

FAST MVAE: Joint Separation and Classification of Mixed Sources based on Multichannel Variational Autoencoder with Auxiliary Classifier

Li Li¹, Hirokazu Kameoka², Shoji Makino¹

1. University of Tsukuba, 2. NTT Communication Science Laboratories, NTT Corporation

1. Introduction

2. Problem Formulation

◆ Multichannel Source Separation

- Separating out individual source signals from microphone array inputs without any prior information
- Important front-end process for numerous applications, e.g., ASR, VC

Research Objective

To reduce computational costs and improve performances of source label classification of MVAE without decreasing performances of source separation

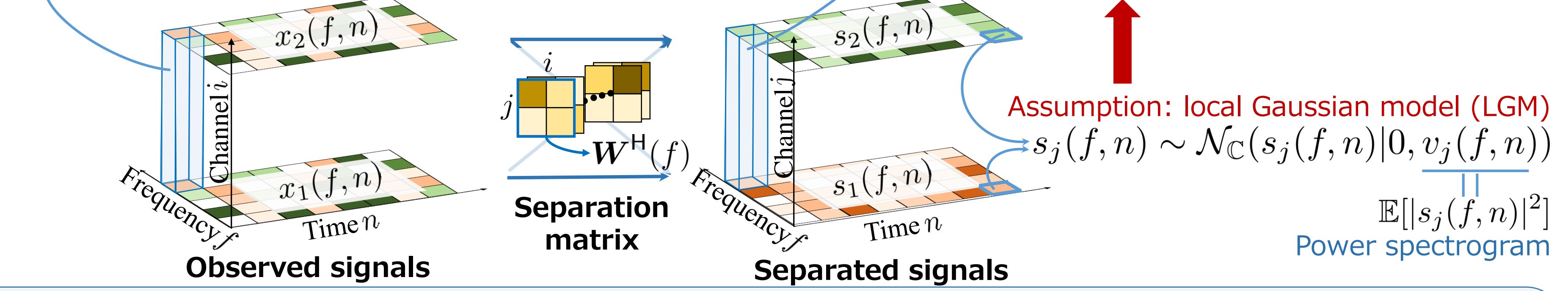
◆ Multichannel Variational Autoencoder (MVAE)

[Kameoka+2018]

- Recently proposed powerful approach to multichannel source separation that incorporates conditional variational autoencoder (CVAE) into the spectrogram modeling to enhance the representation power of spectrograms of source signals
- Advantages:**
 - Significant improvement in source separation performances
 - Simultaneously performing source label classification
 - Convergence-guaranteed optimization algorithm
- Disadvantages:**
 - Time-consuming optimization process
 - Unsatisfactory classification performances

◆ Probability model ($i = j = 2$ case)

$$\mathbf{x}(f, n) \sim \mathcal{N}_C(\mathbf{x}(f, n) | \mathbf{0}, (\mathbf{W}^H(f))^{-1} \mathbf{V}(f, n) \mathbf{W}(f)^{-1}) \quad \mathbf{s}(f, n) \sim \mathcal{N}_C(\mathbf{s}(f, n) | \mathbf{0}, \mathbf{V}(f, n))$$



◆ Objective function (negative log-likelihood):

$$-\log \mathcal{L} \stackrel{c}{=} \sum_{f, n, j} \left(\log v_j(f, n) + \frac{w_j^H(f) \mathbf{x}(f, n)^2}{v_j(f, n)} \right) - 2N \sum_f \log |\det \mathbf{W}^H(f)|$$

