

Introduction

Previous challenges focused solely on high accuracy

➔ **We aim to reduce inference speed while improving recognition accuracy.**

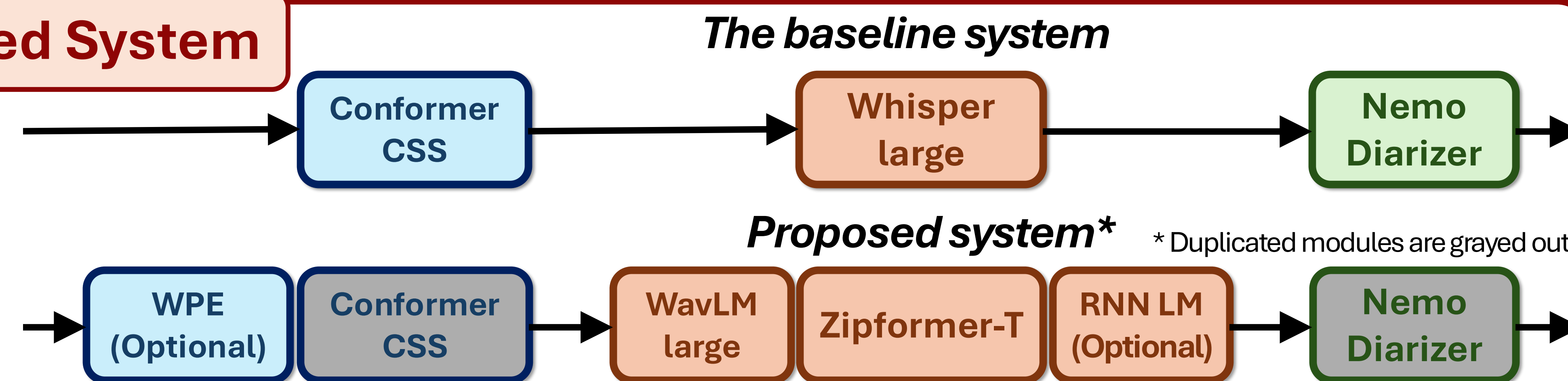
Inference time of the baseline system mostly comes from the ASR module

