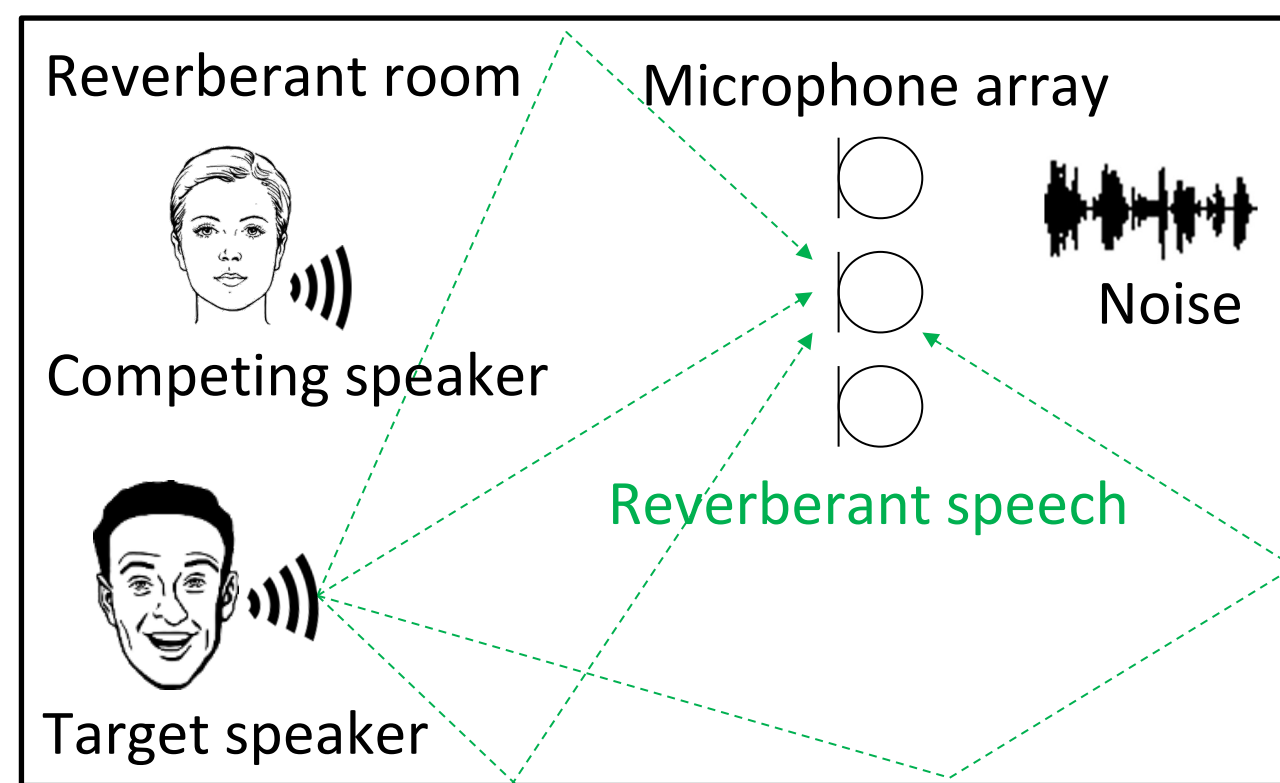


Zhong-Qiu Wang and Shinji Watanabe
Carnegie Mellon University

Unsupervised separation: motivations