Terms w.r.t source model Permutation problem Terms w.r.t separation matrix

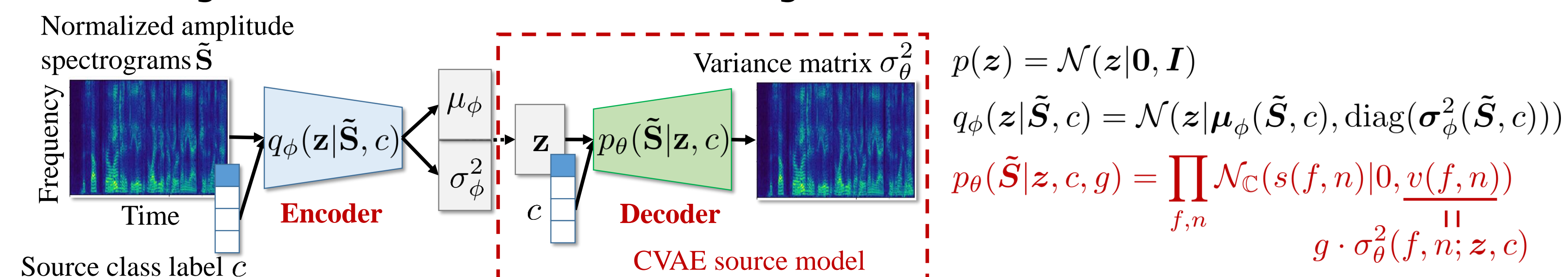
3. Conventional Methods

◆ ILRMA [Kameoka+2010, Kitamura+2016]

- Independent low-rank matrix factorization; determined multichannel nonnegative matrix factorization (NMF)
- Simultaneously solving source separation and permutation problem by capturing spectral structures of sources with NMF: $v_j(f, n) = \sum_k b_{j,k}(f) h_{j,k}(n)$
- Low-rank assumption of sources does not always hold

◆ MVAE [Kameoka+2018] **stronger representation power of spectrograms**

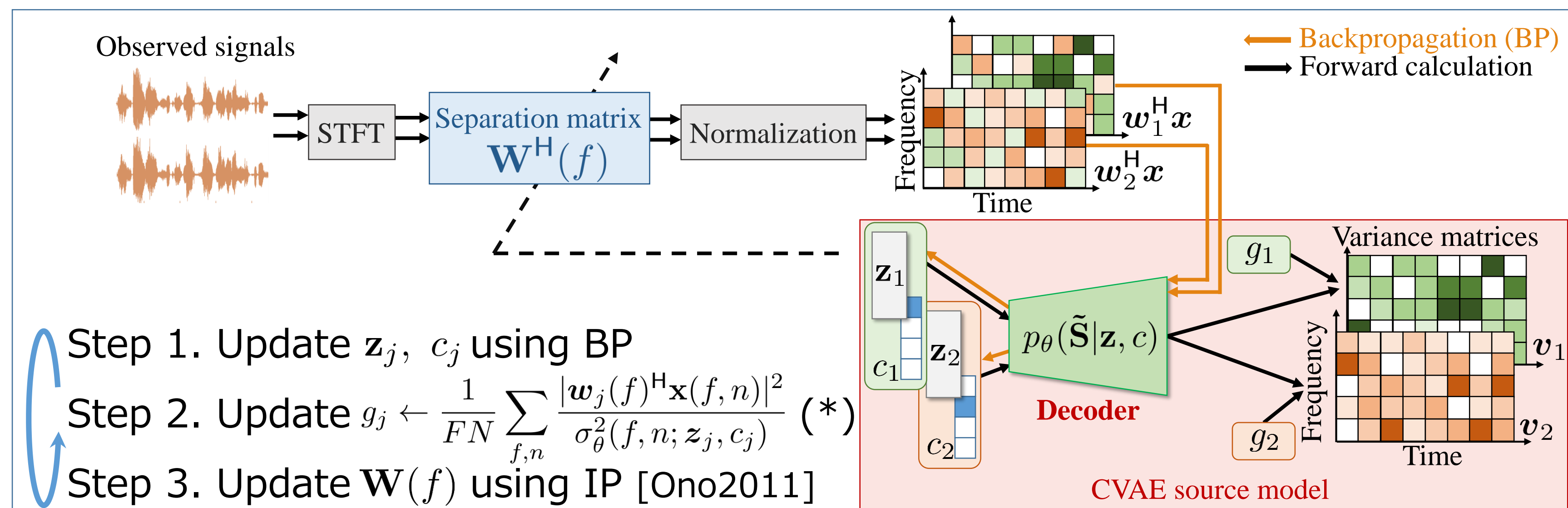
- Adopting CVAE to model complex spectrograms of sources with carefully designed decoder distribution having the same form with LGM



