

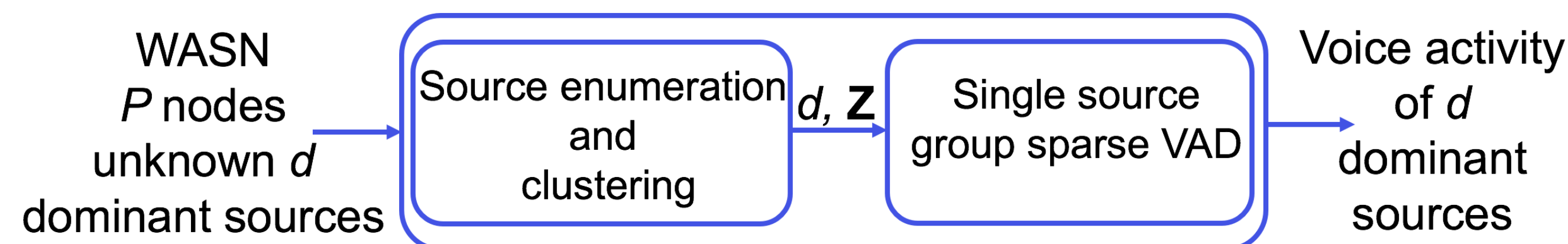


Overview

- Multispeaker voice activity detection (VAD) prerequisite for **speech enhancement** and **noise cancellation** in WASN
- Traditional techniques: Based on energy signatures, e.g., Multiplicative non-negative ICA (**M-NICA**), **degrade** in performance
 - as number of speakers increase
 - in presence of impulsive noise
- Proposed technique:
 - clusters nodes** around each speaker
 - single-source M-NICA with **block-sparse** penalization
 - no knowledge of source/node positions or number of speakers required

