

# CoVoMix: Advancing Zero-Shot Speech Generation for Human-like Multi-talker Conversations



**Presenter: Leying Zhang (Shanghai Jiao Tong University)**

Advisor: Yanmin Qian, Yao Qian

*NeurIPS 2024, Vancouver, Canada*



# Outline

- 1. Background and Motivation**
- 2. Design of CoVoMix**
- 3. Data Processing**
- 4. Evaluation Results**
- 5. Conclusion**

# Outline

**1. Background and Motivation**

2. Design of CoVoMix

3. Data Processing

4. Evaluation Results

5. Conclusion





# Background and Motivation

**Conversation** is the most frequent form of human communication



Daily Talk



Meeting



Podcast



Human-Machine  
Interaction

## Challenges

### 1) Rendering the spontaneous-style speech

filled pauses, interjections, repair, repetition, laughing, etc

### 2) Managing natural speech transition of multiple speakers

the timing of their speech, determining when one speaks and when the other will follow

# Outline

1. Background and Motivation

**2. Design of CoVoMix**

3. Data Processing

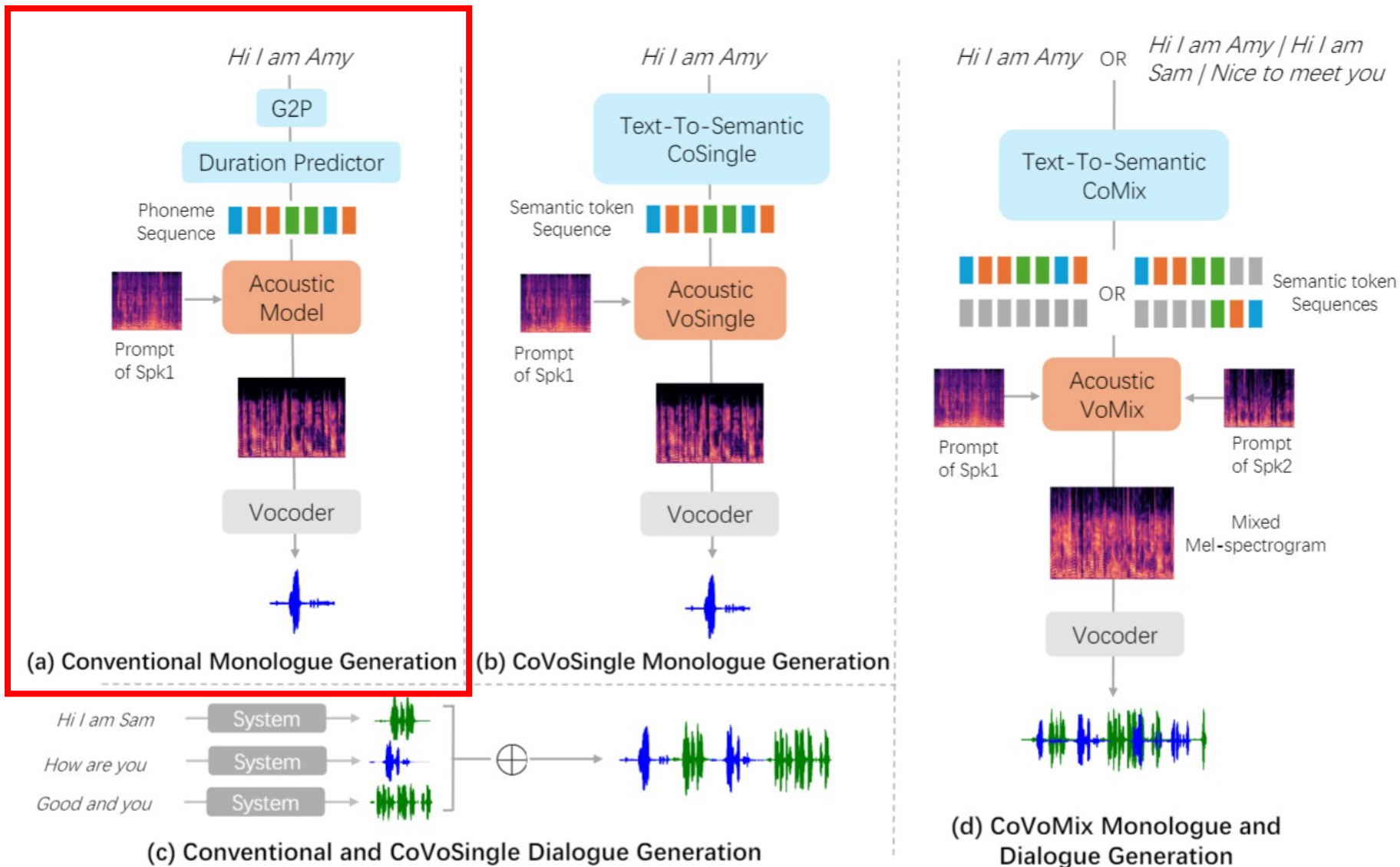
4. Evaluation Results

5. Conclusion



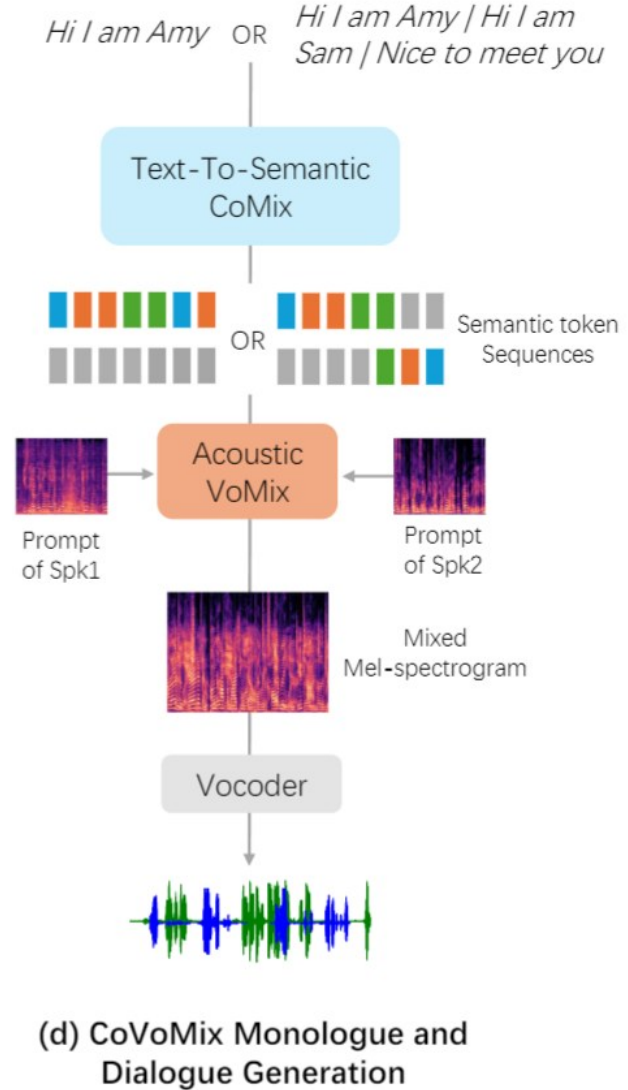
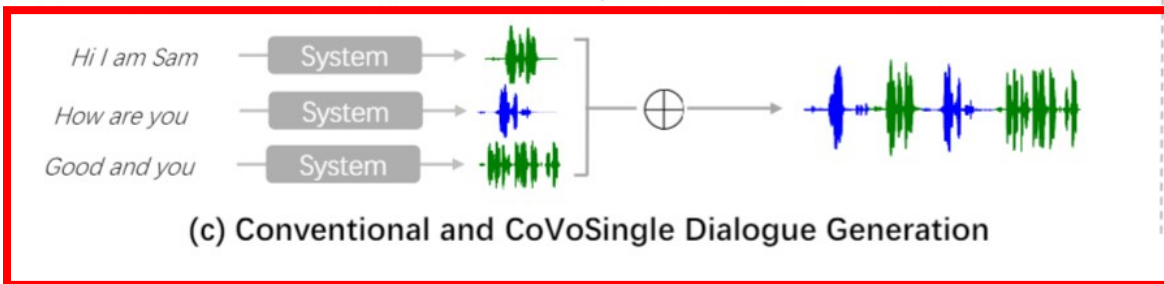
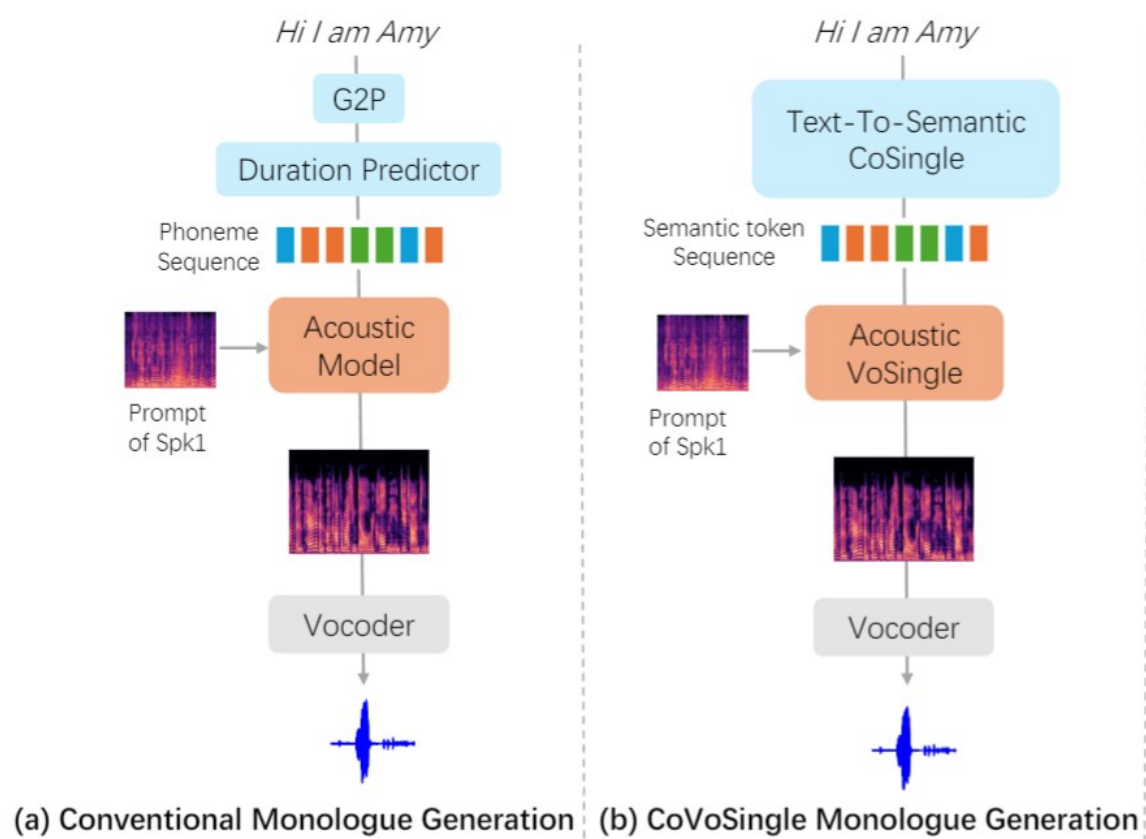


