

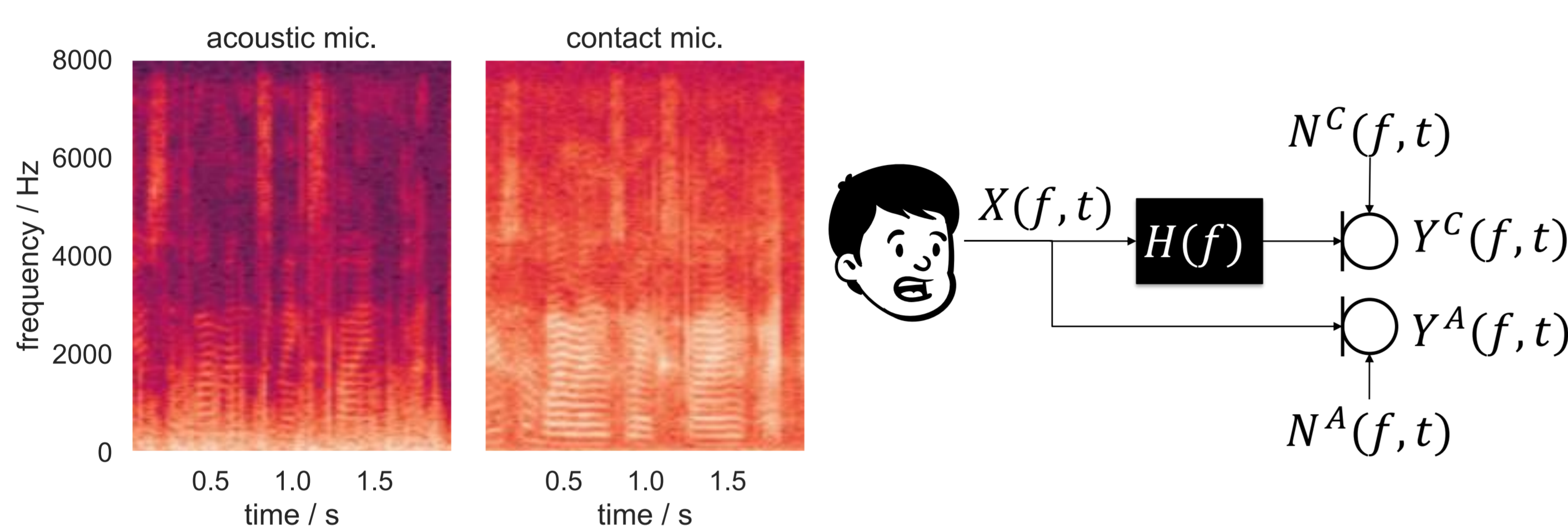
PROBLEM STATEMENT

- wind noise may severely reduce speech intelligibility in mobile communication scenarios
- performance of speech enhancement algorithms tends to degrade in extremely challenging conditions
- auxiliary information** from sensors such as **contact microphones** can help improve performance, but labeled data is scarce

this poster: fuse advantages of acoustic and contact microphones for wind noise reduction

SIGNAL MODEL

- two microphones capturing speech in wind noise
 - acoustic mic. (Y^A): **high wind noise level**, **undistorted speech**
 - contact mic. (Y^C): **low wind noise level**, **distorted speech**, **additional noise**

