

Speech-guided source separation using a pitch-adaptive guide signal model

Romain HENNEQUIN, Juan-José BURRED, Simon MALLER, Pierre LEVEAU

Audionamix, 171 quai de Valmy, 75010 Paris, France

Introduction

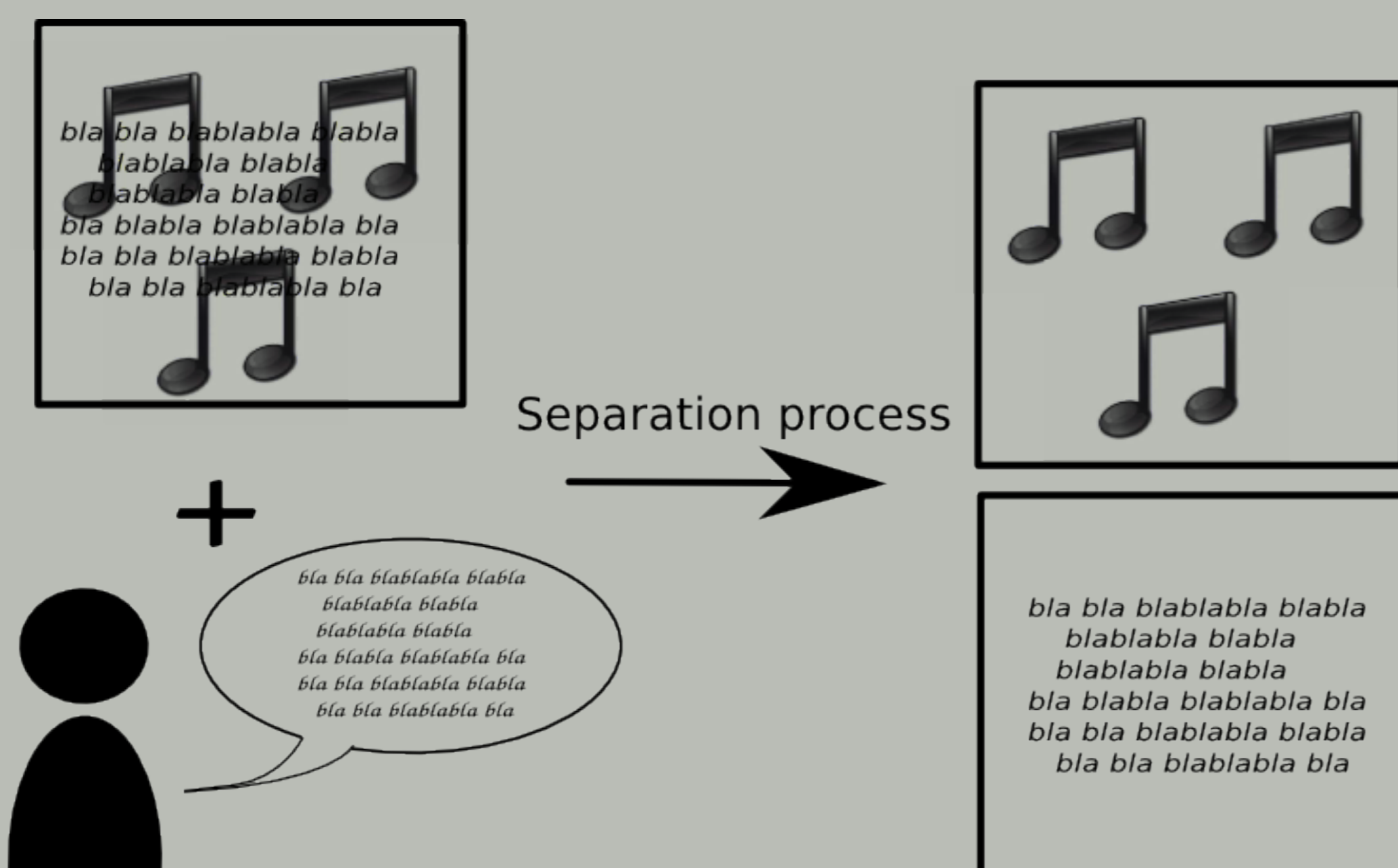
New method to perform underdetermined audio source separation using a spoken or sung reference signal to inform the separation process.

- ▶ Explicitly models differences between reference and target signal.
- ▶ Outperforms state-of-the-art methods.

