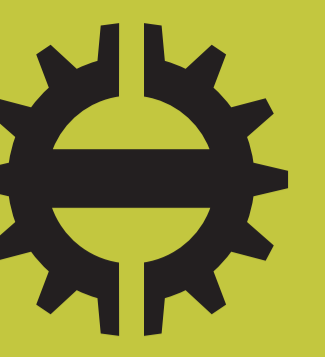


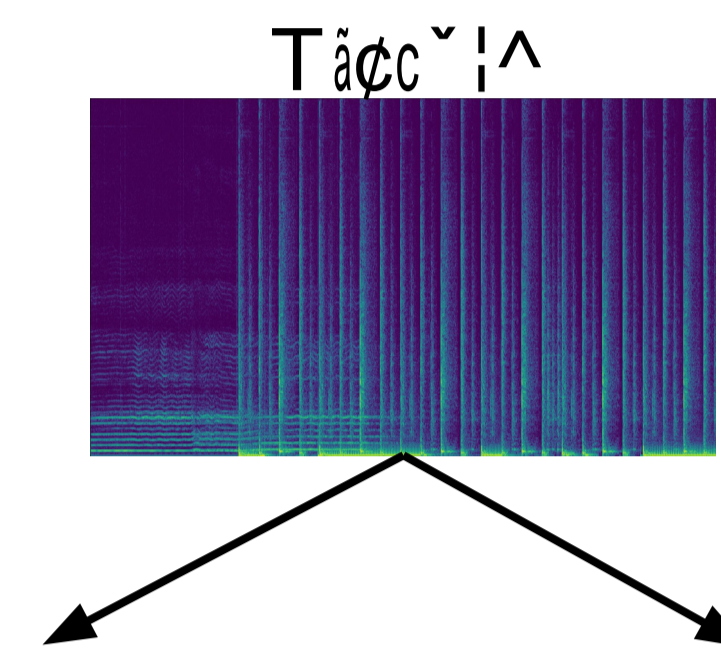
Harmonic-Percussive Source Separation with Deep Neural Networks and Phase Recovery



Konstantinos Drossos¹, Paul Magron¹, Stylianos Ioannis Mimitakis², Tuomas Virtanen¹

¹Laboratory of Signal Processing, Tampere University of Technology, Finland, ²Fraunhofer IDMT, Ilmenau, Germany

Harmonic/Percussive Source Separation (HPSS)



- Separate percussive (e.g. drum, percussion) from harmonic (e.g. guitar, piano, singing voice) components.

$$P_{\hat{m}} \{ \hat{m} \} \quad \hat{U} \wedge \hat{m} \bullet \bullet \hat{m} \wedge$$

- Applications: rhythm analysis, augmented mixing, time-stretching, etc.

Contributions

- We propose a novel HPSS method, based on two components.
 - A recently proposed deep neural network (DNN) method for monaural music source separation [1].
 - A recently introduced algorithm for phase recovery [2]
- Reproducible research** / Source code available, results on freely available dataset.