# Pipeline Comparison



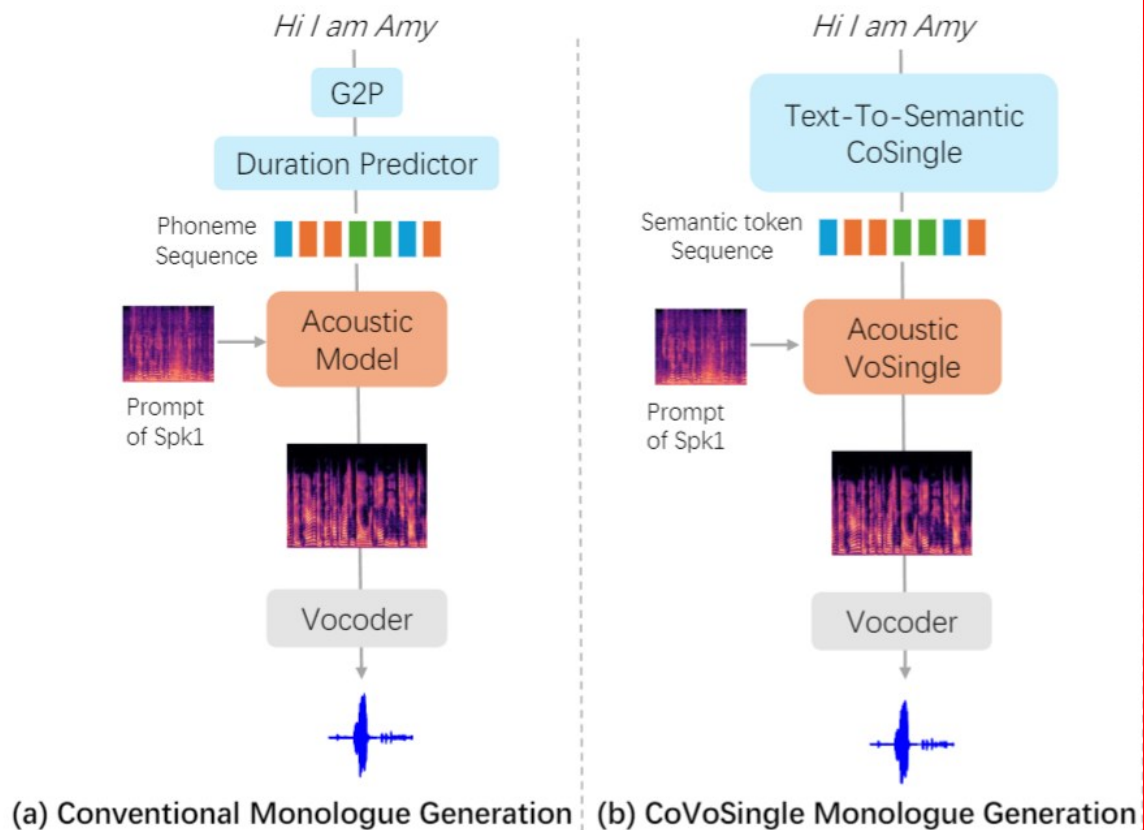


# Pipeline Comparison

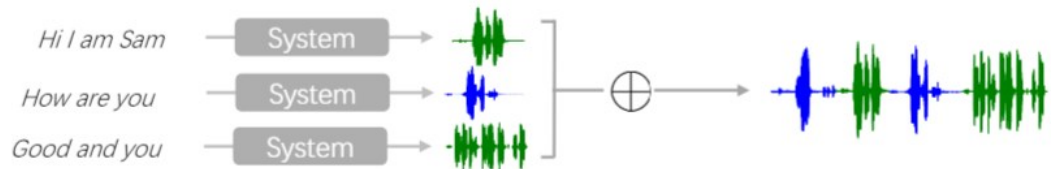




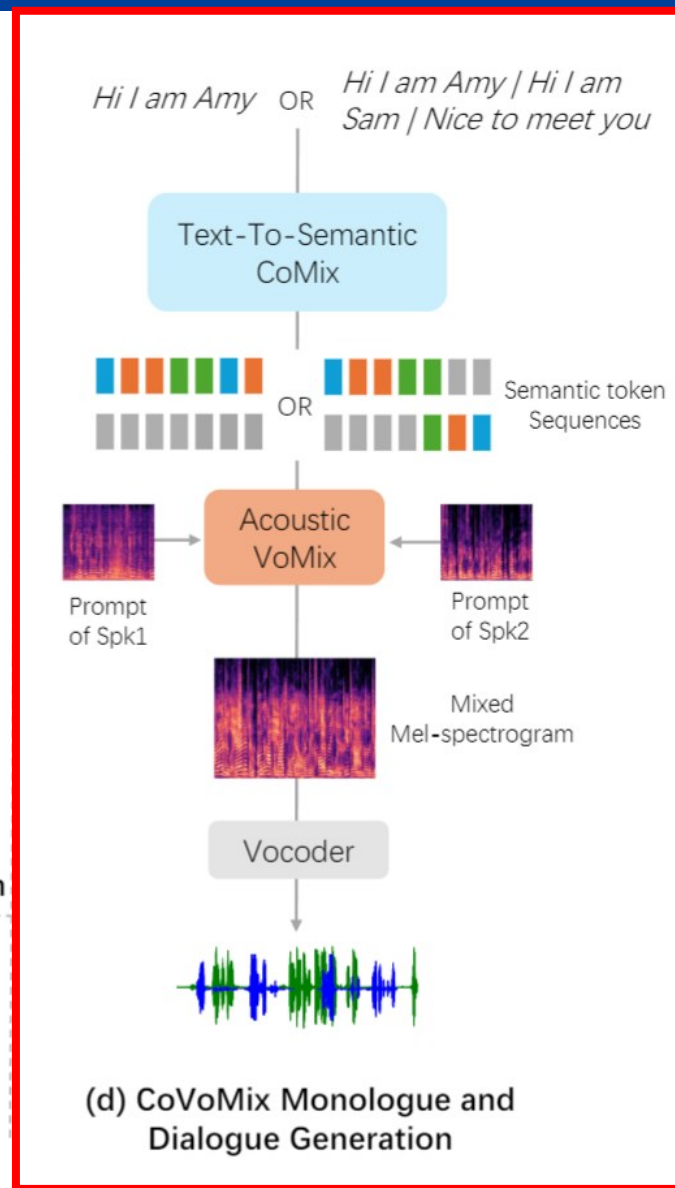
# Pipeline Comparison



(a) Conventional Monologue Generation (b) CoVoSingle Monologue Generation



