



NORTHWESTERN
UNIVERSITY

Combining Modeling of Singing Voice and Background Music for Automatic Separation of Musical Mixtures

Zafar Rafii¹, François G. Germain², Dennis L. Sun^{2,3}, and Gautham J. Mysore⁴

¹Northwestern University, Department of Electrical Engineering & Computer Science

²Stanford University, Center for Computer Research in Music and Acoustics

³Stanford University, Department of Statistics

⁴Adobe Research

zafarrafii@u.northwestern.edu, {fgermain,dlsun}@stanford.edu, gmysore@adobe.com



Adobe

Introduction

Context

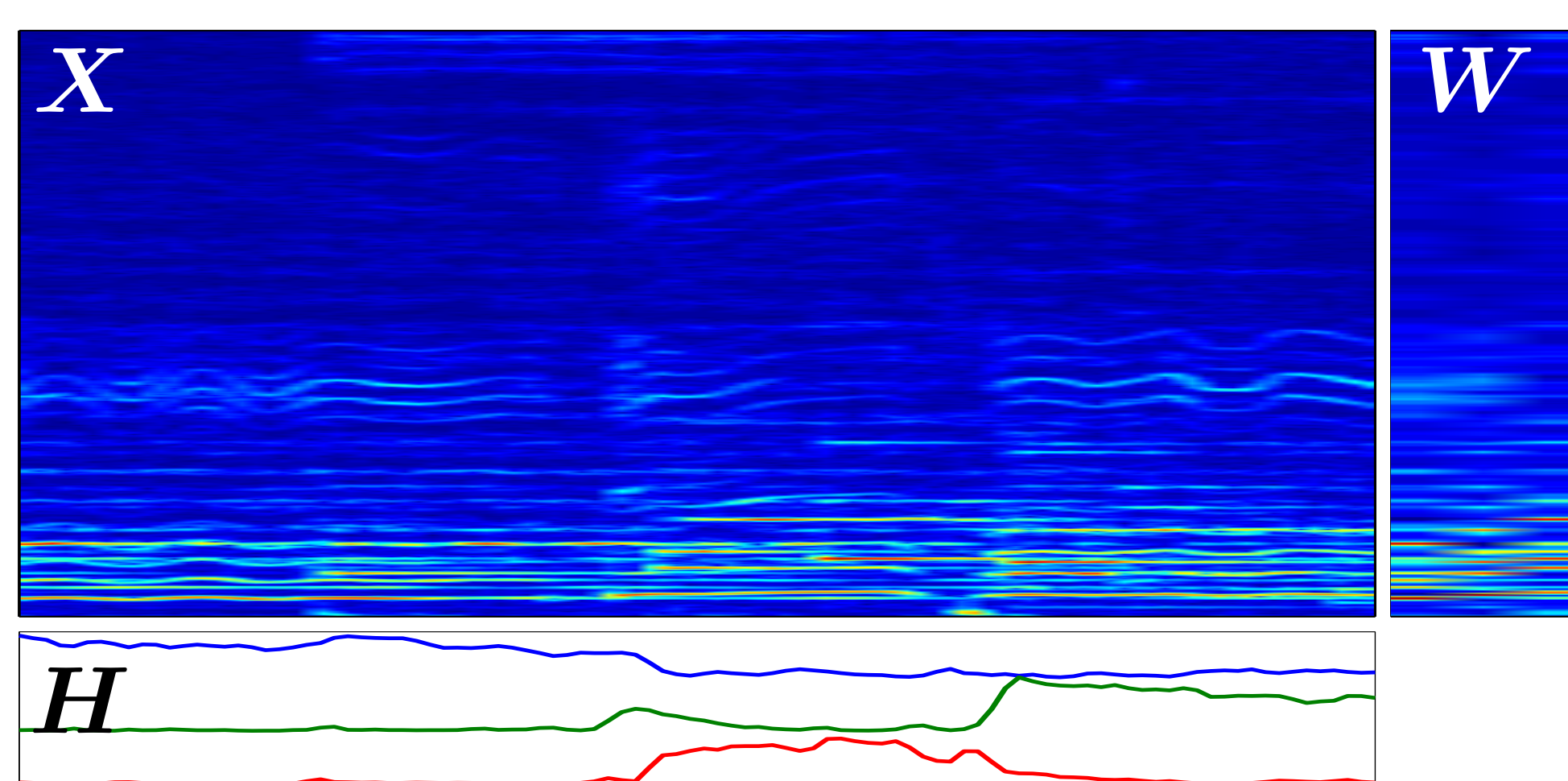
Many music/voice separation techniques focus on modeling one source and explain the other one as the residual, often resulting in degraded performance for the latter.