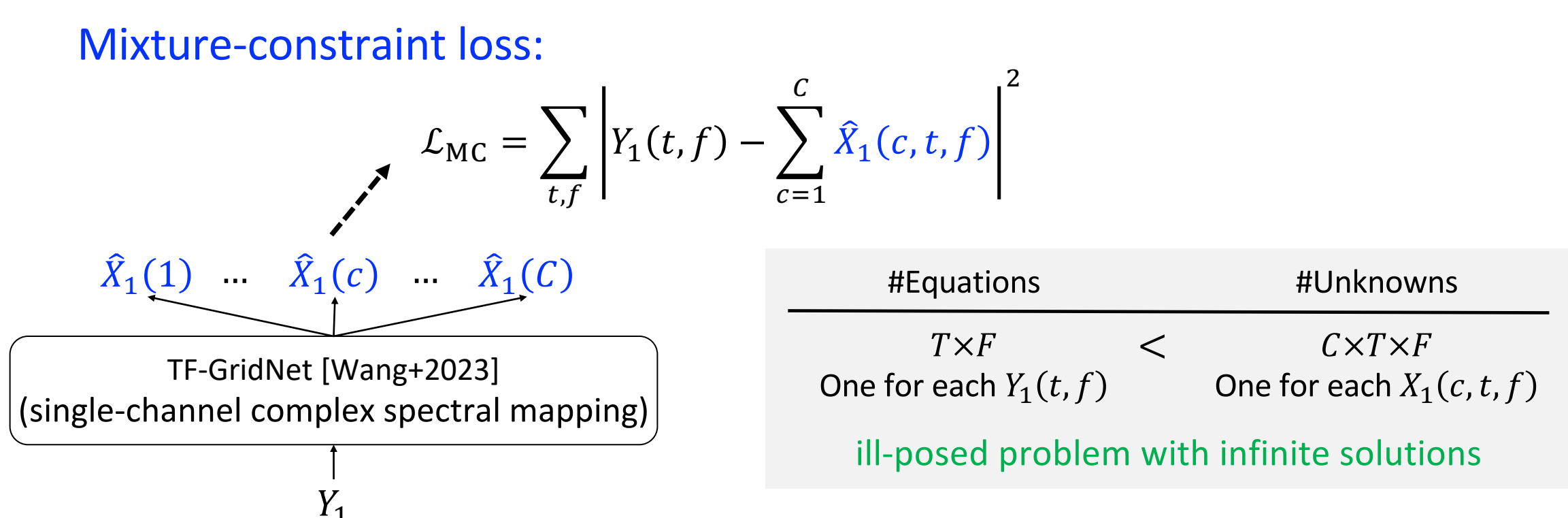
- Speech separation, *a.k.a.* cocktail party problem, aims at separating multi-speaker mixture to individual speaker signals



