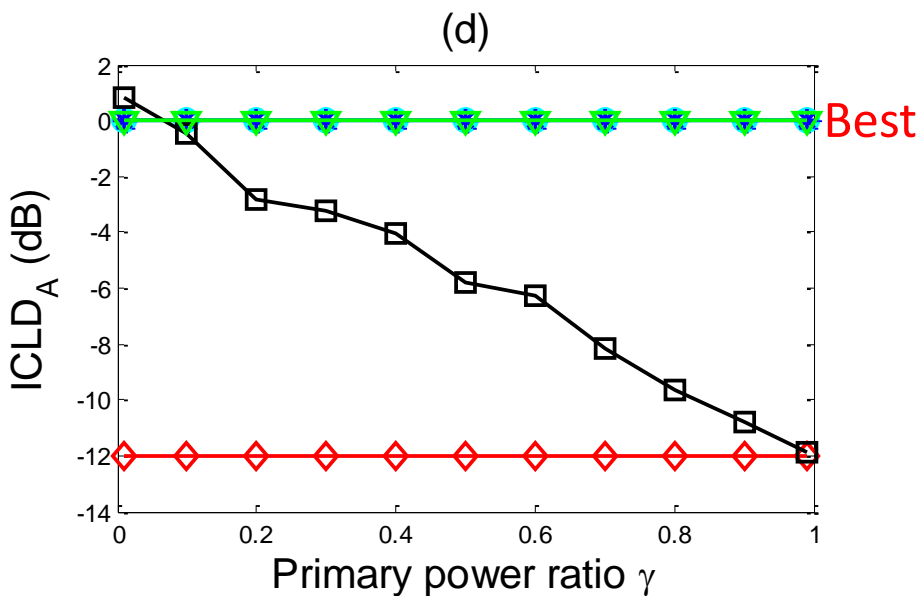
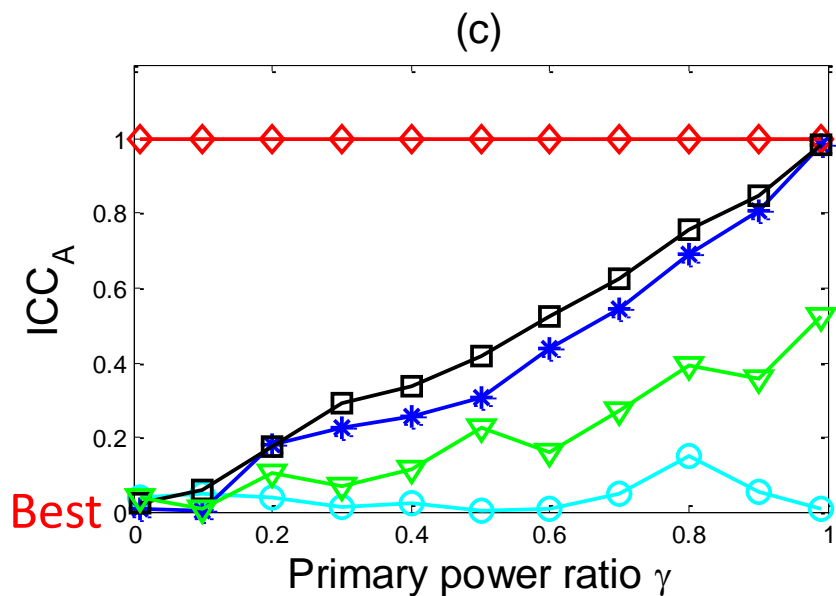
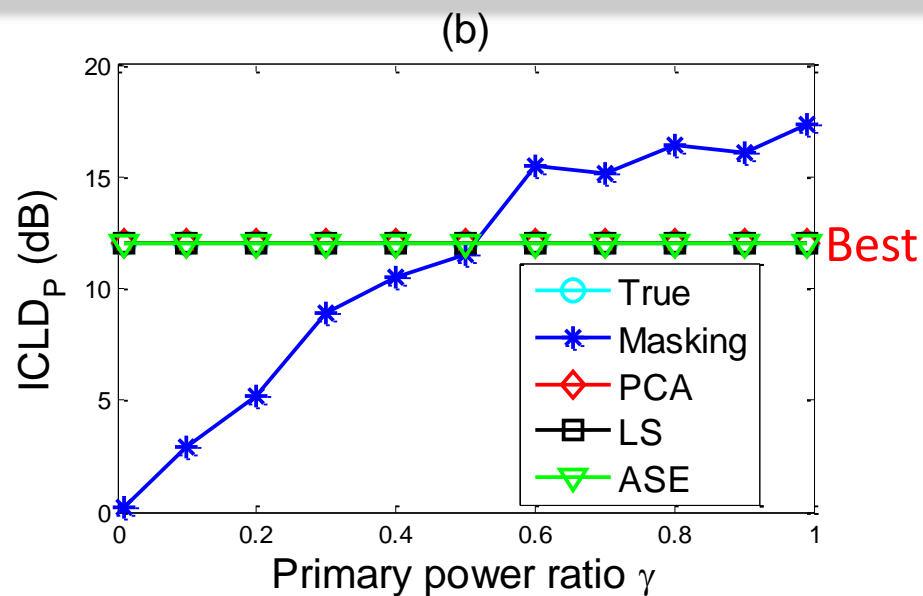
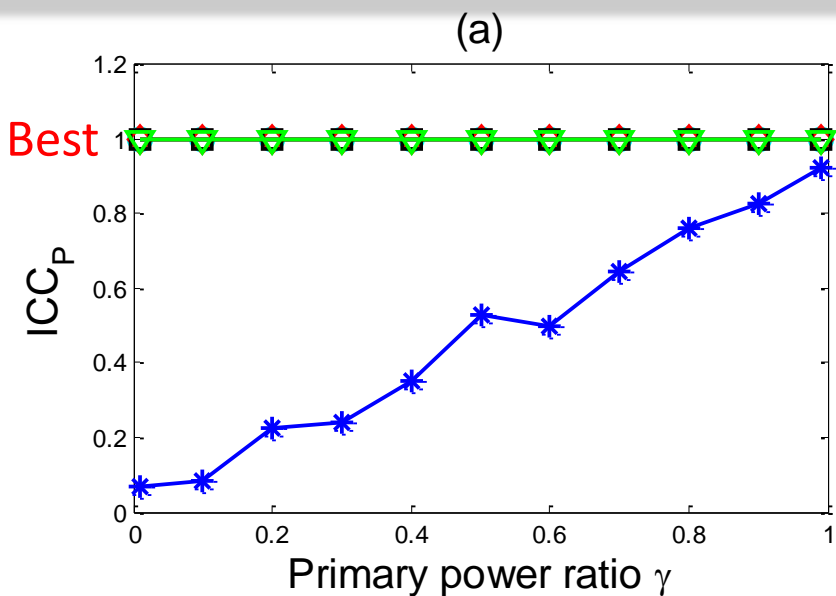
$$ESR_P = 10 \log_{10} \left\{ \frac{1}{2} \sum_{c=0}^1 \frac{\|\hat{\mathbf{p}}_c - \mathbf{p}_c\|_2^2}{\|\mathbf{p}_c\|_2^2} \right\},$$

$$ESR_A = 10 \log_{10} \left\{ \frac{1}{2} \sum_{c=0}^1 \frac{\|\hat{\mathbf{a}}_c - \mathbf{a}_c\|_2^2}{\|\mathbf{a}_c\|_2^2} \right\}.$$

# Extraction accuracy



# Spatial accuracy





# Subjective evaluation

## Stimuli

- Primary component:
  - speech, music, and bee sound,  $k = 2$
- Ambient component:
  - forest, canteen, and waterfall sound
- Primary power ratio (PPR):
  - (0.3, 0.7)
- Duration: 2-4 seconds