Typical scenario

Isolate a target sound, imitated by the user, from a background sound.

- ▶ Isolation/removal of dialogs in a movie soundtrack.
- ▶ Isolation/removal of lead vocal in a song.



Issues: differences between guide and target signal:

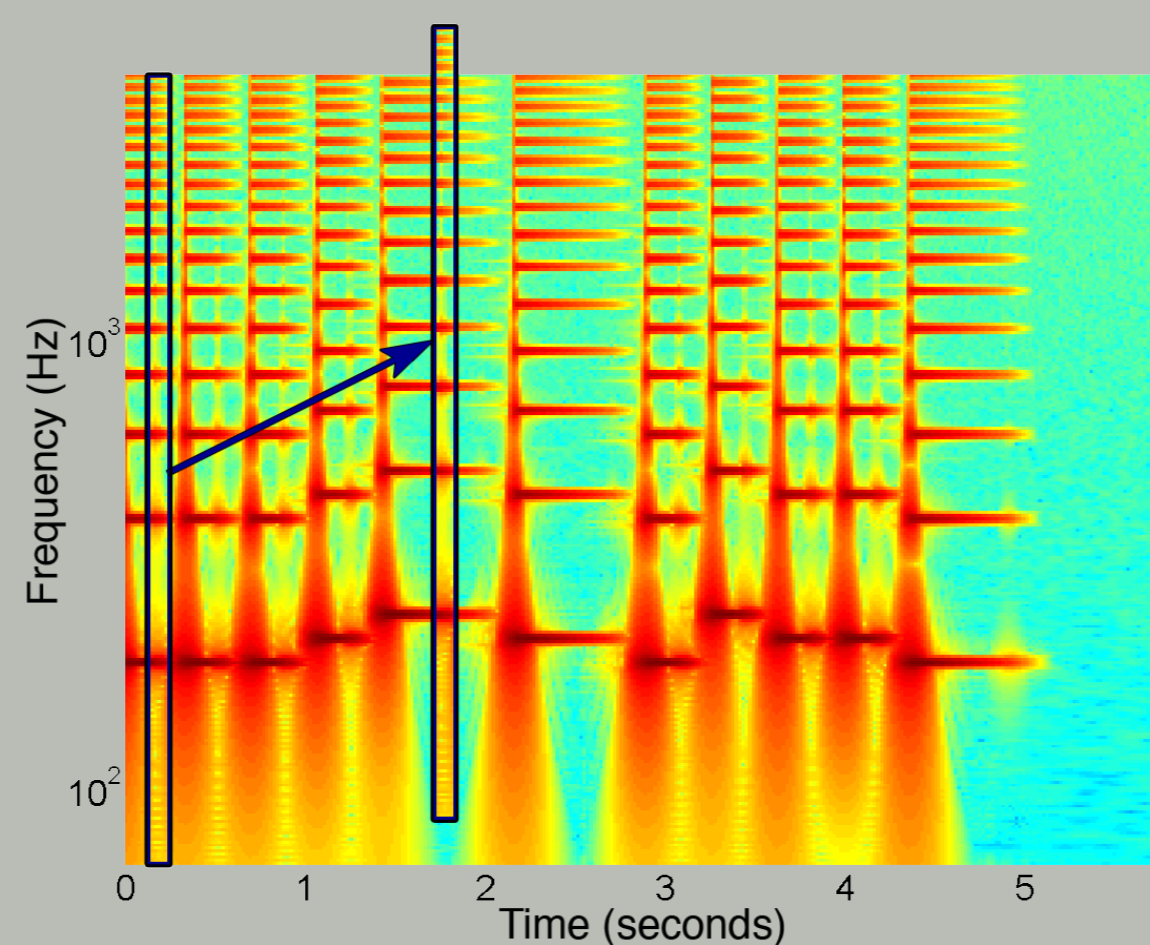
- ▶ pitch differences: absolute pitch, intonation, vibrato...
- ▶ time differences: misalignments.
- ▶ equalization/timbre differences.

Model

Input spectrogram

V : power log-frequency spectrogram (Constant-Q transform) of the mixture signal.

Shift invariance property: a pitch modification can be modeled by a vertical shift.



