

# NON-PARALLEL VOICE CONVERSION USING VARIATIONAL AUTOENCODERS CONDITIONED BY PHONETIC POSTERIORGRAMS AND D-VECTORS

Yuki Saito<sup>† ‡</sup>, Yusuke Ijima<sup>‡</sup>, Kyosuke Nishida<sup>‡</sup>, and Shinnosuke Takamichi<sup>†</sup>

<sup>†</sup>The University of Tokyo, Japan

<sup>‡</sup>NTT Corporation, Japan

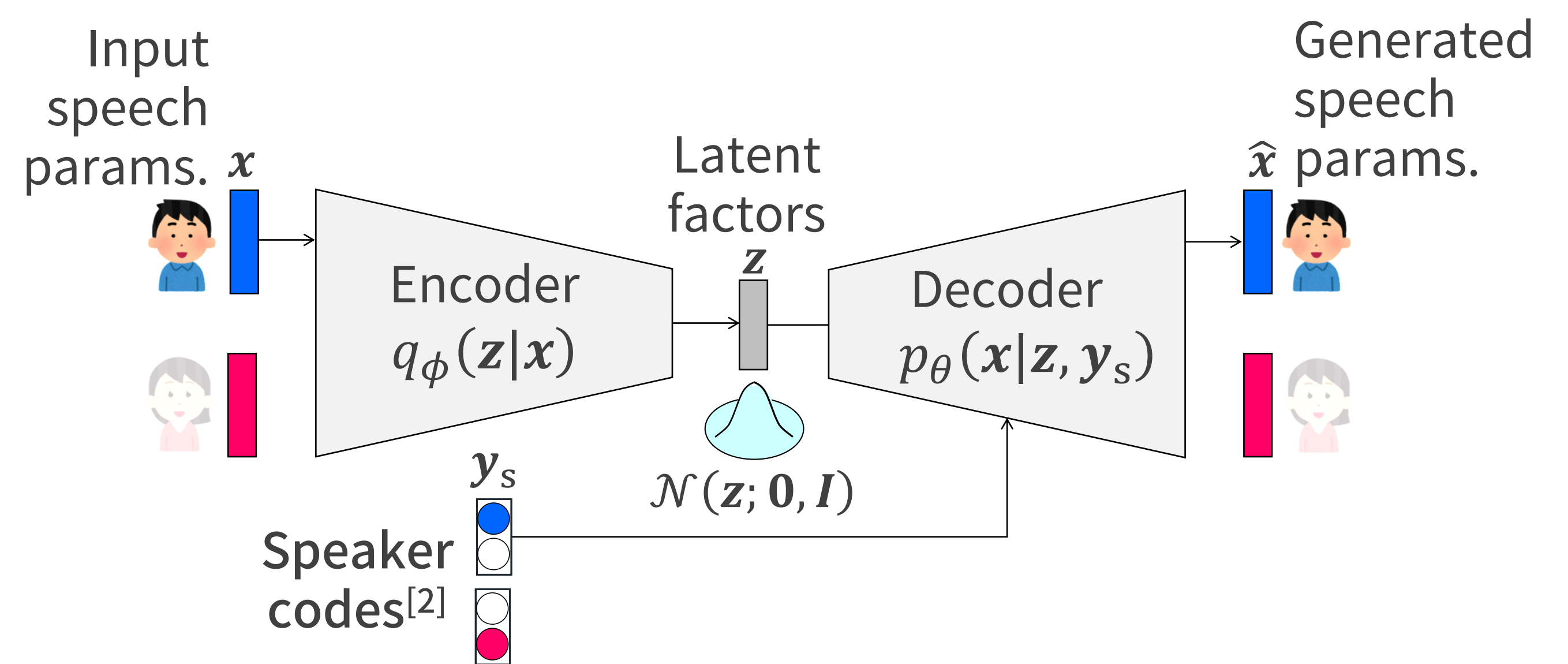
## 1. SYNOPSIS

Our approaches for non-parallel Voice Conversion (VC) using Variational AutoEncoders (VAEs):

- (1) Introduce Phonetic PosteriorGrams (PPGs) for dealing with speech quality degradation.
- (2) Extend conventional one-to-one VC to many-to-many VC (any speakers to any other speakers).