## Performance evaluated

1. Extraction accuracy
2. Ambient diffuseness

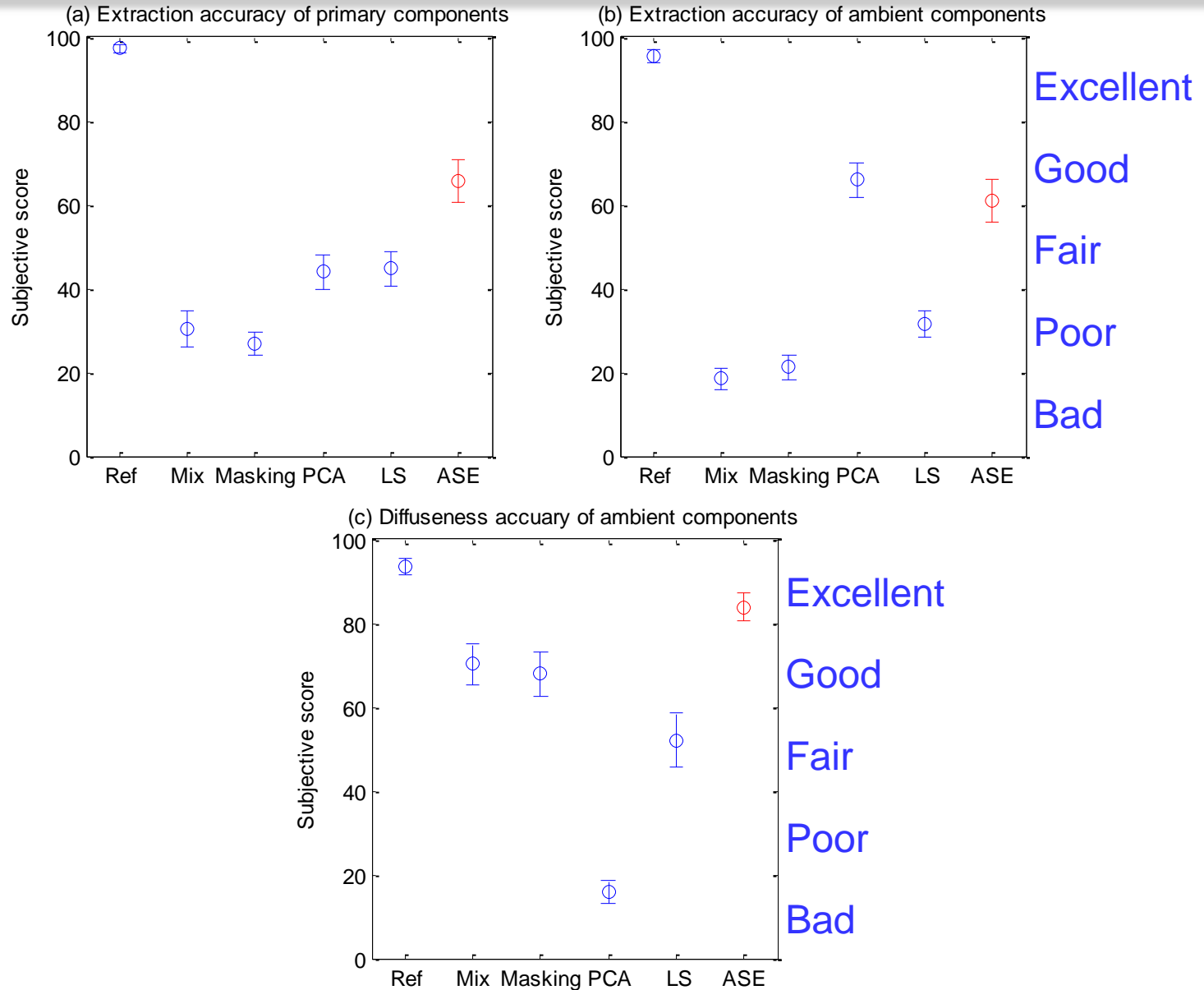
## Approaches compared

- Masking
- PCA
- LS
- ASE (APEX)
- Reference
- Mixture

## Listening tests

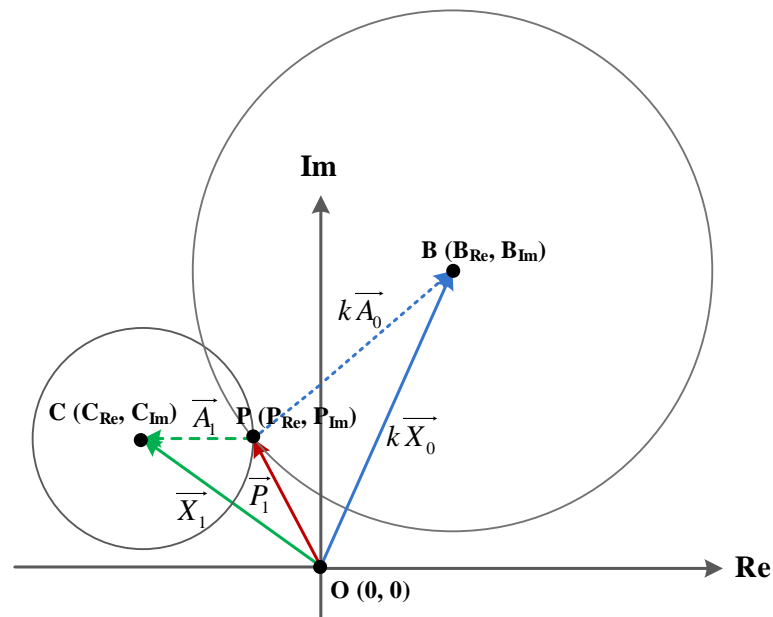
- 17 subjects
- Headphone listening
- Procedure similar to MUSHRA

# Subjective scores



# Contributions on ASE based PAE

1. Reformulated PAE as ASE by exploiting the equal ambient magnitude.
2. Solved ASE using the criterion of sparse primary component, resulting in APES, AMES, and APEX.
3. Proposed a technique to compute error measures for PAE approaches without analytic solutions.
4. Validated the improved performance (ESR reduction: 3-6 dB average) in the objective and subjective experiments.



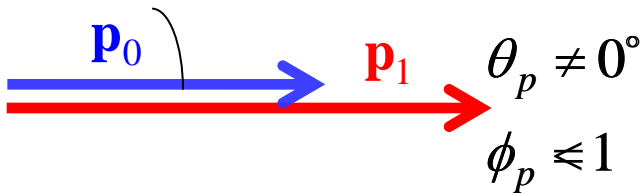
Approximate efficient solution APEX :

$$\hat{\theta}_1^* = \begin{cases} \angle \mathbf{X}_1 & , \forall k > 1 \\ \angle (\mathbf{X}_1 - \mathbf{X}_0) & , \forall k = 1 \end{cases}$$

[J2] J. He, W. S. Gan, and E. L. Tan, "Primary-ambient extraction using ambient phase estimation with a sparsity constraint," *IEEE Signal Process. Letters*, vol. 22, no. 8, pp. 1127-1131, Aug. 2015.

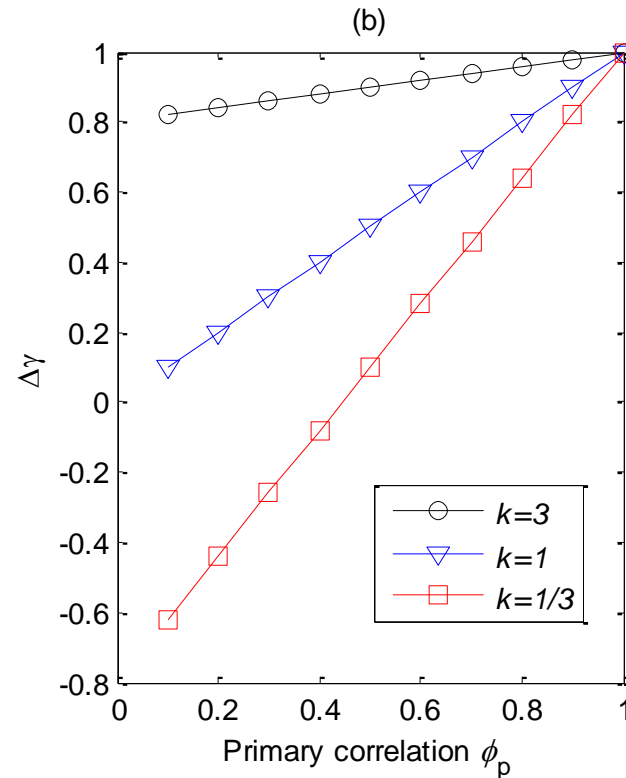
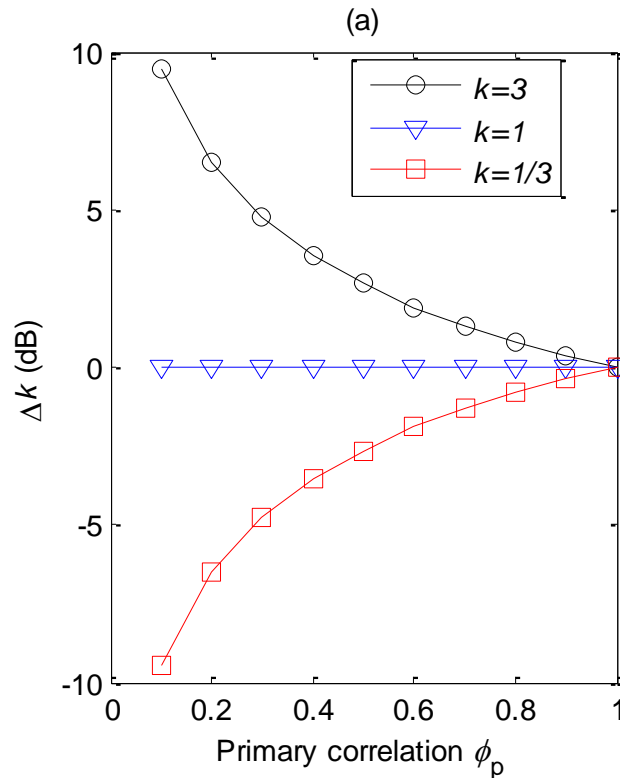
[J3] J. He, E. L. Tan, and W. S. Gan, "Primary-ambient extraction using ambient spectrum estimation for immersive spatial audio reproduction," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 9, pp. 1431-1444, Sept. 2015.

# PAE in ideal case $\rightarrow$ complex case

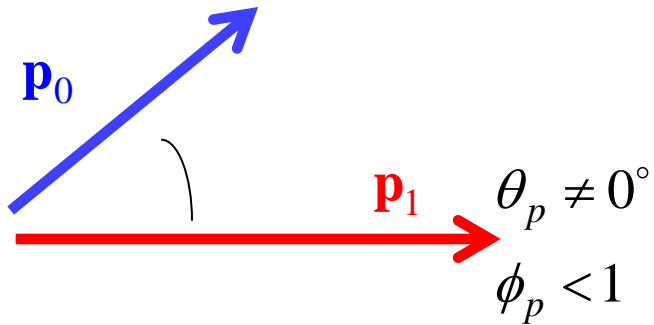


$$\Delta k = \frac{\hat{k}_{ic}}{k} = \frac{k^2 - 1}{2\phi_P k^2} + \sqrt{\left(\frac{k^2 - 1}{2\phi_P k^2}\right)^2 + \frac{1}{k^2}},$$

$$\Delta\gamma = \frac{\hat{\gamma}_{ic}}{\gamma} = \frac{k^2 - 1 + 2\phi_P}{k^2 + 1}.$$



# PAE in ideal case $\rightarrow$ complex case

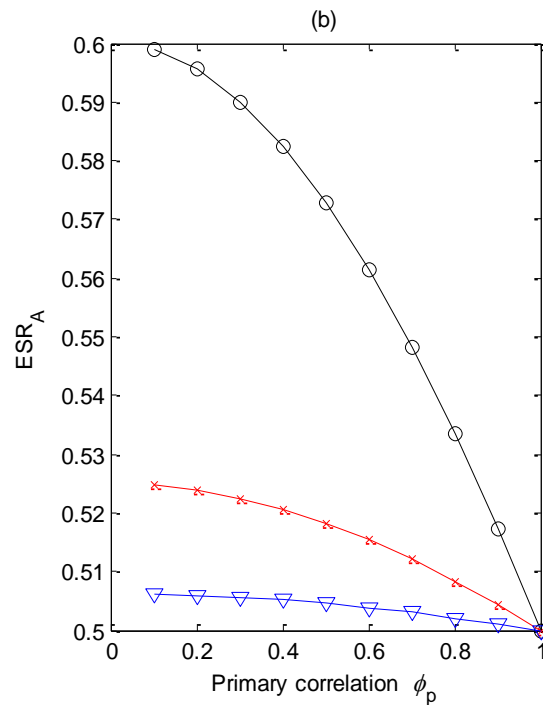
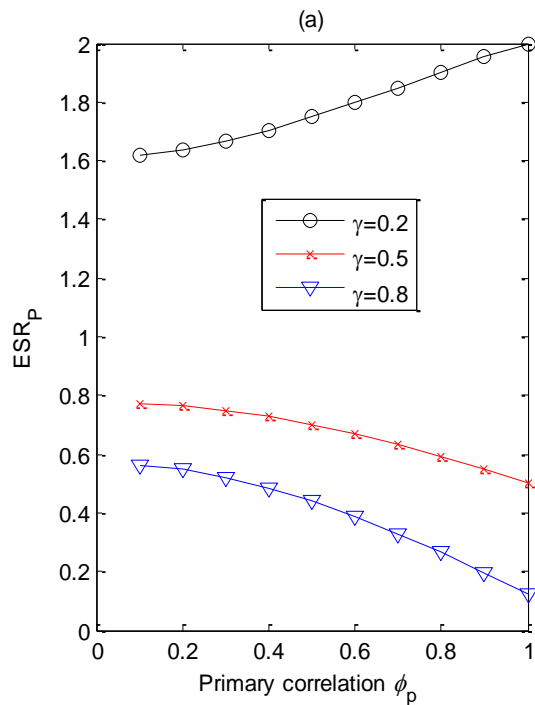