Proposed method

We propose to combine:

- A new Non-negative Matrix Factorization (NMF) based technique that explicitly models the voice without requiring singer-dependent training examples.
- An existing technique that explicitly models the background music.

Modeling the singing voice



We model the non-negative magnitude spectrogram \mathbf{X} as a linear combination of non-negative basis vectors stored as columns of \mathbf{W} , with the rows of a matrix \mathbf{H} representing the activations of the basis vectors over time, such that:

$$\mathbf{X} \approx \mathbf{W}\mathbf{H}$$

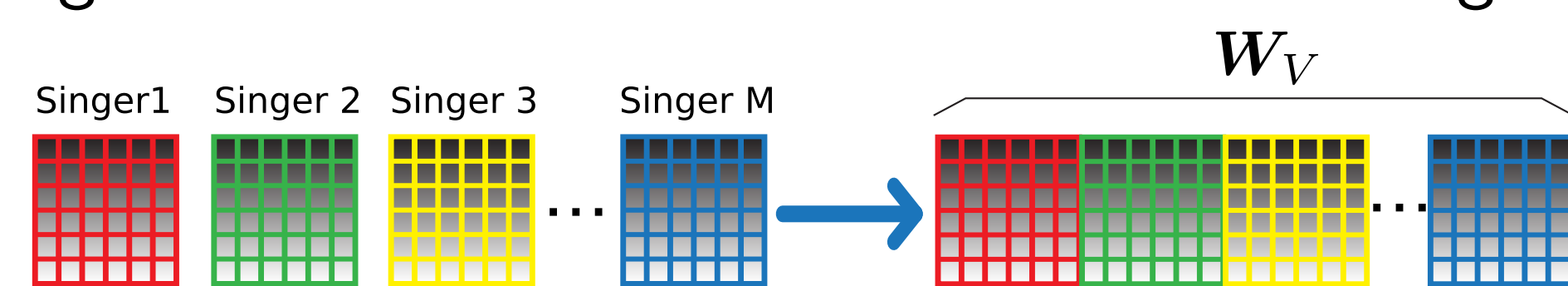
The factorization is learned from \mathbf{X} through the minimization of a cost function $D(\mathbf{X}||\mathbf{W}\mathbf{H})$.

The separation algorithm happens in two steps:

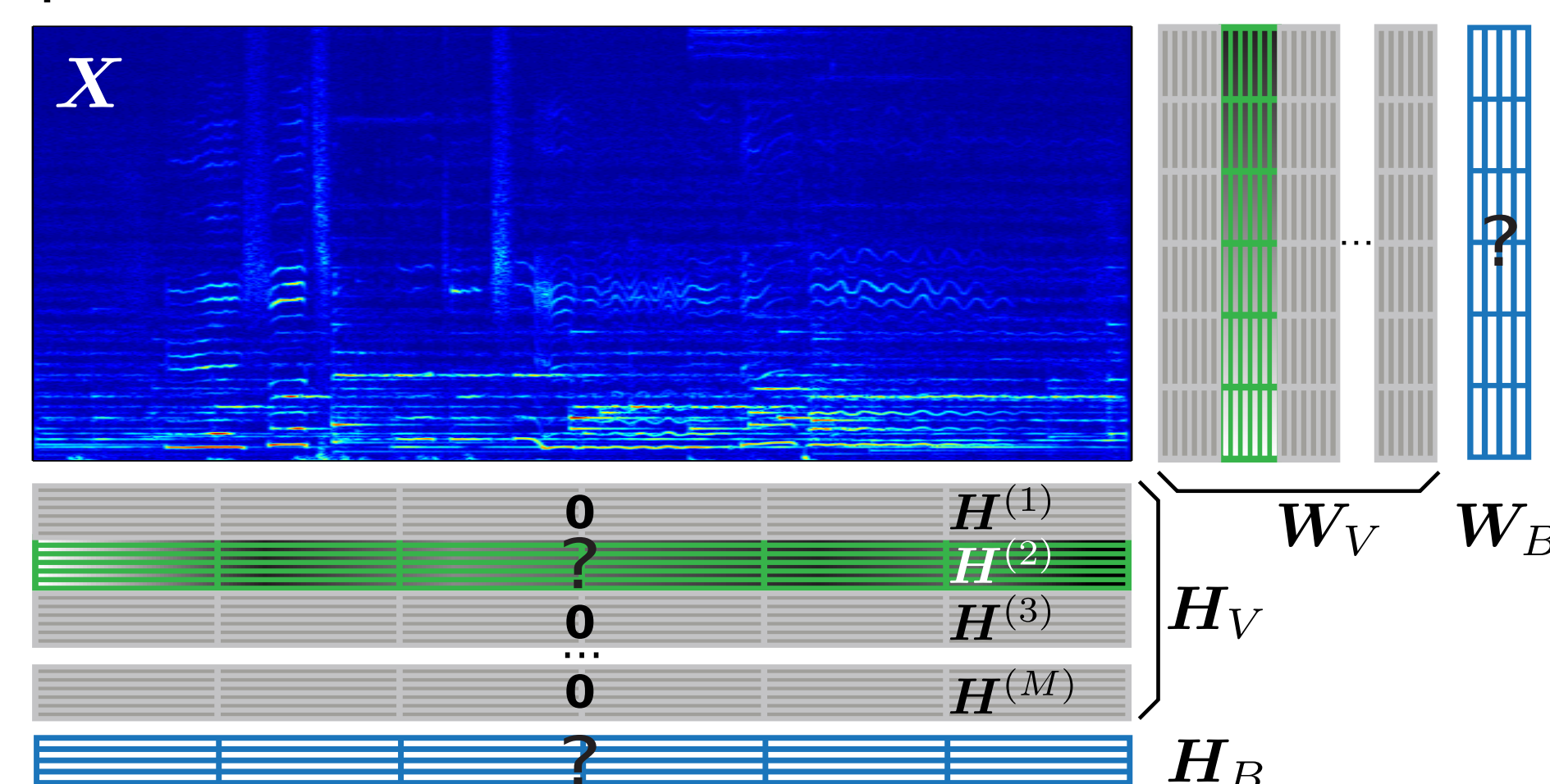
- 1 Train the basis vectors for the source models in \mathbf{W} , using *clean isolated* samples of each source.
- 2 Estimate the activations \mathbf{H} of each source from the test samples, keeping \mathbf{W} fixed.

Universal Voice Models (UVM)

Gathering isolated training samples for either the singer or the music is usually impractical. This limitation can be alleviated in the proposed approach by learning a **universal voice model** concatenating the models associated with M different singers.

