

Generating Speakers by Prompting Listener Impressions for Pre-trained Multi-Speaker Text-to-Speech Systems



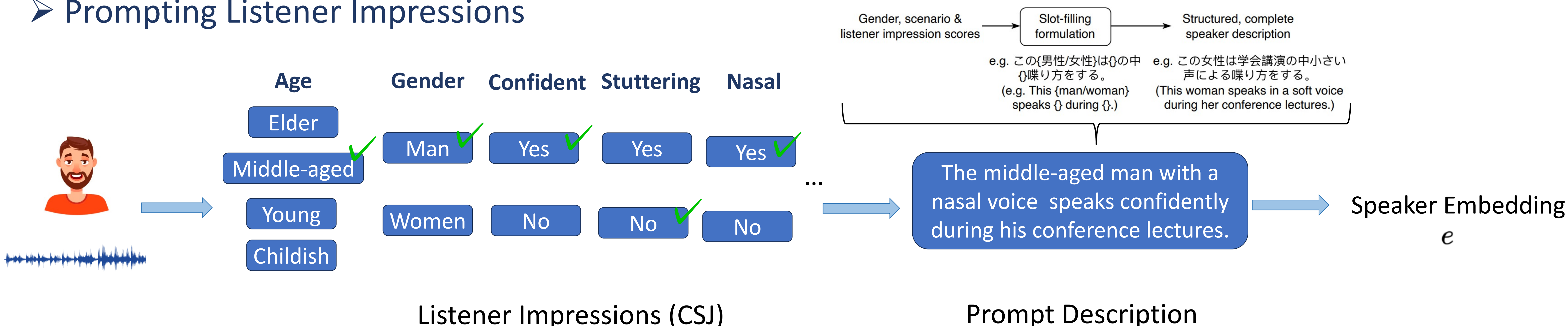
Zhengyang Chen^{1,2}, Xuechen Liu², Erica Cooper²,
Junichi Yamagishi², Yanmin Qian¹



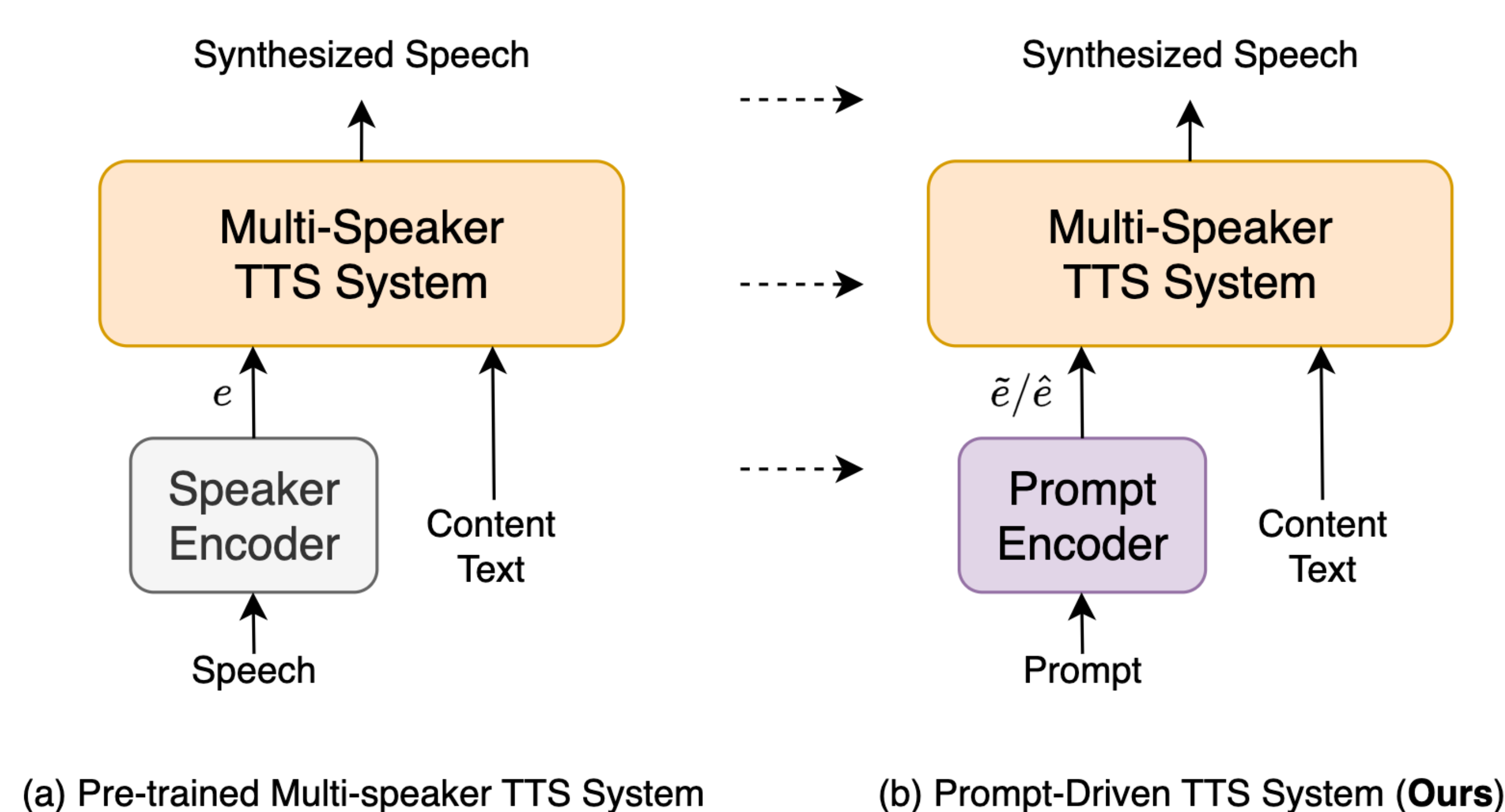
¹AudioCC Lab, CS Dept, Shanghai Jiao Tong University, Shanghai, China