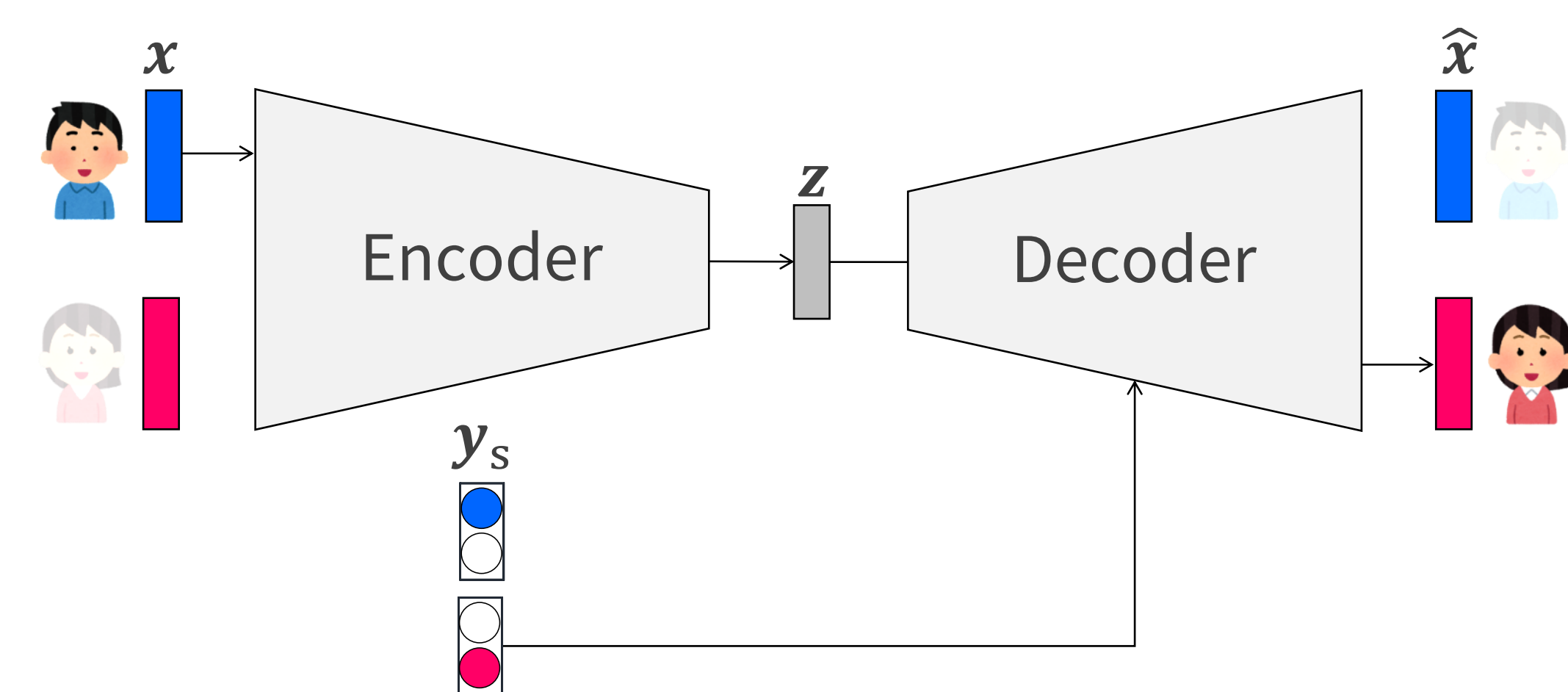
## 2. CONVENTIONAL VAE-BASED VC<sup>[1]</sup>

Training VAEs for VC



$$\mathcal{L}(\theta, \phi; x, y_s) = \underbrace{-D_{\text{KL}}(q_{\phi}(z|x) || \mathcal{N}(z; \mathbf{0}, \mathbf{I}))}_{\text{Regularization term of } z} + \underbrace{\mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z, y_s)]}_{\text{Reconstruction error of } x}$$

Conversion using trained VAEs



Assumption:  $z$  is independent of  $y_s$   
 $\Rightarrow$  Expected to represent phonetic contents.