- Training loss function of CVAE:

$$\text{minimize } \mathcal{J}(\phi, \theta) = \underbrace{-\mathbb{E}_{(\tilde{\mathbf{S}}, c) \sim p(\tilde{\mathbf{S}}, c)} [\mathbb{E}_{\mathbf{z} \sim q_\phi(\mathbf{z} | \tilde{\mathbf{S}}, c)} [\log p_\theta(\tilde{\mathbf{S}} | \mathbf{z}, c)]]}_{\text{Reconstruction error}} + \underbrace{\text{KL}[q_\phi(\mathbf{z} | \tilde{\mathbf{S}}, c) || p(\mathbf{z})]}_{\text{Regularization}}$$

- Convergence-guaranteed optimization algorithm for separation
- BP process in each iteration is highly time-consuming



4. Proposed: Fast MVAE

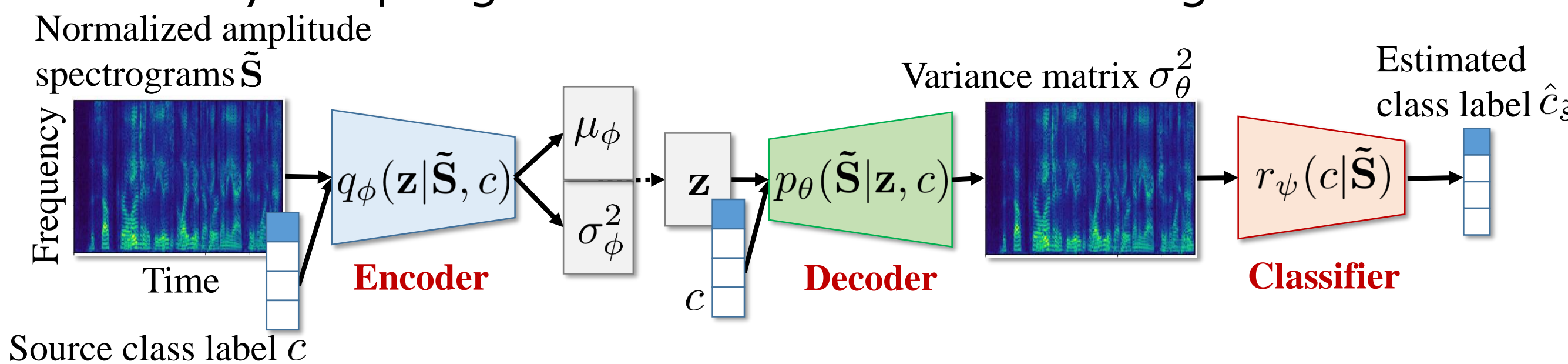
◆ Main idea

- Approximating the maximum of posterior distribution $p(\mathbf{z}_j, c_j | \tilde{\mathbf{S}}_j)$ searched by BP with product of two approximate distributions

$$p(\mathbf{z}_j, c_j | \tilde{\mathbf{S}}_j) = p(\mathbf{z}_j | \tilde{\mathbf{S}}_j, c_j) p(c_j | \tilde{\mathbf{S}}_j) \approx q_\phi(\mathbf{z}_j | \tilde{\mathbf{S}}_j, c_j) r_\psi(c_j | \tilde{\mathbf{S}}_j)$$

◆ Auxiliary classifier VAE (ACVAE) [Chen+2016, Kameoka+2018]

- Enhancing effect of class label on controlling the generative model by adopting an information-theoretic regularization