²National Institute of Informatics, Tokyo, Japan

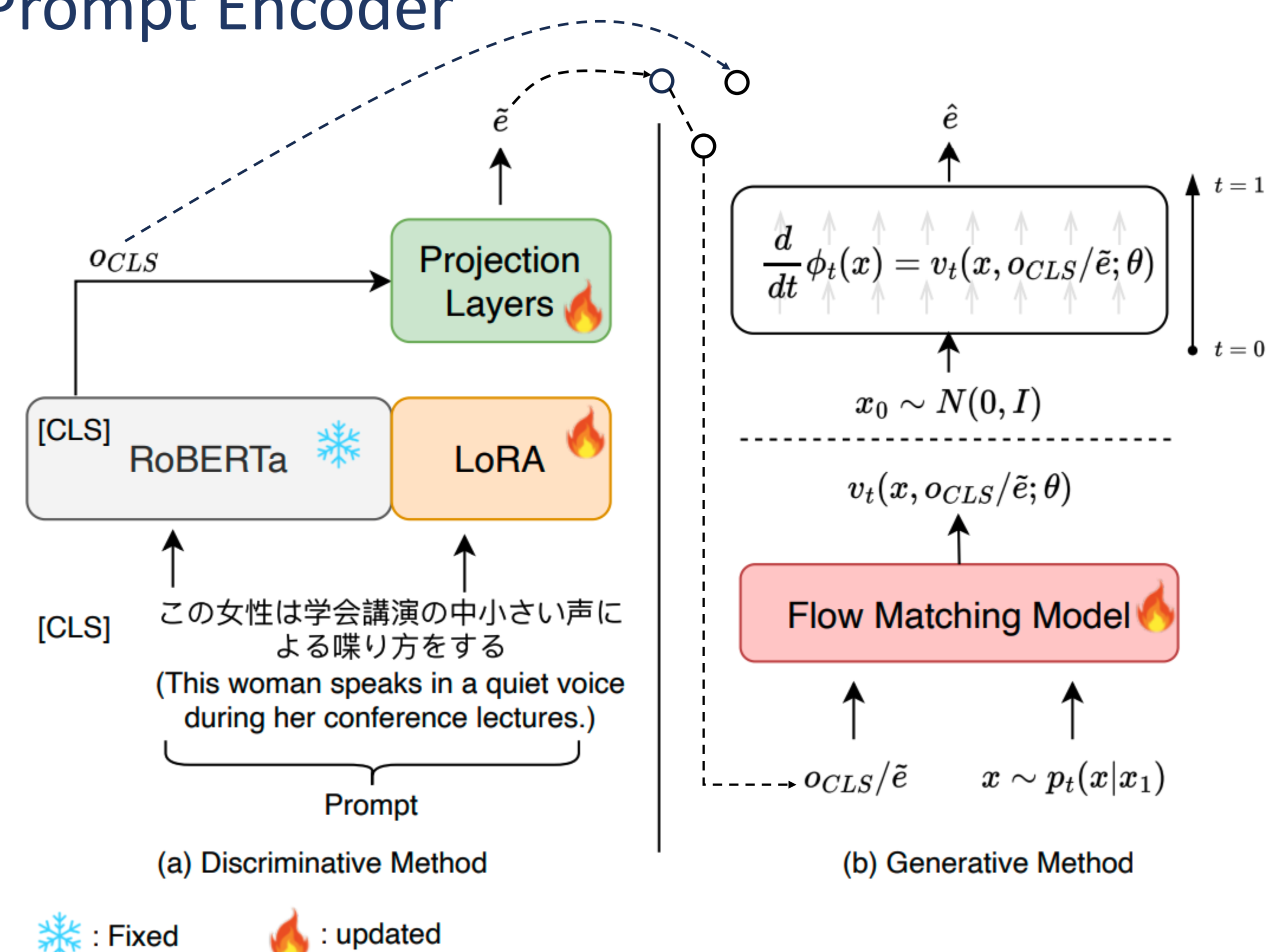
Prompting Listener Impressions



System Overview



Prompt Encoder



Experimental Setup

- Dataset: Japanese CSJ dataset
 - Speaker attributes #: 26
 - Train: 2672 speakers
 - Test: 30 speakers
- Objective Evaluation:
 - FAD score
- Subjective Evaluation:
 - Naturalness MOS
 - Speaker attribute impression MOS on a 5-point rating scale (each sample is rated 8 times by different raters)

Objective Evaluation and Naturalness Moss

Table 2: **FAD score and Naturalness MOS results on the CSJ evaluation set.**

System	FAD Score	Naturalness MOS
ground-truth	-	4.06 ± 0.25
Discriminative (w/o LoRA)	11.217	3.15 ± 0.25
Discriminative (w/ LoRA)	5.244	3.45 ± 0.19
Flow-Matching (w/ LoRA)	3.559	3.52 ± 0.26
Discriminative + Flow-Matching	3.126	3.50 ± 0.24

- The generative method significantly contributes to producing high-fidelity speech

Correlation between the Impression MOS from Reference Speech and Generated Speech

Table 1: **Spearman Rank Correlation Coefficient (SRCC) between MOS scores from reference and synthesized speech.**

Scenario	System	Speaker Attribute									
		expressiveness	confidence	relaxation	voice_depth	age	energy	pitch	speed	clarity	Avg