(c) Conventional and CoVoSingle Dialogue Generation

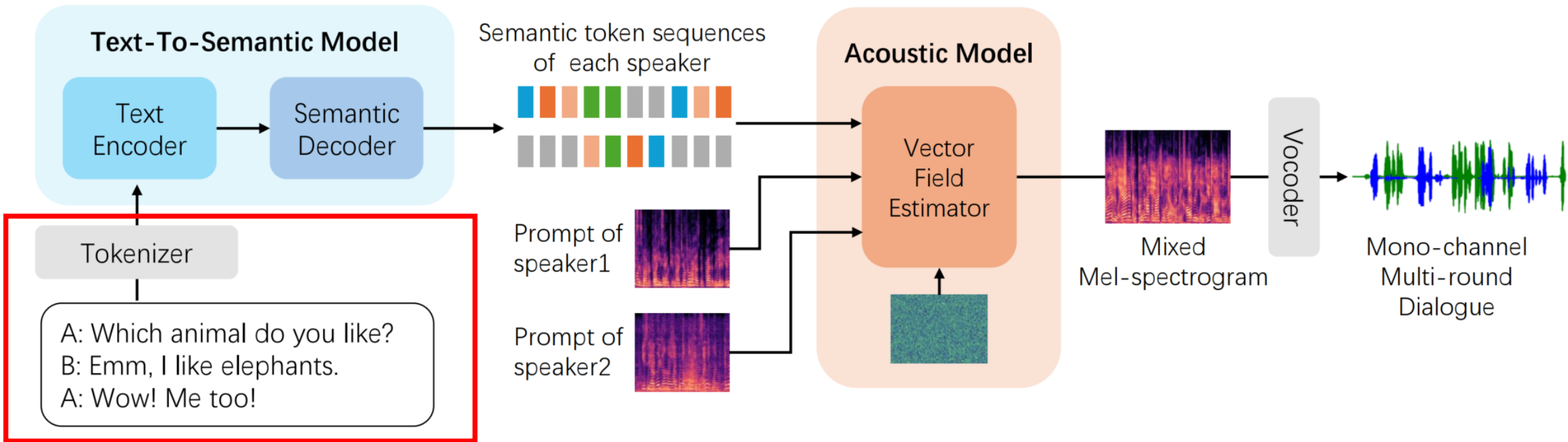


(d) CoVoMix Monologue and Dialogue Generation



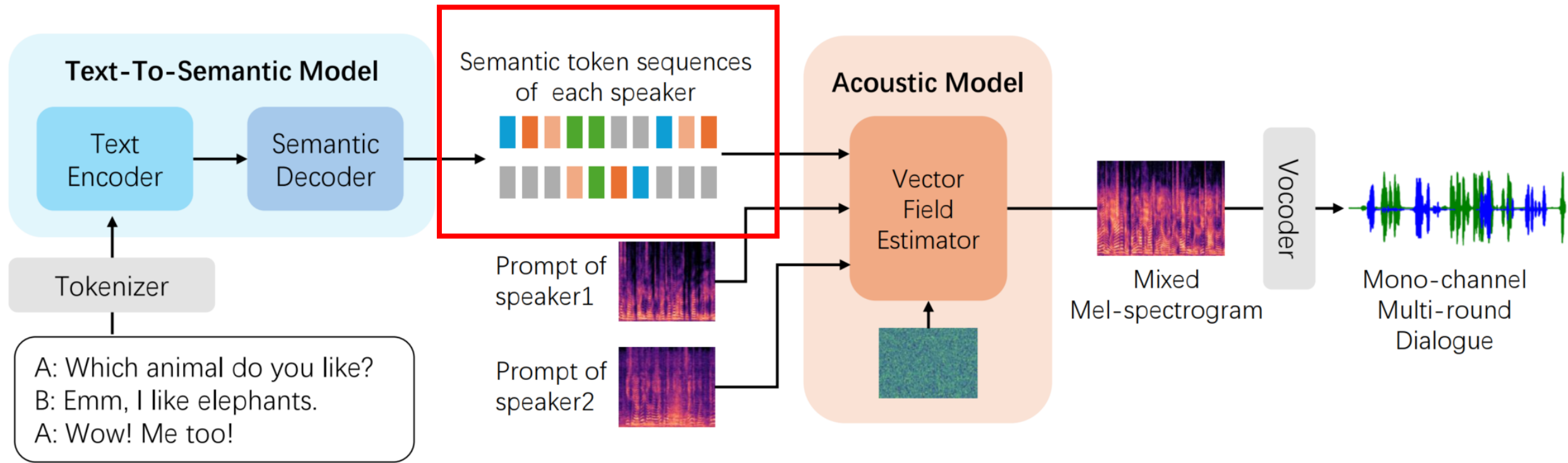


# CoVoMix



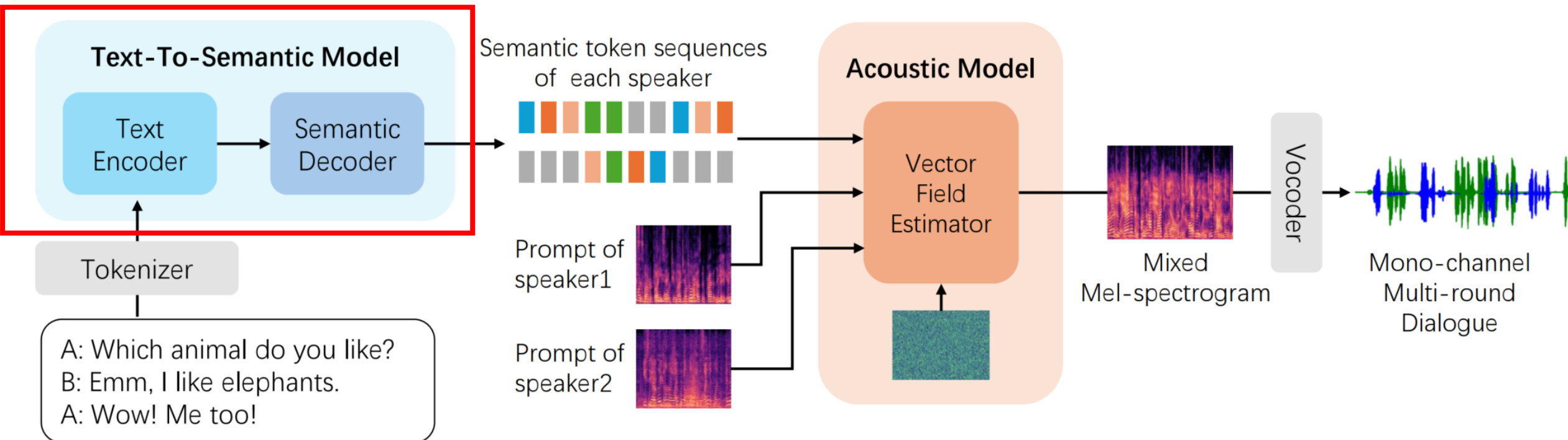


# CoVoMix



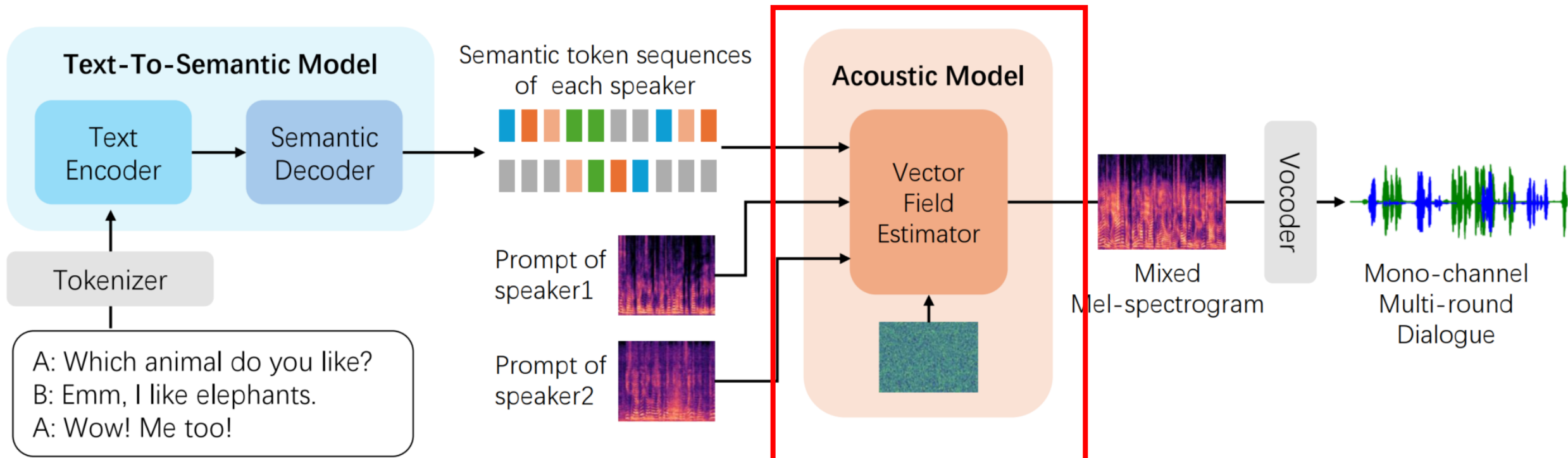