$$\hat{\mathbf{p}}_{\text{PCA},0} = \mathbf{p}_0 - \mathbf{v} + \frac{1}{1 + \hat{k}_{ic}^2} (\mathbf{a}_0 + \hat{k}_{ic} \mathbf{a}_1),$$

$$\hat{\mathbf{p}}_{\text{PCA},1} = \mathbf{p}_1 + \frac{1}{\hat{k}_{ic}} \mathbf{v} + \frac{\hat{k}_{ic}}{1 + \hat{k}_{ic}^2} (\mathbf{a}_0 + \hat{k}_{ic} \mathbf{a}_1),$$

$$\hat{\mathbf{a}}_{\text{PCA},0} = \frac{\hat{k}_{ic}^2}{1 + \hat{k}_{ic}^2} \mathbf{a}_0 + \mathbf{v} + \frac{-\hat{k}_{ic}}{1 + \hat{k}_{ic}^2} \mathbf{a}_1,$$

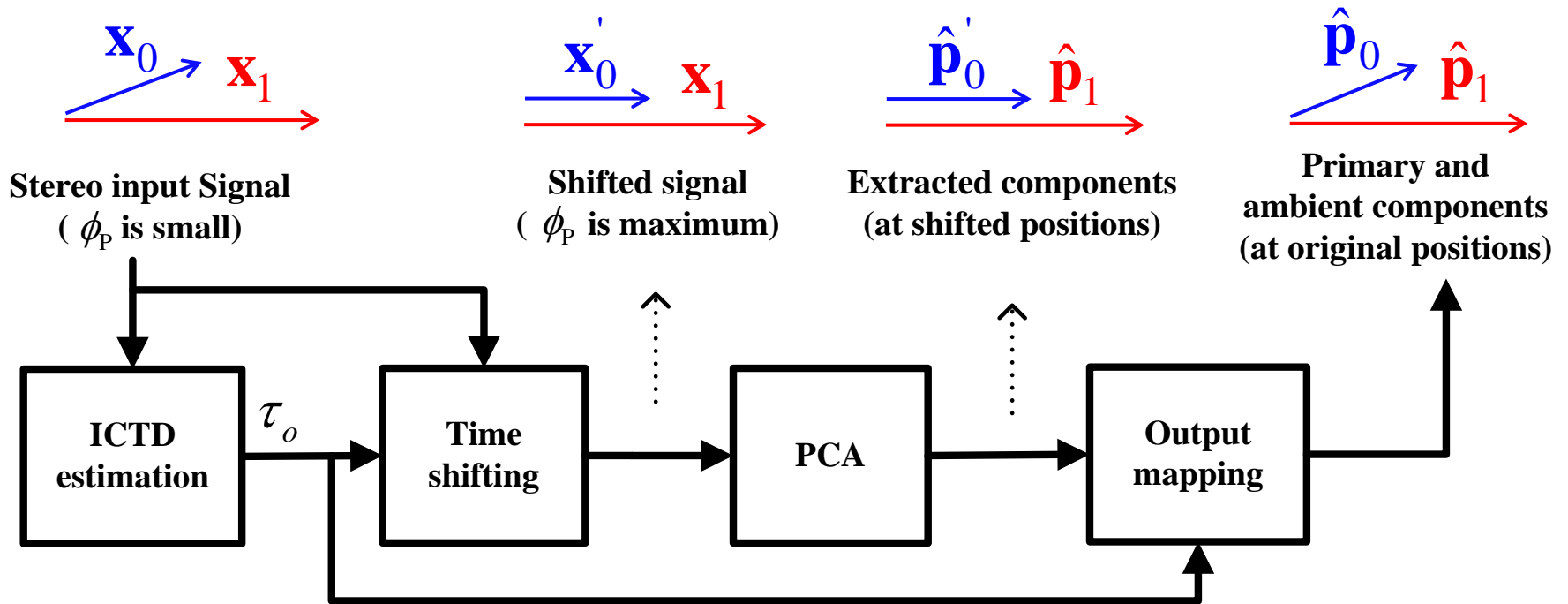
$$\hat{\mathbf{a}}_{\text{PCA},1} = \frac{1}{1 + \hat{k}_{ic}^2} \mathbf{a}_1 - \frac{1}{\hat{k}_{ic}} \mathbf{v} + \frac{-\hat{k}_{ic}}{1 + \hat{k}_{ic}^2} \mathbf{a}_0,$$



$$\mathbf{v} = \frac{\hat{k}_{ic}}{1 + \hat{k}_{ic}^2} (\hat{k}_{ic} \mathbf{p}_0 - \mathbf{p}_1)$$

# Time-shifting for PAE

To account for the partial primary correlation (0-lag) caused by the inter-channel time difference (ICTD).

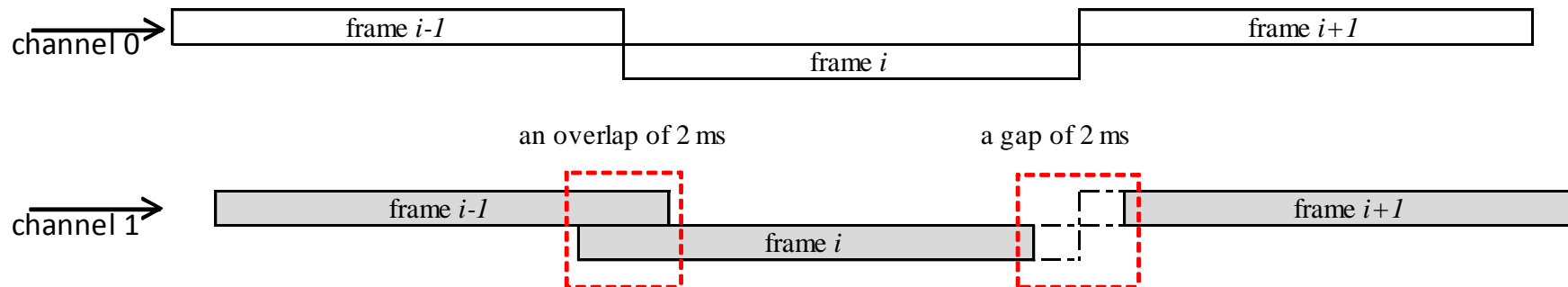


$$\hat{p}_{\text{SPCA},0}(n) = \frac{1}{1 + \hat{k}_{ic}^2} \left[ x_0(n) + \hat{k}_{ic} x_1(n - \tau_o) \right], \quad \hat{p}_{\text{SPCA},1}(n) = \frac{\hat{k}_{ic}}{1 + \hat{k}_{ic}^2} \left[ x_0(n + \tau_o) + \hat{k}_{ic} x_1(n) \right],$$

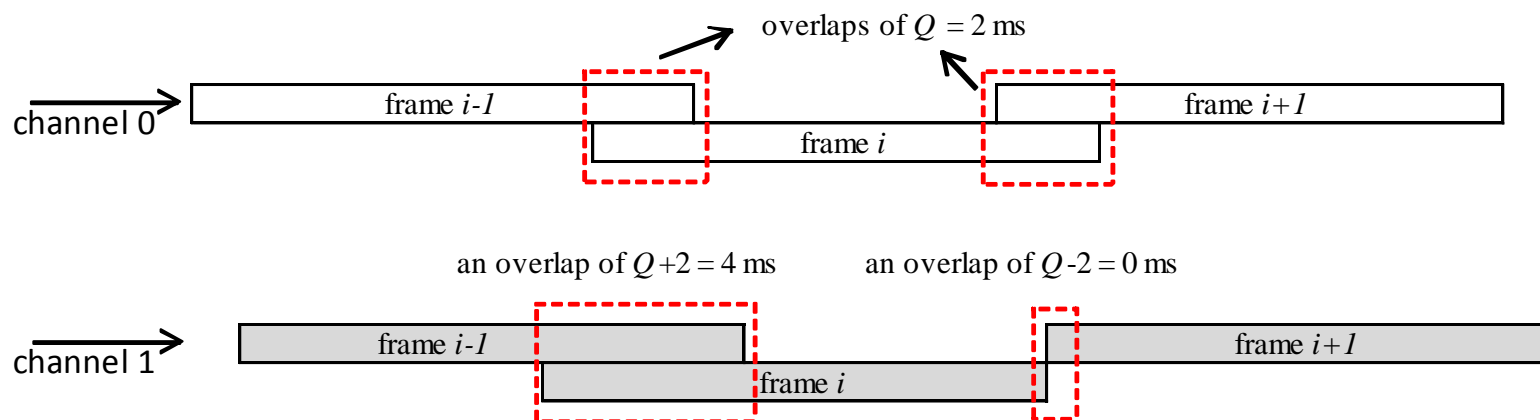
$$\hat{a}_{\text{SPCA},0}(n) = \frac{\hat{k}_{ic}}{1 + \hat{k}_{ic}^2} \left[ \hat{k}_{ic} x_0(n) - x_1(n - \tau_o) \right], \quad \hat{a}_{\text{SPCA},1}(n) = -\frac{1}{1 + \hat{k}_{ic}^2} \left[ \hat{k}_{ic} x_0(n + \tau_o) - x_1(n) \right].$$

