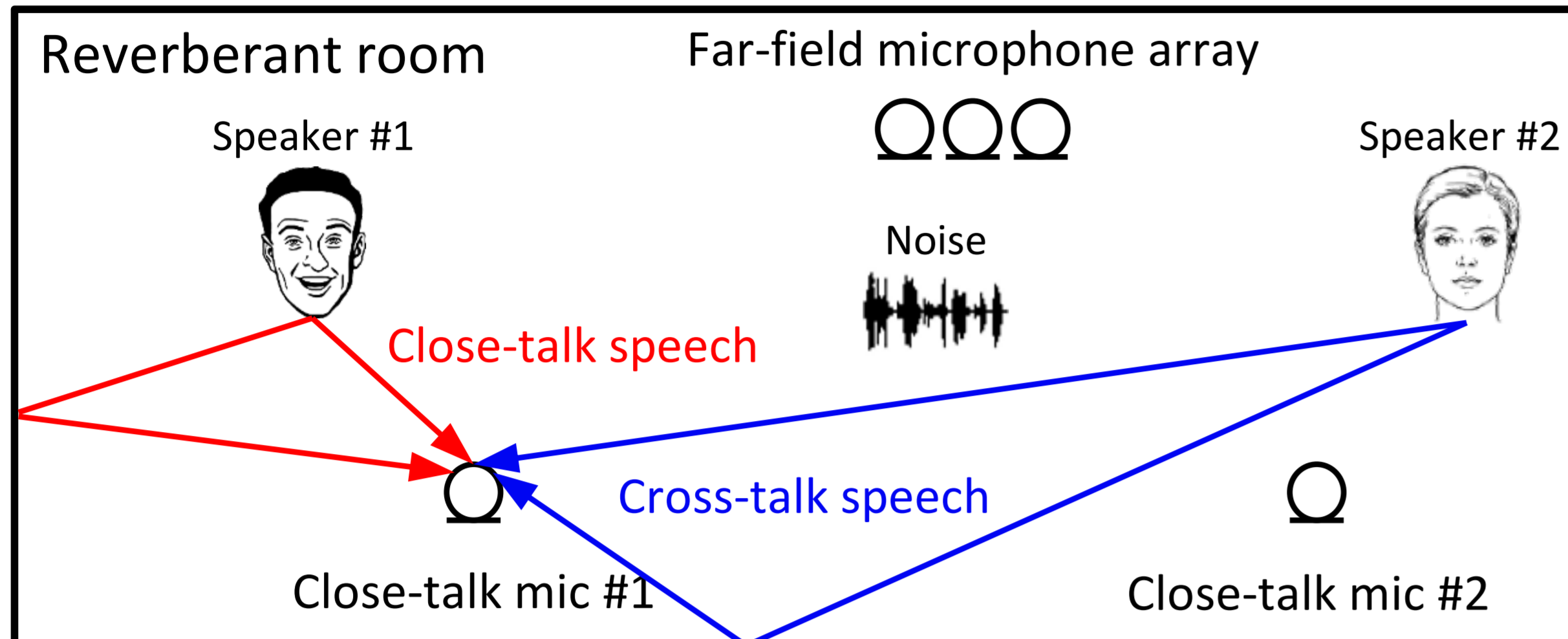




## 1. Motivation

- Close-talk mixture has a high input SNR of target speaker, but often contains significant cross-talk speech



- Cross-talk reduction (CTR) aims at reducing cross-talk speech and enhancing close-talk speech

- Could enable many applications, e.g.,
  - Generate pseudo-labels for real-recorded far-field mixtures
  - Generate pseudo-reference signals for metric computation
  - Reduce labeling efforts of annotators
- Supervised CTRnet on simu. data?
  - Suffers from generalization issues, as simu. data often mismatches real data
- We propose un-/weakly-supervised CTRnet
  - Can be trained directly on real data, realizing better generalizability

## 2. Formulating CTR as blind deconvolution

- Physical model

- Assuming  $P$  far-field mics, and  $C$  speakers, each wearing a close-talk mic

$$\text{Close-talk mixture } c: Y_c(t, f) = \sum_{c'=1}^C X_c(c', t, f) + \varepsilon_c(t, f)$$

Note: subscript indexes mics

$$\text{Far-field mixture } p: Y_p(t, f) = \sum_{c=1}^C X_p(c, t, f) + \varepsilon_p(t, f)$$

