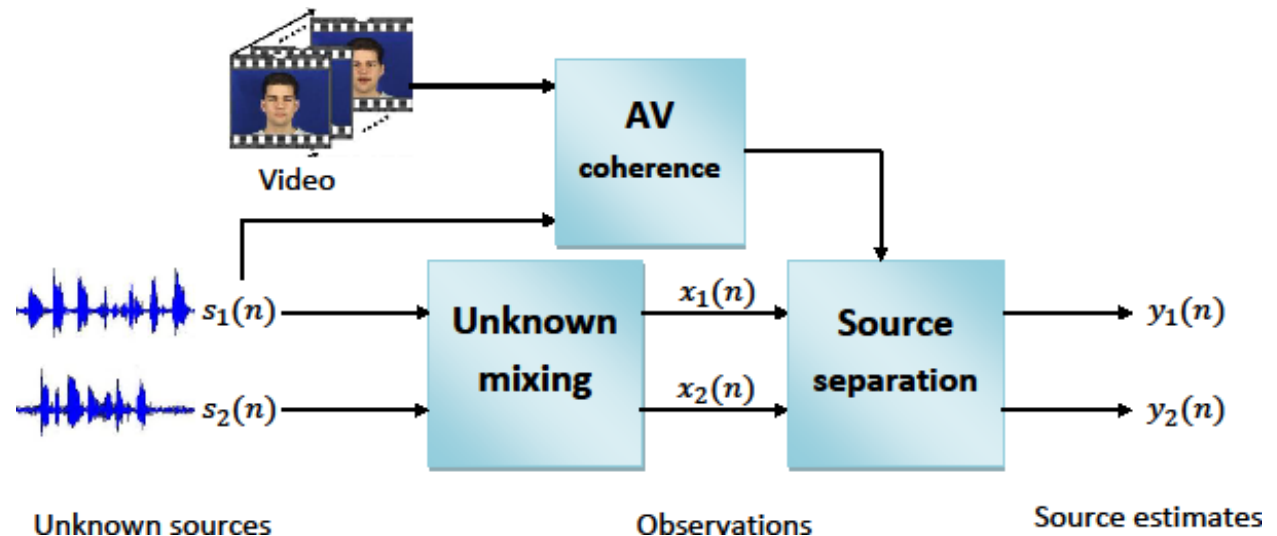# Audio-Visual Dictionary Learning and Probabilistic Time-Frequency Masking in Convolutive and Noisy Source Separation

## Wenwu Wang

**w.wang@surrey.ac.uk**
**Senior Lecturer in Signal Processing**
Centre for Vision, Speech and Signal Processing
Department of Electronic Engineering
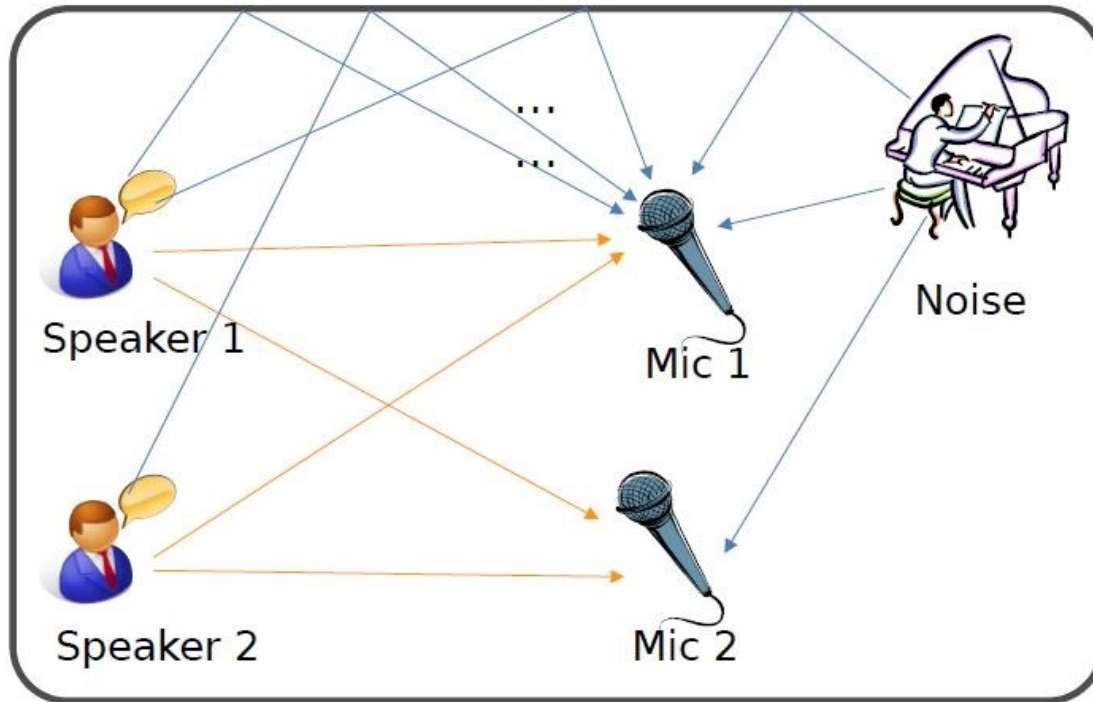University of Surrey, Guildford

# Acknowledgement

➢Joint work with Dr Qingju Liu (former PhD student & current postdoc)

➢Collaborators: Dr Philip Jackson, Dr Mark Barnard, Prof Josef Kittler, Prof Jonathon Chambers (Loughborough University), and Dr Wei Dai (Imperial College London)

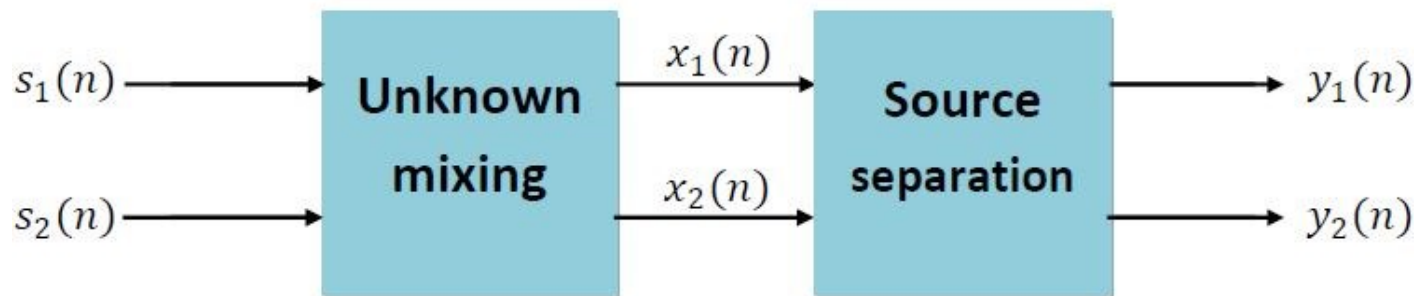➢Financial support: EPSRC & DSTL, UDRC in Signal Processing

# Outline

# Introduction----Cocktail party problem



➤Independent component analysis (ICA)