# CoVoMix





# CoVoMix



# Outline

1. Background and Motivation

2. Design of CoVoMix

**3. Data Processing**

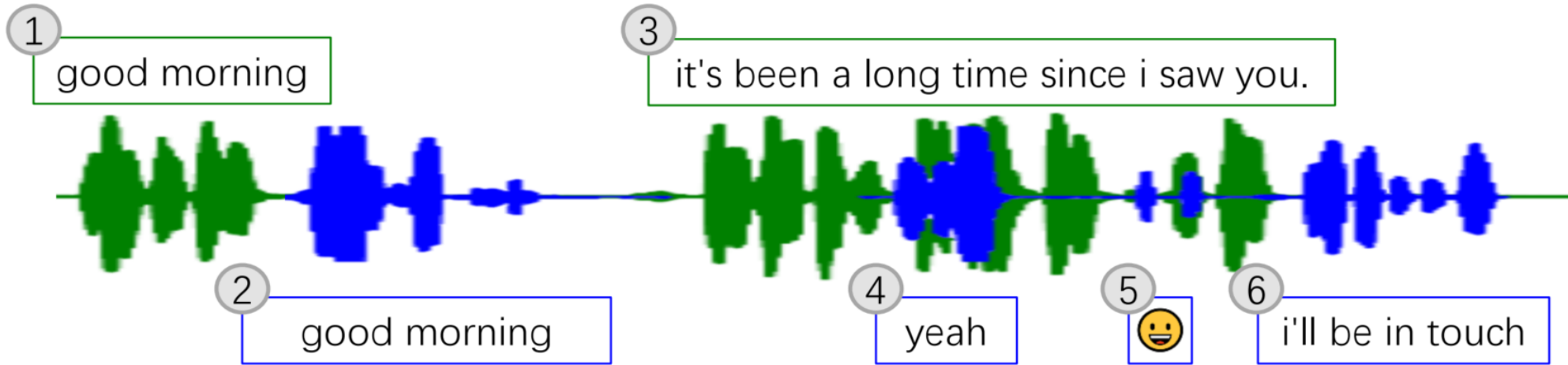
4. Evaluation Results

5. Conclusion





# Data Processing



good morning | good morning | it's been a long time since i saw you | yeah 😊 i'll be in touch

