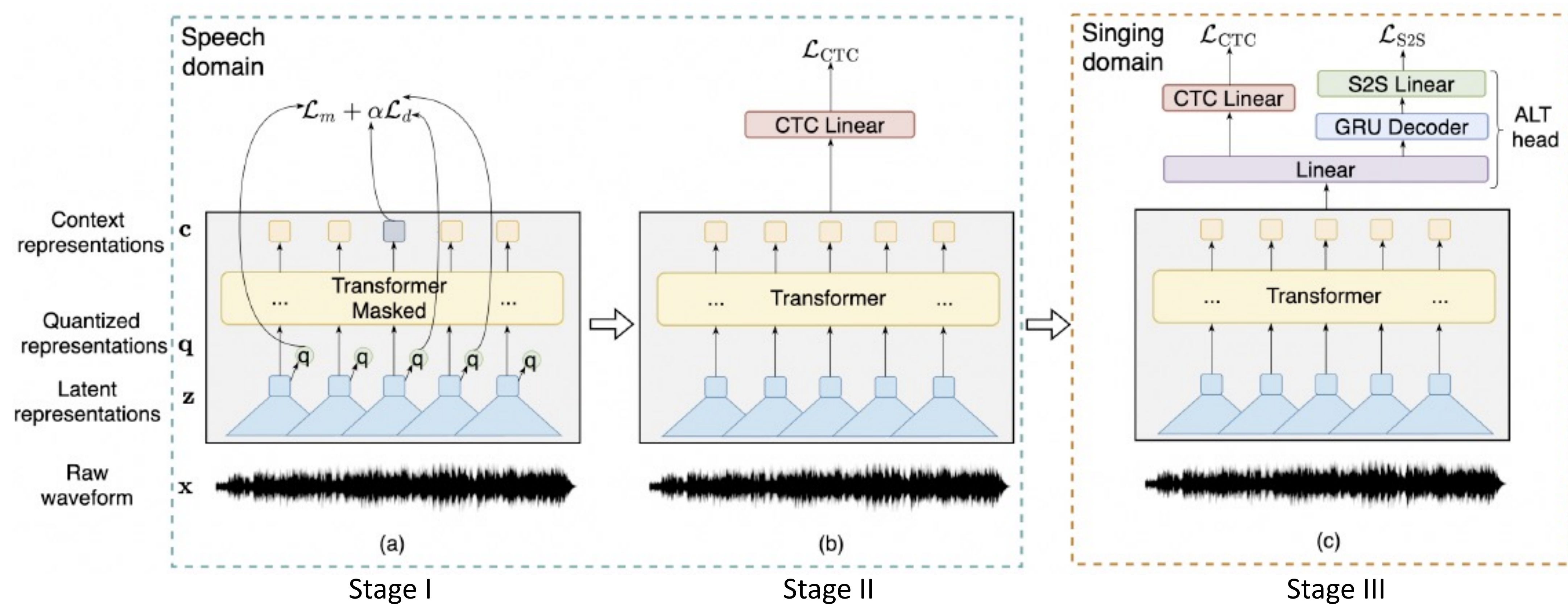


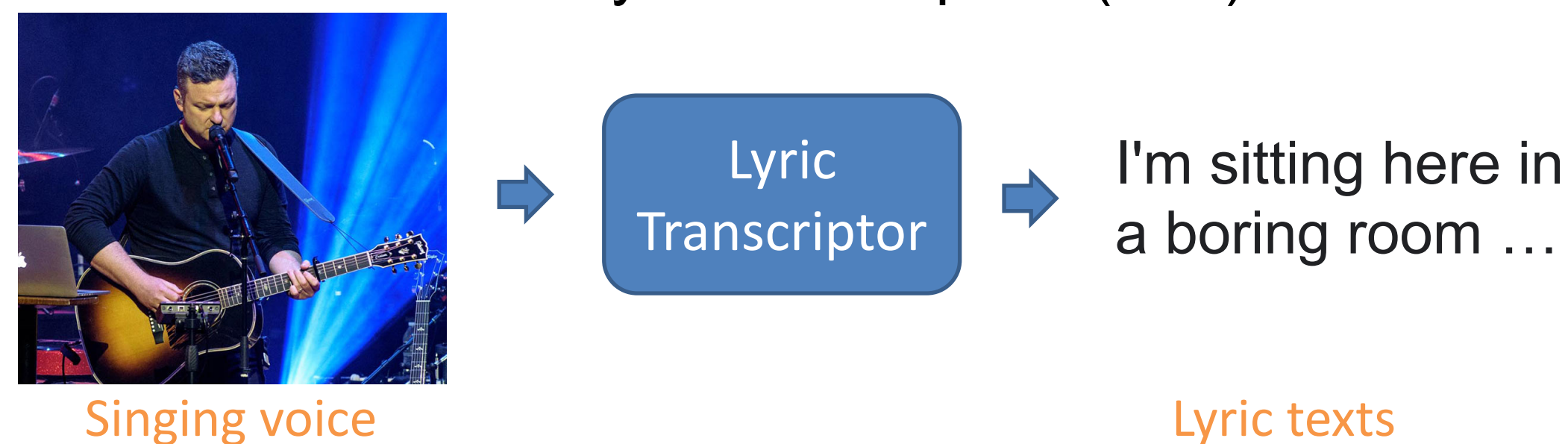
Transfer Learning of wav2vec 2.0 for Automatic Lyric Transcription