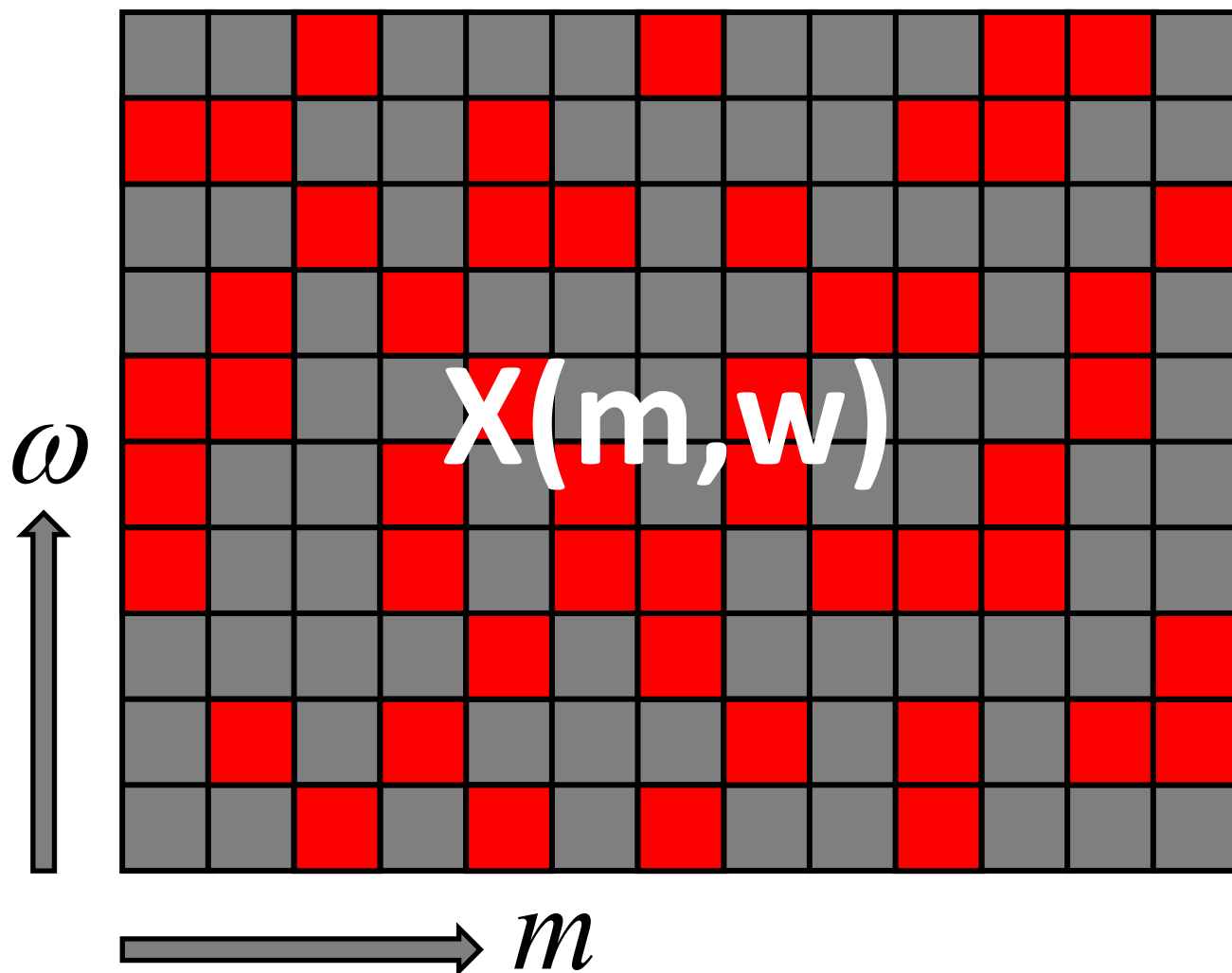➤Time-frequency (TF) masking

**"Blind" source separation BSS**

$s_1(n)$ ——→ **Unknown mixing** —$x_1(n)$→ **Source separation** ——→ $y_1(n)$

$s_2(n)$ ——→ —$x_2(n)$→ ——→ $y_2(n)$

Sources     Observations     Source estimates
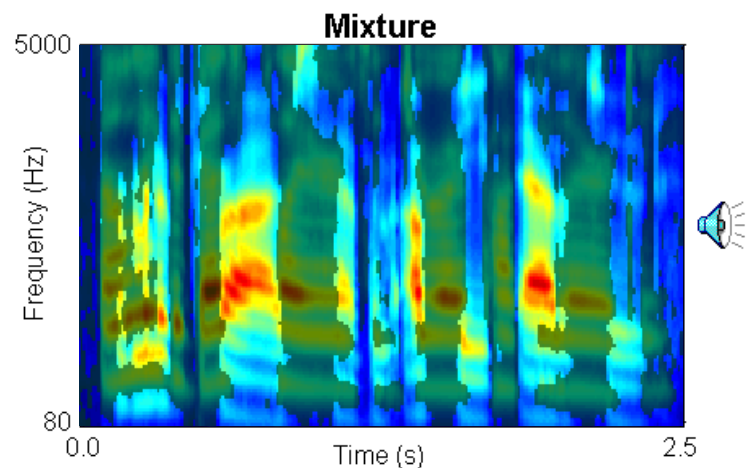
# BSS using TF masking
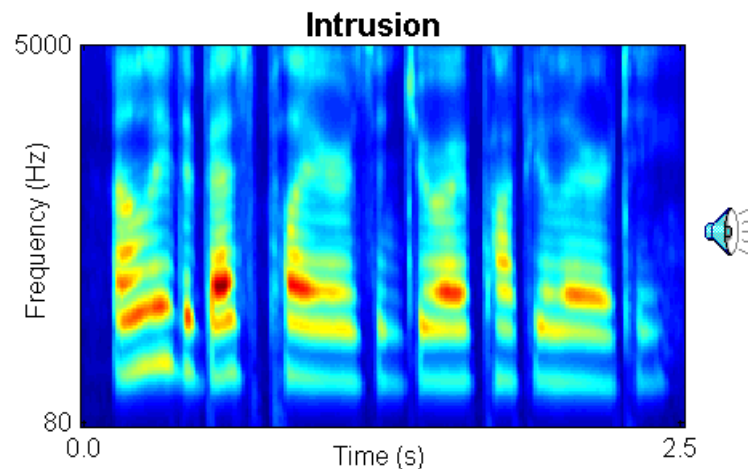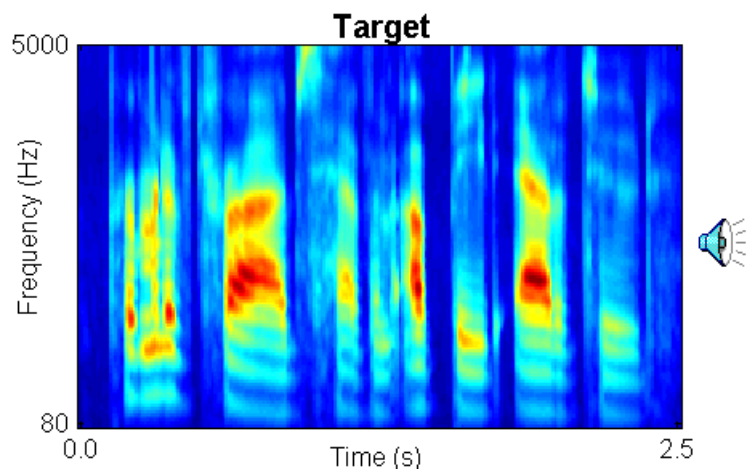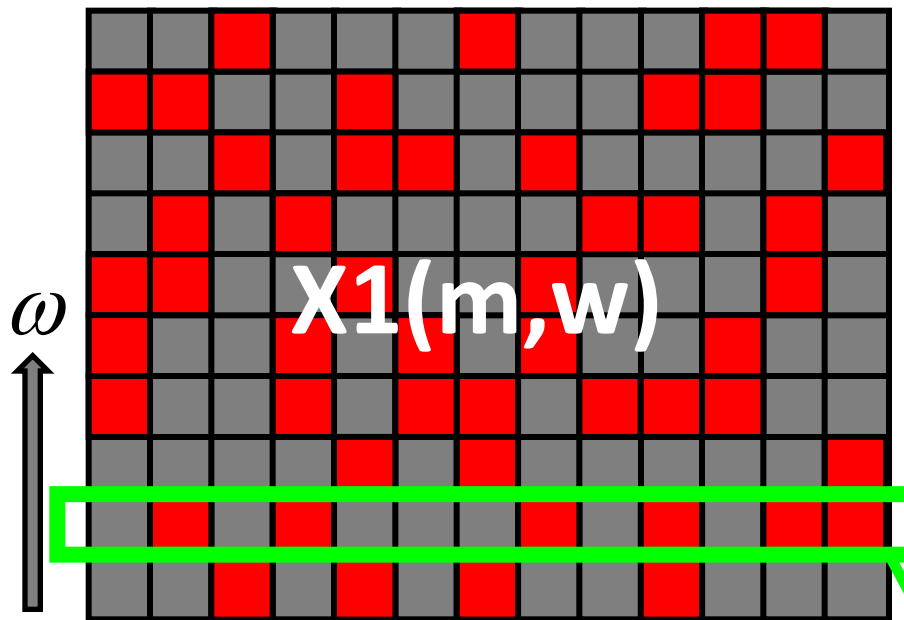