# Outline

1. Background and Motivation
2. Design of CoVoMix
3. Data Processing
- 4. Evaluation Results**
5. Conclusion





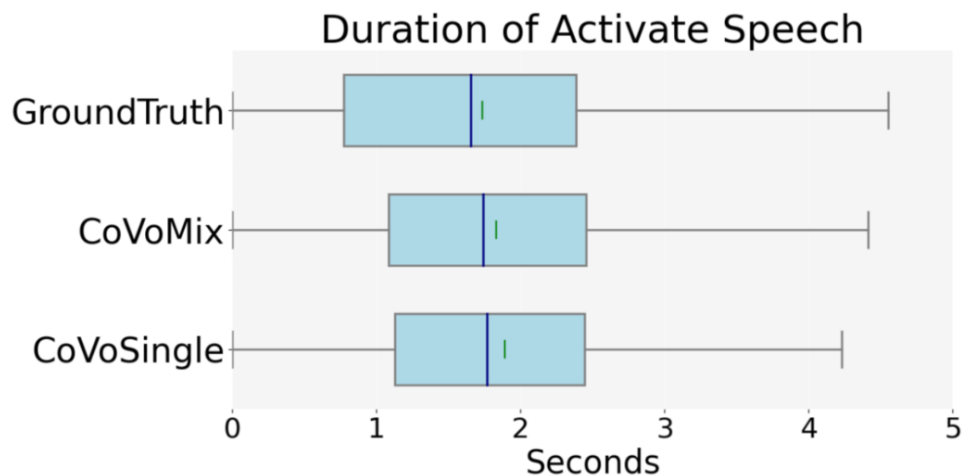
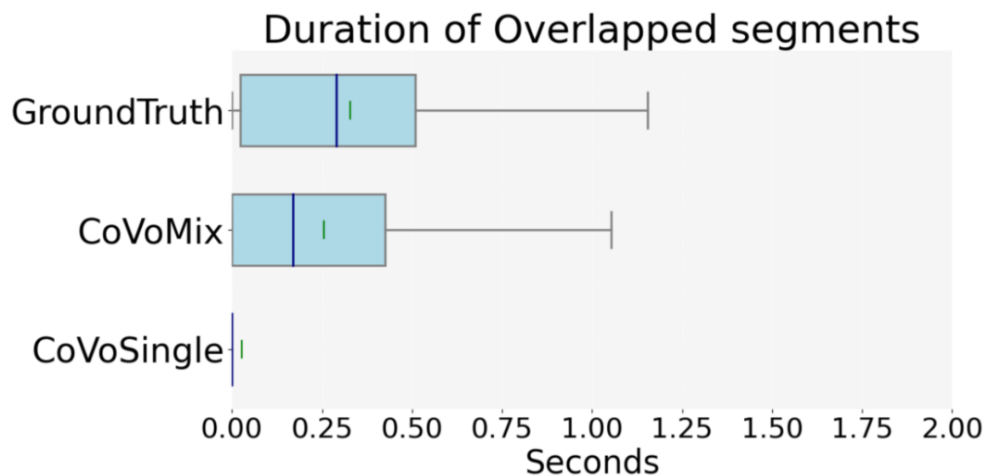
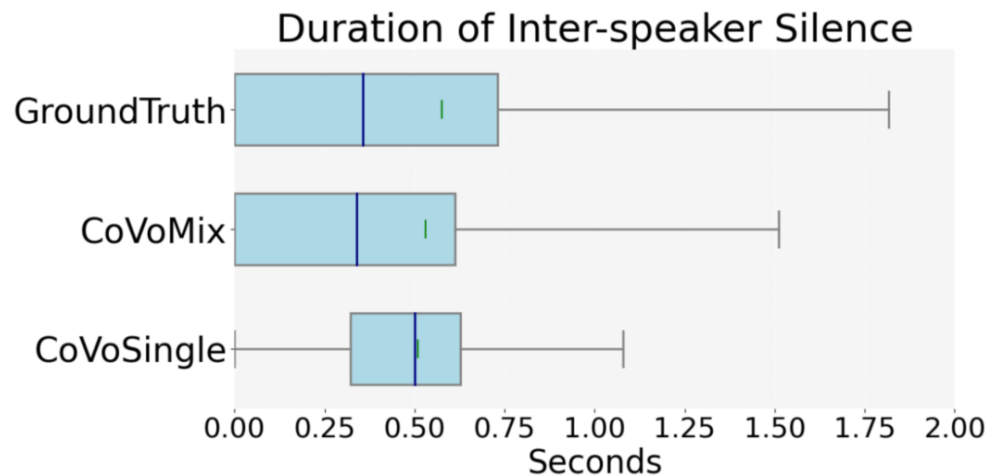
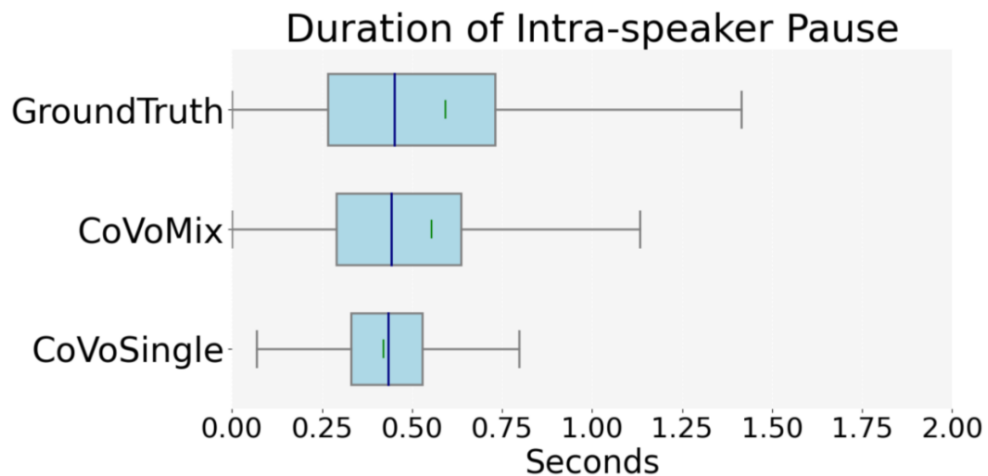
# Results for Monologue and Dialogue Generation

