

TOWARDS ONLINE SOURCE COUNTING IN SPEECH MIXTURES APPLYING A VARIATIONAL EM FOR COMPLEX WATSON MIXTURE MODELS

Lukas Drude, Aleksej Chinaev, Dang Hai Tran Vu, Reinhold Häb-Umbach

University of Paderborn, Germany, {drude, chinaev, tran, haeb}@nt.uni-paderborn.de, <http://www-nt.uni-paderborn.de>

Introduction

- Source counting treated as a model selection problem
- Directions learned by a complex Watson mixture model
- Observation selection based on power-quantile
- Comparison with DoA-based variational EM algorithm
- Proof of concept for an online algorithm

Modeling and feature extraction

Convolutional mixture model:

$$\mathbf{X}(t, f) = \sum_{k=1}^K \mathbf{H}_k(f) S_k(t, f) + \mathbf{N}(t, f)$$

Phase, frequency and unit-norm normalization:

$$\tilde{X}_d(t, f) = |X_d(t, f)| \exp\left(j \frac{\arg(X_d(t, f) X_1^*(t, f))}{2f_{\text{real}} c^{-1} d_{\text{max}}}\right)$$

$$\mathbf{Y}(t, f) = \tilde{\mathbf{X}}(t, f) / \|\tilde{\mathbf{X}}(t, f)\|$$

⇒ Phase solely determined by source position

Statistical model

Complex Watson mixture model:

$$p(\mathbf{Y} | \mathbf{W}_{1:K+1}; \pi_{1:K+1}, \kappa_{1:K+1}) = \sum_{k=1}^{K+1} P(c(t, f) = k; \pi_{1:K+1}) \frac{1}{\mathcal{O}_W(\kappa_k)} e^{\kappa_k |\mathbf{W}_k^H \mathbf{Y}(t, f)|^2}$$

$P(c(t, f) = k; \pi_{1:K+1})$ = Categorical distribution

$p(\mathbf{W}_k; \mathbf{B}_k)$ = Complex Bingham distribution as prior

Arguments for a complex Watson mixture model:

- All spatial information preserved in observations
- A priori distribution available
⇒ Variational EM (VEM) developed
- Distance measure $\mathbf{W}^H \mathbf{Y}$ resembles a spatial correlation in the beamforming concept

Offline algorithm

- First iteration ($\nu = 1$), Quantile criterion:

$$A^{(1)}(t, f) = \begin{cases} \frac{1}{2} + \frac{1}{2}a, & \mathbf{X}^H(t, f) \mathbf{X}(t, f) > P, \\ \frac{1}{2} - \frac{1}{2}a, & \mathbf{X}^H(t, f) \mathbf{X}(t, f) < P, \end{cases} \quad P = \text{quantile}(\mathbf{X}^H(t, f) \mathbf{X}(t, f), q)$$

⇒ Emphasize observations containing a dominant source

- Next iterations, updating observation weights:

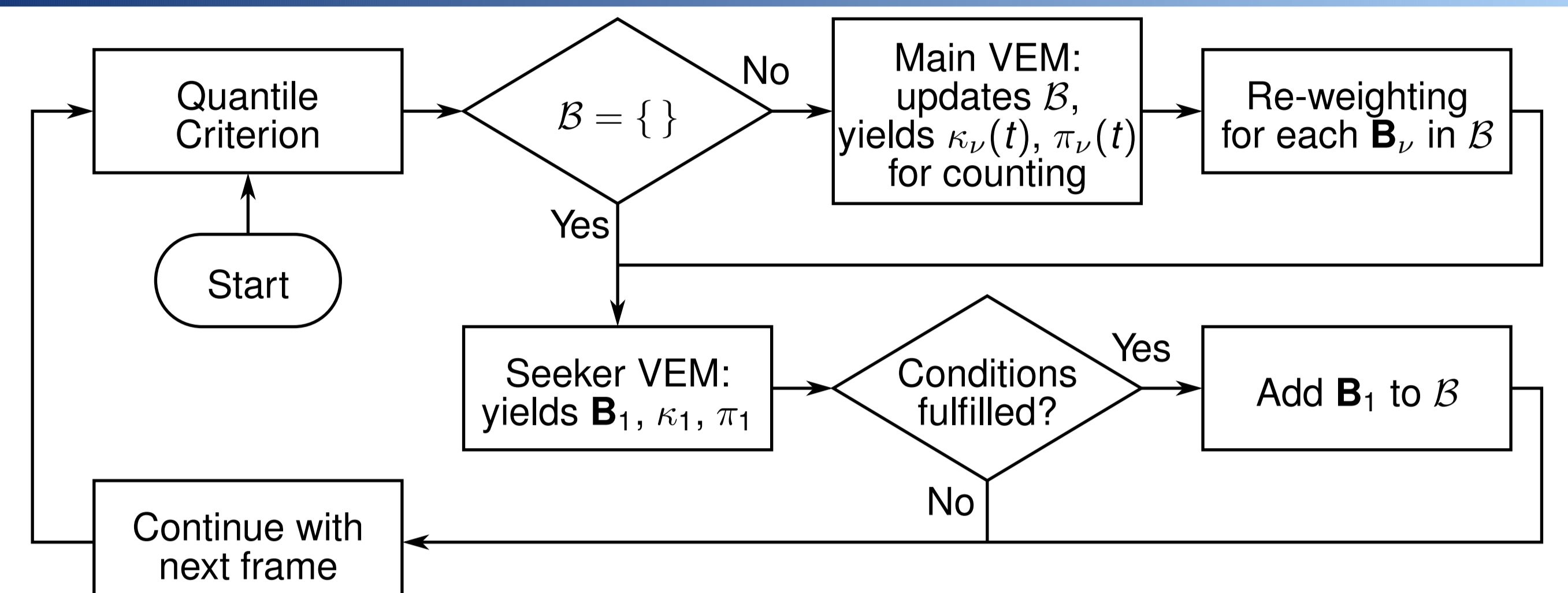
$$A^{(\nu+1)}(t, f) = A^{(\nu)}(t, f) \left(1 - e^{-\kappa_{\text{Re}} (|\hat{\mathbf{W}}_{\nu}^H \mathbf{Y}(t, f)|^2 - 1)}\right)$$