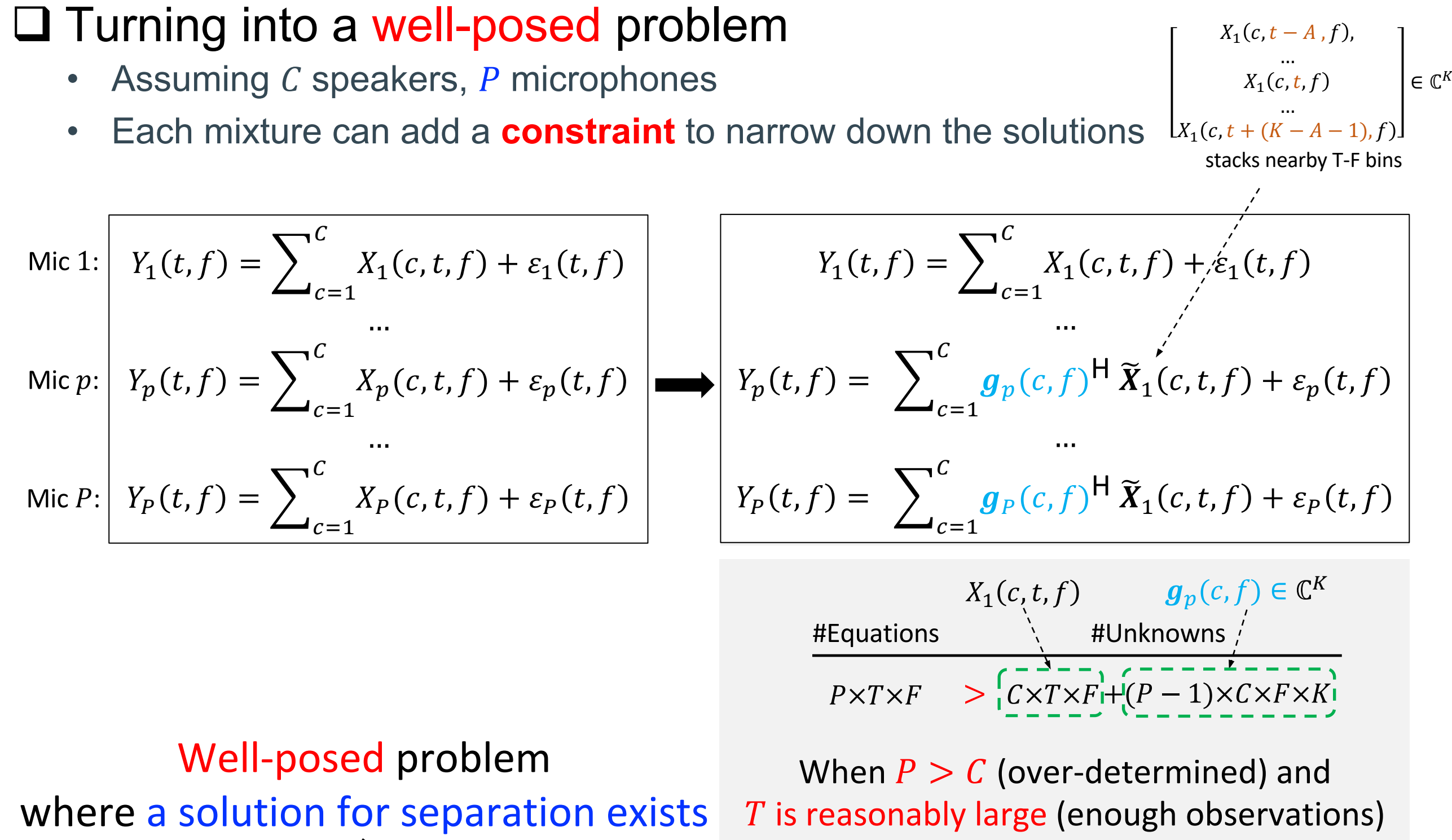
- Supervised separation
 - Use synthetic training data, exhibiting generalization problems
- Unsupervised separation
 - Leverage unlabeled data for training, alleviating generalization issues
 - Earlier studies still require synthesized mixtures (e.g., MixIT), or rely on conventional spatial processing (e.g., unsupervised deep clustering)
- Our solution: **training DNNs directly on mixtures for separation**

Problem formulation

- Unsupervised monaural separation is **ill-posed**
 - Assuming C speakers, **1** microphone
 - Physical model: $Y_1(t, f) = \sum_{c=1}^C X_1(c, t, f) + \varepsilon_1(t, f)$
 - A possible solution



- Turning into a **well-posed** problem
 - Assuming C speakers, P microphones
 - Each mixture can add a **constraint** to narrow down the solutions



Proposed algorithm: UNSSOR

- Solve a blind deconvolution problem [Levin+2021]

$$\arg \min_{\hat{X}_1(c, t, f), g_p(c, f)} \sum_{t,f} \left(\left| Y_1(t, f) - \sum_{c=1}^C X_1(c, t, f) \right|^2 + \sum_{p=2}^P \left| Y_p(t, f) - \sum_{c=1}^C g_p(c, f)^H \tilde{X}_1(c, t, f) \right|^2 \right)$$