➔ **We mainly worked on the ASR module.**

Faster Inference Speed

Higher Accuracy

Proposed System

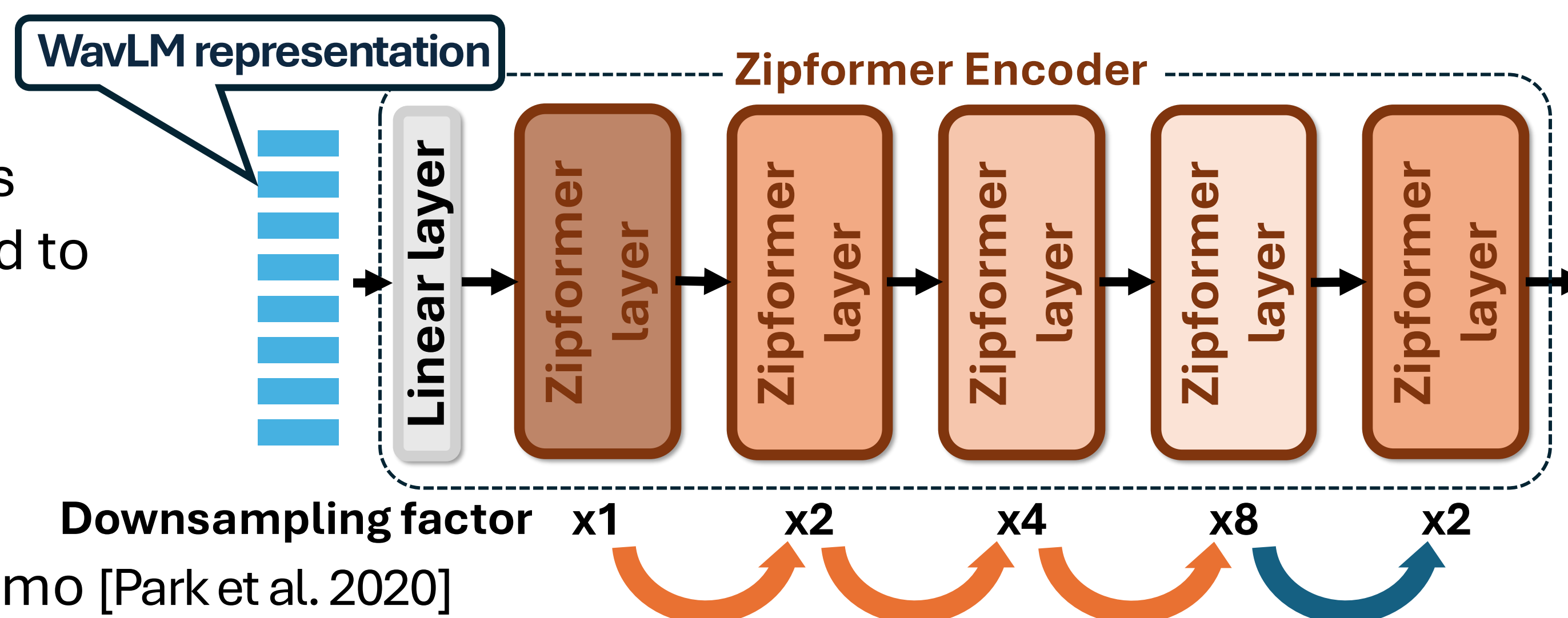


Continuous Speech Separation

- Conformer CSS [Chen, et al. 2020]
- Checkpoints from the baseline system was used
- Block-wise WPE dereverberation**
 - Applied to 3.0 sec of audio chunks
 - Implementation in nara_wpe was used

Automatic Speech Recognition

- WavLM-large feature extractor** [Chen, et al. 2021]
 - The output of the last (24th) layer was used
- Zipformer transducer** [Zengwei, et al. 2023]
 - Fast and memory-efficient because of U-Net-like downsampling and upsampling operations
 - Better performance compared to other transformers
 - Word timestamps
 - Obtained by referring to input audio frame indexes
 - Word duration was limited to 1.0 sec



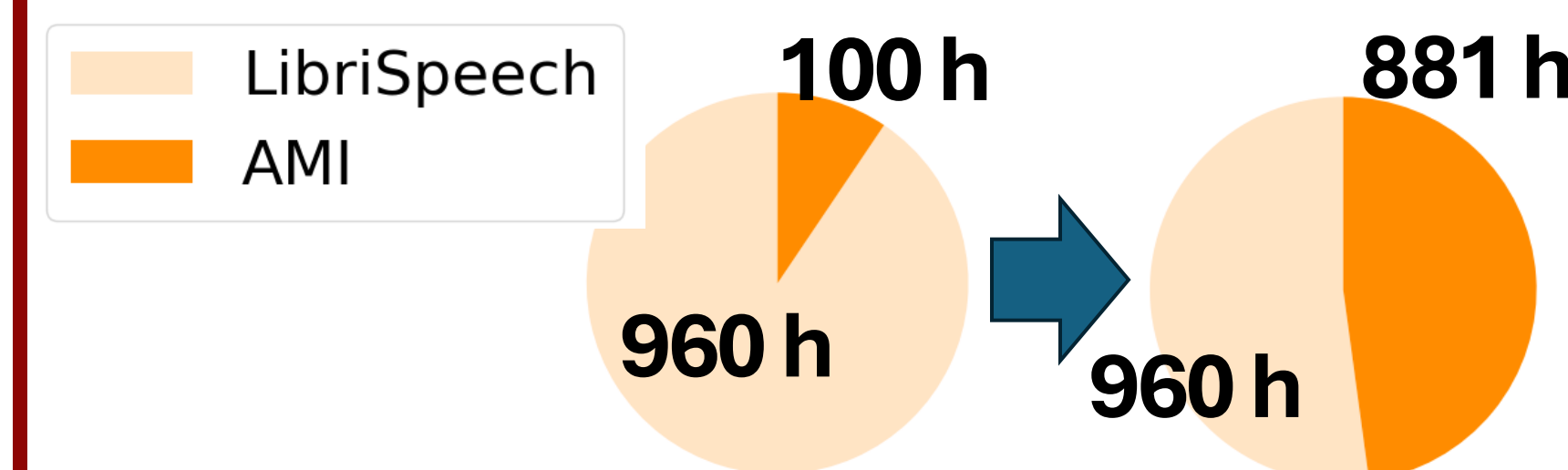
Speaker Diarization

- Word-nmesc dierizer from Nemo [Park et al. 2020]
 - Extract speaker embedding for each word using word timestamps
 - Assign speaker labels by performing spectral clustering