$\omega$

$X(m,w)$

$m$

CASA

Onset
Periodicity
Harmonicity
Locations
Binarual cues

**Sparsity assumption** ------ each TF point is dominated by one source signal.
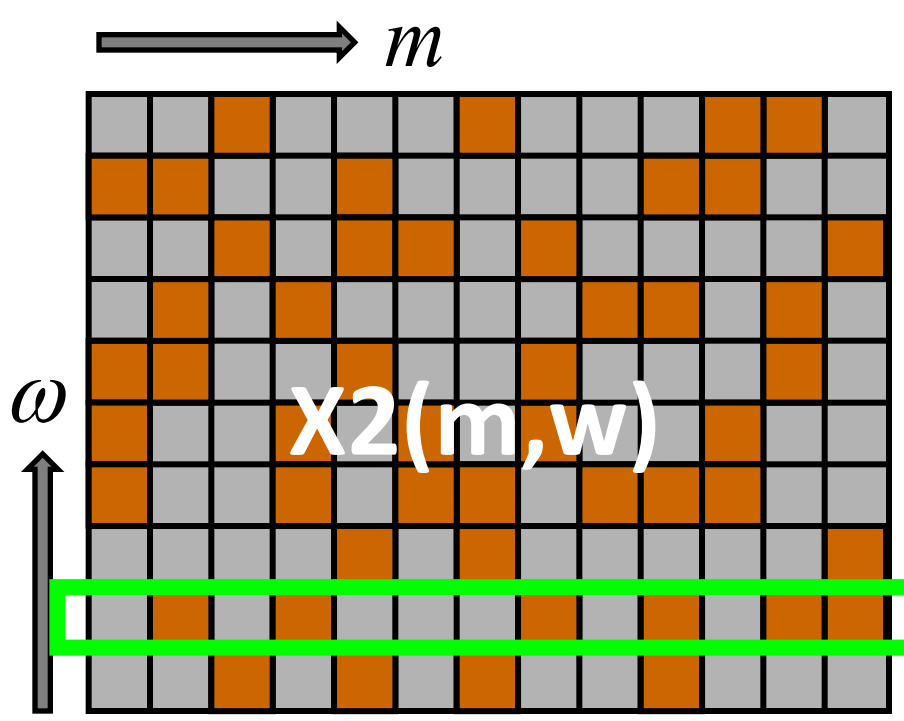
# Benchmark: ideal binary mask (IBM)



Demonstrations by DeLiang Wang, The Ohio State Univ.

$$\frac{X_1(m,\omega)}{X_2(m,\omega)} \Rightarrow \alpha(m,\omega), \beta(m,w)$$

**IPD**

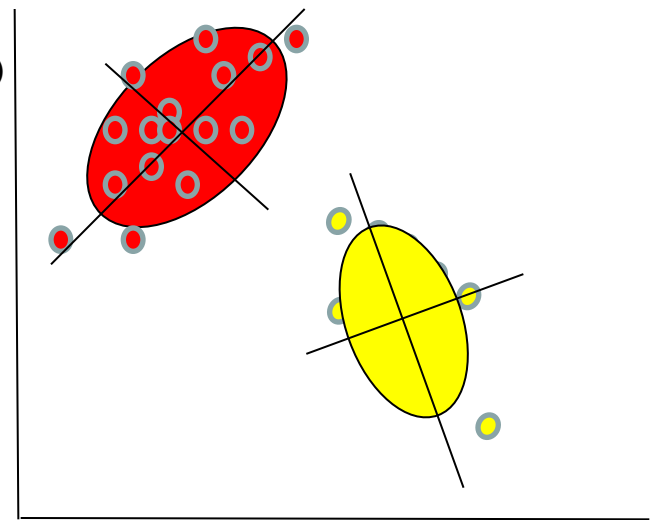**ILD**

UNIVERSITY OF SURREY

www.surrey.ac.uk
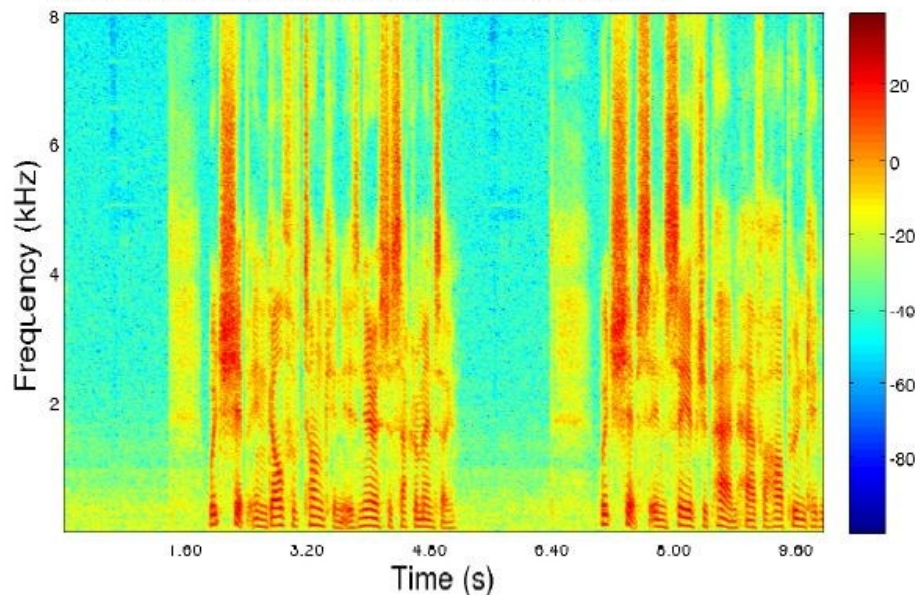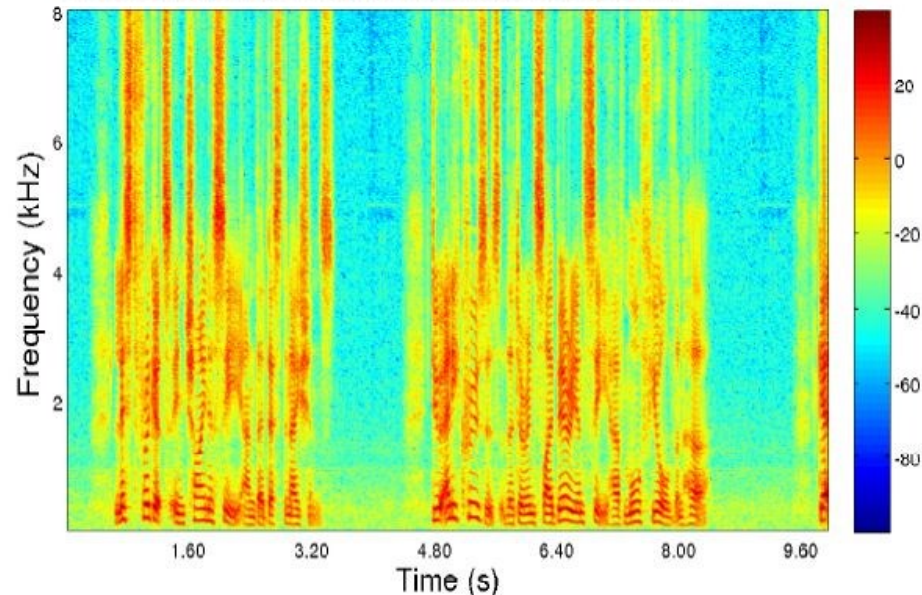
6