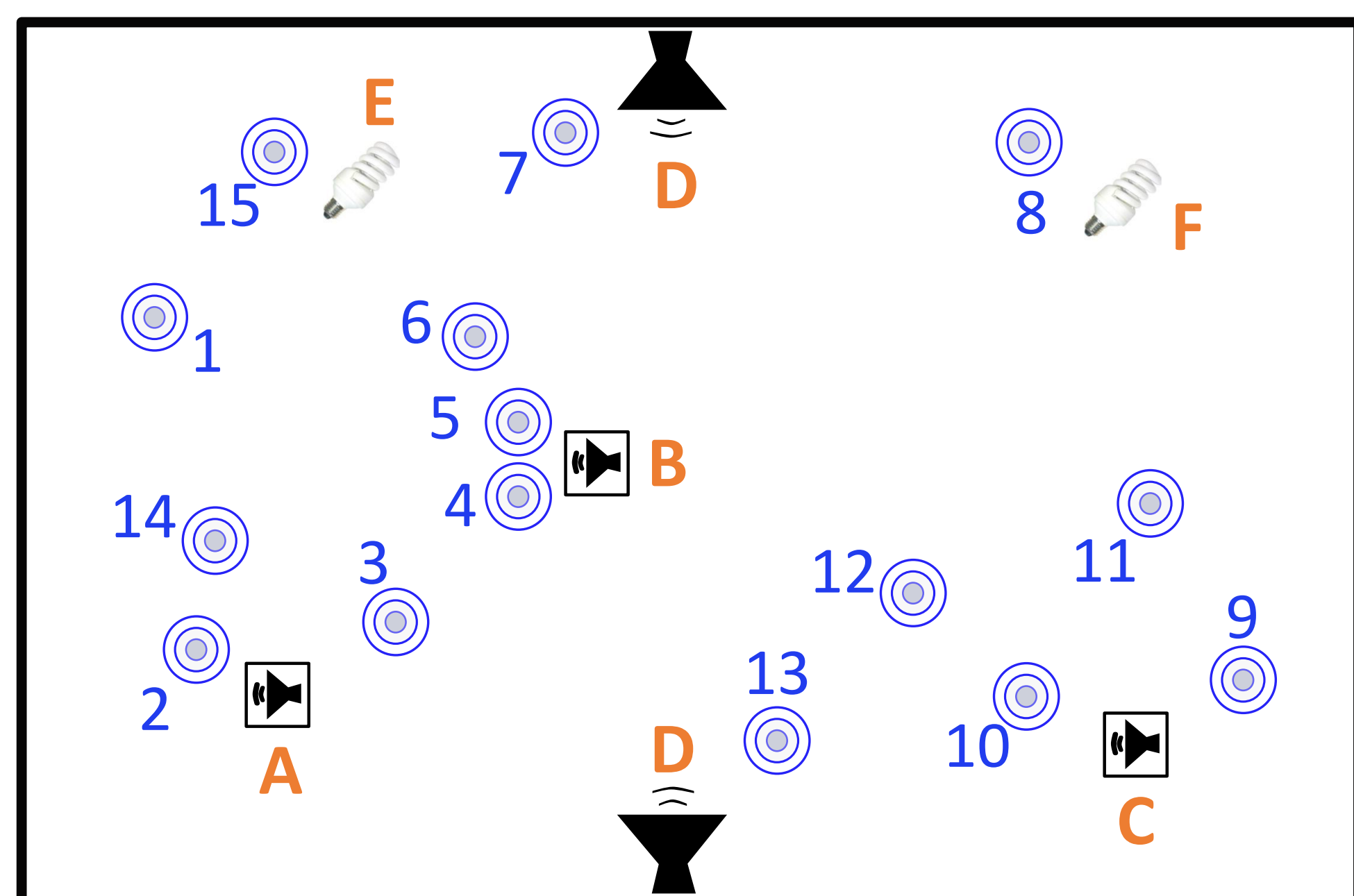
Problem Formulation

- P sensor nodes each with n microphones
- Unknown number of d **dominant sources** among all q sources
- Goal:** Robustly estimate voice activity of d dominant sources



- Cluster information matrix $\mathbf{Z} \in \mathbb{B}^{d \times P}$ whose $\{ij\}$ -th entry is 1 if i th dominant source is observed by j th node and 0 otherwise

WASN Setup



- 20m \times 10m room with $P = 15$ nodes (in blue), $n = 3$ microphones
- $d = 4$ dominant sources (A-D), bulb-flickering noise by sources E and F

Source Enumeration and Node Clustering

- P-channel model:** Observed data at f th frequency index, $\mathbf{x}_p(f) \in \mathbb{C}^n$

$$\mathbf{x}_p(f) = \mathbf{A}_p(f)\mathbf{s}_p(f), \quad p = 1, \dots, P,$$
 $\mathbf{A}_p(f) \in \mathbb{C}^{n \times m_p}$ is acoustic transfer function, $\mathbf{s}_p(f) \in \mathbb{C}^{m_p}$ contains m_p uncorrelated sources
- Strong correlation** among nodes observing **same** dominant source
- Let $\mathbf{R}_{pq} = E[\mathbf{x}_p \mathbf{x}_q^H]$ and $\mathbf{C}_{pq} = \mathbf{R}_{pp}^{-\frac{1}{2}} \mathbf{R}_{pq} \mathbf{R}_{qq}^{-\frac{1}{2}}$

$$\mathbf{C} = \begin{bmatrix} \mathbf{I} & \mathbf{C}_{12} & \dots & \mathbf{C}_{1P} \\ \mathbf{C}_{21} & \mathbf{I} & \dots & \mathbf{C}_{2P} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{C}_{P1} & \mathbf{C}_{P2} & \dots & \mathbf{I} \end{bmatrix}$$
- Composite coherence matrix \mathbf{C}
 - \mathbf{C} has exactly d eigenvalues greater than one
 - Let $\mathbf{u}^{(i)} = [\mathbf{u}_1^{(i)T}, \mathbf{u}_2^{(i)T}, \dots, \mathbf{u}_P^{(i)T}]^T$, be the i th eigenvector associated with the above eigenvalue of \mathbf{C} . The i th source is observed by the p th node iff $\mathbf{u}_p^{(i)} \neq \mathbf{0}$
- Bootstrap-based hypothesis tests for estimating d and \mathbf{Z}