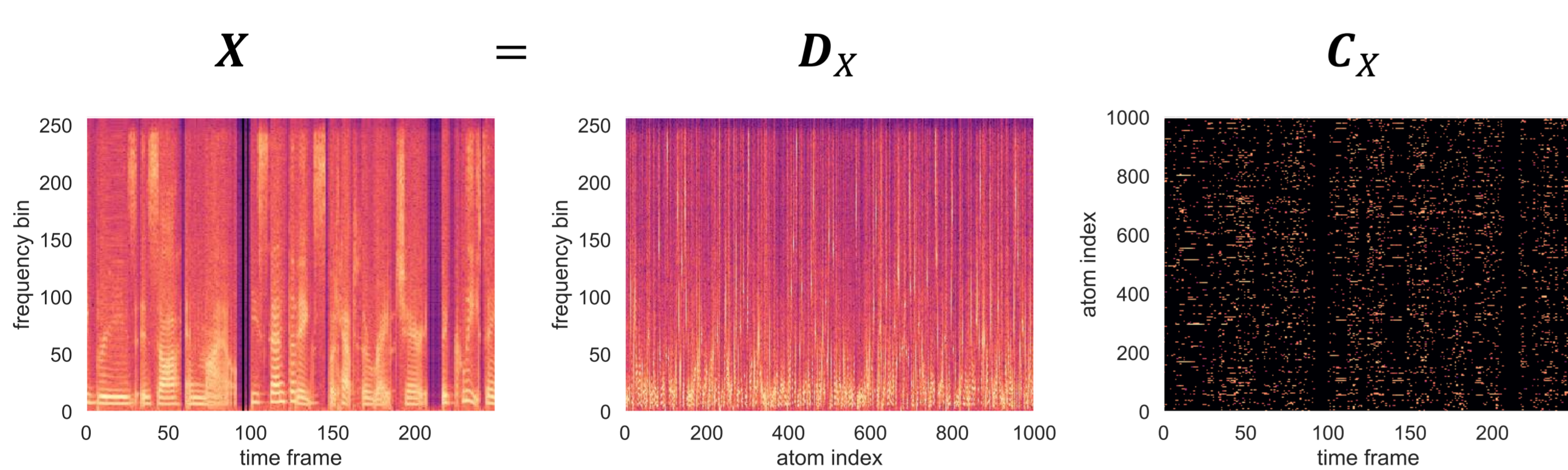


- signal model assumptions:
 - speech components related via **(linear) transfer function** $H(f)$
 - separate additive noise** components

$$\begin{bmatrix} Y^C \\ Y^A \end{bmatrix} = \begin{bmatrix} HX + N^C \\ X + N^A \end{bmatrix} \in \mathbb{C}^{2F \times T}$$

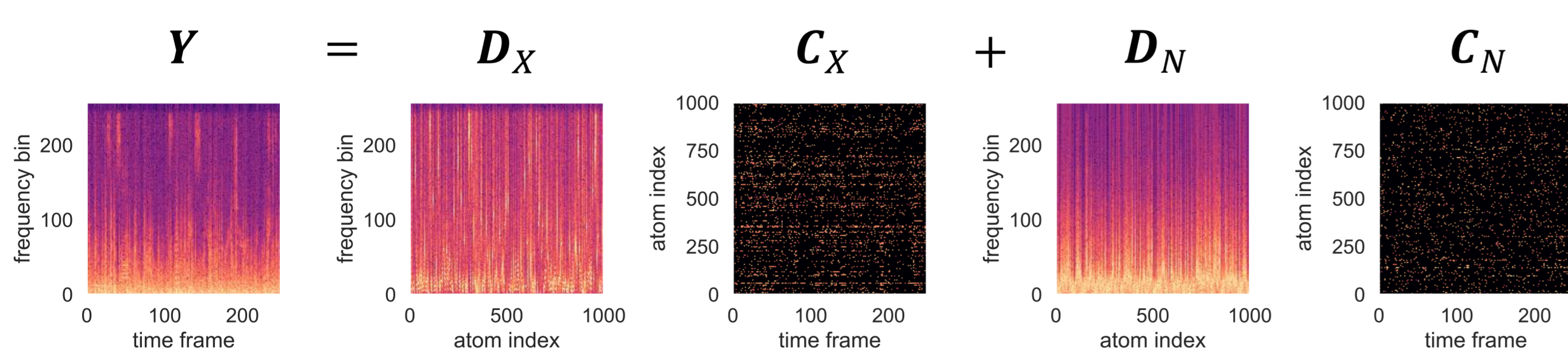
DICTIONARY-BASED SIGNAL REPRESENTATION

- speech contains spectro-temporal structures** (e.g., harmonics)
- these structures can be learned from data ("**dictionary atoms**") [1]
 - represent speech signal X using (complex-valued) learned dictionary D_X and sparse code C_X :



keeping number of activated atoms low yields noise reduction

- extension towards noisy signals:**
 - learn both speech and noise dictionaries D_X and D_N
 - speech can be estimated using speech dictionary D_X and sparse code C_X
 - signal components will be captured according to coherence with dictionaries



- dictionaries can be trained by solving non-convex optimization problem:**

$$D_X, C_X = \underset{D, C}{\operatorname{argmin}} \|X - DC\|_2^2 + \lambda \|C\|$$