# Adverse effects
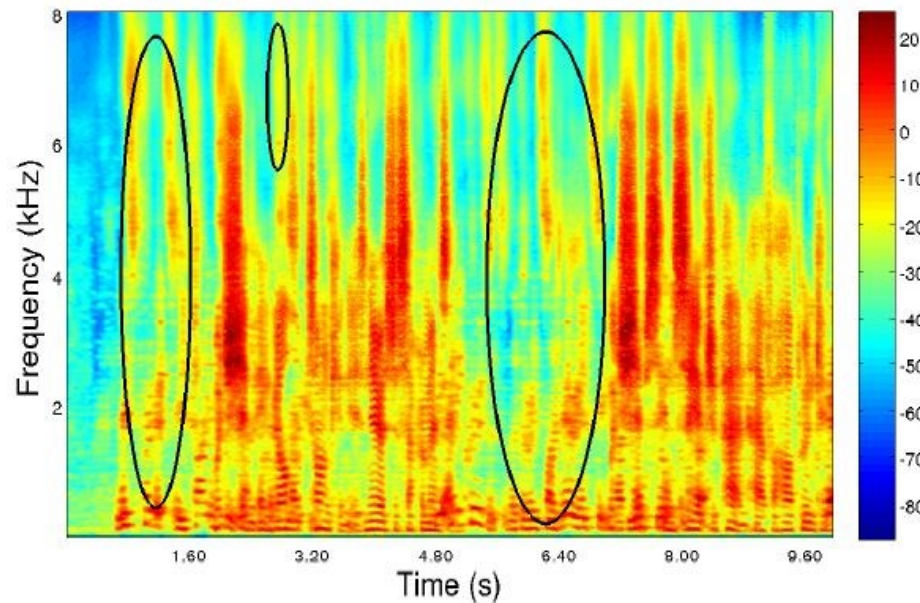
➤ Acoustic noise

➤ Reverberations
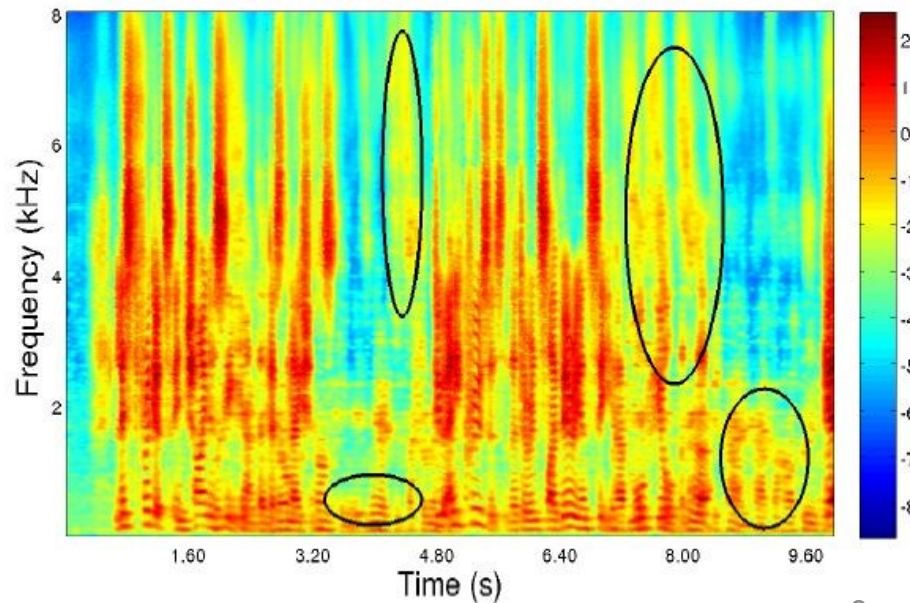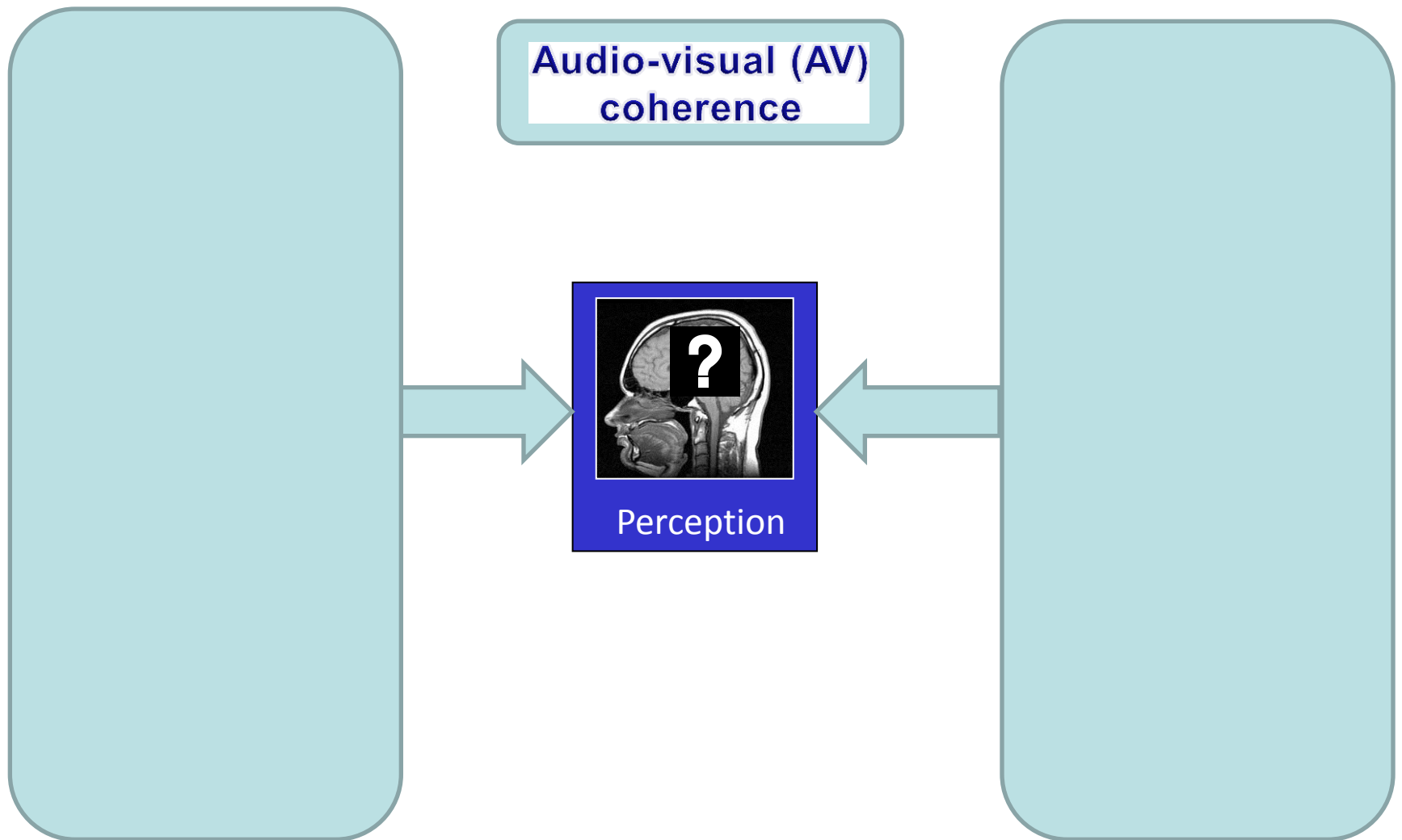
(a) Magnitude spectrum of source 1

(b) Magnitude spectrum of source 2

(c) Magnitude spectrum of source 1 estimate

(d) Magnitude spectrum of source 2 estimate

# Why AV-BSS?----AV coherence

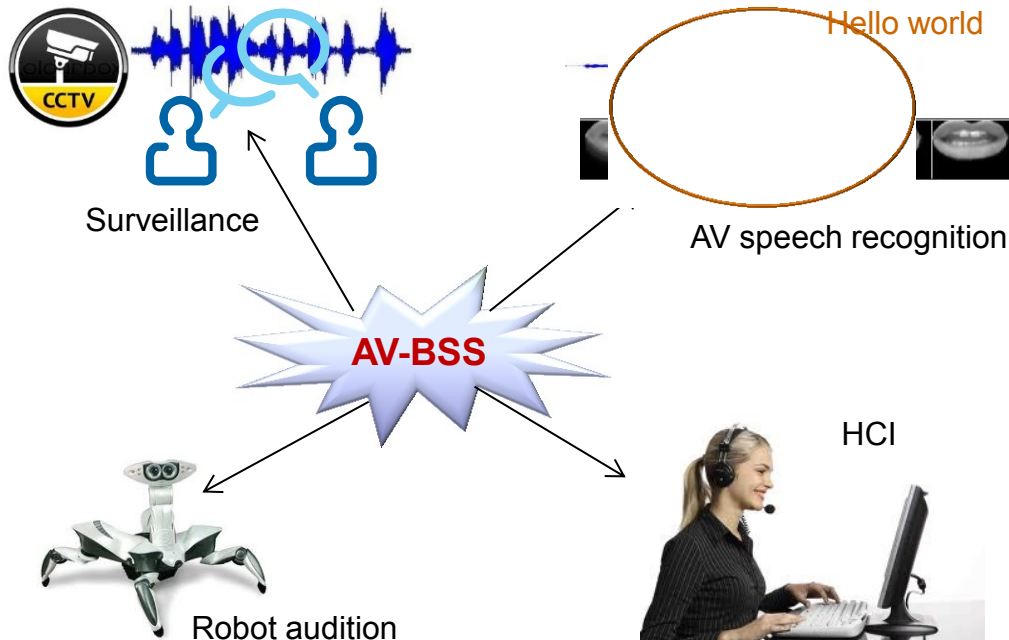**Audio-visual (AV) coherence**



Perception

# Why AV-BSS?

- The audio-domain BSS algorithms **degrade in adverse conditions**.