No heuristic thresholds required

Group Sparse Voice Activity Detection

- Received signals \mathbf{y} composed of N time samples summarized in matrix \mathbf{Y}

$$\mathbf{Y} = \mathbf{a}\mathbf{s} + \mathbf{W}$$

$$\mathbf{Y} \in \mathbb{R}^{n^{(i)} \times N} \quad \mathbf{a} \in \mathbb{R}^{n^{(i)} \times 1} \quad \mathbf{s} \in \mathbb{R}^{1 \times N} \quad \mathbf{W} \in \mathbb{R}^{n^{(i)} \times N}$$
- Singular value decomposition $\text{SVD}(\mathbf{Y}) = \sigma \mathbf{u} \mathbf{v}^T$
- Information on shape of dominant source energy signal contained in \mathbf{v} [2]
- Divide signal matrix \mathbf{Y} and \mathbf{v} into L **groups** of length N_g

$$\mathbf{Y} = [\mathbf{Y}_{g,1}, \dots, \mathbf{Y}_{g,L}], \quad \mathbf{Y}_{g,l} \in \mathbb{R}^{n^{(i)} \times N_g}$$

- Penalize** SVD optimization problem **to enforce group sparsity**

$$\argmin_{\mathbf{v}} \left\| \sum_{l=1}^L \mathbf{Y}_{g,l} - \sigma \mathbf{u} \mathbf{v}_{g,l}^T \right\| + \lambda_{\mathbf{v}} \sum_{l=1}^L \underbrace{\sqrt{\sum_k^{N_g} |v_{g,l}[k]|^2}}_{\text{mixed } \ell_1/\ell_2 \text{ norm}}$$

- Tuning parameter $\lambda_{\mathbf{v}}$ determines degree of sparsity in \mathbf{v}
- Intrinsic Voice Activity Detection** without **heuristic threshold**
- Without sparsity enforcing penalty term, all entries of \mathbf{v} are generally non-zero
- Group sparse structure of \mathbf{v} forces entire groups $\mathbf{v}_{b,l}$ to be zero

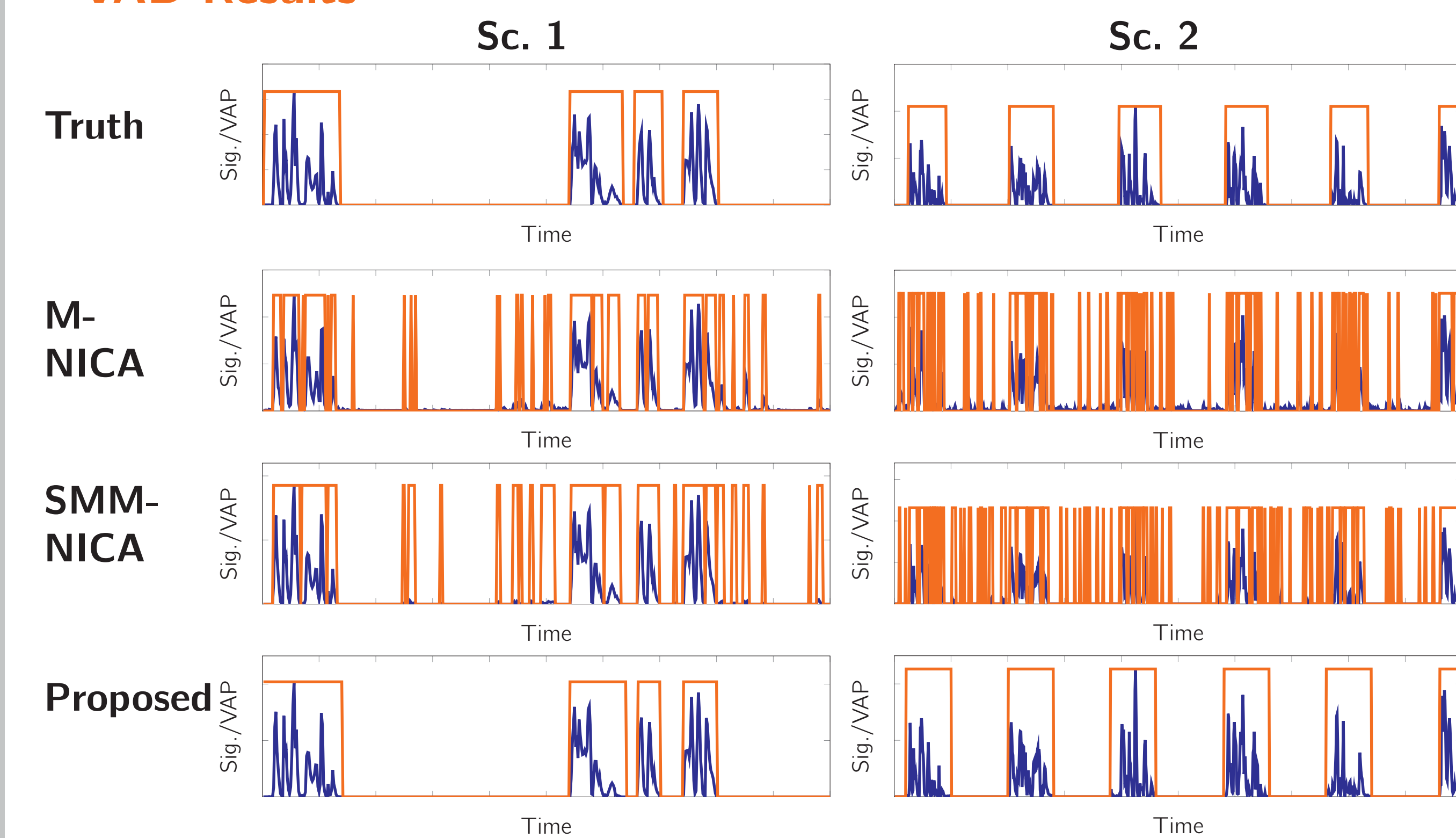
Results

- Each node corrupted by AWGN, sampled at 16kHz
- Hamming window STFT
- Sc. 1: $\hat{d}=3$, 15s
- Sc. 2: $\hat{d}=4$, 30s

Clustering results

Source	Cluster Nodes (Sc. 1)	Cluster Nodes (Sc. 2)
A	2 and 3	2, 3 and 14
B	4, 5 and 6	4, 5 and 6
C	9 and 10	9, 10 and 11
D	Not active	7, 12 and 13

VAD Results



Advantages of Proposed Group-Sparse Method

- + Outperforms standard methods and clustered sparse approach
- + Is scalable to large WASNs with many speakers
- + Handles impulsive noise
- + Does not require heuristic thresholds for node clustering and VAD

Acknowledgements

This research was supported by the German Research Foundation (DFG) under grants SCHR 1384/3-2 and ZO 215/17-2. The WASN speech dataset has been generated within the EU FET-Open Project HANDiCAMS (GA no. 323944).

Bibliography

- [1] T. Hasija, C. Lameiro, T. Marrinan & P.J. Schreier, "Determining the Dimension and Structure of the Subspace Correlated Across Multiple Data Sets," *arXiv*, 2019
- [2] A. Bertrand & M. Moonen, "Energy-based multi-speaker voice activity detection with an ad hoc microphone array," *ICASSP*, 2010
- [3] L. K. Hamaidi, M. Muma & A.M. Zoubir, "Multi-Speaker Voice Activity Detection by an Improved Multiplicative Non-Negative Independent Component Analysis with Sparseness Constraints", *ICASSP*, 2017