- λ trades off more accurate vs. sparser representation
- optimization alternates between dictionary update and sparse coding step

DICTIONARY-BASED MICROPHONE FUSION

- integrate signal model for acoustic and contact microphone signals into dictionary training and sparse coding:**
 - model noise components using separate dictionaries D_N^C, D_N^A
 - model speech components using single dictionary D_X and (unknown) transfer function H

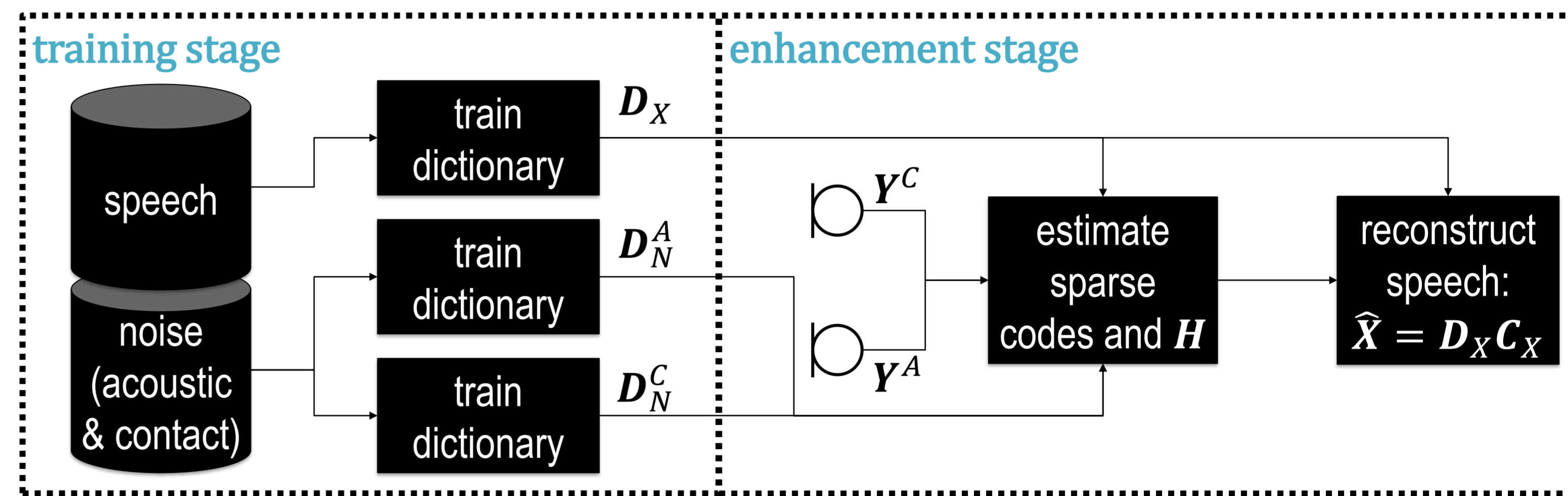
$$\begin{bmatrix} Y^C \\ Y^A \end{bmatrix} = \begin{bmatrix} HD_X C_X + D_N^C C_N^C \\ D_X C_X + D_N^A C_N^A \end{bmatrix}$$

- cost function:**

$$\eta \|Y^C - (HD_X C_X + D_N^C C_N^C)\|_2^2 + (1 - \eta) \|Y^A - (D_X C_X + D_N^A C_N^A)\|_2^2 + \lambda \|\tilde{C}\|$$

- η controls influence of each microphone
- alternating optimization:**

- (iterative) sparse coding algorithm to estimate $\tilde{C} = [C_X^T \ C_N^{C,T} \ C_N^{A,T}]^T$
- closed-form solution to estimate H



SIMULATIONS

Dataset

	training	evaluation
clean speech	LibriTTS (dev-clean subset) [2]	
noise	airflow-induced noise: 3 m/s and 5 m/s (2 min each) recorded by acoustic and contact microphones	
SNR	no mixed signals required	{-10, -5, 0, 5} dB