- The visual stream contains **complementary information** to the coherent audio stream.

**Objective**

**How can the visual modality be used to assist audio-domain BSS algorithms in noisy and reverberant conditions?**

**Potential applications**

Surveillance

Hello world

AV speech recognition

**AV-BSS**

Robot audition
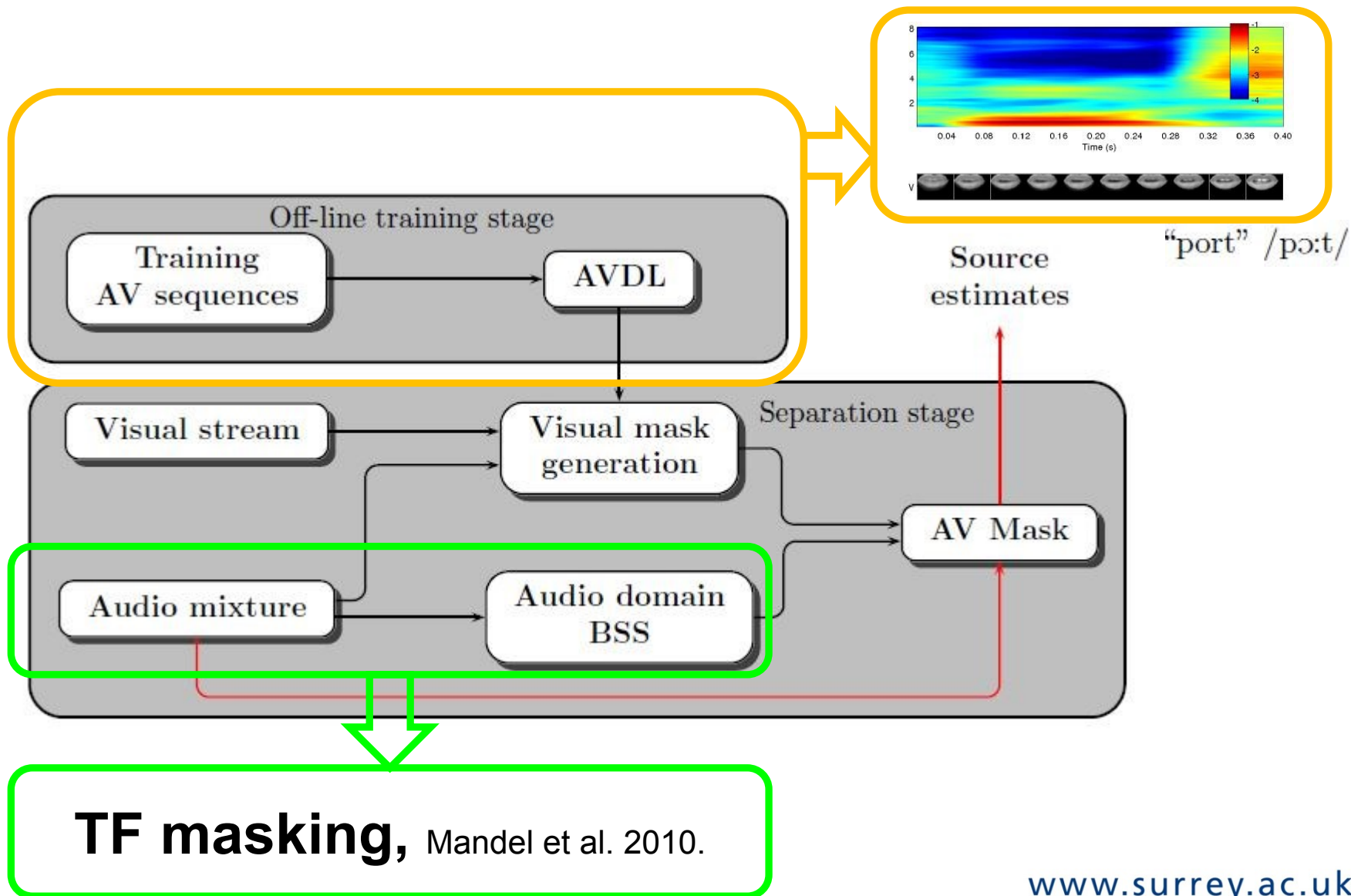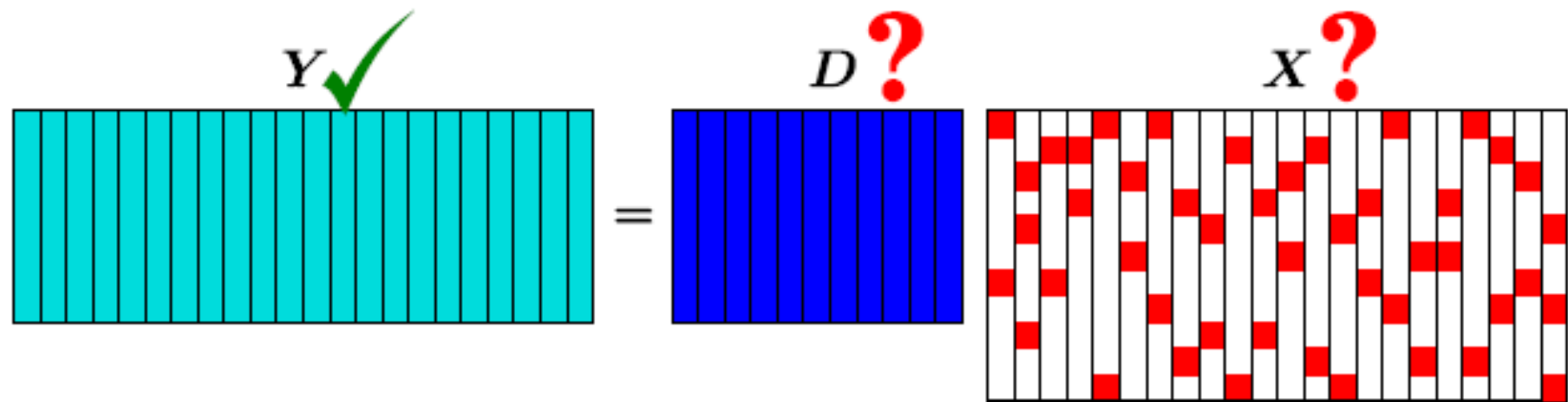
HCI