Proposed method

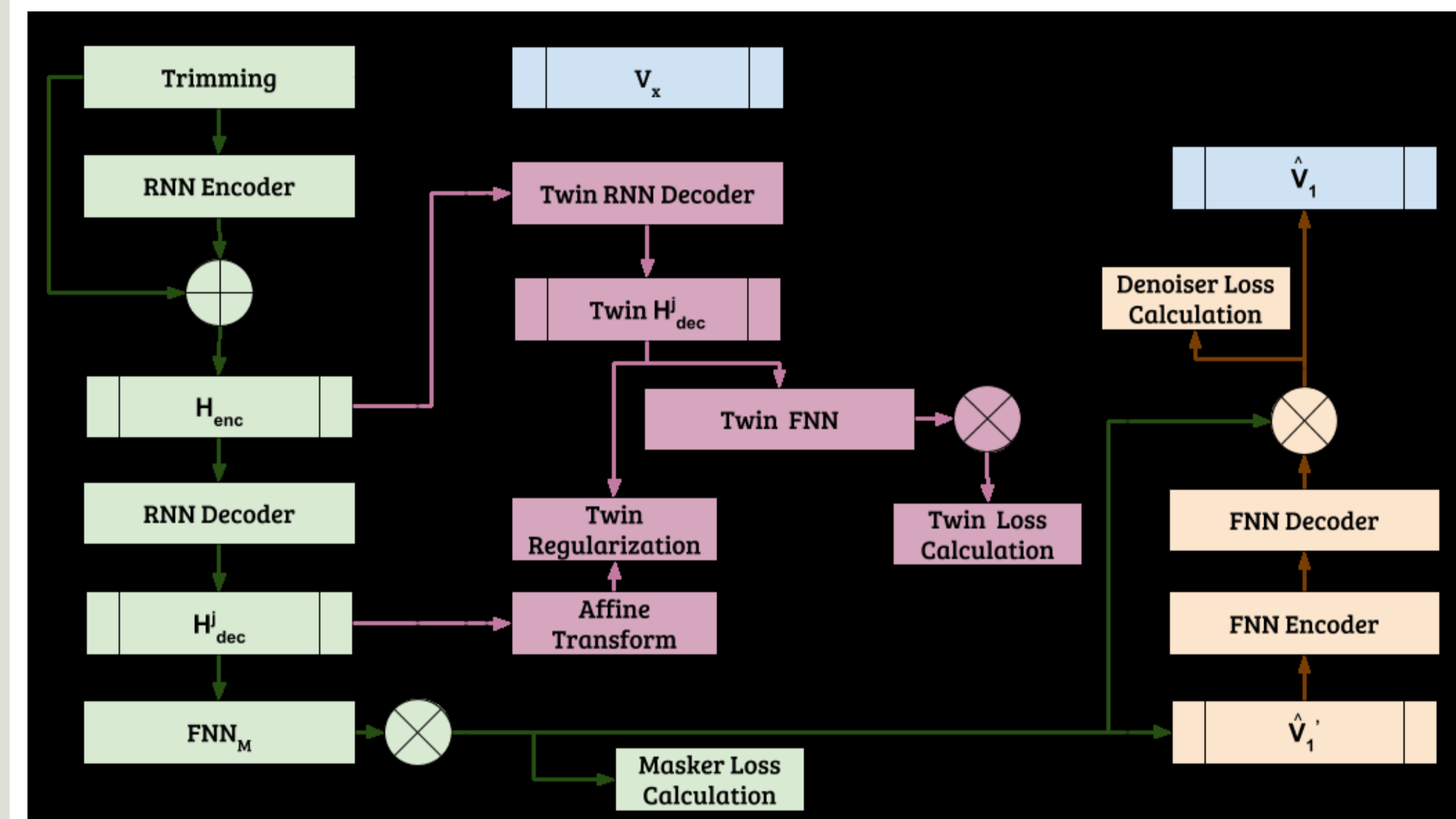
A two-stage approach based on DNNs and phase recovery

- A DNN for estimating the percussive spectrogram [1].
 - Input* / the magnitude spectrogram of the mixture.
 - Output* / the magnitude spectrogram of the percussive component.
 - We estimate harmonic components by spectral subtraction.
- Time-domain signal reconstruction, using either:
 - The phase of the mixture, or
 - An iterative algorithm for improved phase recovery [2].

Magnitude estimation: MaD TwinNet

A two-step monaural source separation system [1].

- Based on denoizing auto encoders framework (DAEs).
- First applies a time-frequency mask, then a time-frequency denoizing filter.
- Takes into account long temporal dependencies through TwinNet regularization.