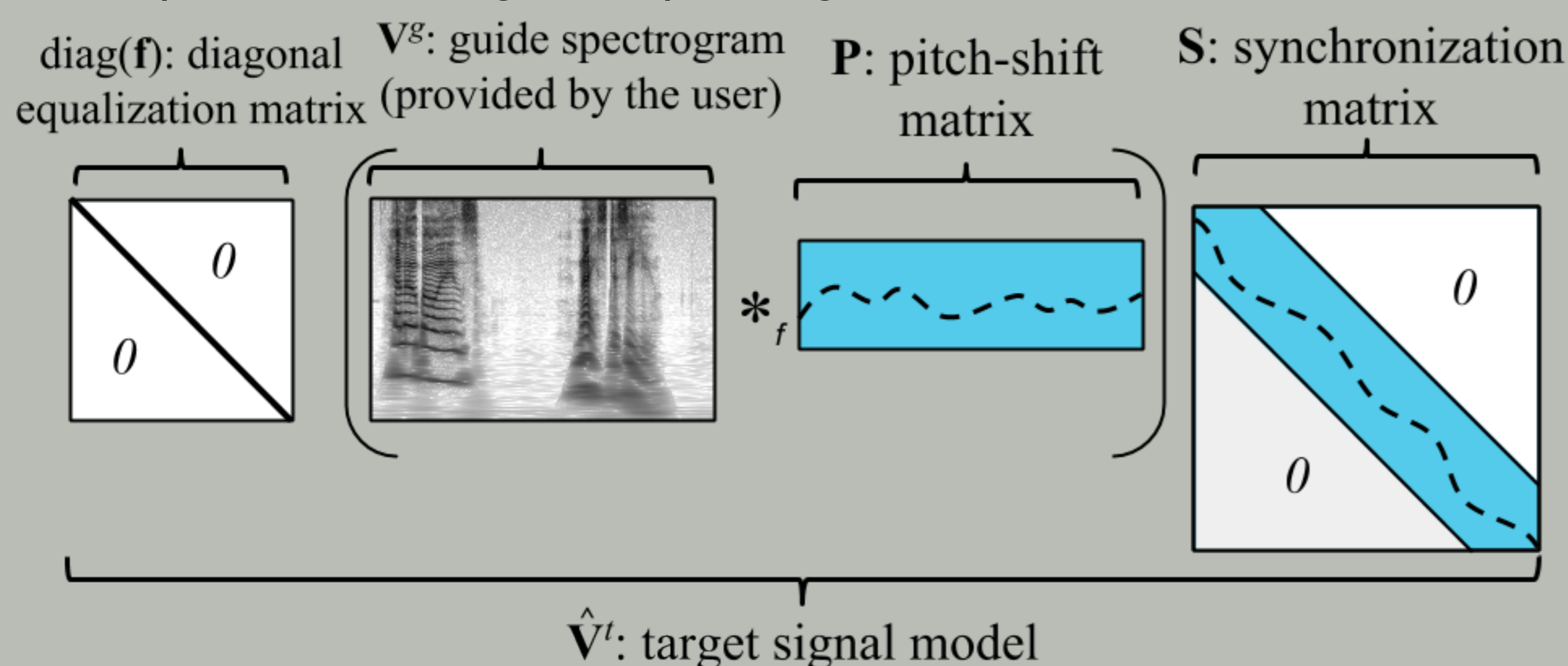
Non-negative spectrogram model

$$V \approx \hat{V} = \hat{V}^t + \hat{V}^b$$

\hat{V}^t : target signal model.
 \hat{V}^b : background signal model.

Target signal model:

\hat{V}^t is adapted from the guide spectrogram V^g :



Pitch shift operator P :

$$V_{\text{shifted}}^g = \left(\sum_{\phi} \downarrow_{\phi} V^g \text{diag}(P_{\phi,:}) \right)$$

Global adaptation filter f and synchronization matrix S :

$$\hat{V}^t = \text{diag}(f) \left(\sum_{\phi} \downarrow_{\phi} V^g \text{diag}(P_{\phi,:}) \right) S$$

Background signal model

Standard Non negative Matrix Factorization (generic model):

$$\hat{V}^b = WH$$

Algorithm

First estimation of parameters

Minimization of Itakura-Saito divergence:

$$\mathcal{C}(P, f, S, W, H) = D_{IS}(V | \hat{V}^t + \hat{V}^b), \text{ (s.t. } P, f, S, W, H \geq 0)$$

using multiplicative update rules.