## Key Challenges

- Reliable **AV coherence modelling**

- **Bimodal differences** in size, dimensionality and sampling rates

- **Fusion of AV coherence** with audio-domain **BSS** methods

UNIVERSITY OF SURREY

www.surrey.ac.uk

10

# AVDL based BSS



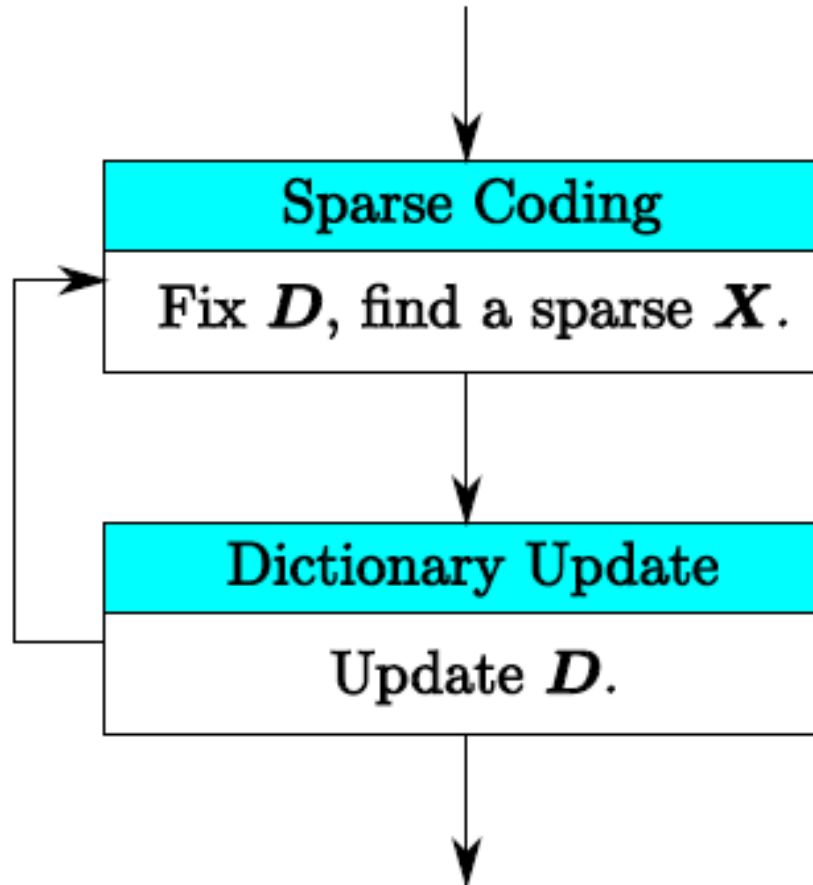"port" /pɔːt/

**TF masking,** Mandel et al. 2010.
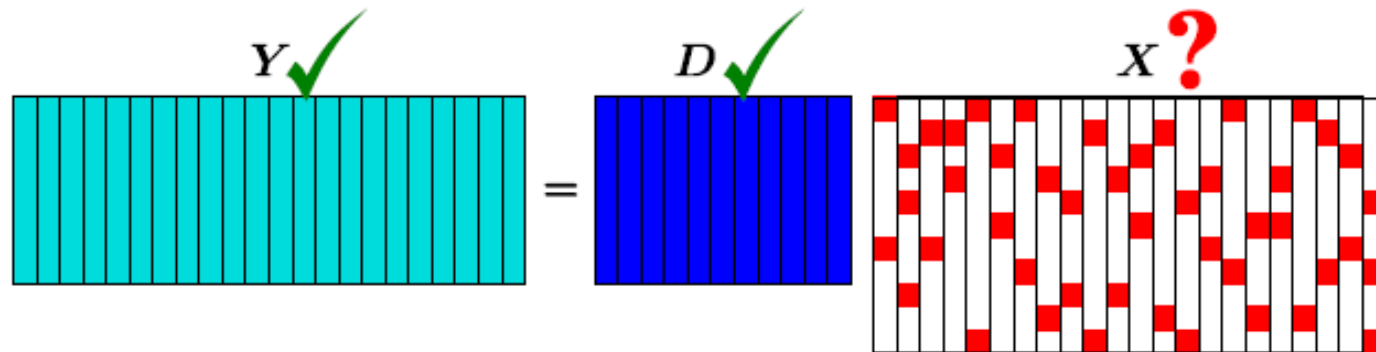
# Dictionary learning



Figures taken from ICASSP 2013 Tutorial 11, by Dai, Maihe and Wang. Likewise for next four pages. Acknowledgement to Wei Dai for making these figures.

# A two-stage procedure

# Sparse coding (approximation)

$$\min \|X\|_0 \text{ s.t. } \|Y - DX\|_F^2 \leq \epsilon.$$

**Greedy algorithms:**

- OMP Y. Pati, et al. 1993; J. Tropp 2004
- Subspace pursuit (SP) W. Dai and O. Milenkovic 2009 CoSaMP D. Needell and J. Tropp 2009
- IHT T. Blumensath and M. Davies 2009

# Dictionary update: the formulation

- **Constraints:**

  ▶ Fixed sparsity pattern

  $$\begin{aligned}\Omega \ &= \{(i,j): \ X_{i,j} \neq 0\}, \\ \mathcal{X}_\Omega \ &= \{X: \ X_{i,j} = 0, \ \forall (i,j) \in \Omega^c\}.\end{aligned}$$
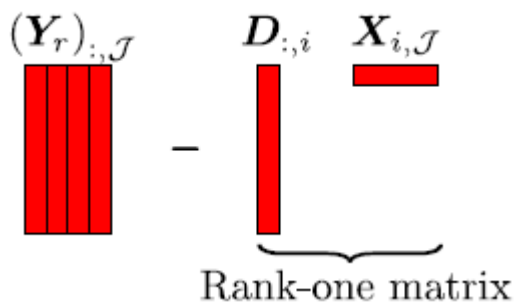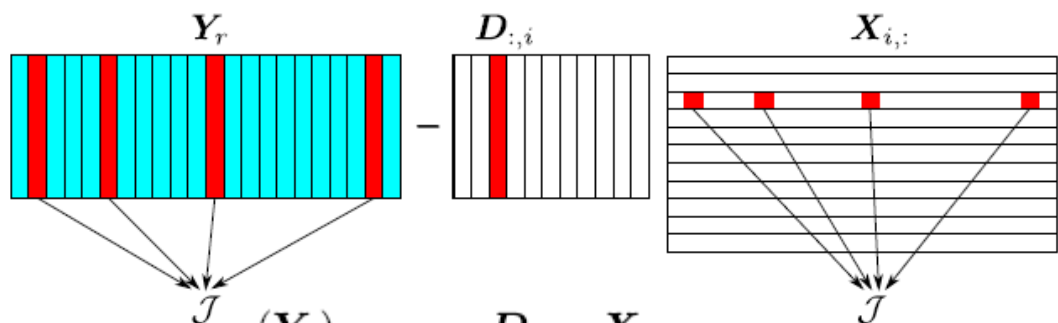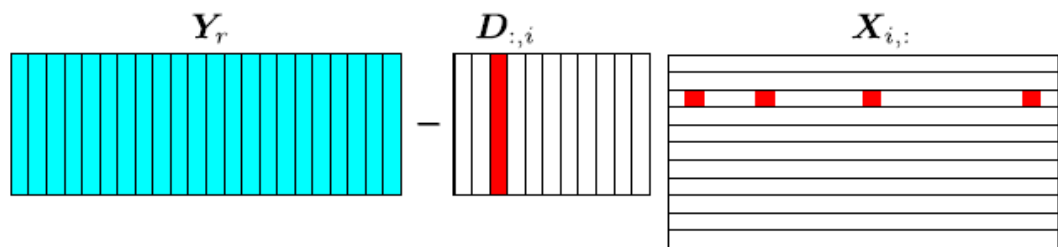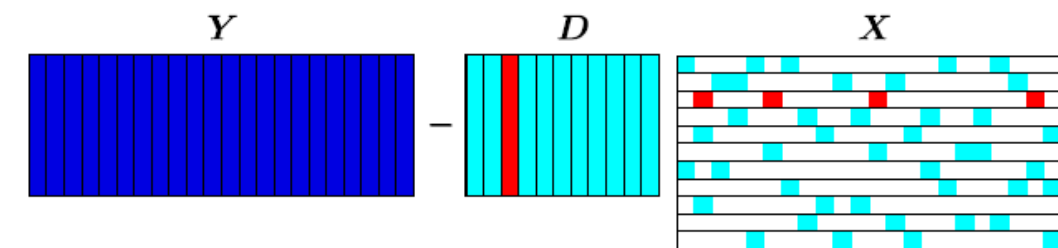
  ▶ Unit norm codewords