Phase recovery: PU-iter

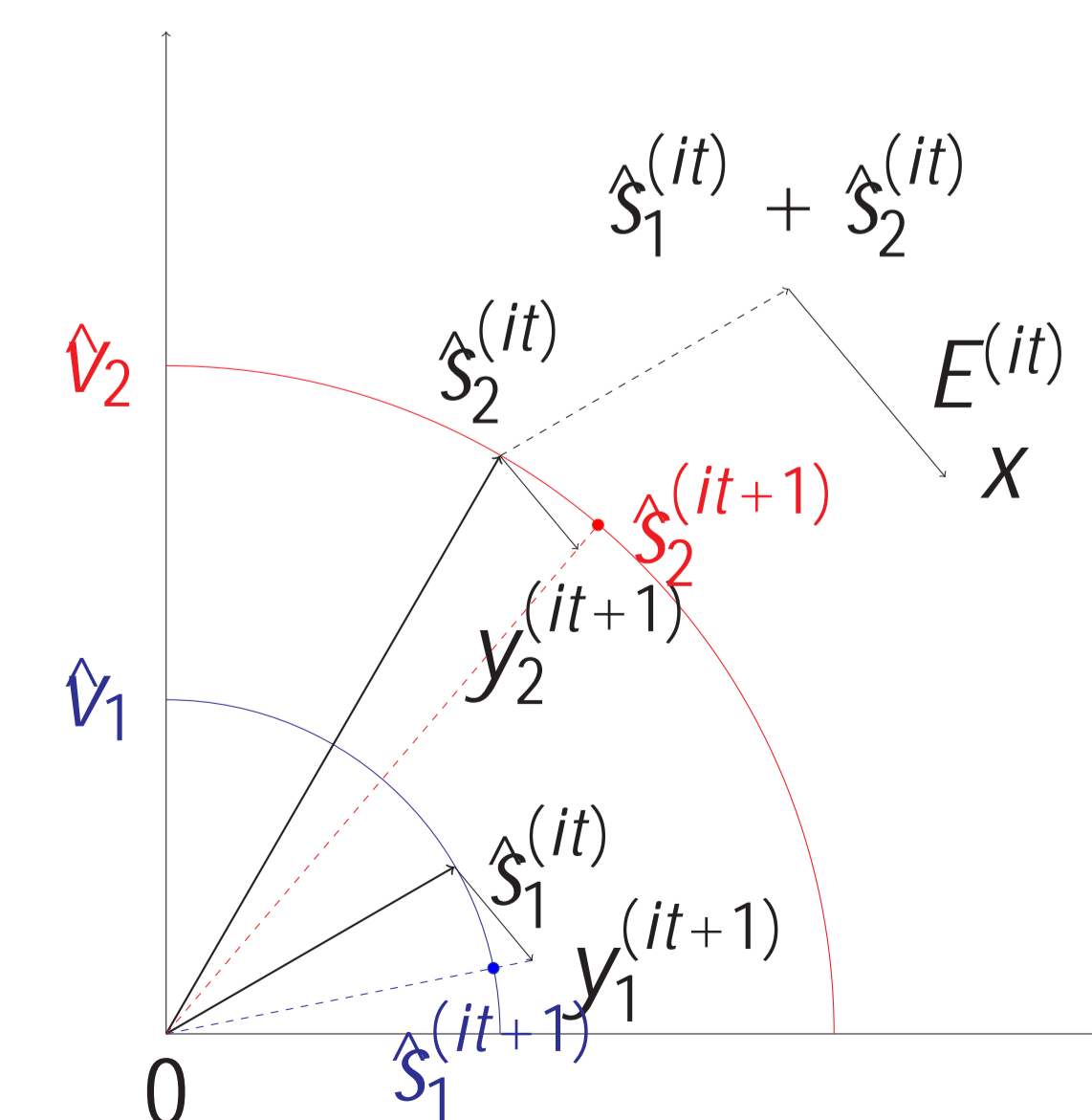
Sinusoidal phase

- The harmonic source is modeled as a sum of sinusoids.
- Explicit phase relationship between successive time frames:

$$\text{harmo}_{f;t} = \text{harmo}_{f;t-1} + 2 \cdot \angle f;t$$

Iterative procedure [2]