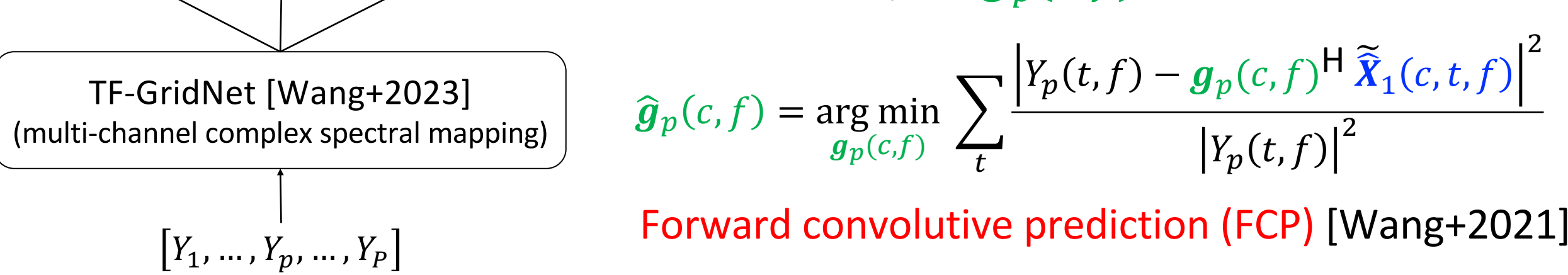
- Not solvable if not assuming prior knowledge about the filter or the source
- Our solution: **model speech patterns via unsupervised deep learning**

- UNSSOR

Mixture constraint at all microphones

$$\mathcal{L}_{MC} = \sum_{t,f} \left(\left| Y_1(t, f) - \sum_{c=1}^C \hat{X}_1(c, t, f) \right|^2 + \sum_{p=2}^P \left| Y_p(t, f) - \sum_{c=1}^C \hat{g}_p(c, f)^H \tilde{X}_1(c, t, f) \right|^2 \right)$$

- \mathcal{L}_{MC} determines whether \hat{X}_1 is good, providing **sample-level supervision**
- How to compute $\hat{g}_p(c, f)$?



- FCP filters \tilde{X}_1 to approximate speaker images at other mics

- When \tilde{X}_1 is reasonably accurate

Let $Y_p = X_p(c) + V_p(c)$

Ideally, $V_p(c)$ is independent from $\tilde{X}_1(c)$ and $X_p(c)$

$$\begin{aligned} \hat{g}_p(c, f) &= \arg \min_{g_p(c, f)} \sum_t \frac{|Y_p(t, f) - g_p(c, f)^H \tilde{X}_1(c, t, f)|^2}{|Y_p(t, f)|^2} \\ &= \arg \min_{g_p(c, f)} \sum_t \frac{|X_p(c, t, f) + V_p(c, t, f) - g_p(c, f)^H \tilde{X}_1(c, t, f)|^2}{|Y_p(t, f)|^2} \\ &= \arg \min_{g_p(c, f)} \sum_t \frac{|X_p(c, t, f) - g_p(c, f)^H \tilde{X}_1(c, t, f)|^2 + |V_p(c, t, f)|^2}{|Y_p(t, f)|^2} \\ &= \arg \min_{g_p(c, f)} \sum_t \frac{|X_p(c, t, f) - g_p(c, f)^H \tilde{X}_1(c, t, f)|^2}{|Y_p(t, f)|^2} \end{aligned}$$

