

LIWhiz: A Non-Intrusive Lyric Intelligibility Prediction System for the Cadenza Challenge

RAM C. M. C. SHEKAR AND IVÁN LÓPEZ-ESPEJO*

*DEPT. OF SIGNAL THEORY, TELEMATICS AND COMMUNICATIONS, UNIVERSITY OF GRANADA, SPAIN
ramcharanmc@gmail.com, iloes@ugr.es



UNIVERSIDAD DE GRANADA

INTRODUCTION

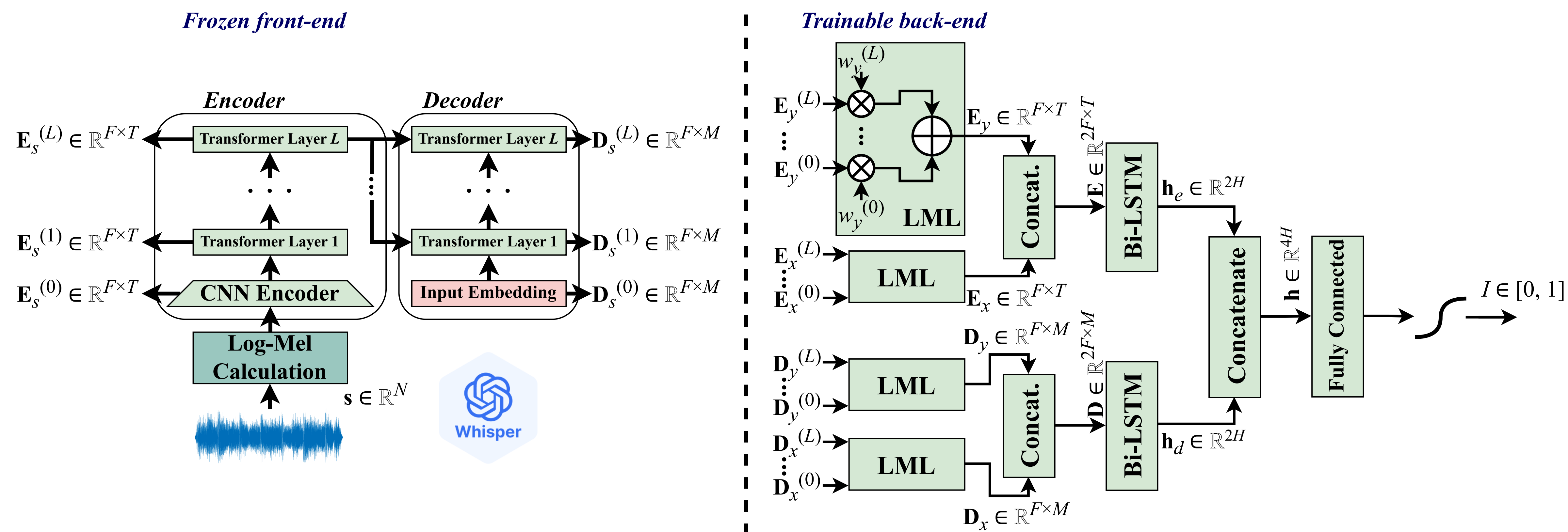
- Understanding lyrics is key to music enjoyment; however, listeners with hearing impairment often struggle to comprehend them
- The development of **lyric intelligibility prediction (LIP)** methods could enable new lyric enhancement technologies
- LIWhiz**: A *non-intrusive* LIP system inspired by our prior work in no-reference speech intelligibility prediction (**wav2vec 2.0 adapted for ASR + lightweight back-end**)¹

¹ H. Wang *et al.*, "No-Reference Speech Intelligibility Prediction Leveraging a Noisy-Speech ASR Pre-Trained Model," in *Proc. of Interspeech 2024*

EXPERIMENTAL SETUP

- To train the back-end, we use *only* the training partition of the **Cadenza Lyric Intelligibility Prediction (CLIP)** dataset
- Ground-truth lyric transcripts are **not** used
- The CLIP dataset is converted to mono and downsampled to 16 kHz
- k -fold cross-validation ($k = 10$) with early stopping (patience = 10 epochs)
- Optimizer**: AdamW ($1r = 10^{-3}$)
- Loss function**: Root mean square error (RMSE)