⇒ Deemphasize observations related to detected sources

- Learn one source and noise component for each ν :

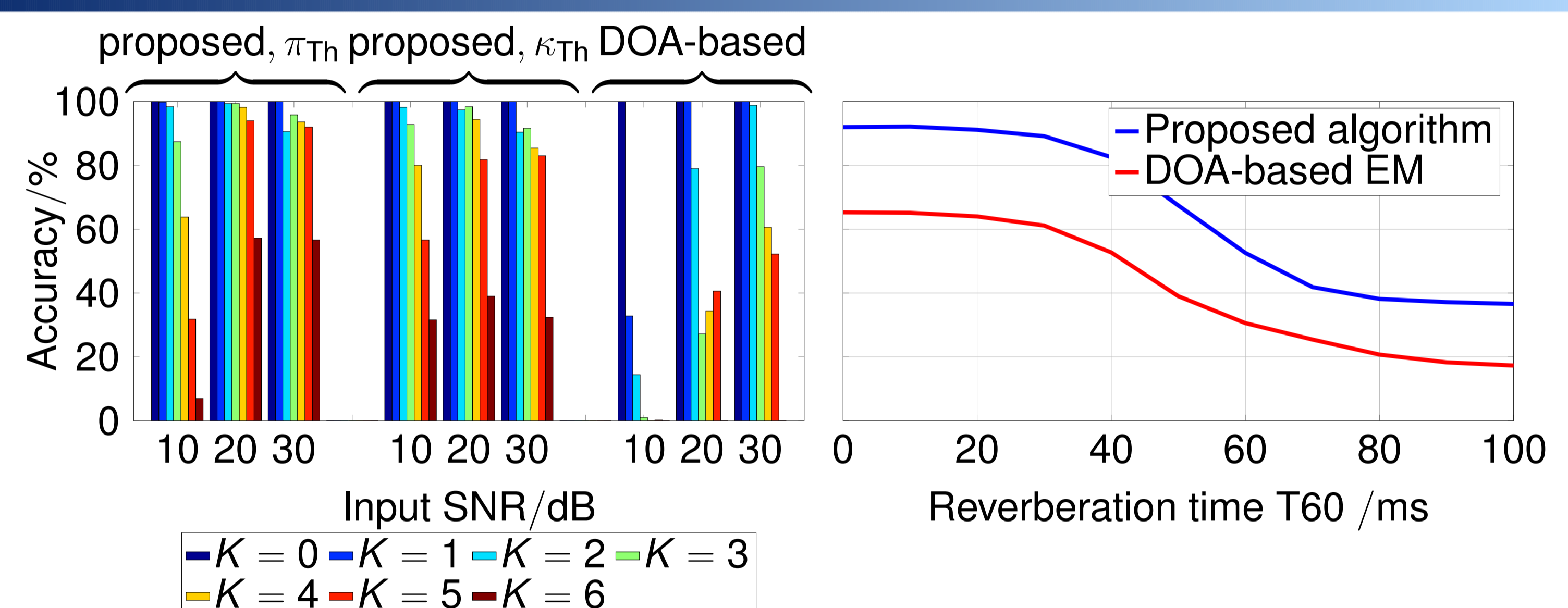
- 1: Calculate $A^{(1)}(t, f)$
- 2: **for** $\nu = 1, \dots, \nu_{\text{max}}$ **do**
- 3: Use VEM with $\mathbf{Y}(t, f)$ and $A^{(\nu)}(t, f)$ to get \mathbf{B}_{ν} and κ_{ν}
- 4: Calculate principal component $\mathbf{W}_{\nu} = \mathcal{P}(\mathbf{B}_{\nu})$
- 5: **if** $\nu < \nu_{\text{max}}$: **then** Reweight observations **end if**
- 6: **end for**
- 7: Calculate $s_1 = 0, s_{\nu} = \max_{\nu'=1, \dots, \nu-1} |\mathbf{W}_{\nu}^H \mathbf{W}_{\nu'}| \forall \nu = 1$
- 8: Count iterations where $\kappa_{\nu} > \kappa_{\text{Th}} \wedge s_{\nu} < s_{\text{Th}}$

