



# A HYBRID MODEL FOR WEAKLY-SUPERVISED SPEECH DEREVERBERATION

Louis Bahrman, Mathieu Fontaine, Gaël Richard

LTCI, Télécom Paris, Institut Polytechnique de Paris, France



## SUMMARY

### Context

- DNN-based approaches for dereverberation often rely on dry/wet pairs
- Anechoic data is rare and expensive
- Metrics-based weak supervision fails to generalize to other metrics

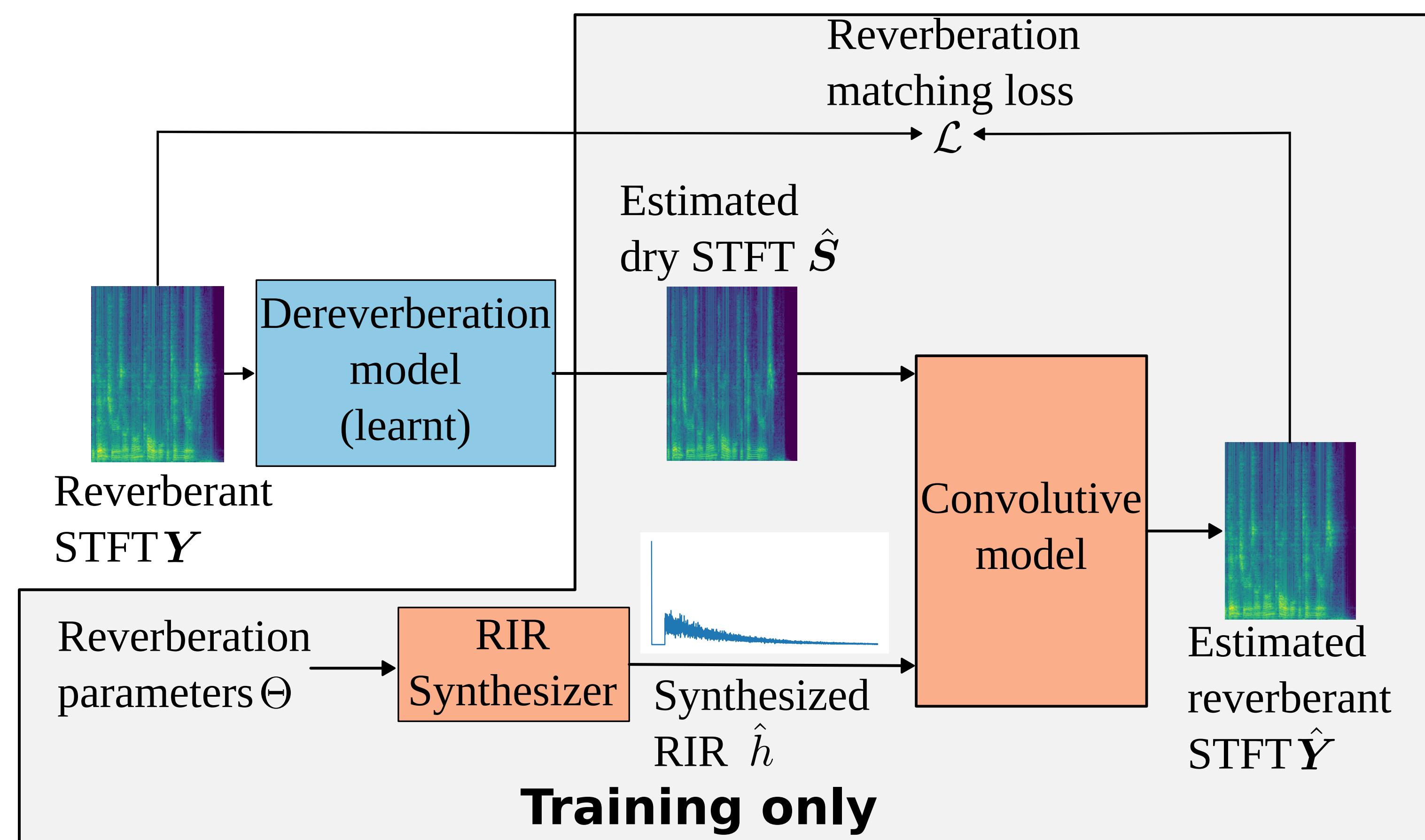
### Main takeways

- Reverberation-based weak supervision for dereverberation
- Dereverberation supervised by only pairs of wet and  $RT_{60}$

