VERY LOW BITRATE SPATIAL AUDIO CODING WITH DIMENSIONALITY REDUCTION

gipsa-lab

Christian Rohlfing¹, Jeremy E.Cohen², Antoine Liutkus³

¹Institut für Nachrichtentechnik, RWTH Aachen University ²Department of Image and Signal-processing, Gipsa-lab, CNRS, Grenoble, France ³Inria, speech processing team, Villers-lès-Nancy, France

Introduction Reference ISS method [1] consists of two steps: mixture mixture • Coder: Sources and mix perfectly known. Nonnegative Tensor Factorization (NTF) of source spectrograms sources **D** $\overline{\mathbf{O}}$



with NTF. Encoding of resulting parameters Θ .

• Decoder: Only mix available. Wiener filtering used to recover sources from mix assisted by transmitted NTF parameters Θ .

Dimension reduction for ISS



Restriction of classical NTF-based ISS methods:

Temporal dimension of source spectrogram can be very large (e.g. for full-length tracks).

 \Rightarrow Growth of temporal dimension of NTF parameters (bitrate) and

Higher-order SVD at coder



Randomized SVD reduces memory costs and numerical complexity: $\mathcal{O}(FT(R_F + R_T))$ (vs. classical NTF with R components: $\mathcal{O}(FTJR)$)

Transmitted parameters do *not* include U_T but only U_F , U_J and \mathcal{G} .

Parameter estimation at decoder

Model of mixture $\boldsymbol{\mathcal{V}}_x \approx (\boldsymbol{U}_F \otimes \boldsymbol{U}_T \otimes \boldsymbol{u}_J) \boldsymbol{\mathcal{G}}$ with $\boldsymbol{u}_J = \sum_{j=1}^J [\boldsymbol{U}_J]_{j:j}$

computational complexity of NTF.

Proposed solution:

- Replace NTF with Higher-order SVD in coder.
- Do not send temporal parameter but estimate it at decoder.
- \Rightarrow Size of side-information independent of track length!

Temporal parameter estimation: Temporal matrix U_T is not transmitted. Estimate U_T given the other transmitted parameters Θ .

$$\widehat{\boldsymbol{U}}_T \leftarrow \underset{\boldsymbol{U}_T^{\mathsf{T}}\boldsymbol{U}_T = \boldsymbol{I}}{\operatorname{argmin}} \|\boldsymbol{\mathcal{V}}_x - \boldsymbol{N}_T\boldsymbol{U}_T^{\mathsf{T}}\|_{\mathrm{F}}^2$$

with $N_T = U_F (u_J \bullet_3 \mathcal{G})$

Iterative re-estimation of parameters U_F and \mathcal{G} initialized with their quantized value obtained from coder.

Experiments



Reference methods:

Performance of NTF method [1] as well as plain HOSVD without quantization and with full transmission of all HOSVD parameters.

Setup:

- Ten full-length tracks (4 sources, taken from DSD100 database)
- HOSVD operates with $R_F, R_T \in [5, 100]$ and $R_J = 4$
- Quantization disabled (q = 0) or enabled (q = 1) with $N_{\mathcal{G}}, N_F \in$ [5, 1000] centroids for quantization of \mathcal{G} and U_F
- Number of iterations for re-estimation $N_{\rm it} \in \{0, 10\}$

• Quality measure: δ_{SDR} over parameter bitrate r and reconstruction scores $\operatorname{RS}_{\widehat{\boldsymbol{\mathcal{V}}}_x} = 10 \log(\|\boldsymbol{\mathcal{V}}_x\|_F^2 / \|\boldsymbol{\mathcal{V}}_x - \widehat{\boldsymbol{\mathcal{V}}}_x\|_F^2)$ and $\operatorname{RS}_{\widehat{\boldsymbol{\mathcal{V}}}_x}$ over iterations

Results:

- Not sending U_T but estimating it at decoder, with other parameters full resolution (—), permits to reach NTF performance (—)
- Coarse quantization of U_F and \mathcal{G} and not sending U_T (—) leads to remarkable decrease of bitrate by a factor of almost 10.
- Iterative re-estimation of quantized parameters (---, ---) is not increasing performance. Score $\mathrm{RS}_{\hat{m{\mathcal{V}}}_x}$ (---) is improved by $0.5\,\mathrm{dB}$. Generalization to $\hat{m{\mathcal{V}}}_s$ does not hold.

[1] A. Liutkus, J. Pinel, R. Badeau, L. Girin, and G. Richard, "Informed source separation through spectrogram coding and data embedding," Signal Processing, vol. 92, no. 8, pp. 1937 – 1949, 2012.



rohlfing@ient.rwth-aachen.de

www.ient.rwth-aachen.de

Institut für Nachrichtentechnik, Melatener Str. 23, 52074 Aachen

ICASSP 2017