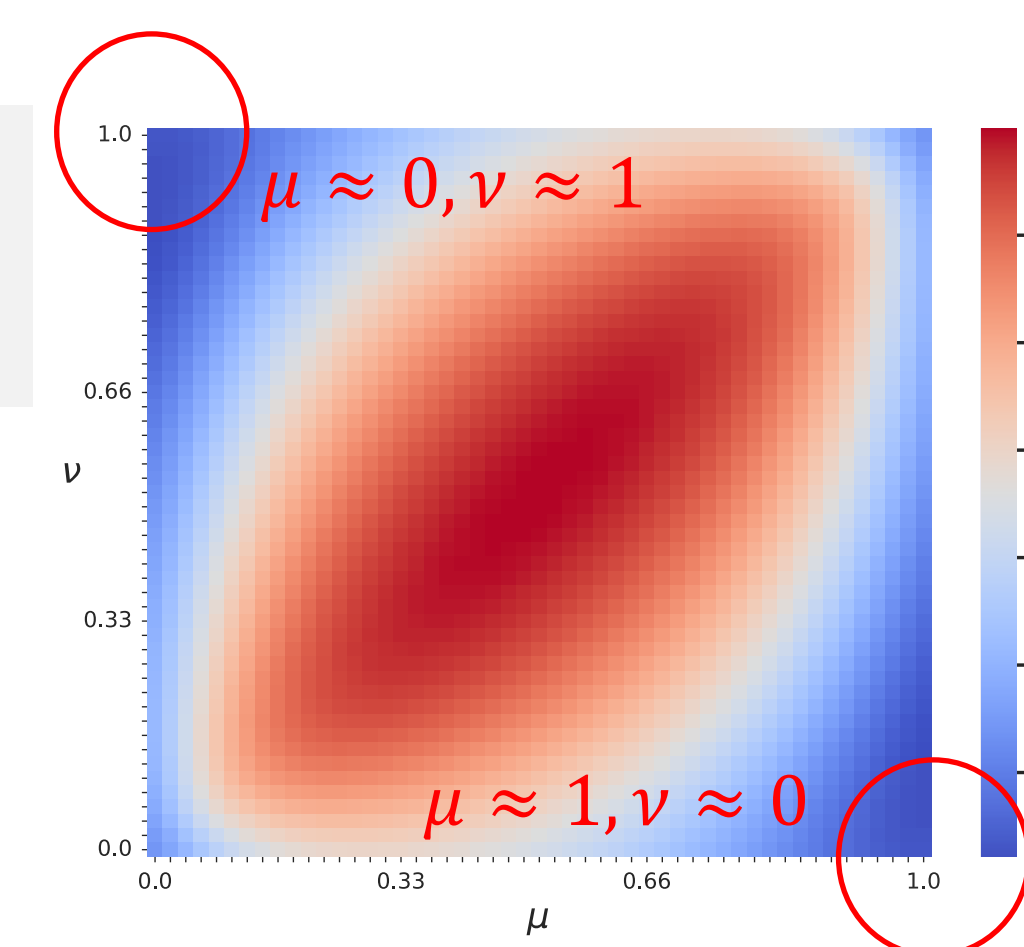
- Minimizing \mathcal{L}_{MC} promotes separation

- Hypothesized separation results

$$\begin{aligned} \hat{X}_1(1) &= \mu \times X_1(1) + \nu \times X_1(2) + \varepsilon_1/2 \\ \hat{X}_1(2) &= (1-\mu) \times X_1(1) + (1-\nu) \times X_1(2) + \varepsilon_1/2 \end{aligned}$$