Seen	Discriminative (w/o LoRA)	<u>0.72</u>	0.53	0.48	<u>0.75</u>	<u>0.86</u>	<u>0.71</u>	<u>0.89</u>	<u>0.89</u>	0.23	0.67
	Discriminative (w/ LoRA)	<u>0.71</u>	<u>0.69</u>	0.65	<u>0.83</u>	<u>0.90</u>	<u>0.76</u>	<u>0.94</u>	<u>0.85</u>	0.37	0.74
	Flow-Matching (w/ LoRA)	<u>0.68</u>	0.53	0.66	<u>0.76</u>	<u>0.79</u>	0.50	<u>0.86</u>	0.38	0.22	0.60
	Discriminative + Flow-Matching	<u>0.74</u>	<u>0.71</u>	<u>0.75</u>	<u>0.87</u>	<u>0.96</u>	<u>0.72</u>	<u>0.90</u>	<u>0.68</u>	0.35	0.74
Unseen	Discriminative (w/o LoRA)	0.04	0.05	0.46	0.38	0.67	0.29	<u>0.73</u>	0.57	-0.37	0.31
	Discriminative (w/ LoRA)	0.54	0.38	0.49	0.48	<u>0.77</u>	0.25	<u>0.81</u>	0.36	0.41	0.50
	Flow-Matching (w/ LoRA)	-0.10	0.12	0.32	0.42	<u>0.82</u>	0.39	<u>0.74</u>	0.14	0.21	0.34
	Discriminative + Flow-Matching	0.36	0.08	0.49	0.35	<u>0.74</u>	0.34	<u>0.75</u>	0.37	0.20	0.41

underline: The statistical significance (p-value) is less than 0.001, indicating the MOS scores of synthetic speech are significantly correlated with the MOS scores of reference speech.

- The Discriminative method enables the model to better follow the prompt's description
- The generative method significantly contributes to producing high-fidelity speech
- Combining both methods can achieve good results in terms of both speech quality and adherence to the prompt.