Experiment

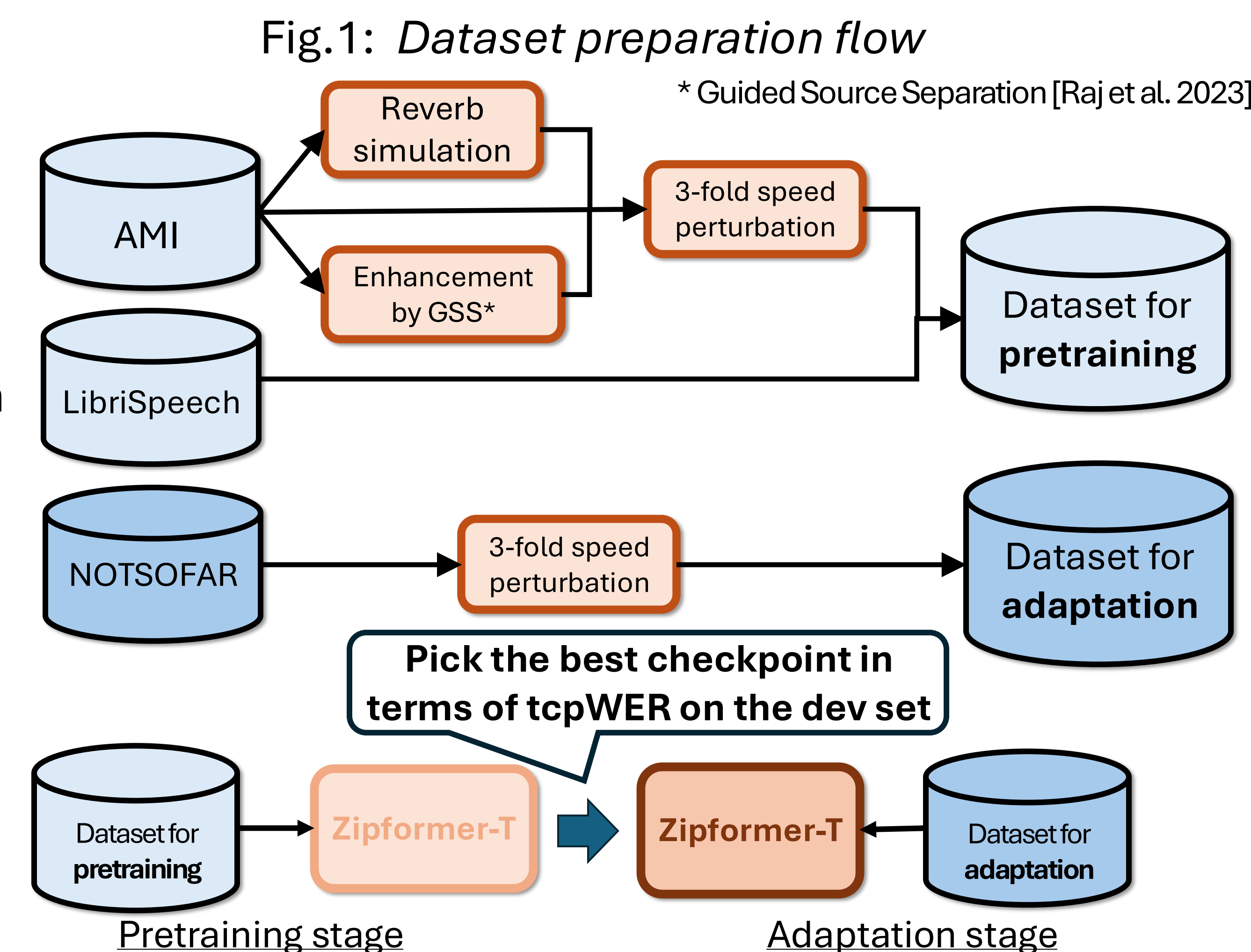
1. Dataset

- Two datasets for model training
- Overall, we tried to **increase the ratio of meeting recordings**