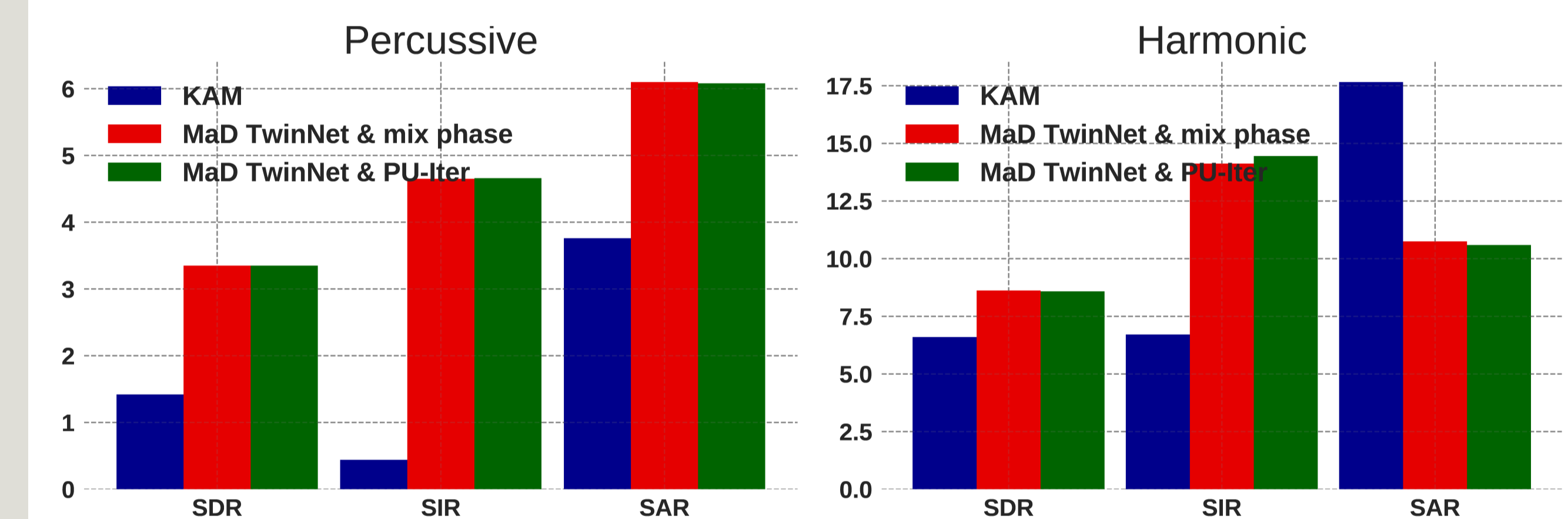
- Minimizes the mixing error;
- Initialized with the mixture's phase (percussive part) or sinusoidal phase (harmonic part):
 - Does not modify the target magnitudes (= MaD TwinNet estimates).



Training & Evaluation

- Demixing secret dataset 100 (DSD100) / 100 audio mixtures and their isolated sources.
- Two different STFT settings:
 - One in favor of MaD TwinNet (worked better).
 - One in favor of the phase recovery algorithm.
- Compared against Kernel Additive Model (KAM) [3].
- Separation quality measured with the signal to: artifacts ratio (SAR), interference ratio (SIR), distortion ratio (SDR).

Objective results



Conclusions & future work

- Supervised HPSS based on deep learning and phase recovery.
- MaD TwinNet and phase recovery improves over KAM.
- Future work
 - Joint magnitude/phase recovery.
 - Phase recovery based on deep learning.

References

- K. Drossos, S.I. Mimitakis, D. Serdyuk, G. Schuller, T. Virtanen, Y. Bengio, "MaD TwinNet: Masker-Denoiser Architecture with Twin Networks for Monaural Sound Source Separation", in Proc. of the IEEE IJCNN, 2018.
- P. Magron, R. Badeau and B. David, "Model-based STFT phase recovery for audio source separation", in the IEEE Trans. on Audio, Speech, and Language Processing, June 2018.
- A. Liutkus, D. Fitzgerald, Z. Rafii, B. Pardo, and L. Daudet, "Kernel Additive Models for Source Separation", in the IEEE Trans. on Signal Processing, Aug. 2014.