# Output mapping

Frame	$i-1$	$i$	$i+1$
ICTD (ms)	1	-1	1



(a) Conventional output mapping



(b) Overlapped output mapping

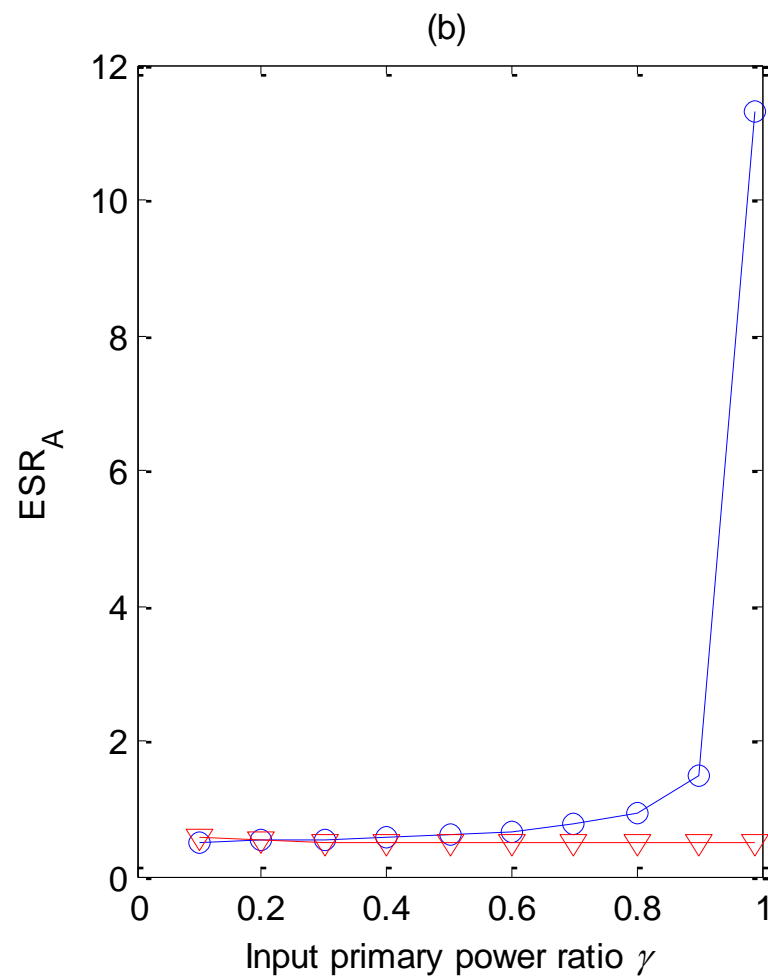
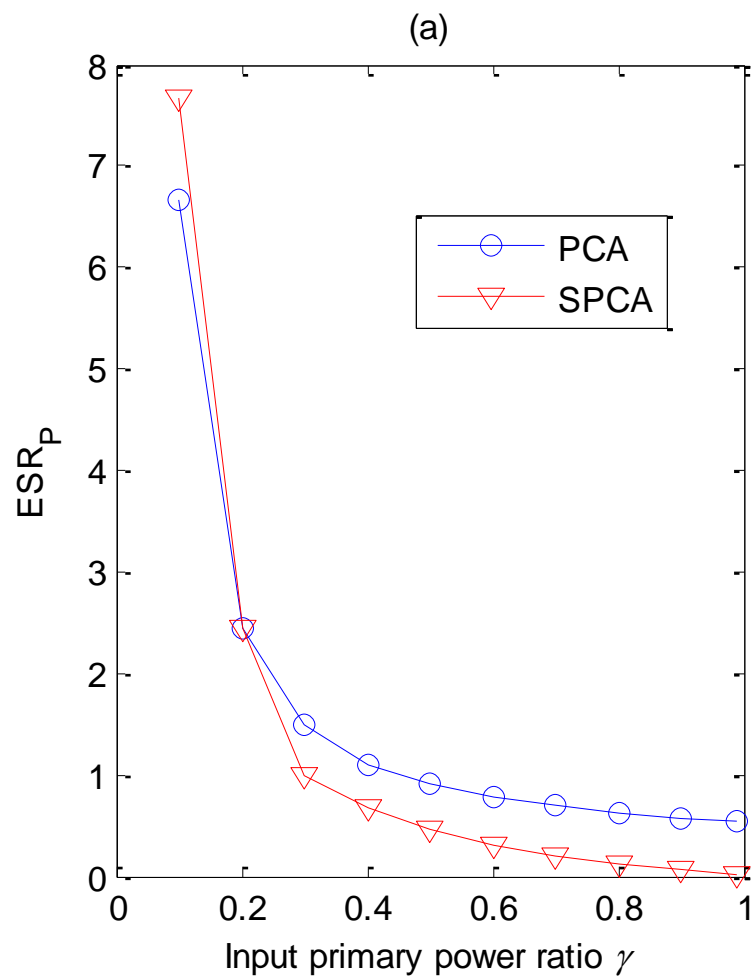
**Overlapping duration:  $Q \geq 2$  ms**

# Experiments

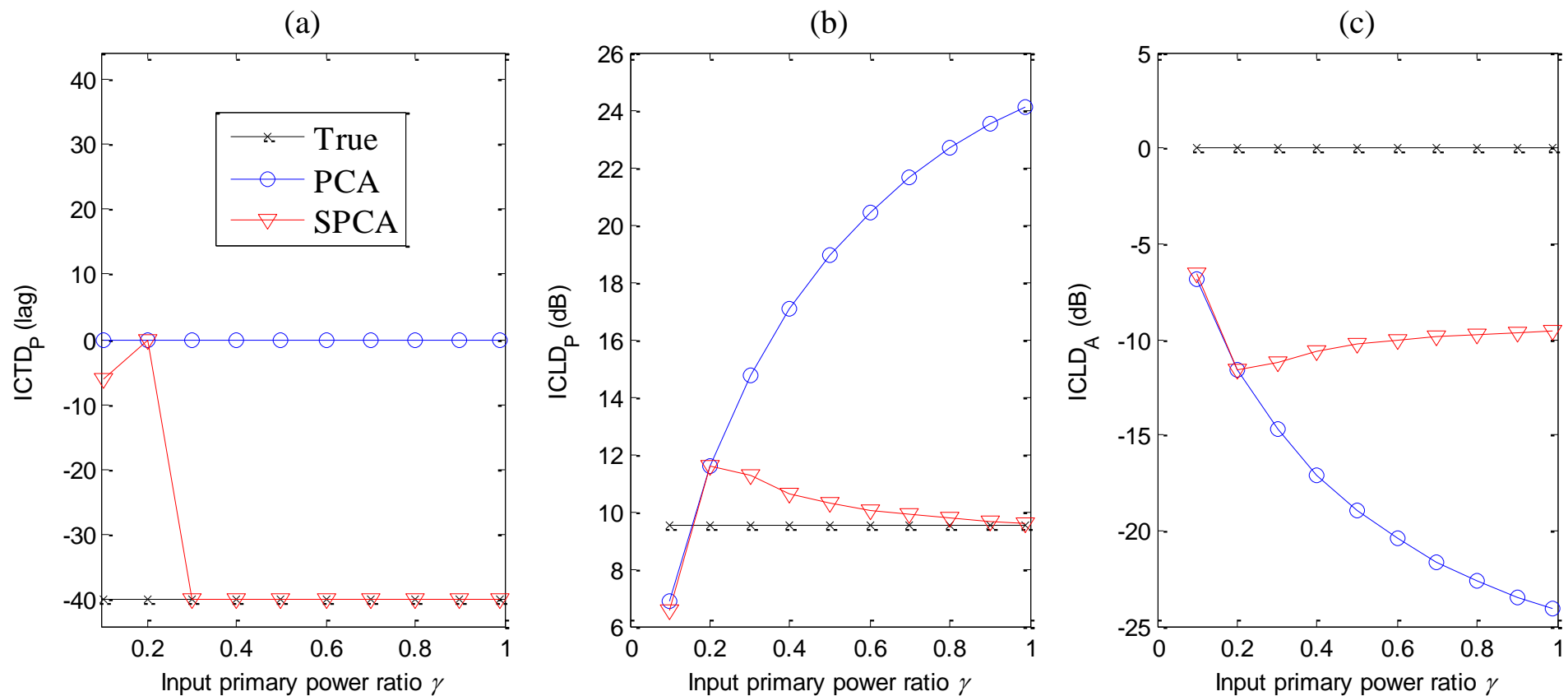
Exp	Input signal	Primary component	Ambient component	Settings
1	Synthesized	Speech	Lapping wave	<b>Fixed</b> direction; different values of $\gamma$
2	Synthesized	Shaking matchbox	Lapping wave	<b>Panning</b> directions with close $\gamma$
3	Synthesized	Direct path of speech	<b>Reverberation</b> of speech	Varying directions with different $\gamma$
4	<b>Recorded</b>	Speech	Canteen recording	Three directions with close $\gamma$



# Experiment 1: Extraction accuracy



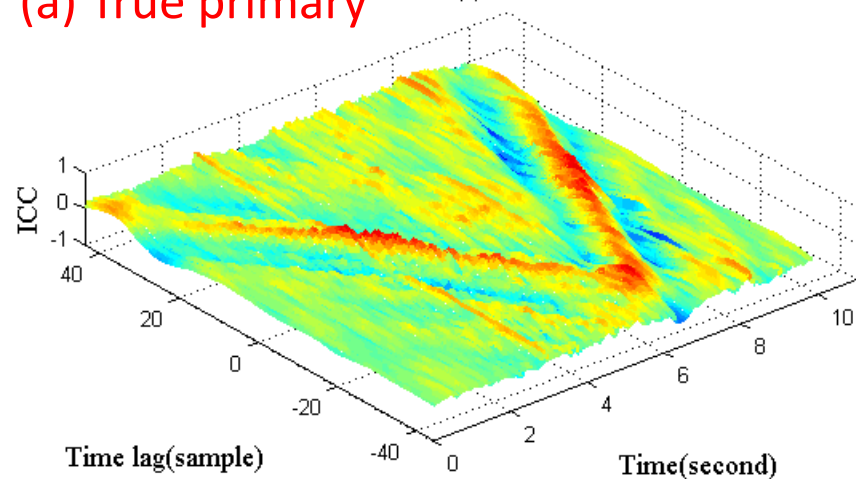
# Experiment 1: Spatial accuracy



# Experiment 2: Tracking of moving source

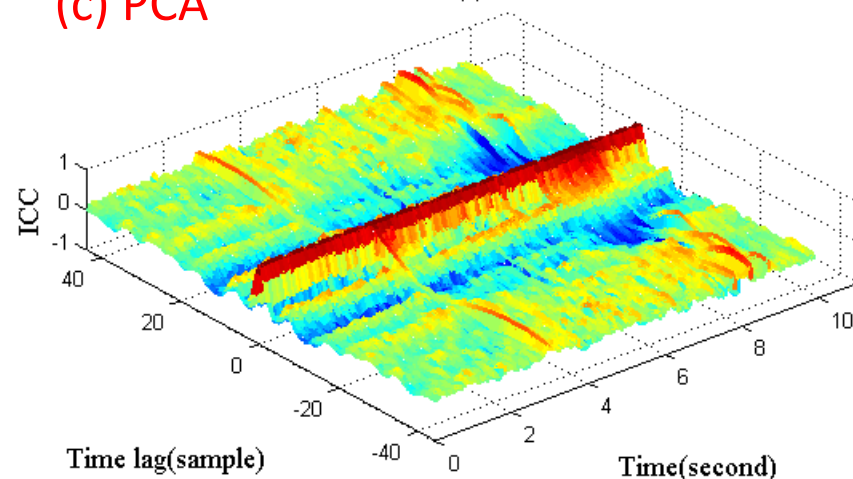
(a) True primary

(a)