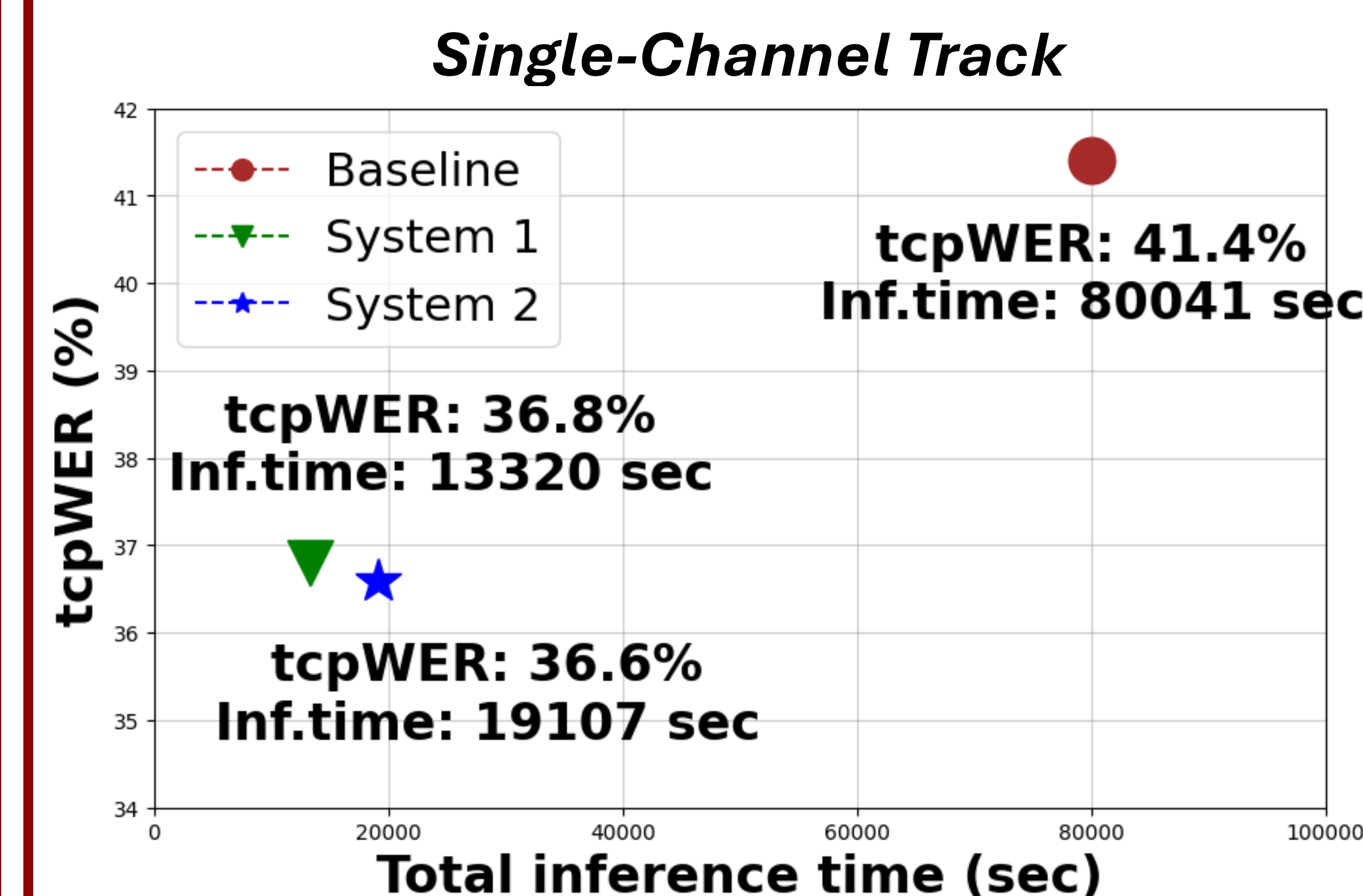
2. Training

- Two NVIDIA A100 40GB GPUs
- Pruned RNN-T loss [Kuang, et al. 2022]
- SpecAugment