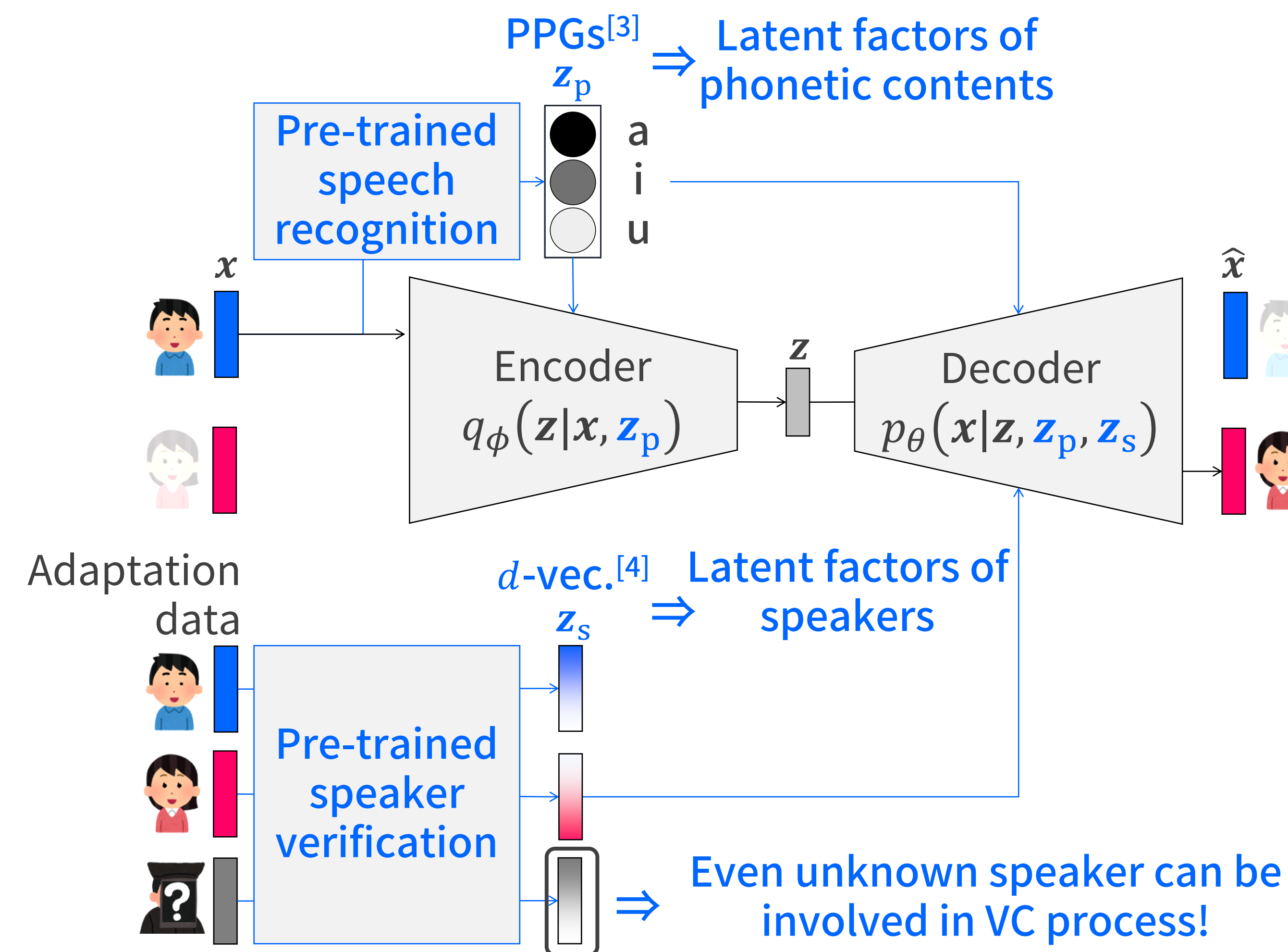
Problems

- (1) Converted speech quality is degraded due to vanishing phonetic properties.
- (2) It supports only one-to-one VC.