Compared Algorithms

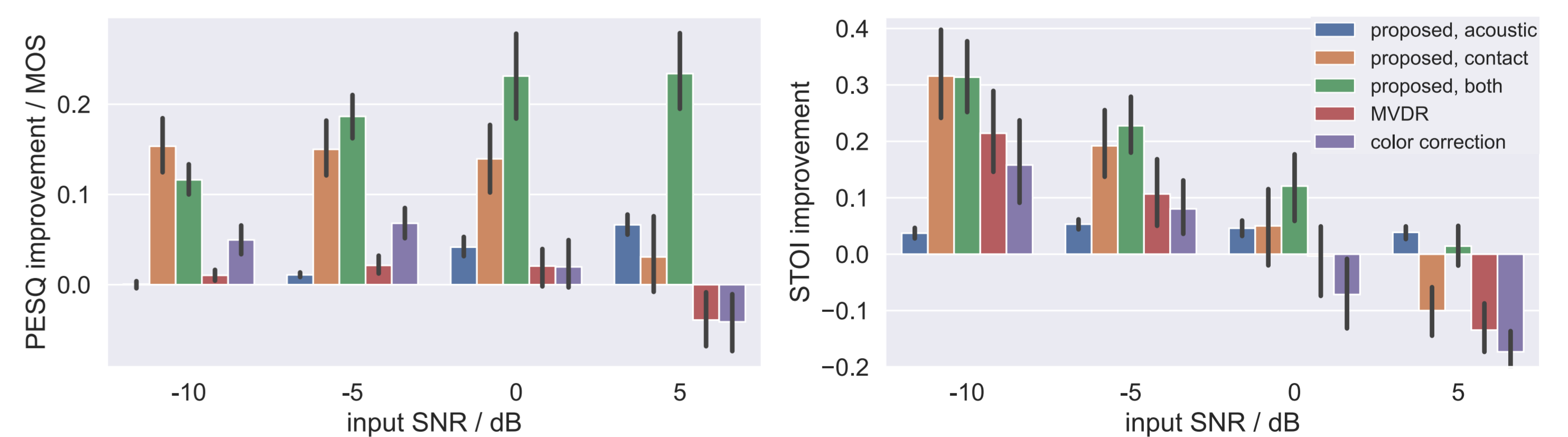
- proposed:** only acoustic microphone, only contact microphone, both microphones (fusion)
- baselines:** MVDR beamformer, color correction (inversion of H)

Settings

- sampling frequency 16 kHz, frame length 32 ms, 50% overlap, $\sqrt{\text{Hann}}$ window
- proposed dictionary-based algorithms:
 - sparse coding:** accelerated proximal gradient method [3] (1000 iterations)
 - dictionary update:** method of optimal directions [4] (25 iterations)
 - dictionary size: 1000 atoms each (3000 atoms in total)
 - $\eta = 0.6$, placing **more weight on contact microphone**
 - 5 iterations alternating between sparse coding and estimating transfer function H
- (initial) transfer function estimation: **covariance whitening method**

Results

- performance measures:**
 - perceptual evaluation of speech quality (**PESQ**)
 - short-time objective intelligibility (**STOI**)



fusing both microphones yields highest improvements in most conditions

- low SNR: **contact mic.** > **acoustic mic.**, high SNR: **contact mic.** < **acoustic mic.**
- benefit of **jointly estimating speech and noise components** as well as transfer function vs. only estimating transfer function

OUTLOOK

- comparison with deep neural network-based approaches** designed for scarce labeled data (e.g., self-supervised pre-training)

REFERENCES

- C. D. Sigg, T. Dikk, and J. M. Buhmann, "Speech Enhancement Using Generative Dictionary Learning," IEEE Trans. Audio Speech Language Processing, vol. 20, no. 6, pp. 1698–1712, Aug. 2012.
- H. Zen, V. Dang, R. Clark, Y. Zhang, R. J. Weiss, Y. Jia, Z. Chen, Y. Wu, "LibriTTS: A Corpus Derived from LibriSpeech for Text-to-Speech," in Proc. Interspeech, Graz, Austria, Sep. 2019, pp. 1526–1530.
- N. Parikh, S. Boyd, "Proximal algorithms," Foundations and trends in Optimization, vol. 1, no. 3, pp. 127–239, 2014.
- K. Engan, S. O. Aase, and J. H. Husøy, "Multi-frame compression: theory and design," Signal Processing, vol. 80, no. 10, pp. 2121–2140, Oct. 2000.