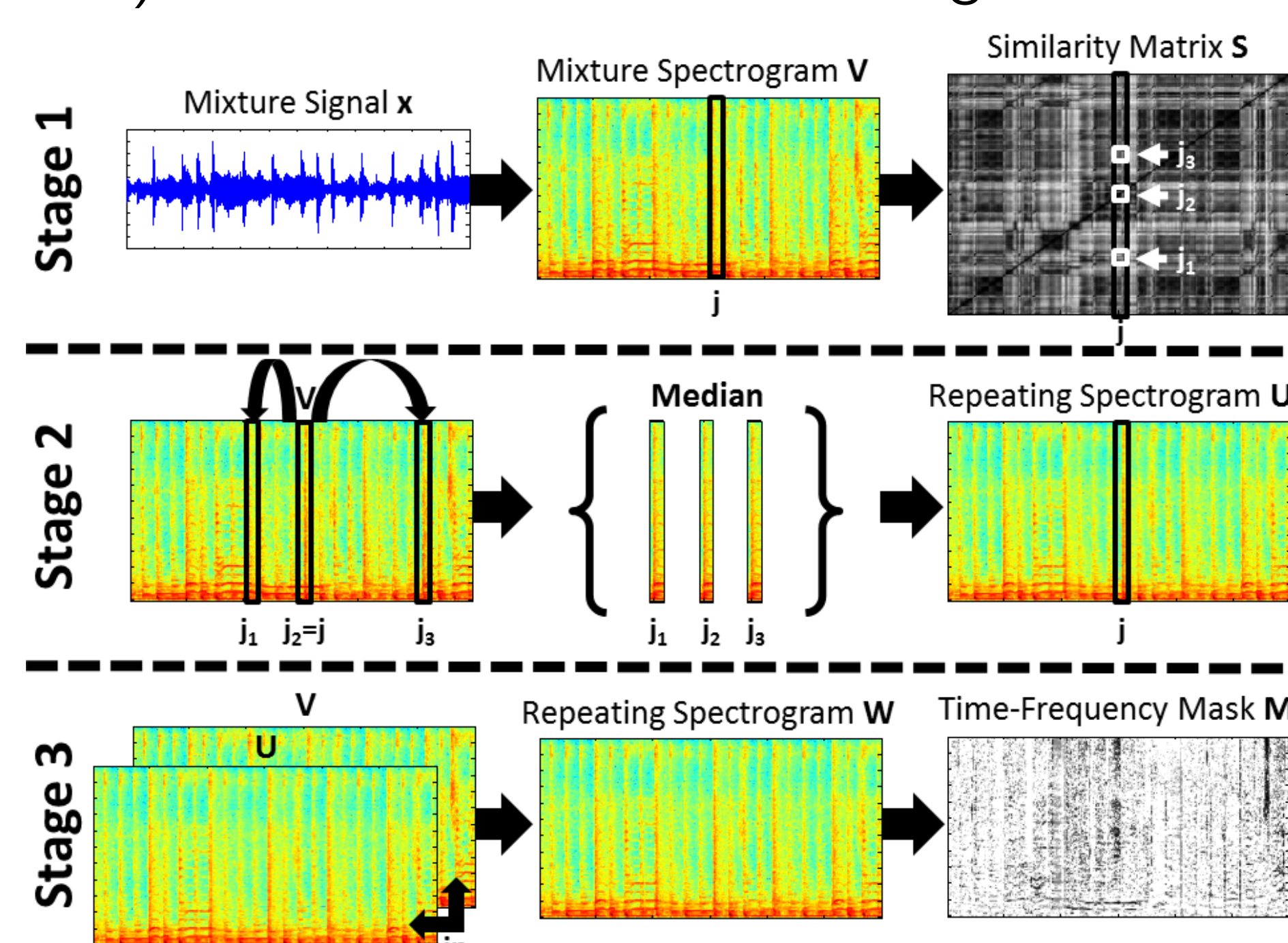


To capture the intuition that only a few singer models should be sufficient to explain an unseen singer, we assume that only a few of them in \mathbf{W} are active by adding block sparsity penalty on \mathbf{H} to the optimization formula.