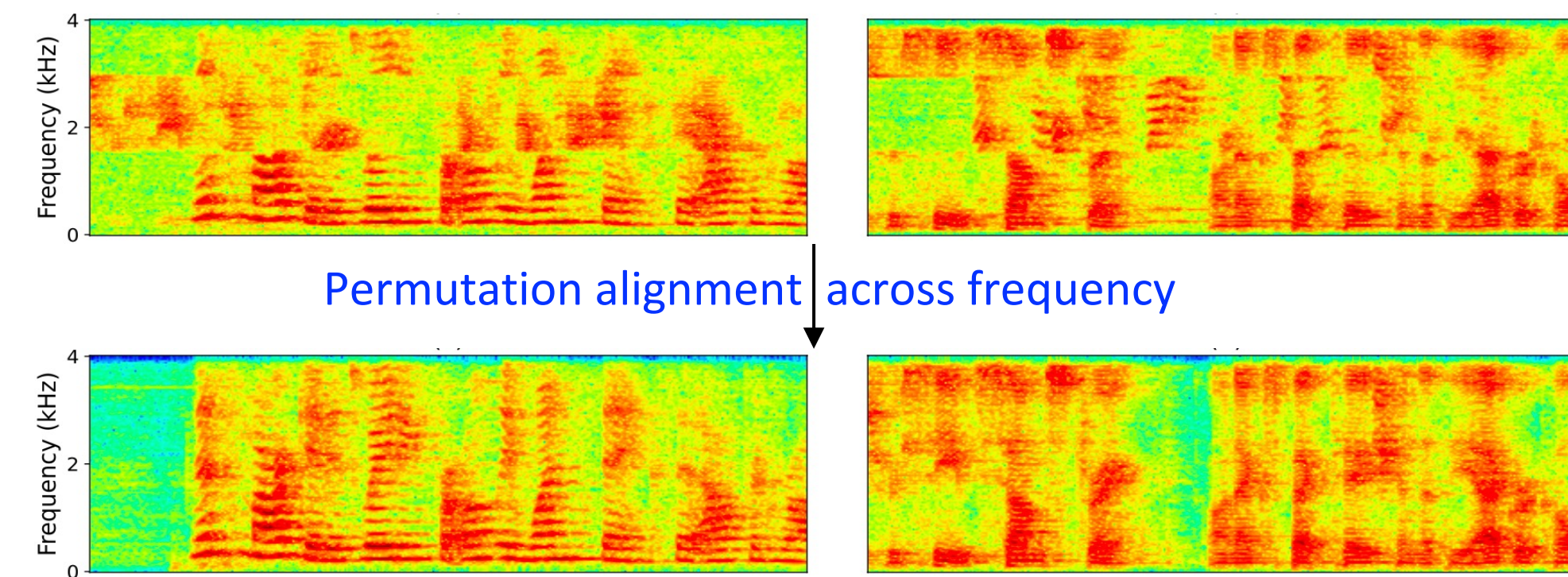
where $0 \leq \mu \leq 1$, $0 \leq \nu \leq 1$, and 2 speakers

- Good separation
 - $\mu \approx 0, \nu \approx 1$ and $\mu \approx 1, \nu \approx 0$
- Bad separation
 - $\mu \approx 0, \nu \approx 0$ and $\mu \approx 1, \nu \approx 1$
 - μ, ν both away from 0 and 1



- Frequency permutation problem

- Happens as FCP is performed at each frequency **independently** from the others



- Propose to addressing frequency permutation during training
- Intra-source magnitude scattering (ISMS) loss
 - Source magnitudes are more scattered when frequency permutation happens

$$\mathcal{L}_{ISMS} = \sum_{p=1}^P \frac{\sum_t \frac{1}{C} \sum_{c=1}^C \text{var}(\log|\hat{X}_p^{FCP}(c, t, \cdot)|)}{\sum_t \text{var}(\log|Y_p(t, \cdot)|)}$$

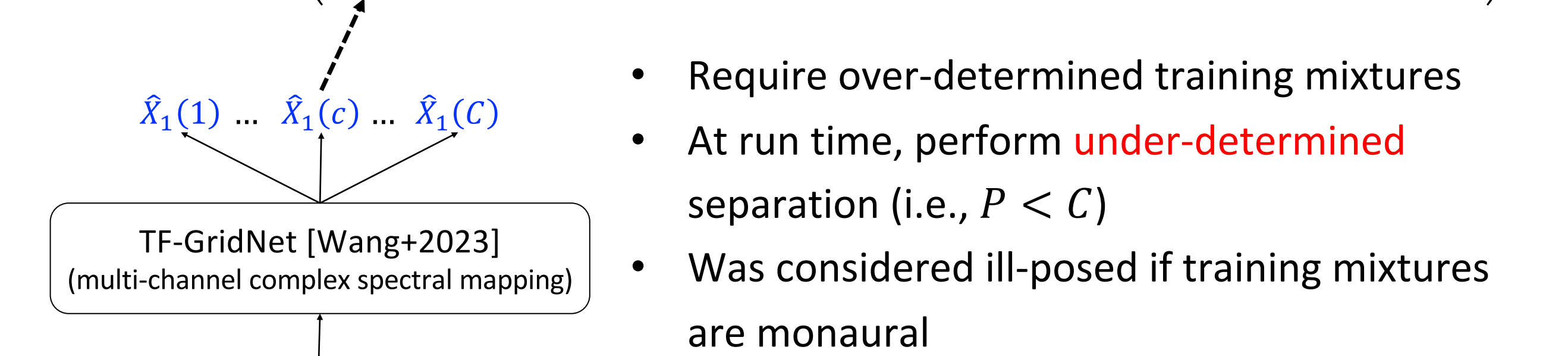
- Combine \mathcal{L}_{MC} and \mathcal{L}_{ISMS} for training