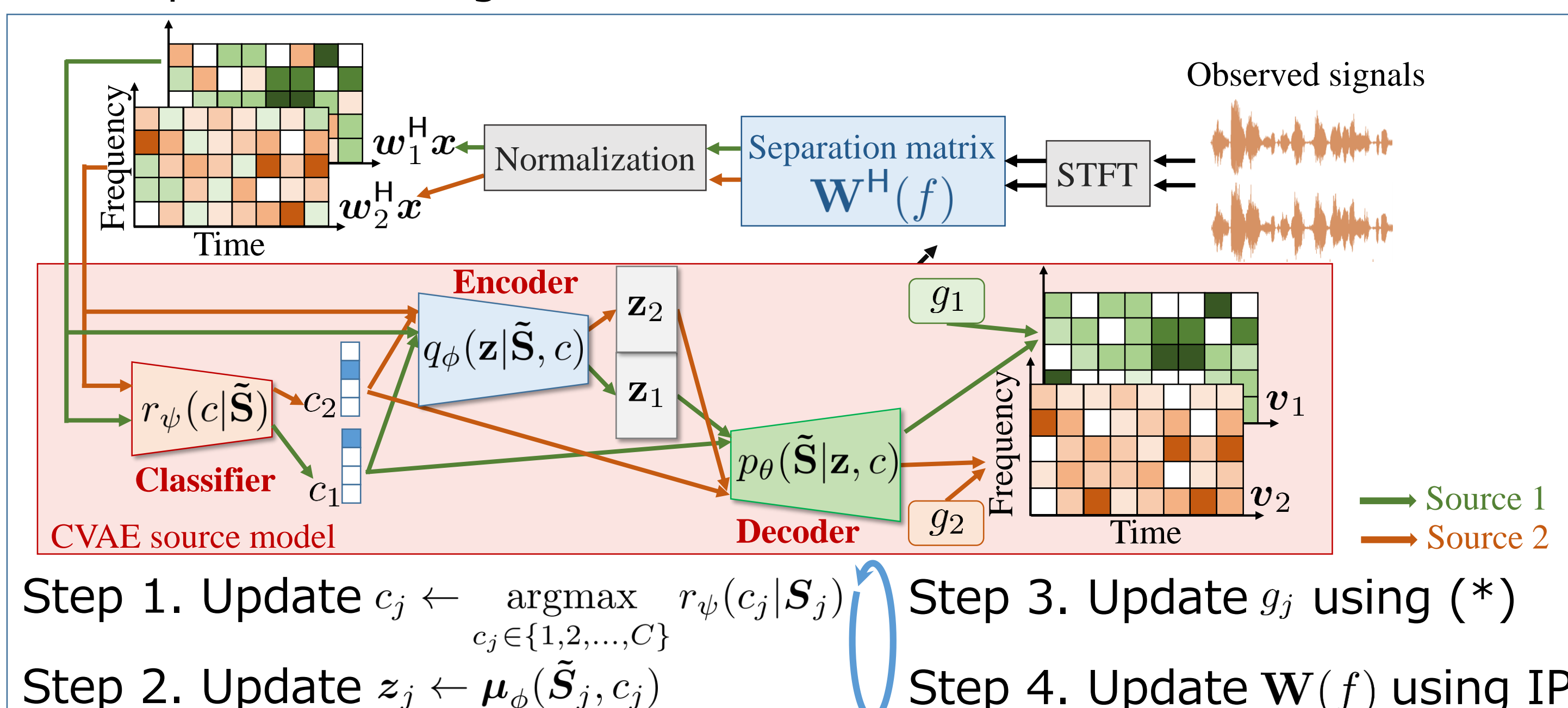


- Loss function of training ACVAE consisting of term of CVAE loss function, regularization on reconstructed sources and regularization on training samples

$$\text{minimize } \mathcal{J}(\phi, \theta) - \lambda_L \mathcal{L}(\phi, \theta, \psi) - \lambda_T \mathcal{I}(\psi) \quad \lambda_L \geq 0 \quad \lambda_T \geq 0$$