  $$\mathcal{D} = \{D: \ \|D_{:,j}\|_2 = 1, \ \forall j \in [d]\}.$$

- **Dictionary Update:**

  $$\min_{D \in \mathcal{D}, \ X \in \mathcal{X}_\Omega} \|Y - DX\|_F^2.$$
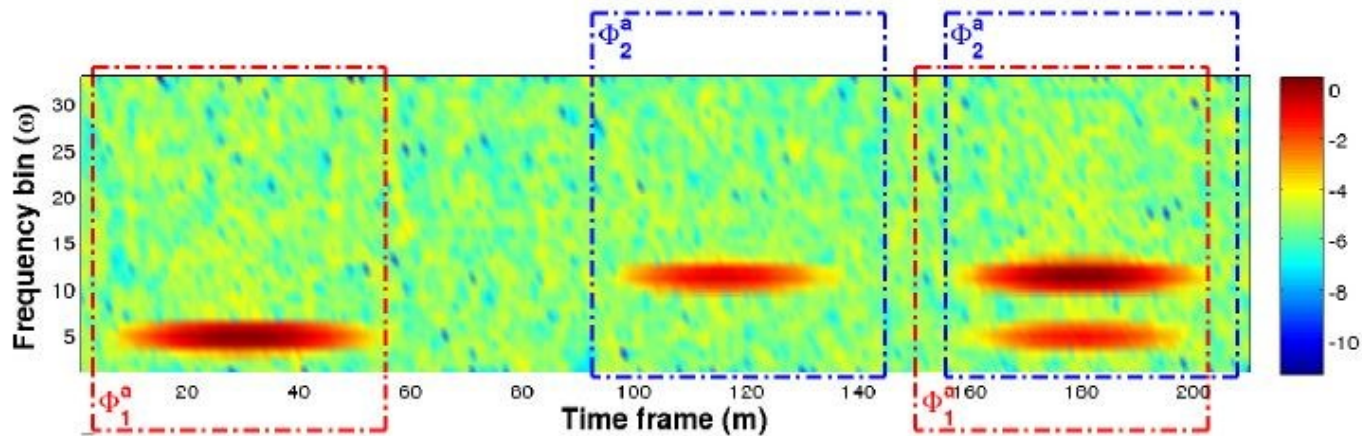
# Dictionary update: K-SVD algorithm



$$\|Y - DX\|^2$$
$$= \|Y - D_{:,j\neq i}X_{j\neq i,:} - D_{:,i}X_{i,:}\|^2$$
$$= \|Y_r - D_{:,i}X_{i,:}\|^2$$
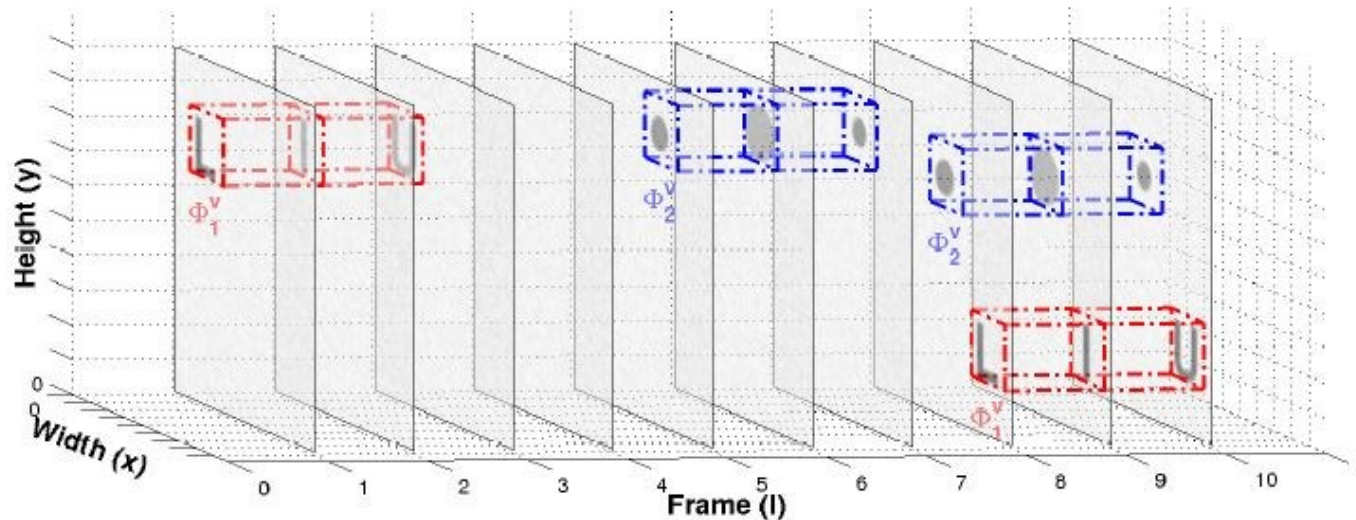$$= \left\|(Y_r)_{:,\mathcal{J}} - D_{:,i}X_{i,\mathcal{J}}\right\|^2 + c$$

Rank-one matrix

# Audio-visual dictionary learning: a generative model

$$\begin{pmatrix} \psi^a(m) \\ \psi^v(y,x,l) \end{pmatrix} \approx \begin{pmatrix} \hat{\psi}^a(m) \\ \hat{\psi}^v(y,x,l) \end{pmatrix} = \sum_{d=1}^{D} \begin{pmatrix} \sum_{\breve{m}=1}^{M_s} c_{d\breve{m}} \phi_d^a(m-\breve{m}) \\ \sum_{\breve{y}=1,\breve{x}=1,\breve{l}=1}^{Y_s,X_s,L_s} b_{d\breve{y}\breve{x}\breve{l}} \phi_d^v(y-\breve{y},x-\breve{x},l-\breve{l}) \end{pmatrix}$$

# Sparse assumption of AVDL



(a) Audio stream $\psi^a$

# Flow of the AVDL



The coding process relies on the matching criterion, how well an atom fits the signal in the MP algorithm

A scanning index is proposed to reduce the computational complexity.

The learning process uses two different update methods, to accommodate different bimodality sparsity constraints.

By mapping the AV sequence to the learned dictionary, a visual mask can be achieved.

# The overall algorithm

**Algorithm 1**: Framework of the Proposed AVDL

**Input**: A training AV sequence $\boldsymbol{\psi} = (\boldsymbol{\psi}^a; \boldsymbol{\psi}^v)$, an initial $\mathcal{D}$ with $K$ atoms, and the number of non-zero coefficients $N$

**Output**: An AV dictionary $\mathcal{D} = \{\boldsymbol{\phi}_k\}_{k=1}^{K}$

1      **Initialization**: $iter = 1, MaxIter$

2      **while** $iter \leq MaxIter$ **do**

3      **%Coding stage**

4      Given $\mathcal{D}$, decompose $\boldsymbol{\psi}$ using **(1)** to obtain $\Omega$.

5      **%Learning stage**

6      Given $\Omega$ and the residual $\boldsymbol{v}$, update $\mathcal{D} = \{\phi_k\}$ for $k = 1, 2, \ldots, K$ to fit model **(1)**.

7      $iter = iter + 1$

# The coding process

$$J^{av}(\bar{\boldsymbol{v}}_{\breve{y}\breve{x}\breve{l}\breve{m}}, \boldsymbol{\phi}_k) = J^a(\bar{\boldsymbol{v}}^a_{\breve{m}}, \boldsymbol{\phi}^a_k) J^v(\bar{\boldsymbol{v}}^v_{\breve{y}\breve{x}\breve{l}}, \boldsymbol{\phi}^v_k),$$

$$J^a_{\mathrm{Mon}} = |\langle \bar{\boldsymbol{v}}^a_{\breve{m}}, \boldsymbol{\phi}^a_k \rangle|$$

$$J^v(\bar{\boldsymbol{v}}^v_{\breve{y}\breve{x}\breve{l}}, \boldsymbol{\phi}^v_k) = \exp\left\{ \frac{-1}{YXL} \left\| \bar{\boldsymbol{v}}^v_{\breve{y}\breve{x}\breve{l}} - \boldsymbol{\phi}^v_k \right\|_1 \right\}.$$

$$[k_n, y_n, x_n, l_n, m_n] = \underset{[k,\breve{y},\breve{x},\breve{l},\breve{m}]}{\arg\max} J^{av}(\bar{\boldsymbol{v}}_{\breve{y}\breve{x}\breve{l}\breve{m}}, \boldsymbol{\phi}_k),$$

$$B(k_n, y_n, x_n, l_n) = 1$$

$$C(k_n, m_n) = J^a(\bar{\boldsymbol{v}}^a_{m_n}, \boldsymbol{\phi}^a_{k_n}).$$

$$\bar{\boldsymbol{v}}^a_{l_n} \leftarrow \bar{\boldsymbol{v}}^a_{l_n} - C(k_n, l_n)\boldsymbol{\phi}^a_{k_n}.$$