## 3. PROPOSED VAE-BASED VC

VC using VAEs w/ PPGs and  $d$ -vectors

- (1) Phonetic contents are given as PPGs<sup>[4]</sup>.  
 $\Rightarrow$  Phonetic contents can be restored!
- (2) Discrete speaker codes are replaced with continuous  $d$ -vectors<sup>[5]</sup>.  
 $\Rightarrow$  Any speakers' characteristics can be converted into any others!



Estimating unknown speaker representation

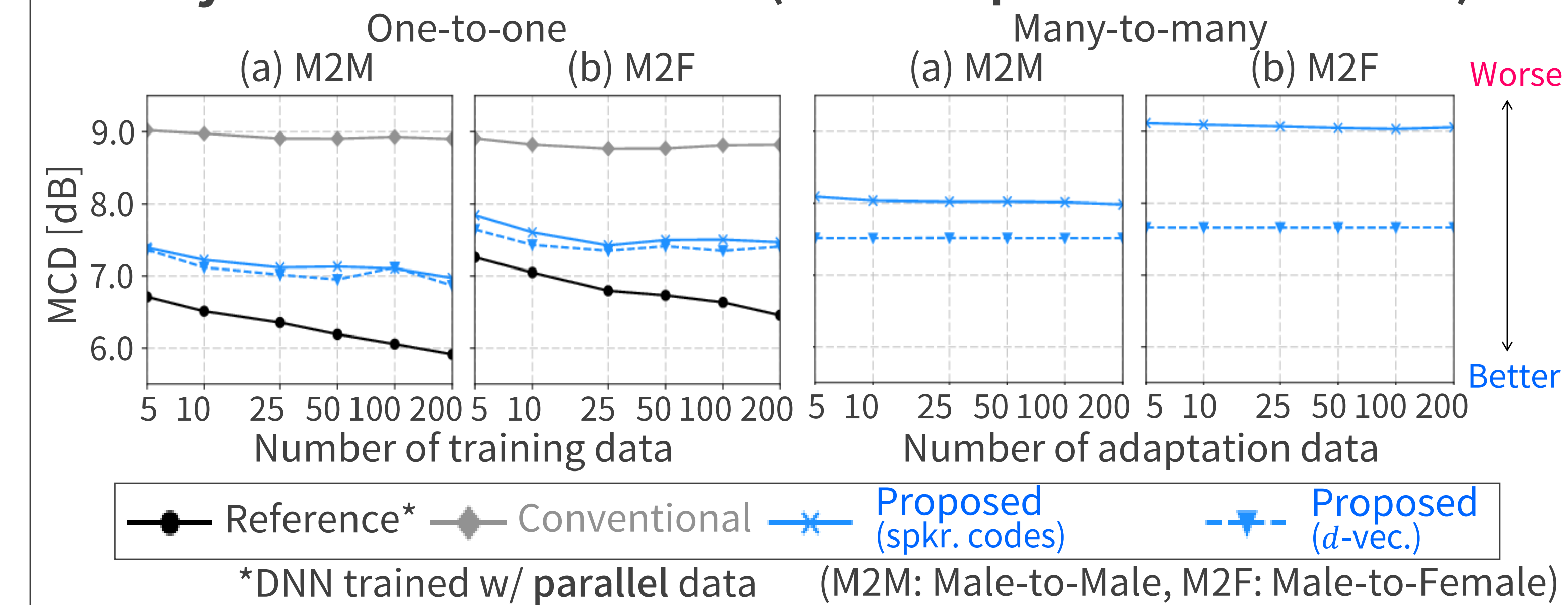
- (1) Speaker code adaptation<sup>[6]</sup>  
 Using backprop. to adapt speaker codes
- (2)  $d$ -vector averaging  
 Using averaged  $d$ -vec. in voiced region of the speaker's adaptation data