UNSSOR for under-determined separation

- Monaural input, but loss on multiple microphones

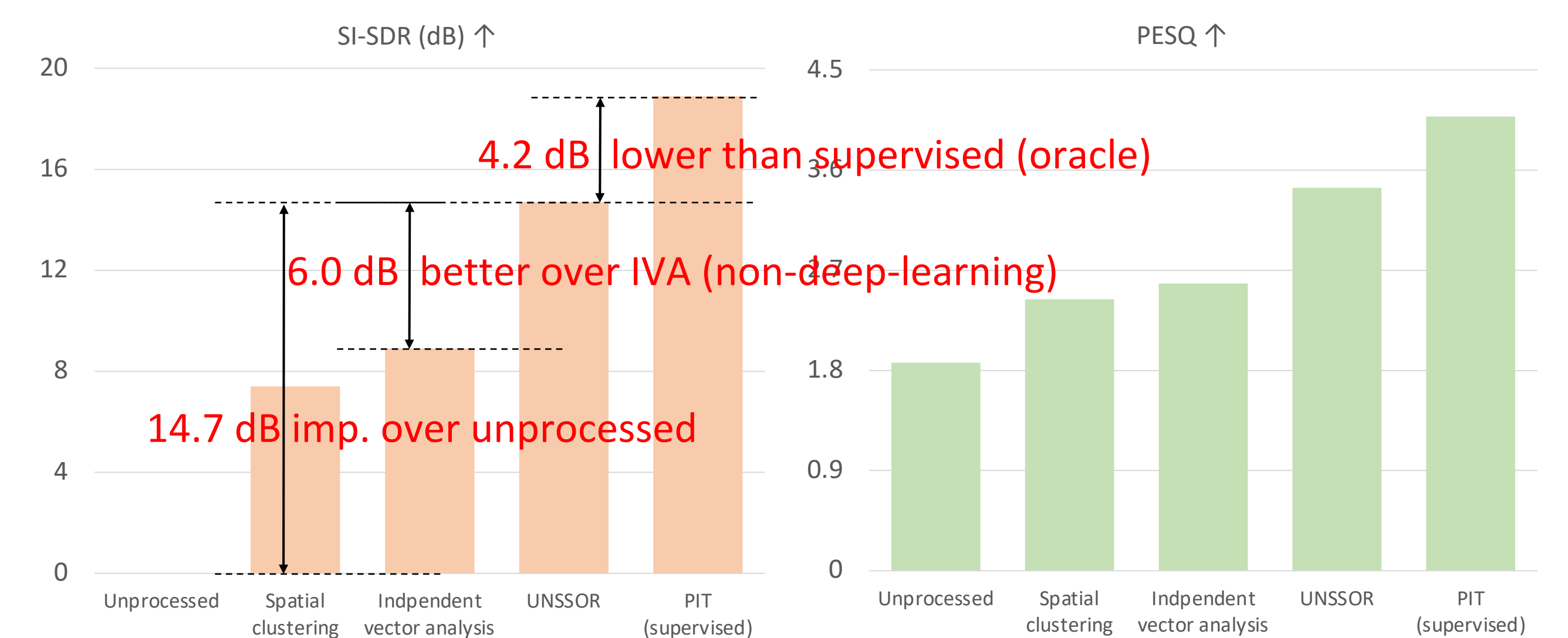
Mixture constraint at all microphones

$$\mathcal{L}_{MC} = \sum_{t,f} \left(\left| Y_1(t, f) - \sum_{c=1}^C \hat{X}_1(c, t, f) \right|^2 + \sum_{p=2}^P \left| Y_p(t, f) - \sum_{c=1}^C \hat{g}_p(c, f)^H \tilde{X}_1(c, t, f) \right|^2 \right)$$