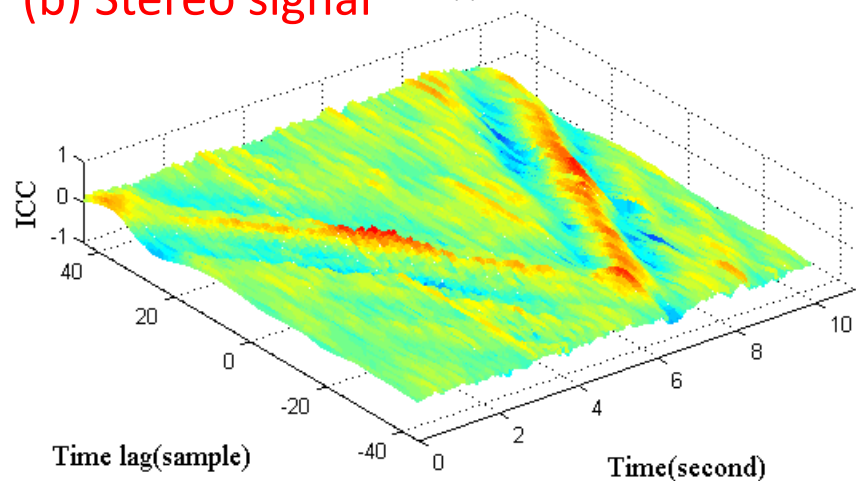
(c) PCA

(c)



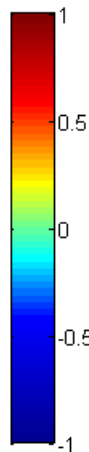
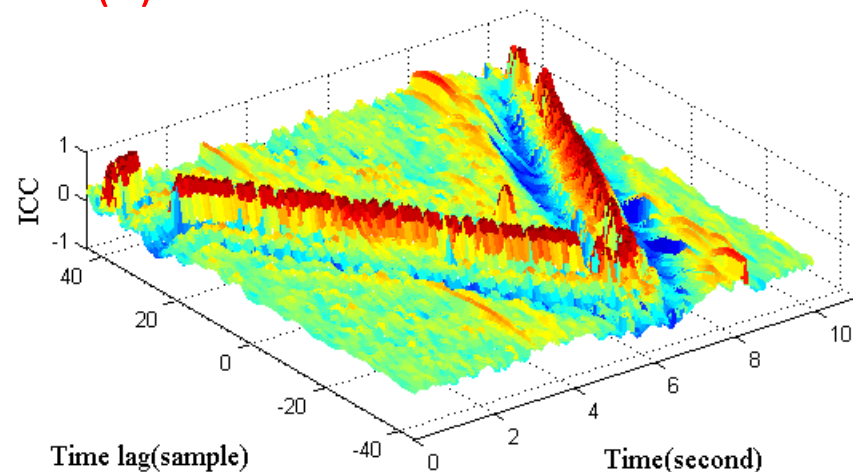
(b) Stereo signal

(b)

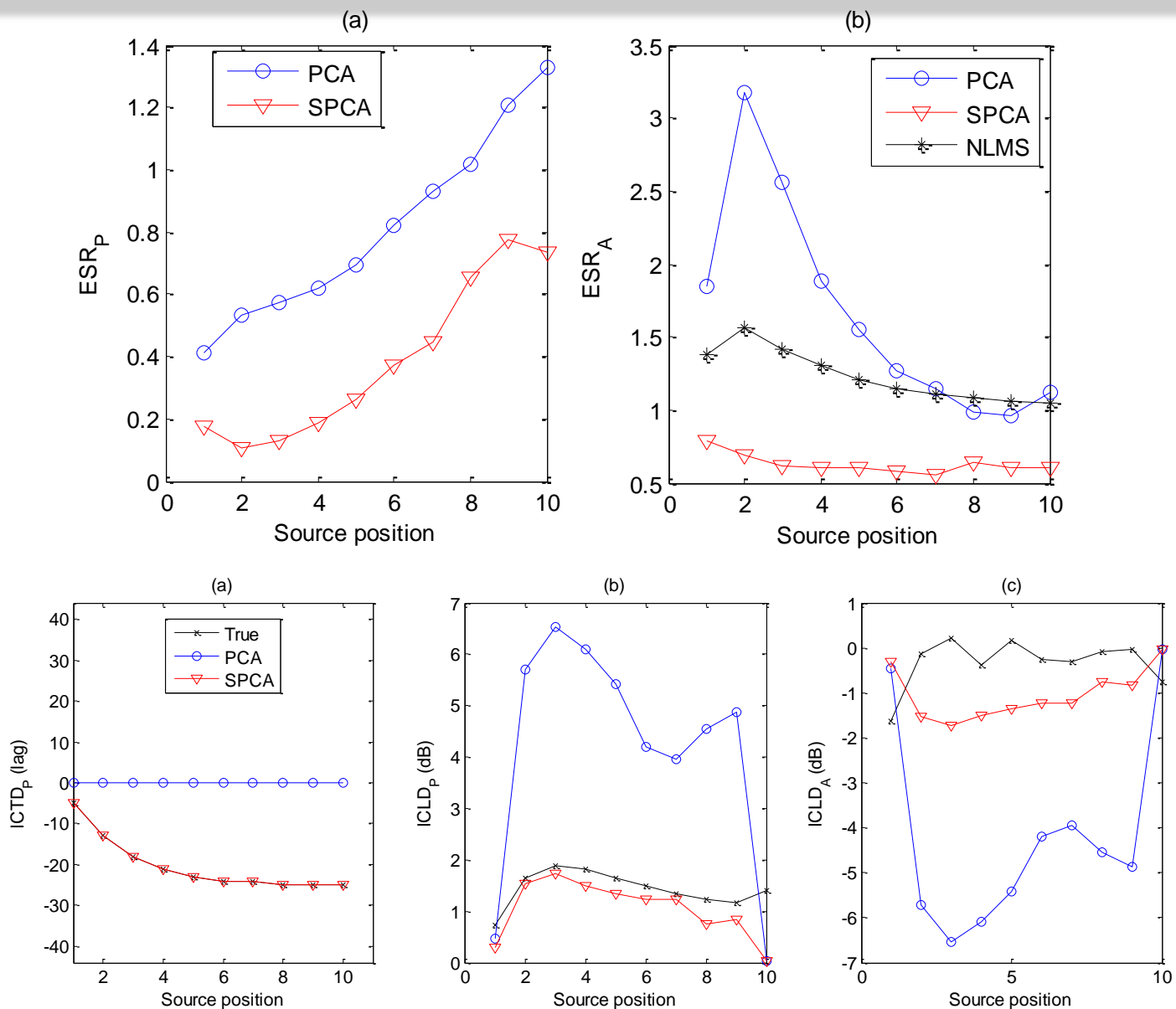


(d) SPCA

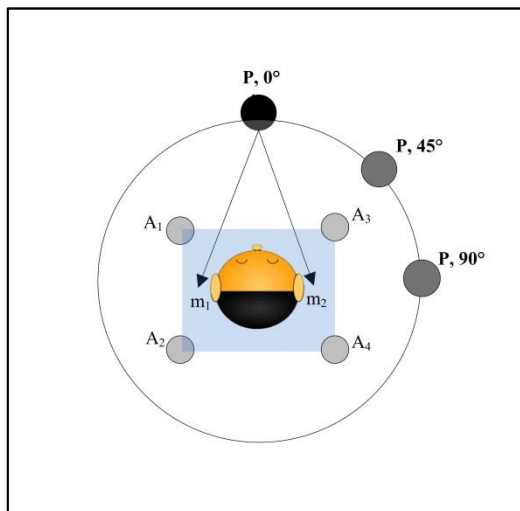
(d)



# Experiment 3: Reverberation experiment



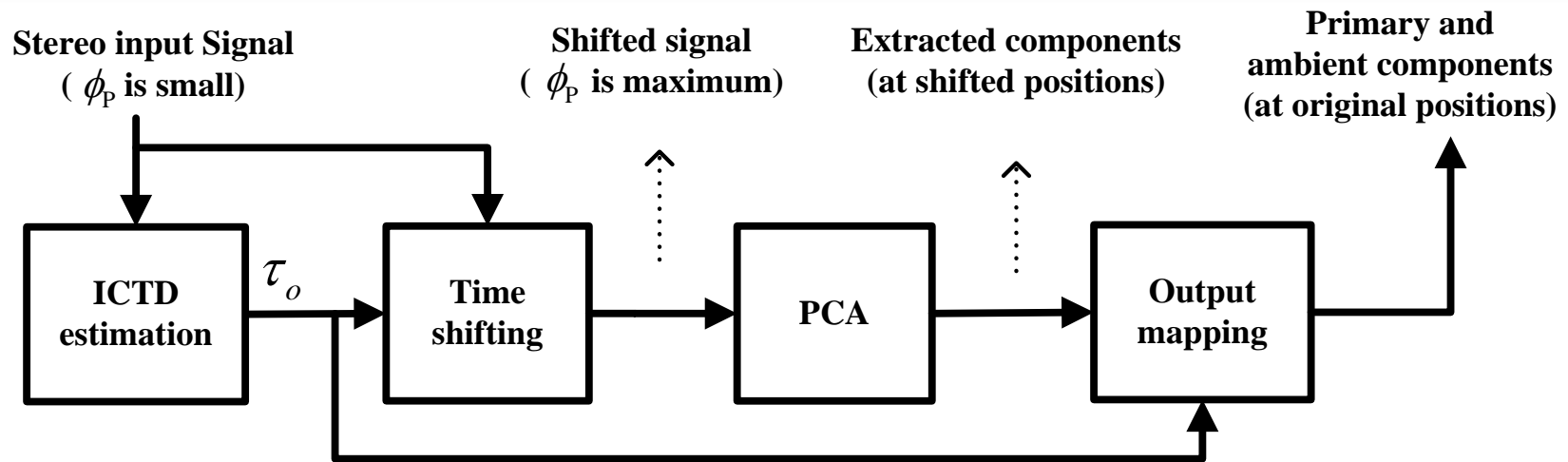
# Experiment 4: Recording experiment



ESR	Primary component			Ambient component			
	$\theta$	$0^\circ$	$45^\circ$	$90^\circ$	$0^\circ$	$45^\circ$	$90^\circ$
PCA		0.27	0.64	0.88	1.08	1.89	2.49
<b>SPCA</b>		<b>0.21</b>	<b>0.31</b>	<b>0.34</b>	<b>0.81</b>	<b>1.02</b>	<b>1.39</b>