## 4. EXPERIMENTAL EVALUATION

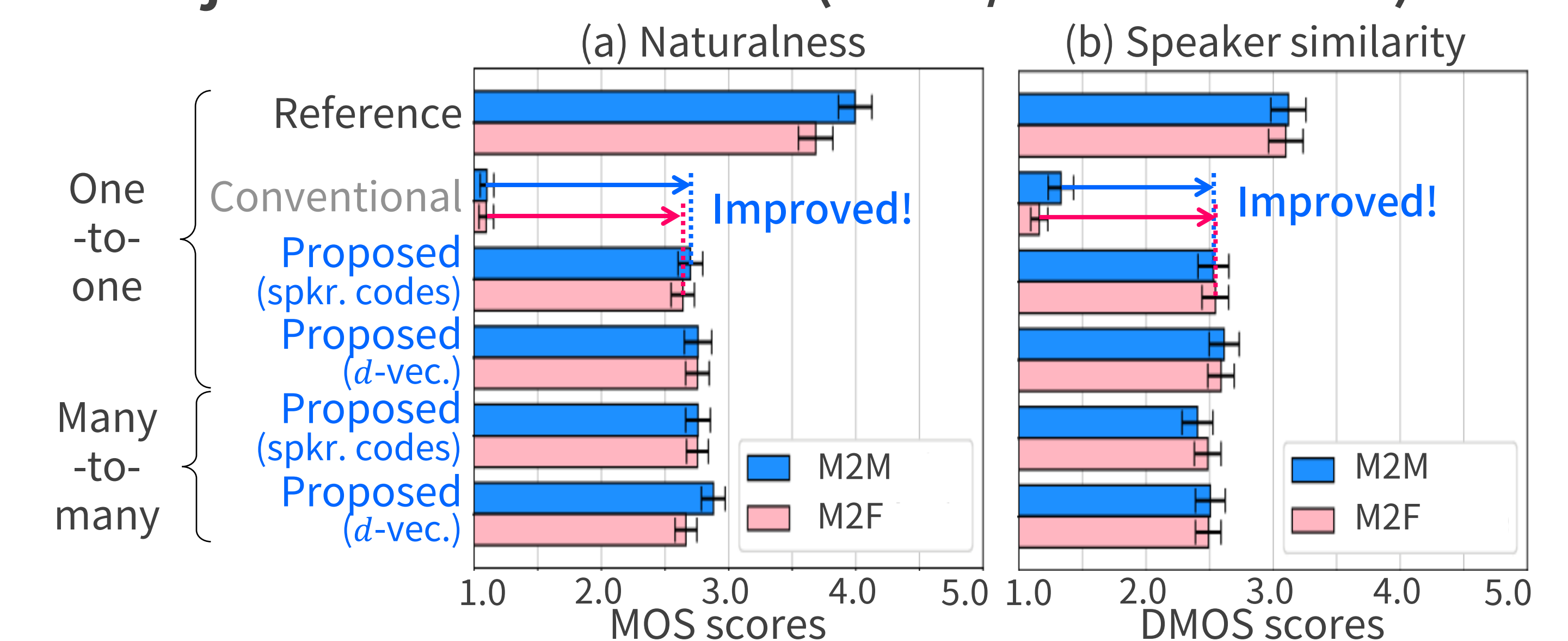
Experimental conditions

Speech corpora	(a) For recognition/verification models: including 260 speakers (130 male and 130 female speakers) (b) For VC models: divided parallel 425 utterances (2 male and 1 female speakers)
Speech params. (including $\Delta$ & $\Delta\Delta$ )	40-dim. mel-cepstral coefficients, Log F0, 10-dim. bap
DNN architecture	Feed-Forward (see our paper)
Dim. of PPGs/ $d$ -vec./ $z$	56 (time variant)/16 (time invariant)/64
VC settings	One-to-one: VC models are trained with corpora (b). Many-to-many: VC models are trained with corpora (a).

Objective evaluation (Mel-Cepstral Distortion)



Subjective evaluation (MOS/DMOS tests)



[References]

[1] Hsu et al., *Proc. APSIPA ASC*, 2016. [2] Hojo et al., *IEICE Trans. on Info. And Syst.*, 2017. [3] Dilokthanakul et al., *arXiv:1611.02648*, 2016. [4] Sun et al., *Proc. ICME*, 2016. [5] Variani et al., *Proc. ICASSP*, 2014. [6] Luong et al., *Proc. ICASSP*, 2017.