Longshen Ou*, Xiangming Gu*, Ye Wang
 (*Both authors contributed equally to this research)



• Introduction:

- The task of Automatic lyric transcription (ALT):



- ALT is difficult, because

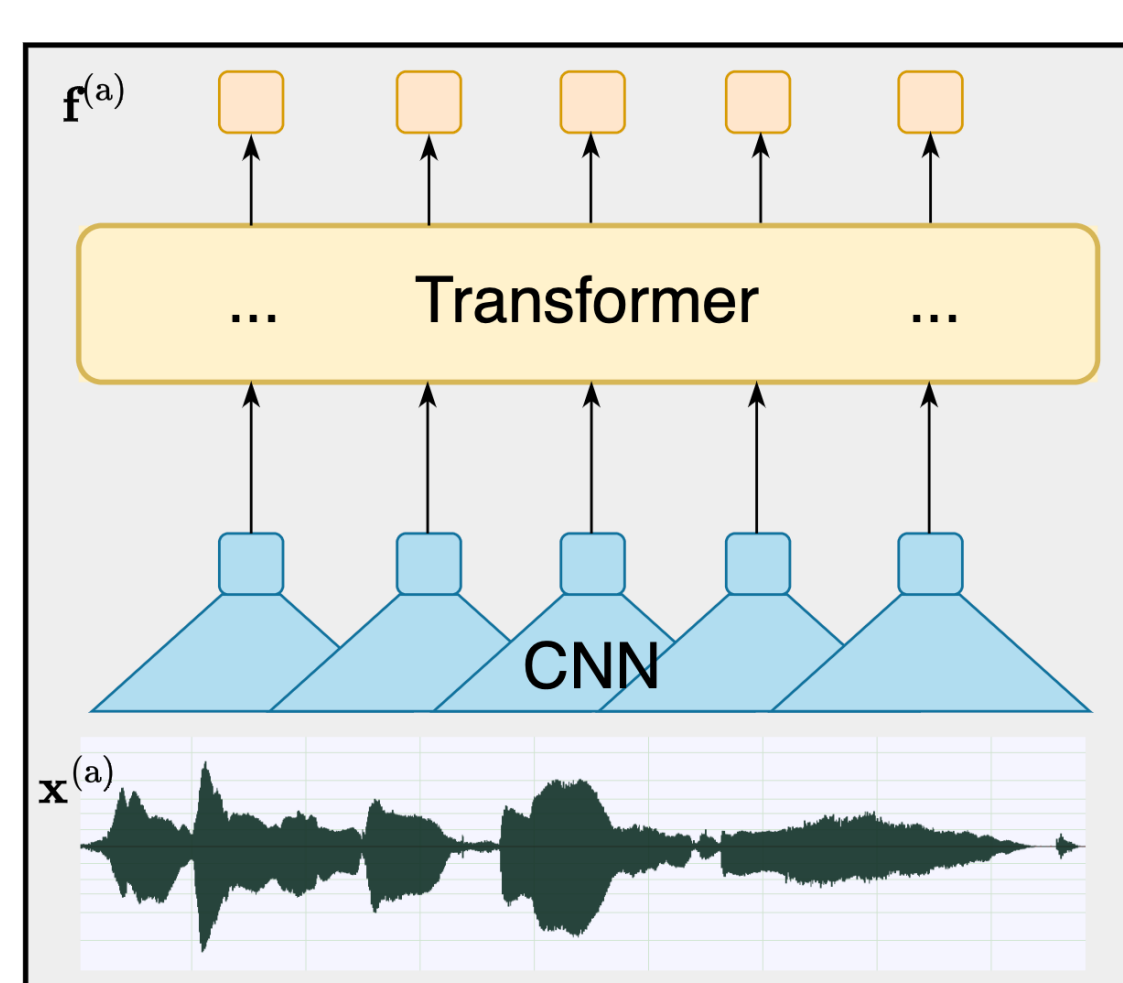
- Degraded intelligibility of sung words
- Perturbations, e.g., musical accompaniments, other ambient noise, etc.
- Small datasets

- We propose to design a transfer-learning-based system to solve this problem

- To utilize the similarity between speech and singing
- To reduce the amount of labeled data needed

• Our System:

- wav2vec 2.0 as the encoder



Structure: CNN + Transformer

Pretraining objectives:

1. **Masked input**, to learn contextualized representation
2. **Contrastive loss**, to encourage contextualized feature to only be similar to local feature of the same frame