Result

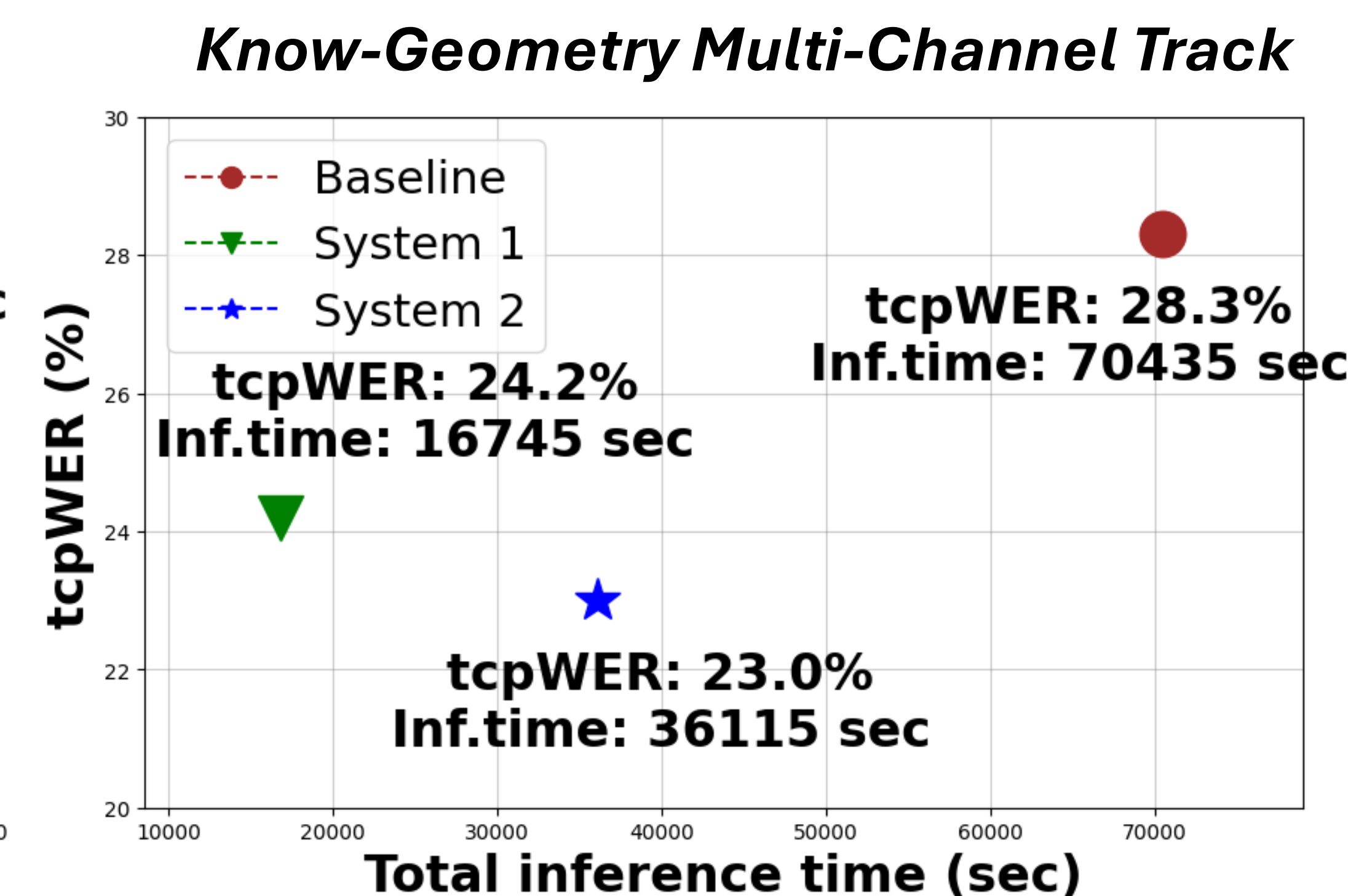
Component	Single-channel track		Multi-channel track	
	System 1	System 2	System 1	System 2
WPE				✓
Conformer CSS	✓	✓	✓	✓
WavLM-large	✓	✓	✓	✓
Zipformer transducer	✓	✓	✓	✓
RNN LM		✓		
Nemo word nmesc diarizer	✓	✓	✓	✓



- We reduced **tcpWER** by **11.1 - 11.6 %** (🥉 3rd place out of 6 teams)
- We reduced **inference time** by **76.1 - 83.4 %**

Brought by our ASR module only

- Shallow fusion by RNNLM resulted in only 0.2 % tcpWER reduction



- We reduced **tcpWER** by **14.5 - 18.7 %** (4th place out of 10 teams)
- We reduced **inference time** by **48.7 - 76.2 %**

- WPE dereverberation reduced tcpWER by 1.2 %, but increased the inference time by nearly 20000 sec.