Online algorithm



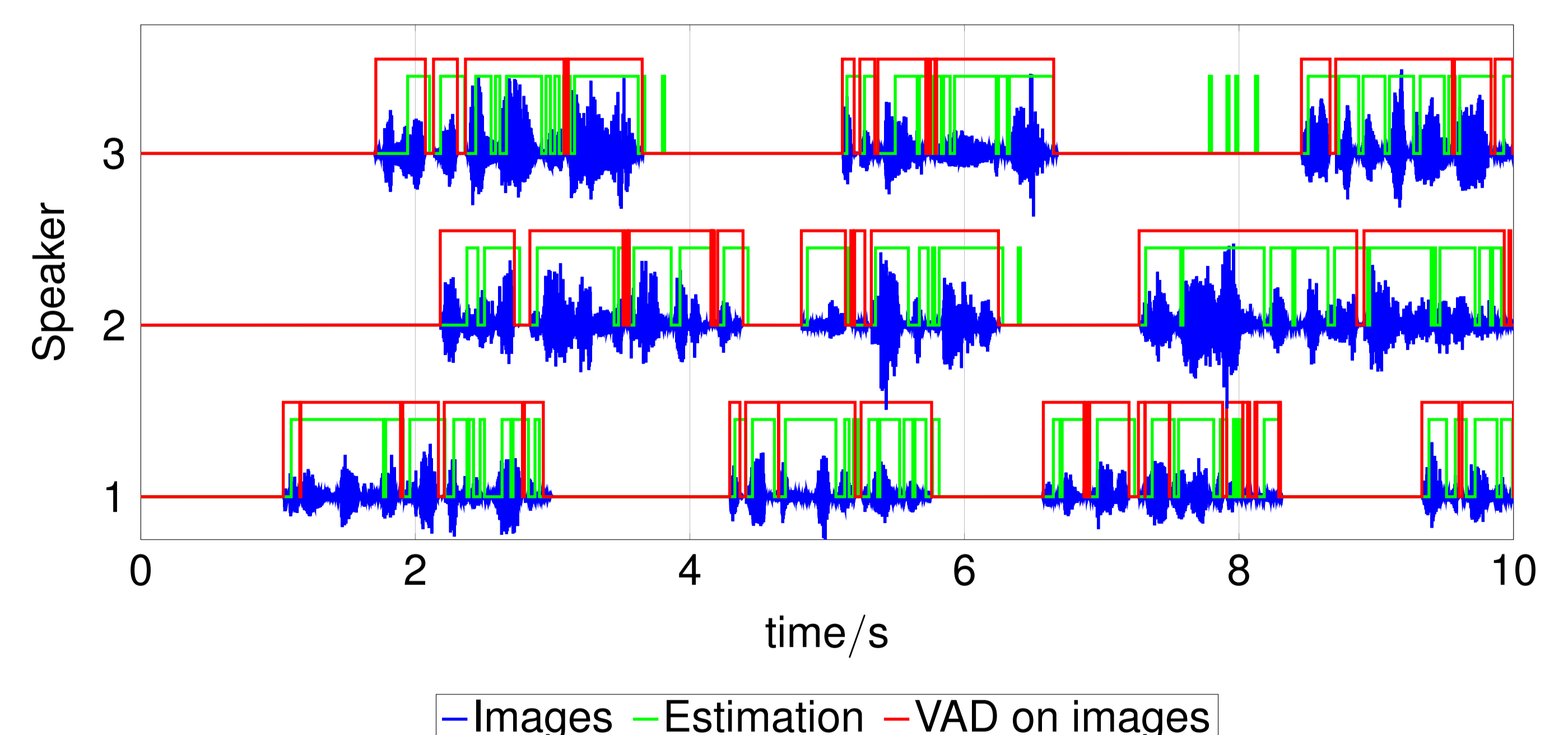
- Frame-wise decision with low latency
- Seeker VEM needs one frame to find a candidate
- Main VEM needs one frame to validate a source
- Conditions to accept a new speaker similar to offline case

Offline results



- Simulated room with image method
- White Gaussian sensor noise
- Uninformative complex Bingham prior
- Comparison: DoA-based VEM
- Proposed algorithm more noise robust
- Both algorithms suffer from reverberation

Online results



- Rediscovery of previously active speakers within one frame
- Diarization error rate 30% of maximum of 187%

Conclusions

- Robust observations emphasized
- Initialization problem relaxed by searching for single speaker at a time
- Low latency online algorithm
- Susceptible to reverberation because of frequency normalization ⇒ Avoid frequency normalization at the cost of introducing the permutation problem.