- Let  $Z(c) = X_c(c)$  denotes close-talk speech of speaker  $c$

$$\begin{aligned} Y_c(t, f) &= Z(c, t, f) + \sum_{c'=1, c' \neq c}^C X_c(c', t, f) + \varepsilon_c(t, f) \\ &= Z(c, t, f) + \sum_{c'=1, c' \neq c}^C g_c(c', f)^H \tilde{Z}(c', t, f) + \varepsilon_c'(t, f) \end{aligned}$$

$$Y_p(t, f) = \sum_{c=1}^C g_p(c, f)^H \tilde{Z}(c, t, f) + \varepsilon_p'(t, f)$$

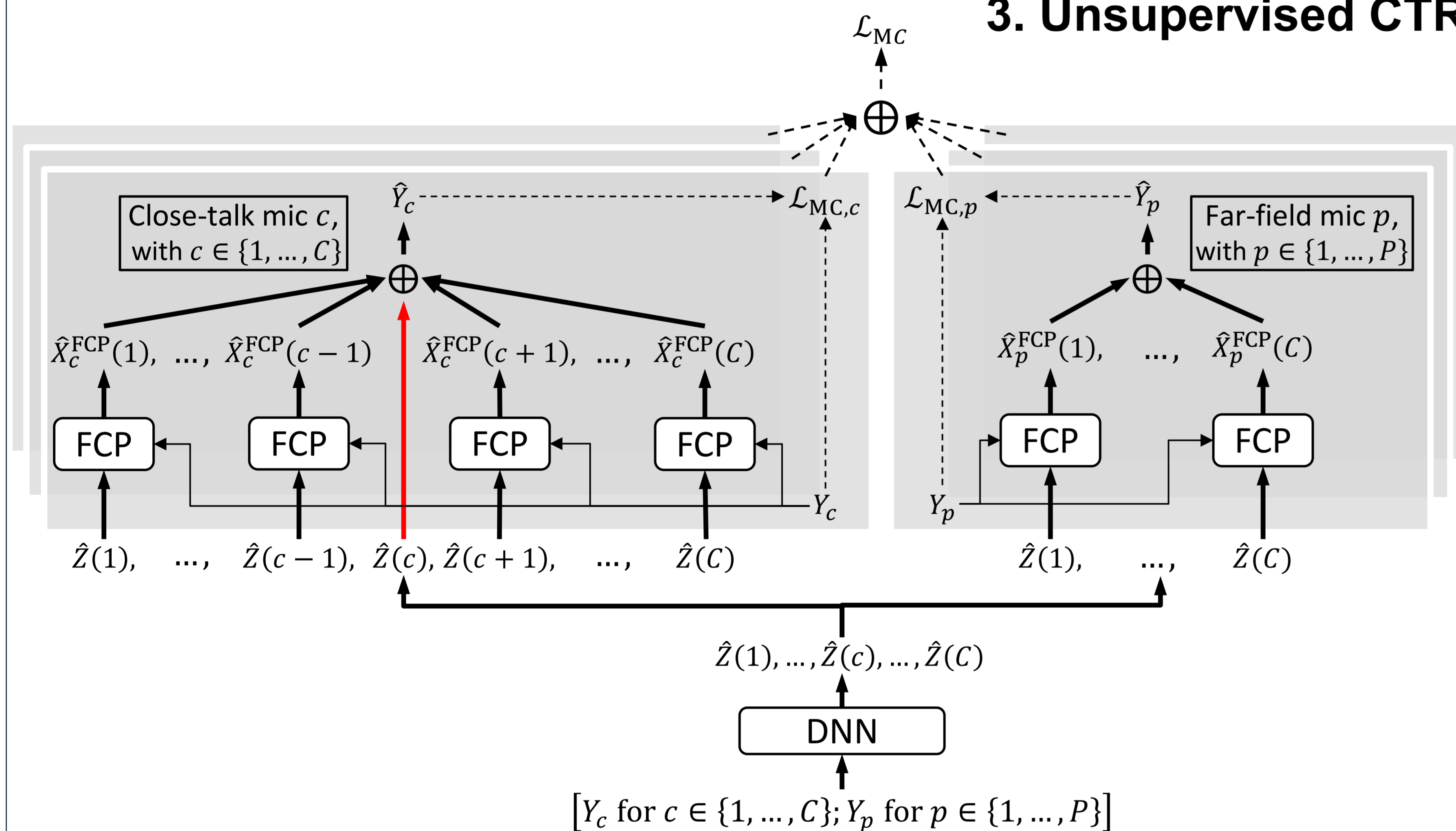
- CTR via blind deconvolution

$$\begin{aligned} \operatorname{argmin}_{Z(\cdot, \cdot), g(\cdot, \cdot)} \sum_{c=1}^C \sum_{t, f} & \left| Y_c(t, f) - Z(c, t, f) - \sum_{c'=1, c' \neq c}^C g_c(c', f)^H \tilde{Z}(c', t, f) \right|^2 \\ & + \sum_{p=1}^P \sum_{t, f} \left| Y_p(t, f) - \sum_{c=1}^C g_p(c, f)^H \tilde{Z}(c, t, f) \right|^2 \end{aligned}$$

- Not solvable, if not assuming prior knowledge about filter or source

- Our solution: model speech patterns via unsupervised deep learning

## 3. Unsupervised CTRnet



- Input: real & imag. of all close-talk and far-field mixtures

- Output: real & imag. of close-talk speech of each speaker

- Loss: mixture-constraint loss

$$\mathcal{L}_{MC} = \sum_{c=1}^C \mathcal{L}_{MC,c} + \alpha \times \sum_{p=1}^P \mathcal{L}_{MC,p}$$

$$\mathcal{L}_{MC,c} = \sum_{t, f} \mathcal{F} \left( Y_c(t, f), \hat{Z}(c, t, f) + \sum_{c'=1, c' \neq c}^C \hat{g}_c(c', f)^H \tilde{Z}(c', t, f) \right)$$

$$\mathcal{L}_{MC,p} = \sum_{t, f} \mathcal{F} \left( Y_p(t, f), \sum_{c=1}^C \hat{g}_p(c, f)^H \tilde{Z}(c, t, f) \right)$$

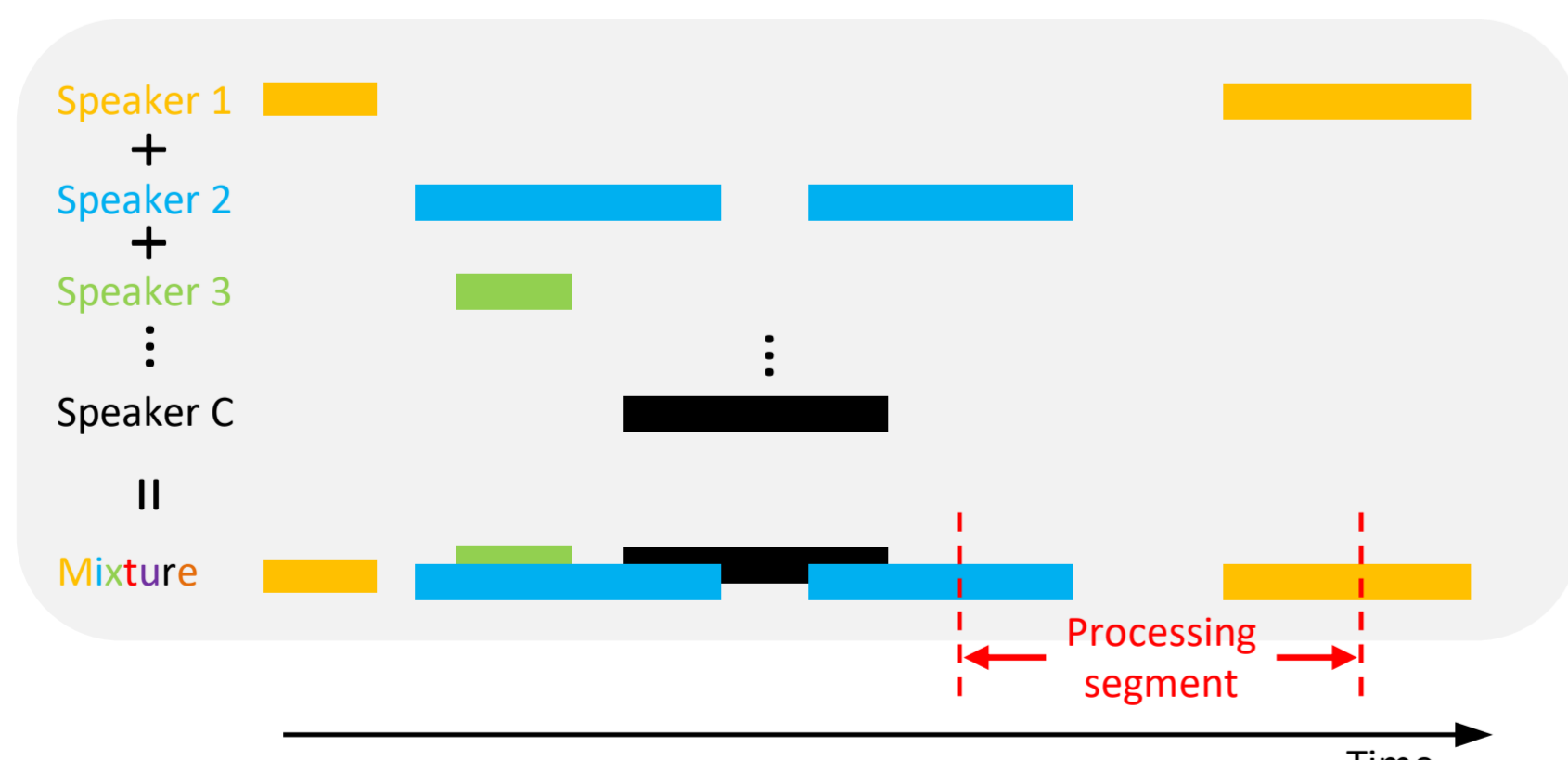
- Linear filter  $\hat{g}_a(\cdot, \cdot)$  is estimated via FCP [Wang+2021]

$$\hat{g}_a(c, f) = \operatorname{argmin}_{g_a(c, f)} \sum_t \frac{|Y_a(t, f) - g_a(c, f)^H \tilde{Z}(c, t, f)|^2}{|Y_a(t, f)|^2}$$

- $a$  indexes all  $P$  far-field and  $C$  close-talk mics

## 4. Weakly-supervised CTRnet

- Realistic speaker overlap is sparse and time-varying



- Unsupervised CTRnet often under-/over-separates mixed speakers

- Like clustering, assuming more clusters  $\rightarrow$  smaller clusters, but some should be merged

- Our solution: leverage speaker-activity timestamps

- Let  $d(c) \in \{0, 1\}^N$  denote timestamps of speaker  $c$ , with  $N$  denoting #samples
- Muting during training: avoid using predictions in silent ranges for FCP

$$\hat{Z}(c, t, f) := \hat{Z}(c, t, f) \times D(c, t) \times E(c)$$

Speaker  $c$  active at frame  $t$ ? Speaker  $c$  active in the training segment?

- Speaker-activity loss: predictions in silent ranges should be zero

$$\mathcal{L}_{SA,c} = \frac{\|\hat{Z}(c) \odot (1 - d(c))\|_1}{\|y_c \odot (1 - d(c))\|_1} \times \frac{N - \|d(c)\|_1}{N}$$

## 5. Experiments

- On a simu. dataset (2-speaker, reverb, weak noise, fully-overlapped)

- Unsupervised CTRnet works, better than spatial clustering (SC) and IVA

Systems	$I$	$J$	$C$	$P$	Masking/Mapping	$\alpha$	$H/L$	SI-SDR (dB) $\uparrow$	SDR (dB) $\uparrow$	PESQ $\uparrow$	eSTOI $\uparrow$
Unprocessed mixture	-	-	-	-	-	-	-	14.7	14.7	2.92	0.875
Unsupervised CTRnet	30	0	2	6	Mapping	$1/P$	$1/-$	<b>26.5</b>	<b>26.8</b>	3.88	<b>0.973</b>
SC [Boeddeker, 2019]	-	-	-	6	-	-	-	-1.9	7.1	2.27	0.561
IVA [Scheibler and Saijo, 2022]	-	-	-	6	-	-	-	22.6	23.7	3.66	0.948

Table 1: Averaged separation results of unsupervised CTRnet on SMS-WSJ-FF-CT.

- On real-recorded CHiME-7 (4-speaker, reverb, noisy, sparse overlap)

Row	Systems	Muting?	DA-WER (%) $\downarrow$				
			$I$	$J$	$C$	$P$	
0	Unprocessed mixture	-	-	-	4	28.3	27.8
1	Unsupervised CTRnet	-	19	1	4	22.5	25.1
2	Weakly-supervised CTRnet	$\times$	19	1	4	79.1	73.0
3	Weakly-supervised CTRnet	$\checkmark$	19	1	4	20.5	22.6
4	GSS [Boeddeker et al., 2018]	-	-	-	4	26.2	26.6

Table 3: ASR results of CTRnet on CHiME-7 close-talk mixtures.

## 6. Conclusion

- CTRnet can be trained on real data and can effectively reduce cross-talk speech on real data

- The proposed un-/weakly-supervised learning based methodology for blind deconvolution works on challenging real data such as CHiME-7