Spatial	ICTD <sub>P</sub>			ICLD <sub>P</sub> (dB)			ICLD <sub>A</sub> (dB)			
	$\theta$	$0^\circ$	$45^\circ$	$90^\circ$	$0^\circ$	$45^\circ$	$90^\circ$	$0^\circ$	$45^\circ$	$90^\circ$
True		1	-17	-31	-1.02	7.74	11.90	1.03	1.18	1.03
PCA		0	0	0	-1.46	36.03	23.11	1.46	-36.03	-23.11
<b>SPCA</b>		<b>1</b>	<b>-17</b>	<b>-31</b>	<b>-1.26</b>	<b>8.65</b>	<b>15.60</b>	<b>1.26</b>	<b>-8.65</b>	<b>-15.60</b>

# Contributions on Time Shifting PAE



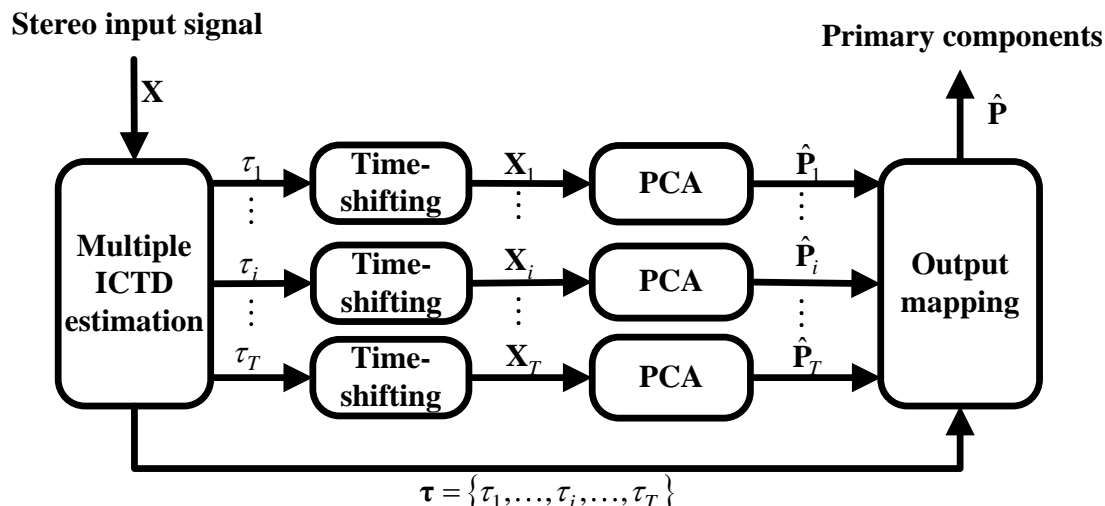
**For mixture signals with partially correlated primary components,**

- Analyzed the performance of conventional PAE;
- Proposed time shifting technique to improve the PAE performance;
- Achieved lower extraction error and more accurate spatial cues.

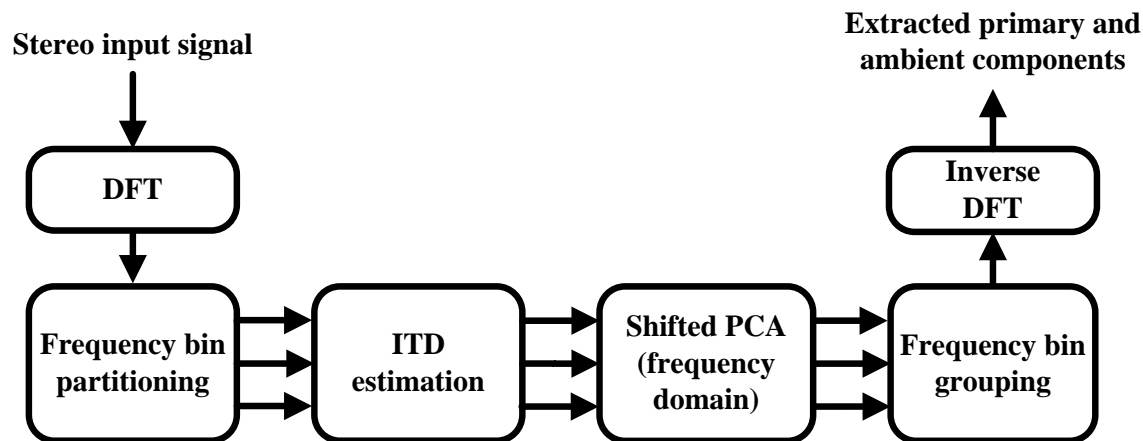
[J4] J. He, W. S. Gan, and E. L. Tan, "Time-shifting based primary-ambient extraction for spatial audio reproduction," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 10, pp. 1576-1588, Oct. 2015.

# Extensions for Multiple Sources

**Multi-shift PAE  
with ICC based  
output weighting**



**Subband PAE  
with frequency  
bin partitioning**

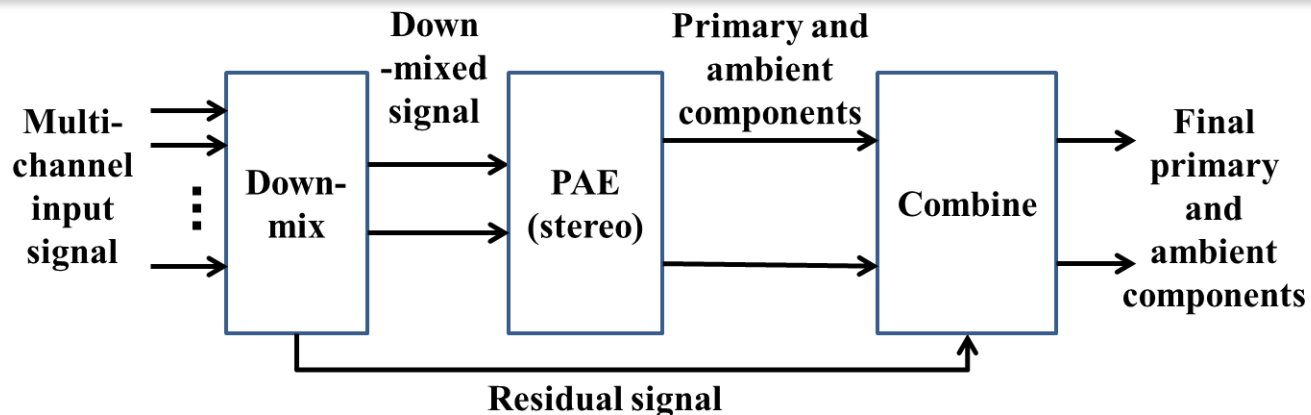


[C3] J. He, and W. S. Gan, "Multi-shift principal component analysis based primary component extraction for spatial audio reproduction," in *Proc. ICASSP*, Brisbane, Australia, Apr. 2015, pp. 350-354.

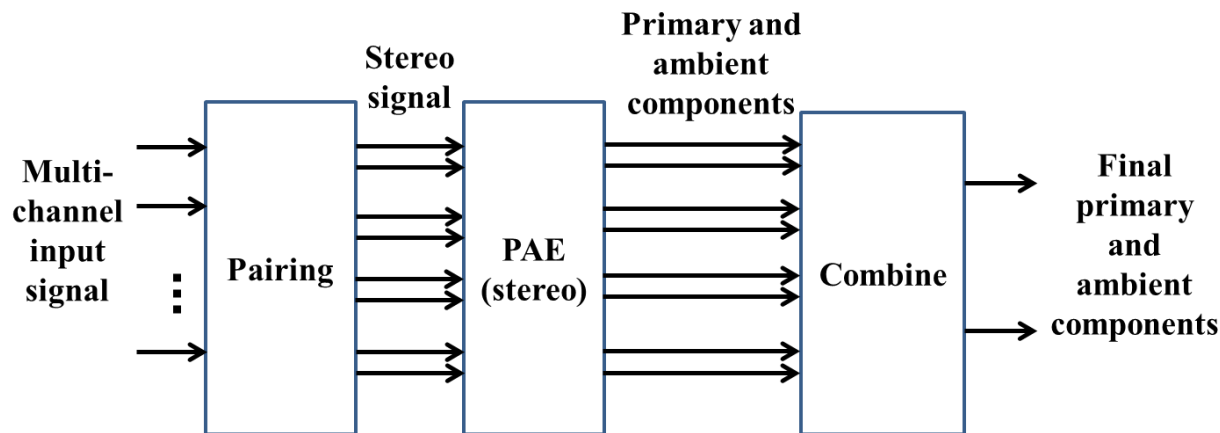
[C2] J. He, E. L. Tan, and W. S. Gan, "A study on the frequency-domain primary-ambient extraction for stereo audio signals," in *Proc. ICASSP*, Florence, Italy, 2014, pp. 2892-2896.

# PAE: from stereo to multichannel signals

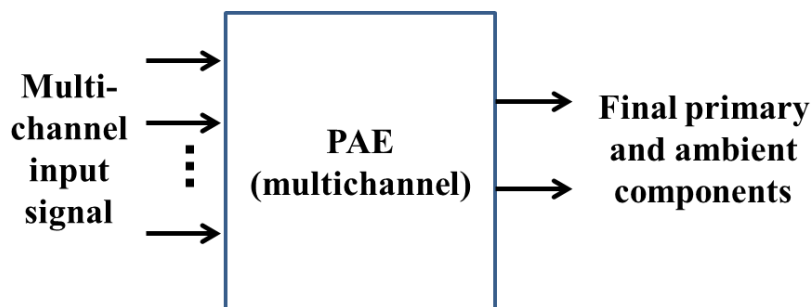
## 1. Using down-mix



## 2. Using pairing



## 3. Direct



[C5] J. He, and W. S. Gan, "Applying primary ambient extraction for immersive spatial audio reproduction," in Proc. *APSIPA Annual Summit and Conference*, Hong Kong, Dec. 2015.