SYSTEM DESCRIPTION



- Motivated by **LyricWhiz**², a state-of-the-art automatic lyric transcription system (*an intrinsically related task*), we use a frozen **Whisper Large v3** model for feature extraction
- Including x alongside y may help LIWhiz better adapt to the listener's degree of hearing loss

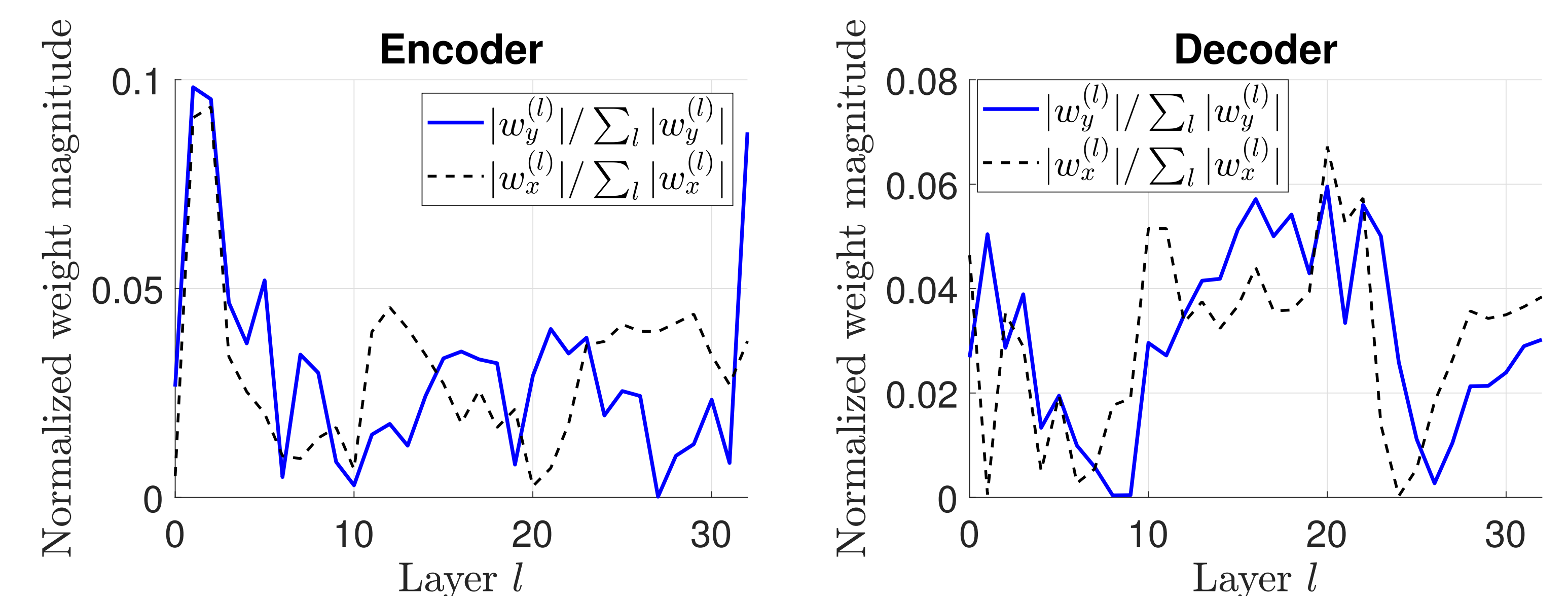
² L. Zhuo *et al.*, "LyricWhiz: Robust Multilingual Zero-Shot Lyrics Transcription by Whispering to ChatGPT," in *Proc. of ISMIR 2023*

RESULTS AND CONCLUSIONS

System	CLIP validation set		CLIP evaluation set	
	RMSE (%) ↓	NCC ↑	RMSE (%) ↓	NCC ↑
Baseline STOI	36.11	0.14	34.89	0.21
Baseline Whisper	29.32	0.59	29.08	0.58
LIWhiz w/o x	27.36	0.66	27.34	0.65
LIWhiz	27.13	0.67	27.07	0.65

- LIWhiz substantially outperforms the baselines on both metrics and sets
- Including the original song excerpt x alongside y leads to a slight improvement in LIP performance

Normalized absolute values of the learned LML weights



- Encoder:** Shallowest layers \rightarrow low-level acoustic encoding
- Decoder:** Middle layers \rightarrow encoder information integrated