### Code



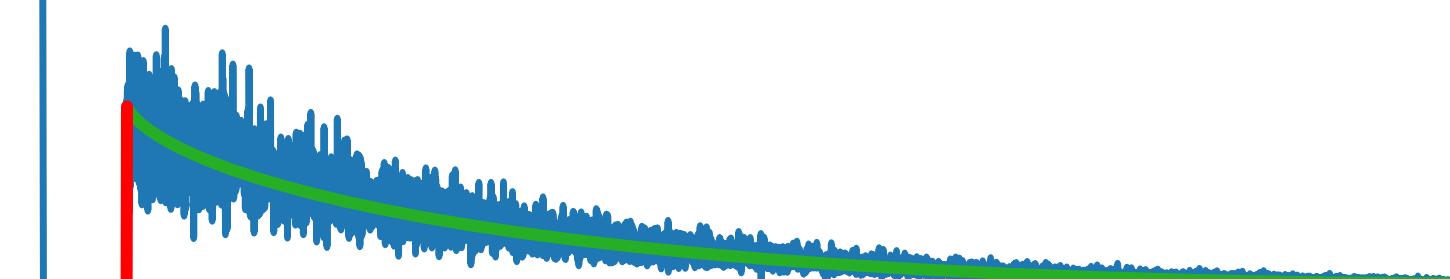
## METHOD



### RIR synthesizer

Inspired from Polack's model

$$\hat{h}_\Theta(n) = \begin{cases} |b(n)| e^{-\frac{3 \ln(10)}{RT_{60} f_s} n} & \text{if } n > 2n_m \\ 1 & \text{if } n = 0 \\ 0 & \text{otherwise} \end{cases}$$



### Convulsive model

Equivalent to  $\hat{y} = \hat{s} * \hat{h}_\Theta$  in time-frequency

$$\mathcal{C}(\hat{S}, \hat{h}_\Theta) = \sum_{f'=f-F'}^{f+F'} \sum_{t'=0}^{\min(t; T_h)} \hat{\mathcal{H}}_{f,f',t'} \hat{S}_{f',t-t'},$$

with  $\hat{\mathcal{H}}_{f,f',t'}$  constructed from  $\hat{h}_\Theta(n)$ .  
(Avargel and Cohen 07)

### Parameters

- $RT_{60}$ : Reverberation time
- $b(n) \sim \mathcal{N}(0, \sigma^2)$
- $n_m = \frac{4Vf_s}{cA}$ : Mixing time

### Weak supervision variants

- $\Theta \triangleq \{RT_{60}, \sigma, V, A\}$ : all the parameters
- $\{RT_{60}, \sigma\}$ : fixed mixing time at 20 ms after the peak (mean over training dataset)
- $\{RT_{60}\}$ : fixed mixing time at 20 ms, fixed  $\sigma = 0.02$  (median over the training dataset).

### Reverberation matching loss

- Compare  $\hat{Y}_{f,t} \triangleq \mathcal{C}(\hat{S}, \hat{h}_\Theta)$  and  $Y$
- Sum of two losses:  $\mathcal{L} = \mathcal{L}_C + \lambda \mathcal{L}_M$ :

$$\mathcal{L}_C = \sum_{f,t} [|\hat{Y}_{f,t} - Y_{f,t}|^2]$$

$$\mathcal{L}_M = \sum_{f,t} \left[ \left| \log \left( \frac{1 + \gamma |\hat{Y}_{f,t}|}{1 + \gamma |Y_{f,t}|} \right) \right|^2 \right]$$

With  $\gamma = \lambda = 1$

## RESULTS

Model	WS?	Supervision	SISDR	ESTOI	WB-PESQ	SRMR
FSN	✗	cRM	5.6 ± 3.9	0.84 ± 0.09	2.55 ± 0.68	8.2 ± 3.5
		$h$	4.3 ± 4.0	0.77 ± 0.12	2.03 ± 0.69	7.8 ± 3.1
	✓	{ $RT_{60}, \sigma, V, A$ }	1.0 ± 2.5	<b>0.71 ± 0.14</b>	<b>1.80 ± 0.70</b>	6.9 ± 2.8
BiLSTM	✗	{ $RT_{60}, \sigma$ }	1.1 ± 2.5	0.70 ± 0.14	1.78 ± 0.69	7.0 ± 2.8
		{ $RT_{60}$ }	<b>2.9 ± 3.4</b>	0.71 ± 0.15	1.78 ± 0.71	6.9 ± 2.8
BiLSTM	✗	$ S_{f,t} ^2, \forall f, t$	1.3 ± 4.2	0.78 ± 0.12	2.25 ± 0.79	7.9 ± 3.0
		$h$	0.1 ± 4.1	0.70 ± 0.15	1.80 ± 0.70	7.2 ± 2.7
	✓	{ $RT_{60}, \sigma, V, A$ }	0.8 ± 4.0	0.70 ± 0.15	1.81 ± 0.74	6.9 ± 2.7
BiLSTM	✓	{ $RT_{60}, \sigma$ }	0.7 ± 4.0	0.70 ± 0.15	1.78 ± 0.72	6.8 ± 2.7
		{ $RT_{60}$ }	<b>1.6 ± 3.7</b>	<b>0.71 ± 0.15</b>	<b>1.84 ± 0.75</b>	6.9 ± 2.8
BiLSTM	✓	SRMR (baseline)	-1.5 ± 3.4	0.64 ± 0.18	1.78 ± 0.74	<b>10.9 ± 4.2</b>
		Reverberant	-1.3 ± 3.4	0.68 ± 0.16	1.75 ± 0.74	6.9 ± 2.9