$$\mathbb{E}_{(\tilde{\mathbf{S}}, c) \sim p_D(\tilde{\mathbf{S}}, c), q_\phi(\mathbf{z} | \tilde{\mathbf{S}}, c)} [\mathbb{E}_{c \sim p(c), \tilde{\mathbf{S}} \sim p_\theta(\tilde{\mathbf{S}} | \mathbf{z}, c)} [\log r_\psi(c | \tilde{\mathbf{S}})]] - \mathbb{E}_{(\tilde{\mathbf{S}}, c) \sim p_D(\tilde{\mathbf{S}}, c)} [\log r_\psi(c | \tilde{\mathbf{S}})]$$

- Optimization algorithm with forward calculations:



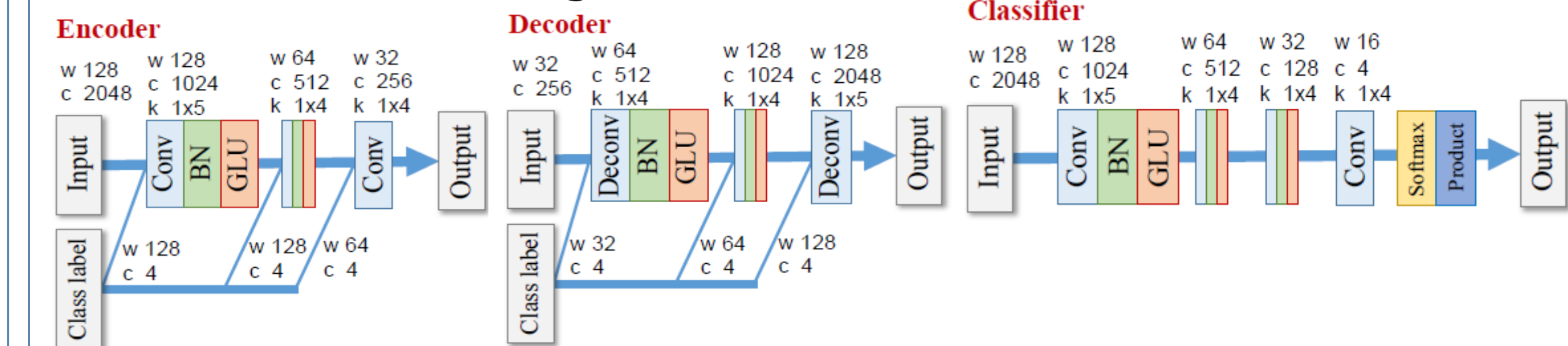
5. Experiments

◆ Data

- 2 male and 2 female excerpted from VCC2018
- Training dataset: 5 mins/speaker
- Test dataset: 2 mins/speaker, 4 speaker combinations, totally 40 utterances
- RT60: 78 ms, 351 ms (simulated RIR) 173 ms, 225 ms (measured RIR)
- Sampling rate: 16 kHz
- Window length/shift: 256 ms / 128 ms

◆ Network architectures

- 1-dimensional gated convolutional neural network



◆ Results

method	SDR	SIR	SAR	method	iter	total
ILRMA	9.01	15.50	12.00	MVAE (GPU)	6.08	202.12
MVAE	13.30	21.22	15.44	fMVAE (CPU)	0.36	20.57
fMVAE	13.98	22.14	15.65	fMVAE (GPU)	0.07	10.90
				ILRMA (CPU)	0.12	12.68

↑ SDR, SIR, SAR [dB] obtained with each method

method	accuracy rate
MVAE	40.63%
fMVAE	78.75%

← Accuracy rate of source classification

- fMVAE achieved comparative source separation performance to MVAE and significantly outperformed ILRMA, which showed the effect of nonlinear CVAE source model
- fMVAE was about 20 times faster than MVAE and achieved source classification accuracy rate of about 80%. (Audio files are available)