Tracking

- ▶ Tracking of the best pitch shift in P using Viterbi algorithm.
 - ▶ Tracking of the best synchronization in S using DTW.
- Thresholding of P and S then parameters reestimation.

Separation

Separation is achieved using Wiener Filtering and CQT inversion.

Experimental results

Experimental setup

Dialog isolation/removal in movie soundtracks.

Compared with other separation algorithms:

- ▶ Main melody extraction method proposed in [1].
- ▶ PLCA-based speech-informed separation algorithm [2].
- ▶ First one but informed with instantaneous pitch of the target signal [3].
- ▶ CQT Wiener oracle as an upper performance bound.

Database

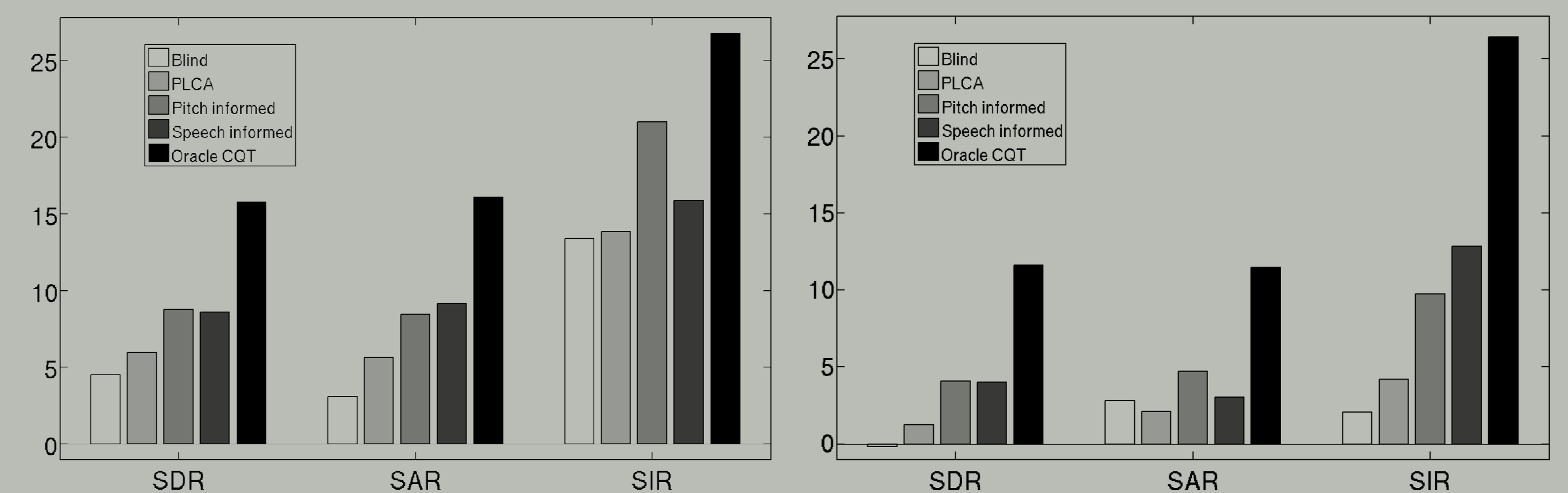
Mix signals:

- ▶ 10 excerpts from synthetic movie soundtracks.
- ▶ 5 different movies in english.
- ▶ Excerpt = dialog only part + music and effects only part mixed down to mono.

Dubs:

- ▶ using the mix signal as a reference.
- ▶ Done by the same male native english speaker.
- ▶ Same dubs for both speech informed algorithms.