Modeling the background music

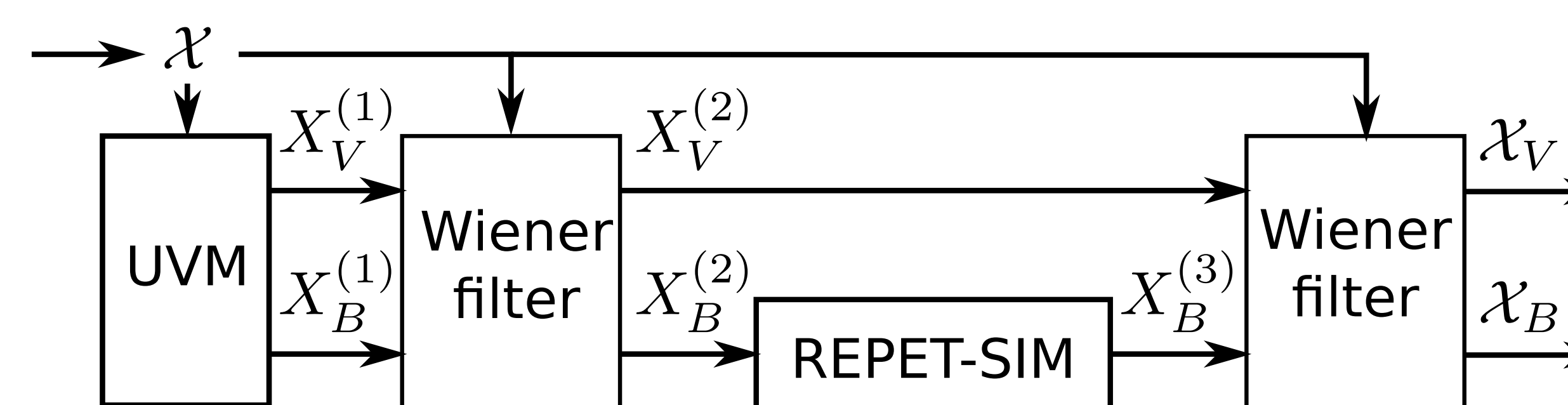
We use the REPET-SIM method, which focuses on extracting the background music. The separation is done by building a similarity matrix and smoothing the non-repeating patterns (assumed to be the voice) to create a mask for the background music.



Combined approach

The two presented methods are such that:

- The UVM approach explicitly models the singing voice.
- REPET-SIM explicitly models the background music.



We propose to combine those methods in order to leverage their ability to model each source in the mixture as follows:

- 1 Estimate $\mathbf{X}_V^{(1)} = \mathbf{W}_V \mathbf{H}_V$ and $\mathbf{X}_B^{(1)} = \mathbf{W}_B \mathbf{H}_B$ from the magnitude spectrogram \mathbf{X} .

- 2 Build new magnitude spectrogram estimates through Wiener filtering of the complex spectrogram \mathcal{X} :

$$\mathbf{X}_V^{(2)} = \left| \frac{\mathbf{X}_V^{(1)}}{\mathbf{X}_V^{(1)} + \mathbf{X}_B^{(1)}} \odot \mathcal{X} \right| \quad \mathbf{X}_B^{(2)} = \left| \frac{\mathbf{X}_B^{(1)}}{\mathbf{X}_V^{(1)} + \mathbf{X}_B^{(1)}} \odot \mathcal{X} \right|$$

- 3 Refine the music estimate using REPET-SIM as $\mathbf{X}_B^{(3)}$.

- 4 Build the final complex estimates through Wiener filtering:

$$\mathcal{X}_V = \frac{\mathbf{X}_V^{(2)}}{\mathbf{X}_V^{(2)} + \mathbf{X}_B^{(3)}} \odot \mathcal{X} \quad \mathcal{X}_B = \frac{\mathbf{X}_B^{(3)}}{\mathbf{X}_V^{(2)} + \mathbf{X}_B^{(3)}} \odot \mathcal{X}$$

Experimental setup

Dataset

We evaluate our algorithm on the 1,000 song clips (at 16kHz) of the MIR-1K dataset, from 19 individual singers. The performance of the algorithm is done using a leave-one-out cross-validation:

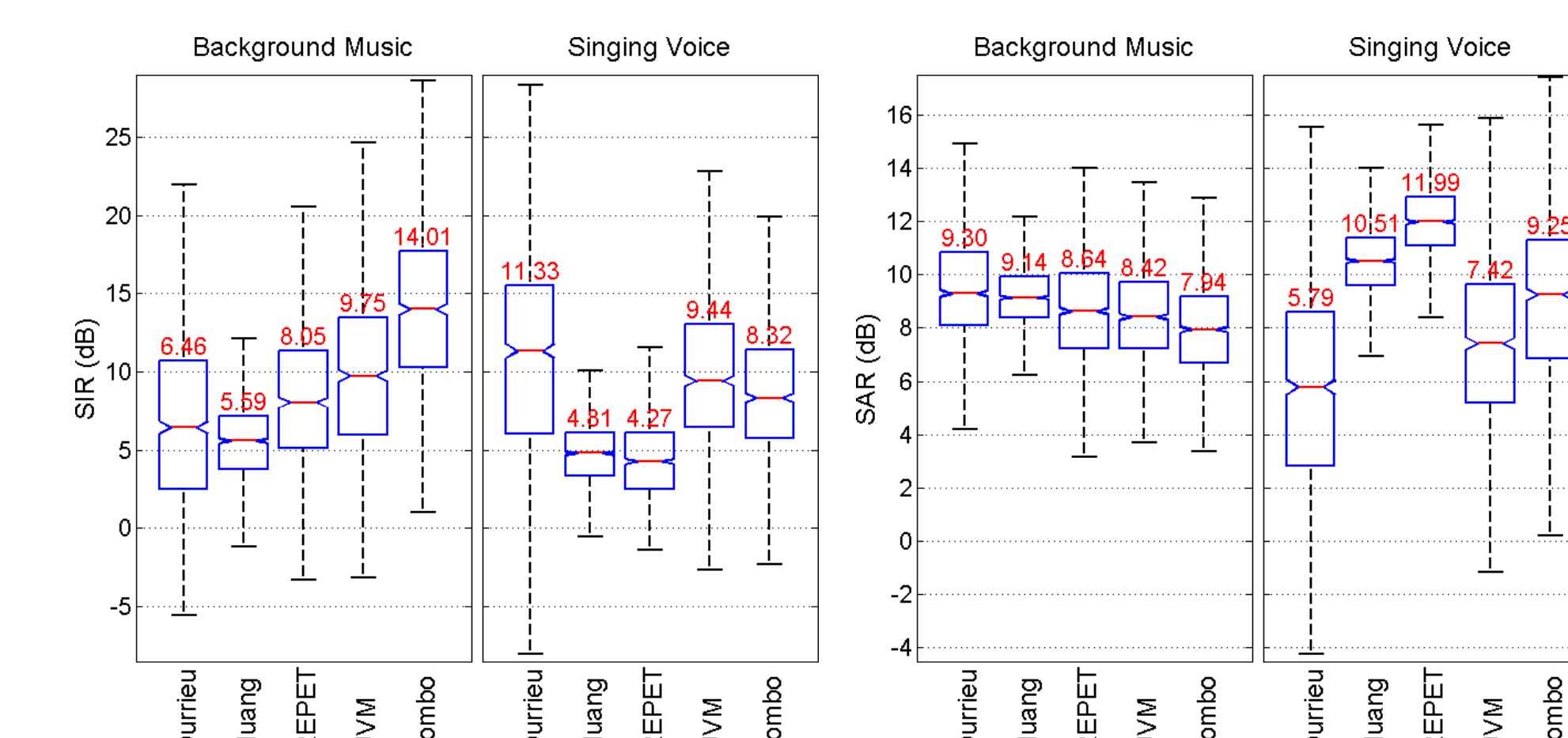
- The clips from one given singer are withheld as test set.
- The clips from the other singers are used to train a universal voice models with 18 dictionaries.

Competitor methods

- *Durrieu*: Source/filter voice model associated with an NMF-based background model learned iteratively from \mathbf{X} .
- *Huang*: Robust Principal Component Analysis assuming music as a low-rank component and voice as a sparse component.
- *REPET*: REPET-SIM alone.

Comparative results

To evaluate the performance, we use the 3 metrics from the BSS Eval toolbox measuring overall (SDR), interference (SIR) and artifact (SAR) related distortion.



Statistical analysis of the results shows that:

- The UVM approach achieves state-of-the-art results for music/voice separation in terms of Signal-to-Distortion Ratio.
- The combined approach further improves the Signal-to-Distortion Ratio scores.

Conclusion

We present here:

- The Universal Voice Models, a new NMF approach for music/voice separation modeling the voice without the need for specific training samples.
- A combined approach using the UVM and the REPET-SIM methods in order to leverage their complementary properties.
- Experimental results demonstrate that the two proposed methods achieve state-of-the-art results for overall music/voice separation performance.

This work was supported in part by NSF grant number IIS-0812314.