# The coding process (algorithm)

---

**Algorithm 2**: The Coding State of the Proposed AVDL

---

**Input**: An AV sequence $\boldsymbol{\psi}$, the dictionary $\mathcal{D} = \{\boldsymbol{\phi}_k\}_{k=1}^K$, the threshold $\delta$, the number of non-zero coefficients $N$

**Output**: The coding parameter set $\Omega = \{\mathbf{B}, \mathbf{C}\}$ and residual $\boldsymbol{v}$

1 **Initialization:** Set $\Omega$ with zero tensors,
$\boldsymbol{v} = \boldsymbol{\psi}, n = 1, J_{opt} = J_{\max} = 0$

2 Calculate $\mathcal{S}^{av}$ using **(10)** to **(13)**.

3 **while** $n \leq N$ *and* $J_{opt} \geq \delta J_{\max}$ **do**

4    % Projection

5    $\mathcal{L} = \begin{cases} \{1 : L_s\}, & n=1 \\ l_{n-1} + \{1 - L : L - 1\}, & \text{otherwise} \end{cases}$

6    **for** $k \leftarrow 1$ **to** $K$ **do**

7       **foreach** $\breve{l} \in \mathcal{L}$ **do**

8          Calculate $J^a(\bar{\boldsymbol{v}}_{\breve{m}}^a, \boldsymbol{\phi}_k^a)$, where $\breve{m}$ is tied with $\breve{l}$ via set **(2)**.

9         **foreach** $(\breve{y}, \breve{x}), \breve{y} \in \{1 : Y_s\}, \breve{x} \in \{1 : X_s\}$ **do**

10            **if** $\mathcal{S}^{av}(\breve{y}, \breve{x}, \breve{l}) = 1$ **then**

11              Obtain $J^v(\bar{\boldsymbol{v}}_{\breve{y}\breve{x}\breve{l}}^v, \boldsymbol{\phi}_k^v)$ via **(6)** and $J^{av}(\bar{\boldsymbol{v}}_{\breve{y}\breve{x}\breve{l}\breve{m}}, \boldsymbol{\phi}_k)$ via **(5)**.

12    % Selection

13    Obtain $[y_n, x_n, l_n, k_n, m_n]$ via **(7)**.

14    Update $\Omega$ via **(8)**.

15    Residual calculation via **(9)**.

16    $J_{opt} = J^{av}(\bar{\boldsymbol{v}}_{y_n x_n l_n m_n}, \boldsymbol{\phi}_{k_n})$

17    **if** $n = 1$ **then**

18    $J_{\max} = J^{av}(\bar{\boldsymbol{v}}_{y_1 x_1 l_1 m_1}, \boldsymbol{\phi}_{k_1})$

19    $n = n + 1$

---

# The learning stage

**Algorithm 3**: The Learning Stage of the Proposed AVDL.

**Input**: The parameter set $\Omega = \{\mathbf{B}, \mathbf{C}\}$, the residual $\boldsymbol{v}$, the old dictionary $\mathcal{D} = \{\boldsymbol{\phi}_k\}_{k=1}^{K}$

**Output**: A new dictionary $\mathcal{D}$

1 **Initialization**: $k = 1$

2 **while** $k \leq K$ **do**

3      Update $\boldsymbol{\phi}_k^a$, $\mathbf{C}$ and $\boldsymbol{v}$ via K-SVD using **(14)** to **(17)**.

4      Update $\boldsymbol{\phi}_k^v$ via the K-means algorithm

5      $\boldsymbol{\phi}_k^v = \text{Mean}\,(b_{k\breve{y}\breve{x}\breve{l}} \bar{\boldsymbol{v}}_{k\breve{y}\breve{x}\breve{l}}^v)$, subject to $b_{k\breve{y}\breve{x}\breve{l}} \neq$

0, $\forall(\breve{y}, \breve{x}, \breve{l})$

6      $k = k + 1$

$$\bar{\boldsymbol{v}}_{\breve{m}}^a \leftarrow \bar{\boldsymbol{v}}_{\breve{m}}^a + c_{k\breve{m}}\boldsymbol{\phi}_k^a, \; \forall\breve{m}. \qquad \boldsymbol{\phi}_k^a \leftarrow \mathbf{ivec}(\mathbf{u}_k|\boldsymbol{\phi}_k^a).$$

$$\Upsilon_k \approx \lambda_k \mathbf{u}_k \mathbf{v}_k^T, \qquad \bar{\boldsymbol{v}}_{\breve{m}}^a \leftarrow \bar{\boldsymbol{v}}_{\breve{m}}^a - c_{k\breve{m}}\boldsymbol{\phi}_k^a, \; \forall\breve{m}.$$

**Synthetic data**



(a) AV: /a/　　(b) AV: /i/　　(c) AV: /o/　　(d) Visual only　　(e)　Audio　　only:



(f) The generated AV synthetic sequence (only one second data is shown)

# AVDL evaluations

**Additive noise added**



(a) AVDL: /a/

(b) AVDL: /i/

(c) AVDL: /o/

(d) Monaci: /a/

(e) Monaci: /i/

(f) Monaci: /o/

# AVDL evaluations

(a) AVDL1

(b) AVDL2

(c) AVDL3

(d) Monaci1

(e) Monaci2

(f) Monaci3

(g) Monaci4

# AVDL evaluations

The approximation error metrics comparison of AVDL and Monaci's method over 50 independent tests over the synthetic data



The proposed AVDL outperforms the baseline approach, giving an average of 33% improvement for the audio modality, together with a 26% improvement for the visual modality.

# AVDL evaluations



(a)

(b)

(c)

# AV mask fusion for AVDL-BSS

$$\mathcal{M}^{av}(m,\omega) = \boxed{\mathcal{M}^a(m,\omega)} \boxed{(\mathcal{M}^v(m,\omega))}$$

**Audio mask**
Statistically generated by evaluating the IPD and ILD of each TF point.

**Visual mask**
Mapping the observation to the learned AV dictionary via the coding stage in AVDL.

# Visual mask generation

$$\mathcal{M}^v(m,\omega) = \begin{cases} 1, & \text{if } \hat{\psi}^a(m,\omega) > \psi^a(m,\omega) \\ \hat{\psi}^a(m,\omega)/\psi^a(m,\omega), & \\ & \text{otherwise.} \end{cases}$$

# AVDL evaluations

**Long Speech**

Sheerman-Chase et al.
LILiR Twotalk database
2011

Lip tracking,
Ong et al. 2008



The first AV atom represents the utterance "**marine**" /mᵊri:n/ while the second one denotes the utterance "**port**" /pᵒ:t/.

www.surrey.ac.uk

31

Demonstration of TF mask fusion in AVDL-BSS

Why do we choose the power law combination, instead of, e.g., a linear combination?

# AVDL-BSS evaluations----SDR

# AVDL-BSS evaluations----OPS-PEASS

Noise-free

10 dB Gaussian noise

# Some examples

| | Mixture | Ideal | Mandel | AV-LIU | **AVDL-BSS** | Rivet | AVMP-BSS |
|---|---|---|---|---|---|---|---|
| A | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 |
| B | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 |
| C | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 |
| D | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 | 🔊 |

## Conclusions

➢AVDL offers an alternative and effective method for modelling the AV coherence within the audio-visual data.
➢The mask derived from AVDL can be used to improve the BSS performance for separating reverberant and noisy speech mixtures

## Future work

➢To achieve dictionary adaptation and source separation simultaneously

Thank you

Q & A

w.wang@surrey.ac.uk

# References

Q. Liu, W. Wang, P. Jackson, M. Barnard, J. Kittler, and J.A. Chambers, "Source separation of convolutive and noisy mixtures using audio-visual dictionary learning and probabilistic time-frequency masking", IEEE Transactions on Signal Processing, vol. 61, no. 22, pp. 5520-5535, 2013.
Q. Liu, W. Wang, and P. Jackson, "Use of bimodal coherence to resolve spectral indeterminacy in convolutive BSS", Signal Processing, 92(8):1916-1927, 2012.
Q. Liu, W. Wang, P. Jackson, and M. Barnard, "Reverberant speech separation based on audio-visual dictionary learning and binaural cues", in Proc. SSP, 2012.