Eval Set	System	SIM $\uparrow$	WER $\downarrow$	MCD $\downarrow$	NISQA $\uparrow$	CMOS $\uparrow$	SMOS $\uparrow$
Monologue	GroundTruth	0.59	6.10	/	3.03	/	/
	Baseline	0.42	15.85	9.45	2.93	-1.60 $\dagger$	-1.18 $\dagger$
	CoVoSingle	0.49	9.99	6.15	3.04	0.00	0.00
	CoVoMix	0.49	8.95	6.04	3.01	0.83 $\dagger$	0.11
Dialogue	GroundTruth	/	14.91	/	2.73	/	/
	CoVoSingle	/	11.77	6.91	2.90	0.00	0.00
	CoVoMix	/	19.84	6.82	2.87	0.81 $\dagger$	0.60 $\dagger$



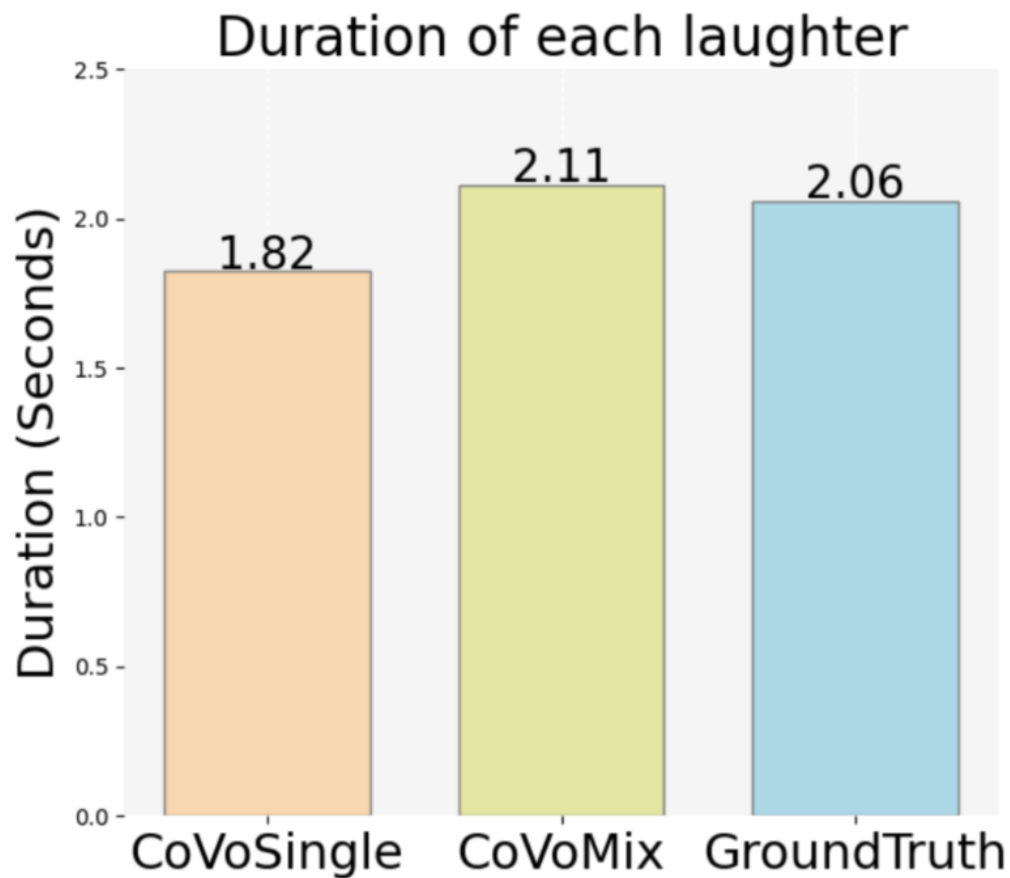
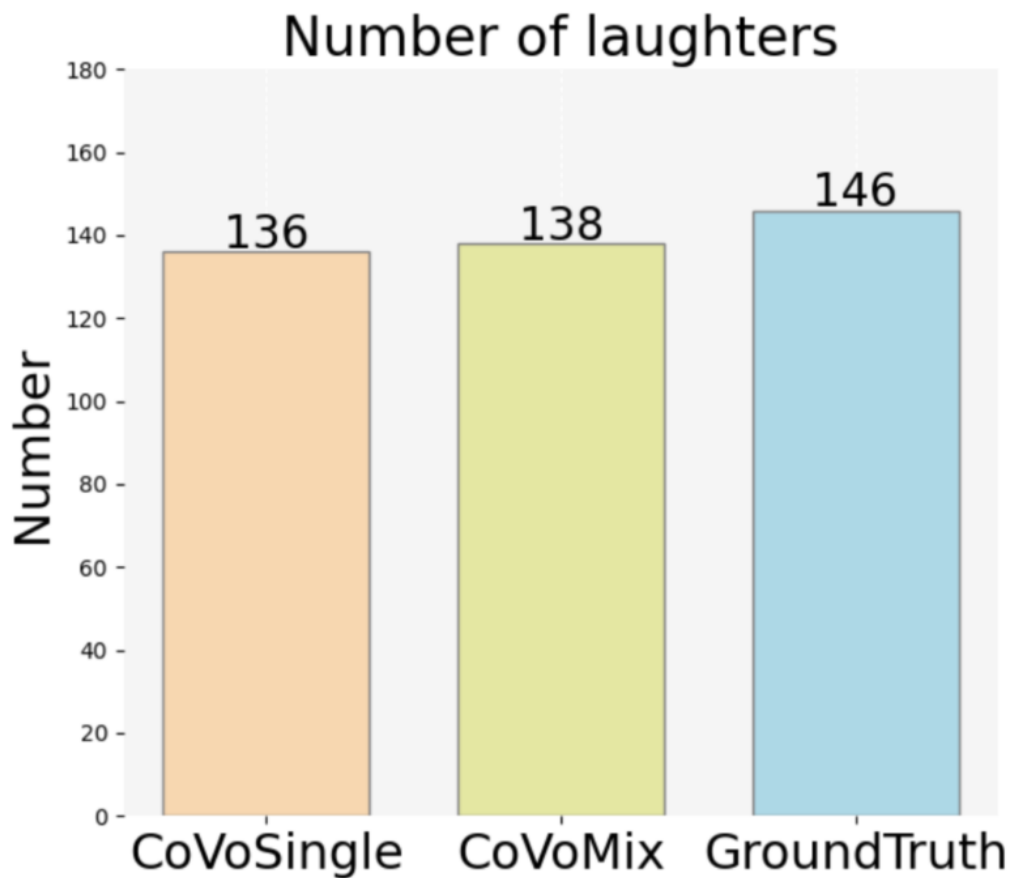


# Turn-taking Statistics for Dialogue Generation





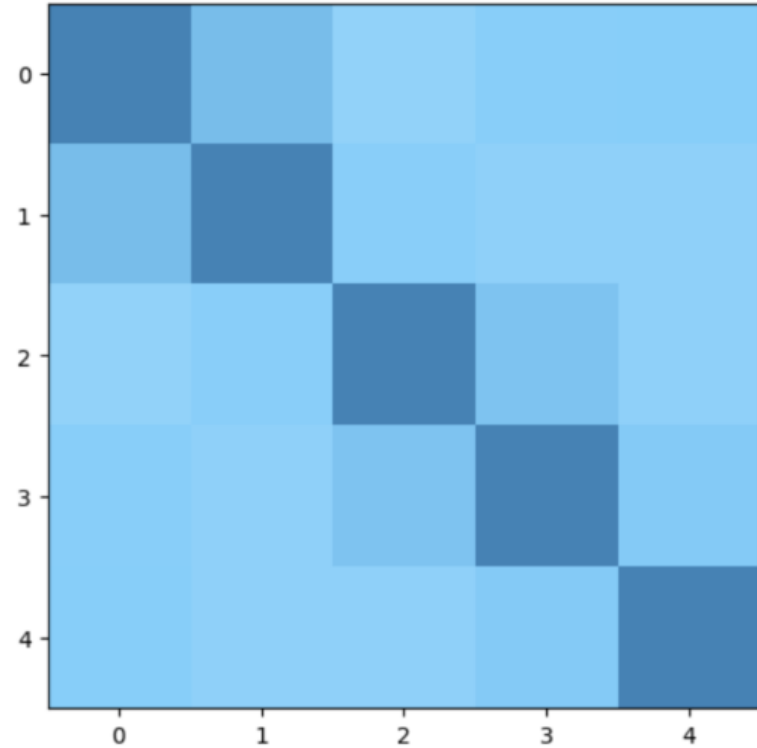
# Laughter Statistics for Dialogue Generation



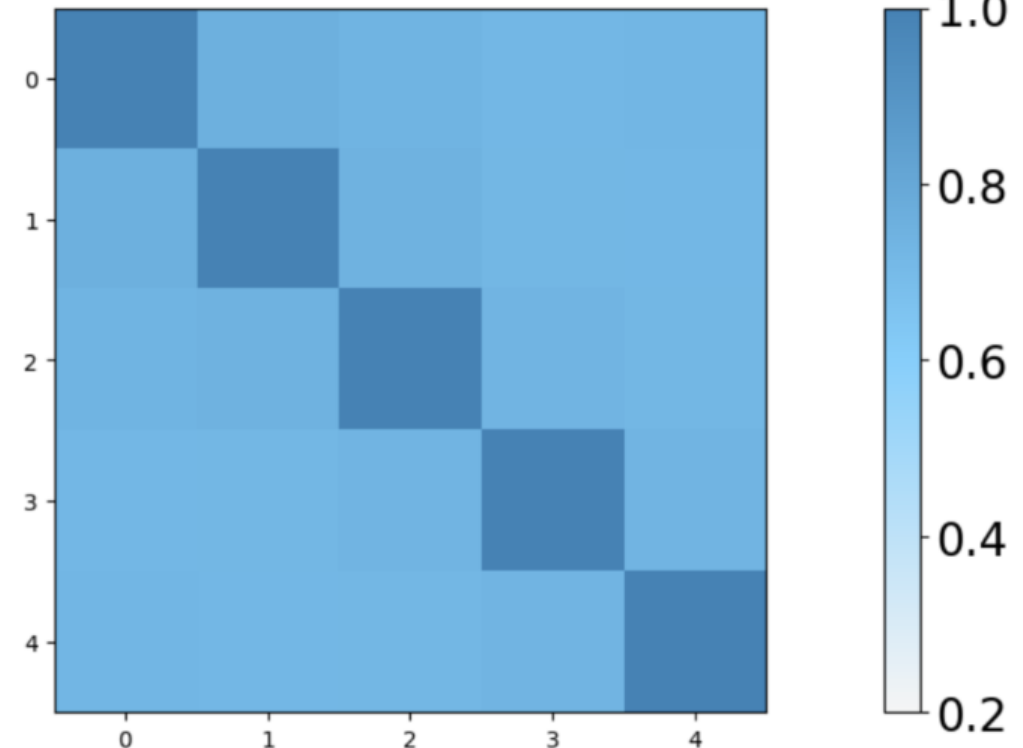


# Speech Consistency for Multi-round Dialogue

Utterance Level Concatenation



No Concatenation



# Outline

1. Background and Motivation

2. Design of CoVoMix

3. Data Processing

4. Evaluation Results

**5. Conclusion**





# Conclusion

- CoVoMix is composed of an auto-regressive text-to-semantic model and a flow-matching based acoustic model, with semantic token sequence as an intermediate representation
- CoVoMix achieves high naturalness and zero-shot speaker similarity in both monologue and dialogue generations
- CoVoMix demonstrates its proficiency in the fluency of dialogue turn-taking and spontaneous behavior generation

# *Thank You!*

## CoVoMix: Advancing Zero-Shot Speech Generation for Human-like Multi-talker Conversations

Discussion: Leying Zhang (Presenter), [zhangleying@sjtu.edu.cn](mailto:zhangleying@sjtu.edu.cn)



上海交通大學

SHANGHAI JIAO TONG UNIVERSITY