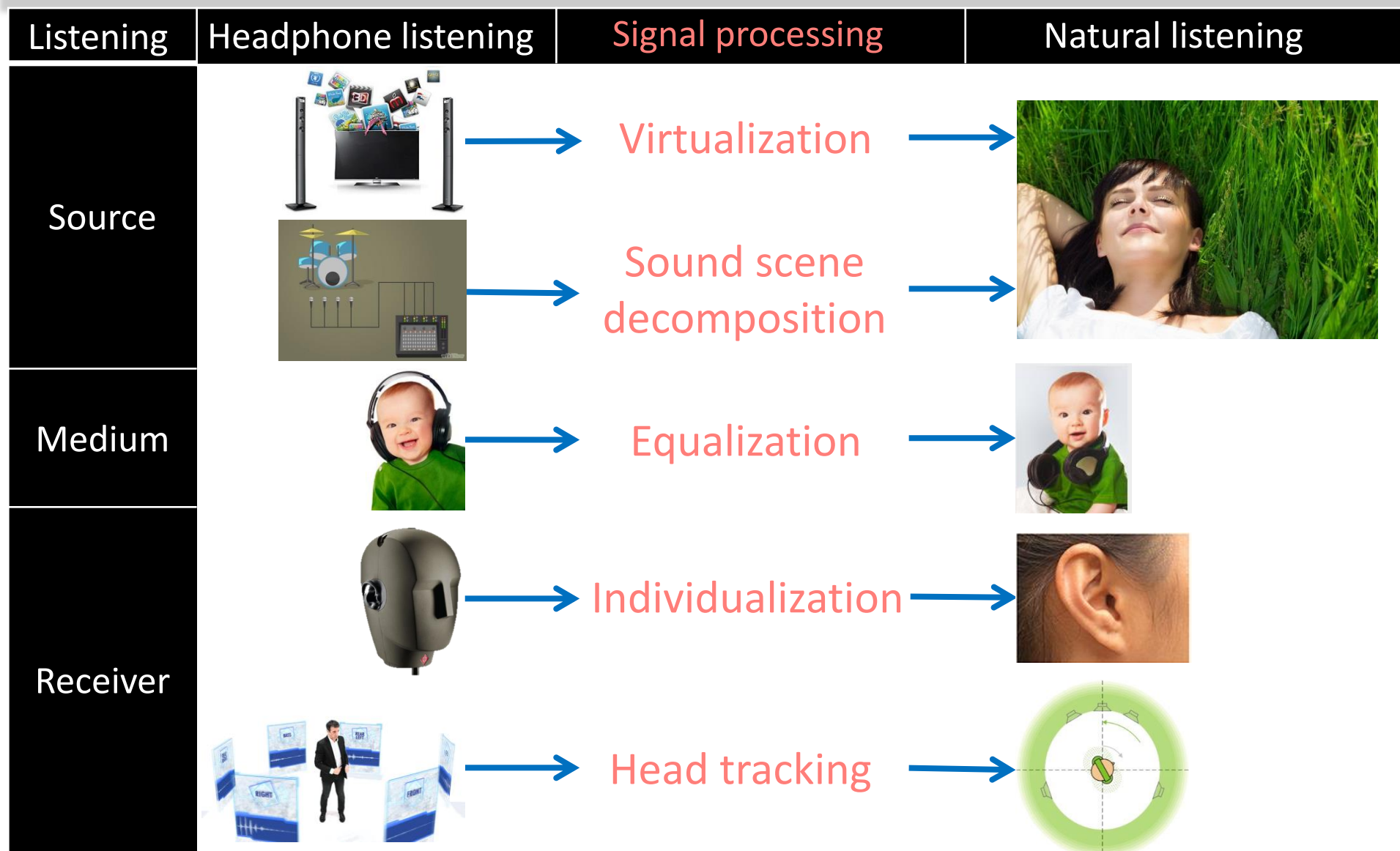


# Natural Sound Rendering for Headphones

K. Sunder, J. He, E. L. Tan, and W. S. Gan, “Natural sound rendering for headphones,” *IEEE Signal Processing Magazine*, vol. 32, no. 2, pp. 100-113, Mar. 2015.