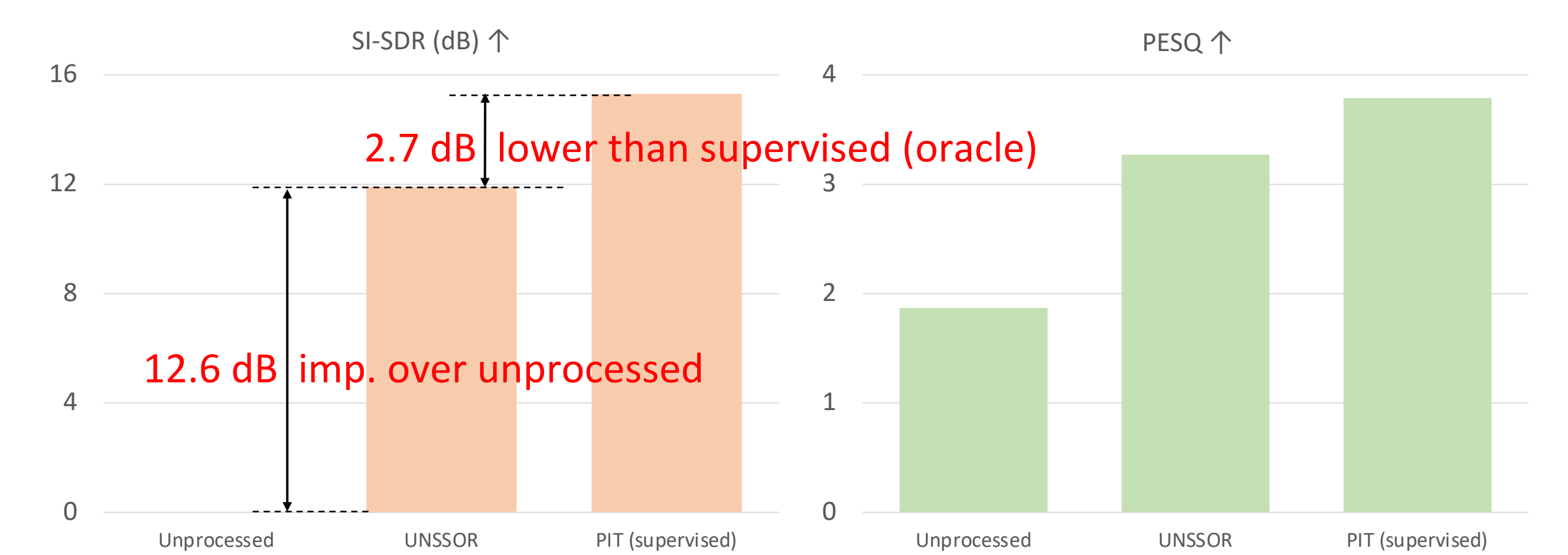


Evaluation results

- SMS-WSJ dataset [Drude+19]: reverb 2-speaker mixture with weak noise
- Results on 2-speaker separation (**6-channel input and loss**)



- Results on 2-speaker separation (**1-channel input, 6-channel loss**)



References

Levin et al. (2011), "Understanding Blind Deconvolution Algorithms," In: IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 33, no. 12, pp. 2354–2367.

Wang et al. (2021), "Convolutional Prediction for Monaural Speech Dereverberation and Noisy-Reverberant Speaker Separation," In: IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 29, pp. 3476–3490.

Wang et al. (2023), "TF-GridNet: Integrating Full- and Sub-Band Modeling for Speech Separation," In: IEEE/ACM Transactions on Audio, Speech, and Language Processing, vol. 31, pp. 3221–3236.

Drude et al. (2019), "SMS-WSJ: Database, Performance Measures, and Baseline Recipe for Multi-Channel Source Separation and Recognition," In: arXiv preprint arXiv:1910.13934.

Sound demo