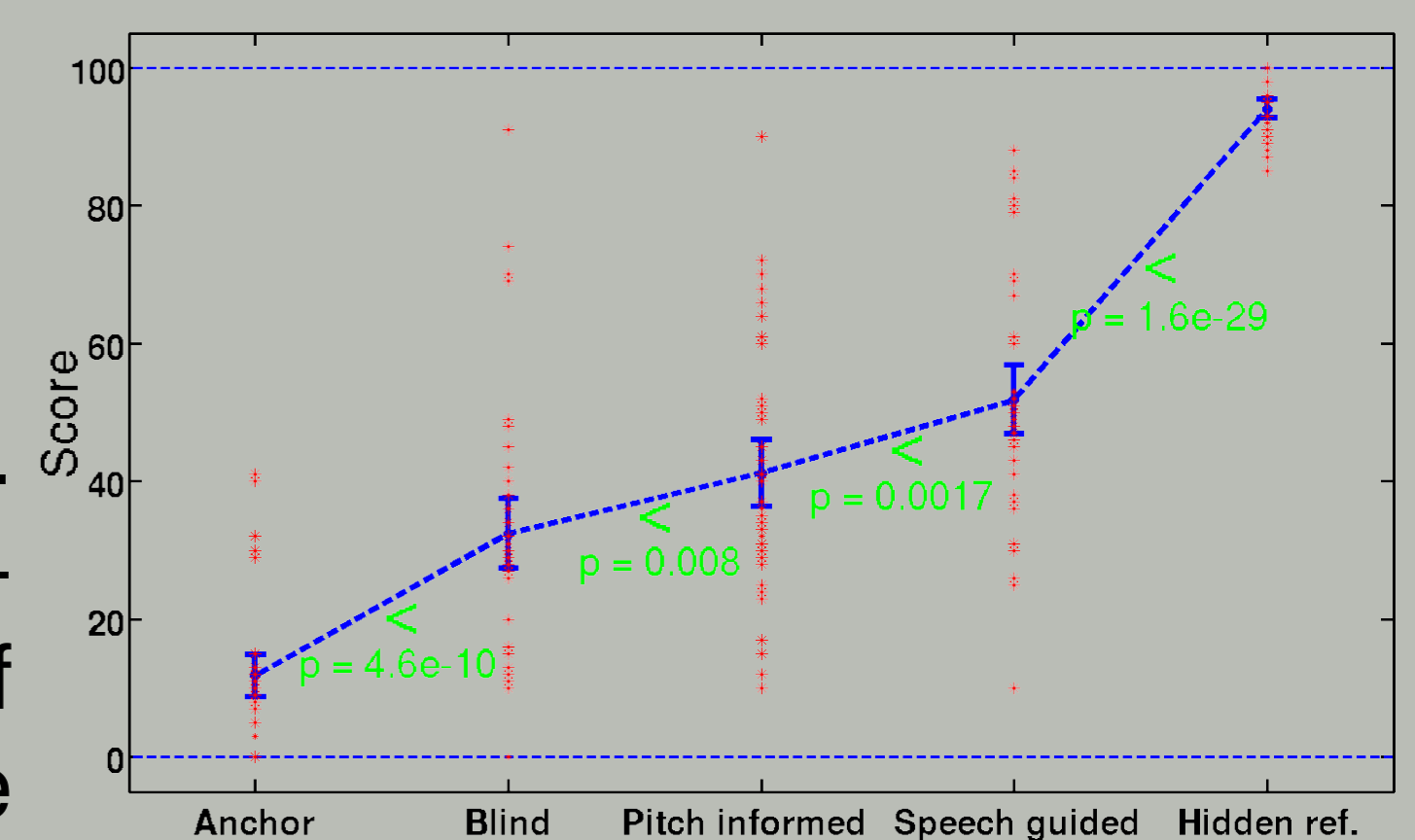
Objective results



Listening tests

Internal blind listening test:

- ▶ Dialog isolation only.
 - ▶ 5 (few!) participants.
 - ▶ MUSHRA protocol.
 - ▶ Rating according to the "usability".
- Results of our algorithm globally preferred over the results of the pitch-informed algorithm (to be taken with care)



Conclusion

- ▶ New method to perform source separation providing a spoken guide.
- ▶ Outperforms a state-of-the-art one and methods performing same task.
- ▶ Future work:
 - ▶ No voice model on the guide signal: any kind of signal can be used. Many other applications.
 - ▶ Other kinds of adaptation (formant adaptation...).
 - ▶ Speeding up the algorithm.

References

- ▶ J.-L. D. et al., "An iterative approach to monaural musical mixture de-soloing," in *ICASSP*, 2009.
- ▶ P. Smaragdis and G. J. Mysore, "Separation by humming: User-guided sound extraction from monophonic mixtures," in *WASPAA*, 2009.
- ▶ J.-L. Durrieu and J.-P. Thiran, "Musical audio source separation based on user-selected f0 track," in *LVA/ICA*, 2012.