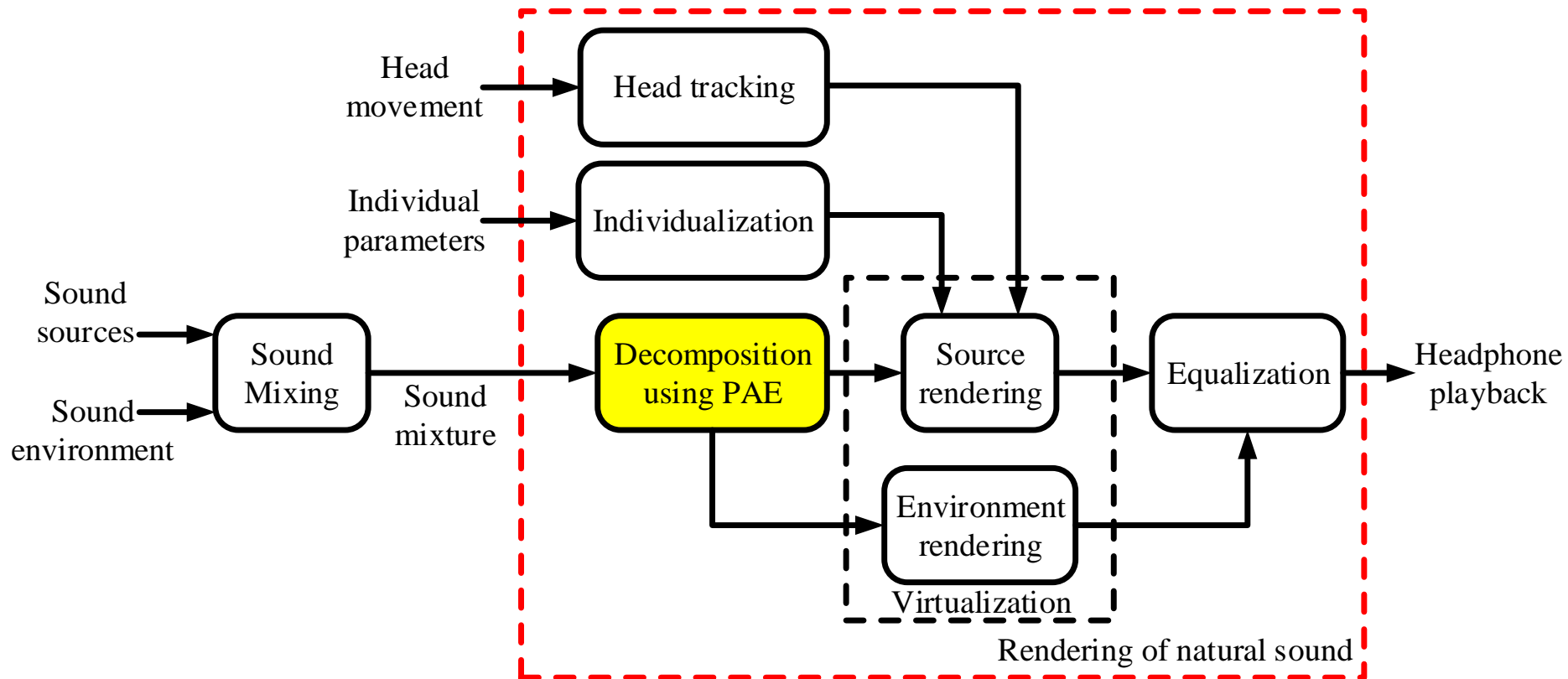


# Challenges and solutions

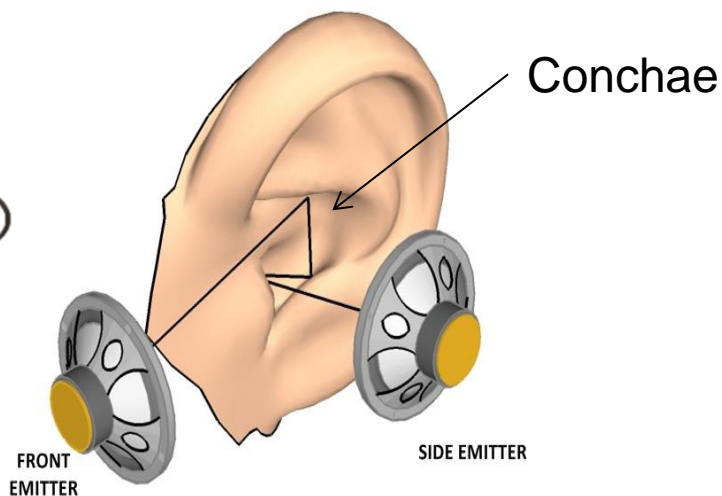


D. R. Begault, *3-D sound for virtual reality and multimedia*: AP Professional, 2000.

# Integration

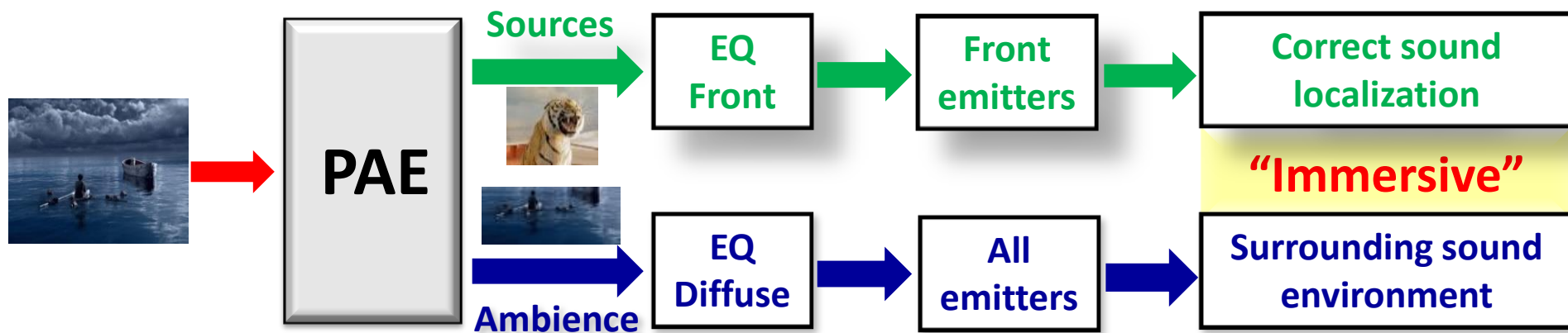


# 3D Audio Headphone: an example



## Key features

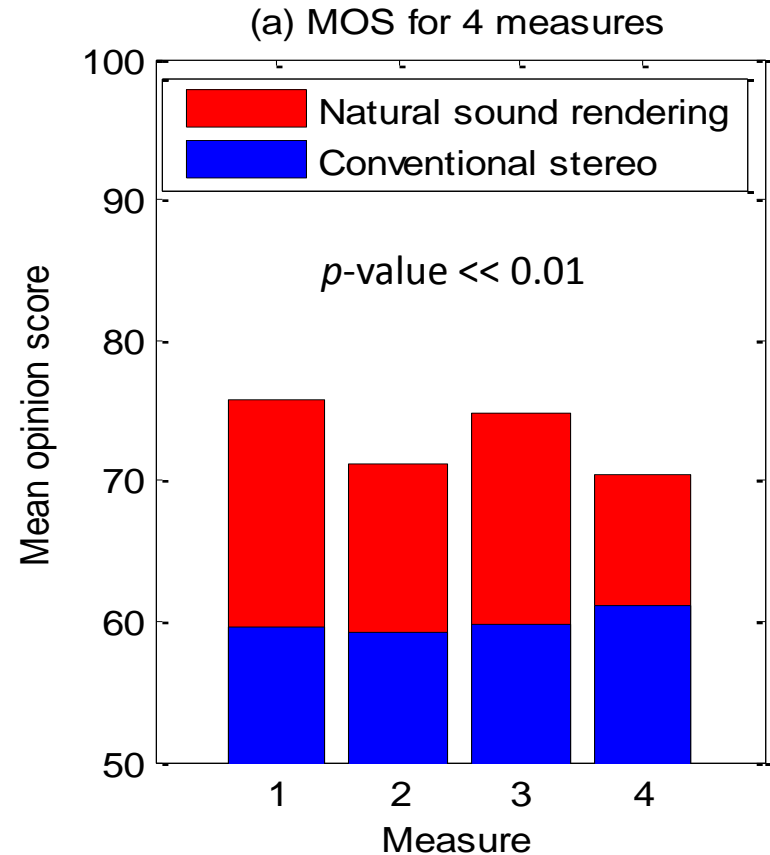
- ✓ Patented structure with strategic-positioned emitters;
- ✓ Individualization via frontal projection; no measurements or training required;
- ✓ Recreate an immersive perception of sound objects with surrounding ambience;
- ✓ Compatible with all existing sound formats.



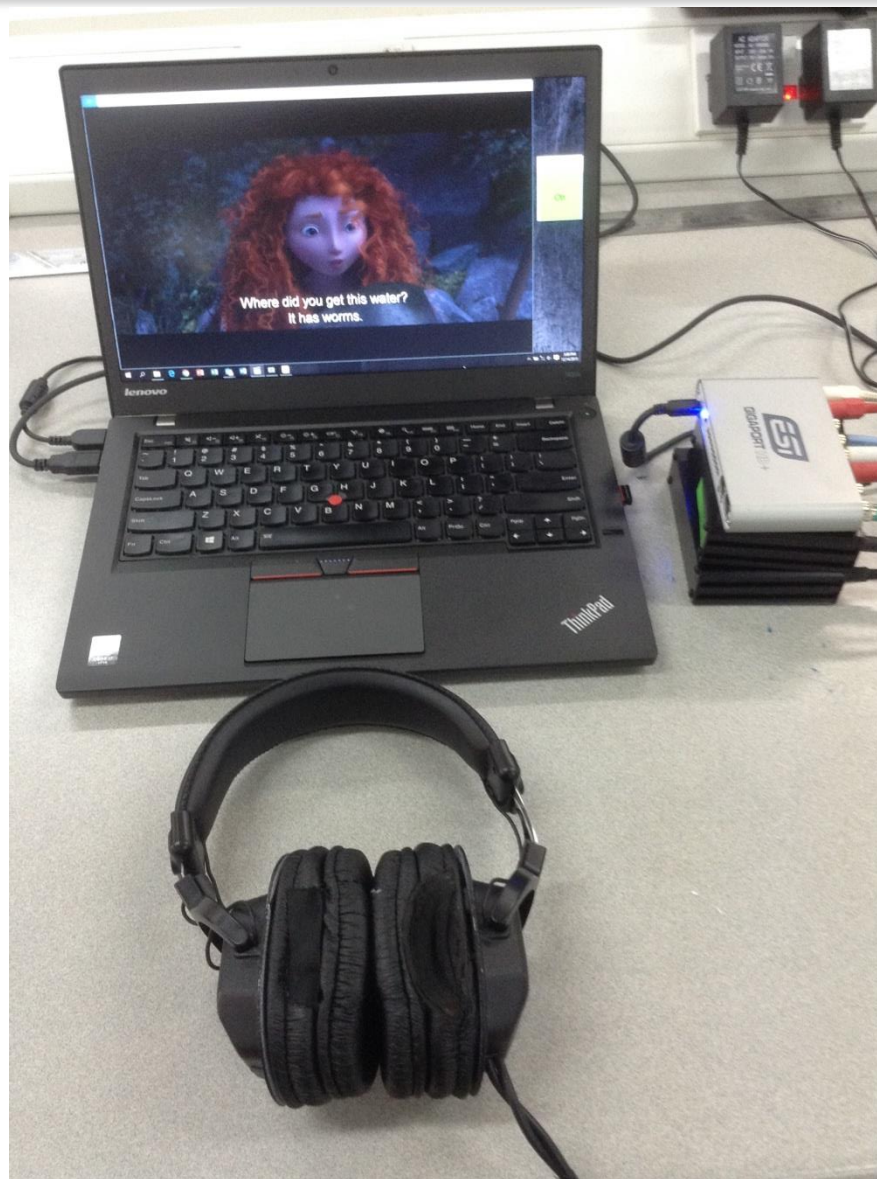
W. S. Gan and E. L. Tan, "Listening device and accompanying signal processing method," US Patent 2014/0153765 A1, 2014.

# Subjective evaluation

- 18 subjects, score of 0-100;
- **Stimuli:** binaural, movie and gaming tracks;
- **Four measures:**
  1. Sense of direction,
  2. Externalization,
  3. Ambience,
  4. Timbral quality.



# Real-time 3D audio headphones



**Show & Tell**  
at  
**ICASSP 2016**

# Conclusions of this Thesis

Primary Ambient Extraction facilitates flexible, efficient, and immersive spatial audio reproduction of channel-based signals for any playback systems

- Comprehensive study on linear estimation based PAE approaches lay the foundation;
- Novel ambient spectrum estimation framework significantly improves PAE performance;
- Time shifting and subband decomposition techniques enhance the robustness of PAE performance in complex scenarios;
- The novel natural sound rendering headphone system validates the advantage of PAE.

# Future work

- Extended study of PAE in complex cases for all approaches;
- To achieve tradeoff between timbre and spatial quality in PAE, and objective evaluation system that can mimic subjective performance;
- Evaluate PAE in specific spatial audio reproduction applications (loudspeakers, headphones);
- Incorporate probabilistic approaches, and make use of Big Data and Machine Learning to improve the robustness in complex audio scenes.



# Author's full publication list

## Journal papers

- [J1] **J. He**, E. L. Tan, and W. S. Gan, "Linear estimation based primary-ambient extraction for stereo audio signals," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 22, no. 2, pp. 505-517, Feb. 2014.
- [J2] K. Sunder, **J. He**, E. L. Tan, and W. S. Gan, "Natural sound rendering for headphones: Integration of signal processing techniques," *IEEE Sig. Process. Mag.*, vol. 32, no. 2, Mar 2015, pp. 100-113.
- [J3] **J. He**, W. S. Gan, and E. L. Tan, "Primary-ambient extraction using ambient phase estimation with a sparsity constraint," *IEEE Signal Process. Letters*, vol. 22, no. 8, pp. 1127-1131, Aug. 2015.
- [J4] **J. He**, E. L. Tan, and W. S. Gan, "Primary-ambient extraction using ambient spectrum estimation for immersive spatial audio reproduction," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 9, pp. 1430-1443, Sept. 2015.
- [J5] **J. He**, W. S. Gan, and E. L. Tan, "Time-shifting based primary-ambient extraction for spatial audio reproduction," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 10, pp. 1576-1588, Oct. 2015.
- [J6] **J. He**, R. Ranjan, W. S. Gan, and K. Sunder, "Scalable data reusing normalized LMS for acoustic system identification with short-duration signals," *IEEE Signal Process. Letters*, under review.

## Conference papers

- [C1] **J. He**, E. L. Tan, and W. S. Gan, "Time-shifted principal component analysis based cue extraction for stereo audio signals," in *Proc. ICASSP*, Vancouver, Canada, 2013, pp. 266-270.
- [C2] **J. He**, W. S. Gan, and E. L. Tan, "A study on the frequency-domain primary-ambient extraction for stereo audio signals," in *Proc. ICASSP*, Florence, Italy, 2014, pp. 2892-2896. (Awarded SPS travel grant)
- [C3] **J. He**, W. S. Gan and Y. K. Chong, "Study on the use of error term in parallel-form narrowband feedback active noise control systems," in *Proc. 2014 Asia Pacific Signal and Information Processing Association*

(*APSIPA*) Annual Summit and Conference (invited), Cambodia, Dec. 2014.

- [C4] **J. He**, W. S. Gan, and E. L. Tan, "On the preprocessing and postprocessing of HRTF individualization based on sparse representation of anthropometry features," in *Proc. ICASSP*, Brisbane, Australia, Apr. 2015, pp. 639-643. (Awarded SPS travel grant)
- [C5] **J. He**, and W. S. Gan, "Multi-shift principal component analysis based primary component extraction for spatial audio reproduction," in *Proc. ICASSP*, Brisbane, Australia, Apr. 2015, pp. 350-354. (Awarded SPS travel grant)
- [C6] S. Fasciani, **J. He**, B. Lam, T. Murao, and W. S. Gan, "Comparative study of cone-shaped versus flat-panel speakers for active noise control of multi-tonal signals in open windows," in *Proc. Internoise 2015 (invited)*, San Francisco, Aug. 2015.
- [C7] **J. He**, and W. S. Gan, "Applying primary ambient extraction for immersive spatial audio reproduction," *2015 Asia Pacific Signal and Information Processing Association (APSIPA) Annual Summit and Conference (invited)*, Hong Kong, Dec. 2015.
- [C8] **J. He**, R. Ranjan, and W. S. Gan, "Fast continuous HRTF acquisition with unconstrained movements of human subjects," in *Proc. ICASSP*, Shanghai, China, Mar. 2016, pp.
- [Book] **J. He**, "Spatial audio reproduction with primary ambient extraction," Monograph in preparation to submit to [SpringerBriefs](#) in Signal Processing.
- [Tutorial] W. S. Gan, and **J. He**, "Assisted listening for headphones and hearing aids: signal processing techniques," Tutorial at *APSIPA ASC 2015*, Hong Kong, Dec. 2015.
- [Show & Tell] D. H. Nguyen, **J. He**, K. K. Phyo, and W. S. Gan, "Real-time audio signal processing platform for natural 3D sound rendering," Show & Tell at *ICASSP 2016*, Shanghai, China, Mar. 2016.