- Training procedures:

1. Stage I: Pretraining on speech
2. Stage II: Finetuning on ASR task, on speech domain
3. Stage III: Finetuning on ALT task, on singing domain

- Extend to hybrid CTC/attention

- Training objective: weighted CTC and sequence-to-sequence loss

$$\mathcal{L}_w = \lambda_a \mathcal{L}_{CTC} + (1 - \lambda_a) \mathcal{L}_{S2S}$$

- Inference: weighted likelihood of CTC, decoder, and language model

$$\begin{aligned} w' = \arg \max_w & \lambda_b \log \sum_{\pi \in \mathcal{B}^{-1}(w)} \prod_{t=1}^T p(\pi_t | f_t) \\ & + (1 - \lambda_b) \log \prod_{s=1}^S p(w_s | w_{<s}, f_{1:T}) \\ & + \lambda_c \log p_{LM}(w) \end{aligned}$$

• Experiments:

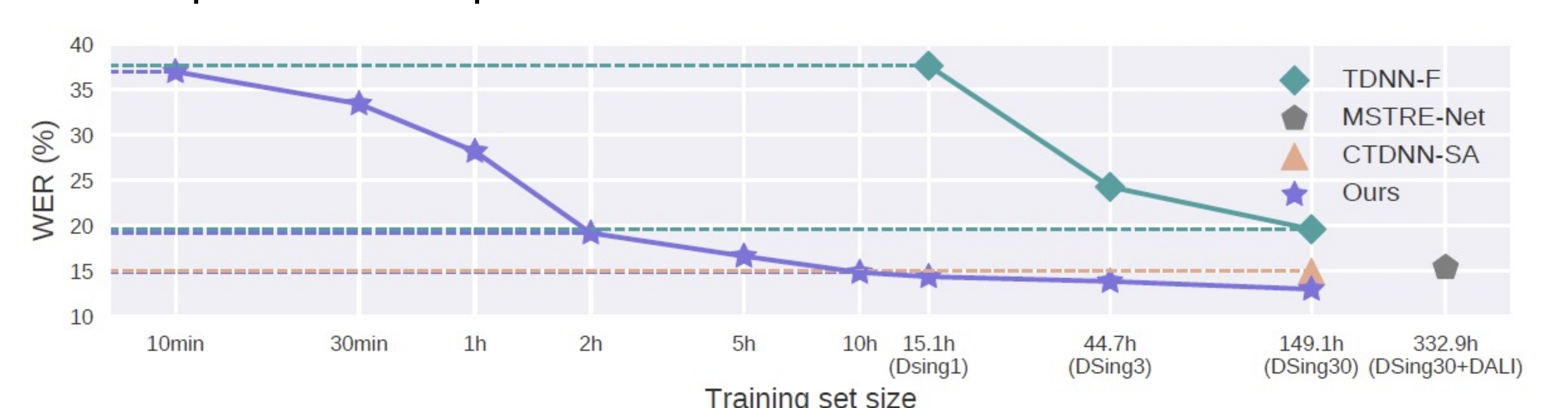
- We achieve state-of-the-art performance on various benchmark datasets (metric: word error rate, WER, lower is better)

Method	DSing ^{dev}	DSing ^{test}	DAL ^{test}	Jamendo	Hansen	Mauch
TDNN-F [5]	23.33	19.60	67.12	76.37	77.59	76.98
CTDNN-SA [6]	<u>17.70</u>	<u>14.96</u>	76.72	66.96	78.53	78.50
Genre-informed AM [3]	-	56.90	-	50.64	39.00	40.43
MSTRE-Net [7]	-	15.38	<u>42.11</u>	<u>34.94</u>	<u>36.78</u>	<u>37.33</u>
DE2 - segmented [4]	-	-	-	44.52	49.92	-
Ours	<u>12.34</u>	<u>12.99</u>	<u>30.85</u>	<u>33.13</u>	<u>18.71</u>	<u>28.48</u>

- Ablation studies show the necessity of all three training stages, and the effectiveness of hybrid CTC/attention

Method	DSing ^{dev}	DSing ^{test}	Method	DSing ^{dev}	DSing ^{test}
Ours	12.34	12.99	CTC	19.86	20.99
- Finetuning	12.64 (+ 0.30)	14.58 (+ 2.59)	+ S2S	15.63 (-4.23)	16.95 (-4.04)
- Pretraining	35.61 (+24.27)	39.13 (+26.14)	+ LM	12.34 (-7.52)	12.99 (-8.00)

- Our system demonstrates effectiveness in low-resource experiment setups:



• Take-aways

- We presented an ALT solution that takes advantage of similarity between speech and singing, by transfer learning of wav2vec 2.0 on singing data
- Our method surpasses previous ones on various ALT datasets, and still demonstrates competitive results with very limited labeled data.
- Check out our project website github.com/guxm2021/ALT_SpeechBrain, or scan this QR code:

