

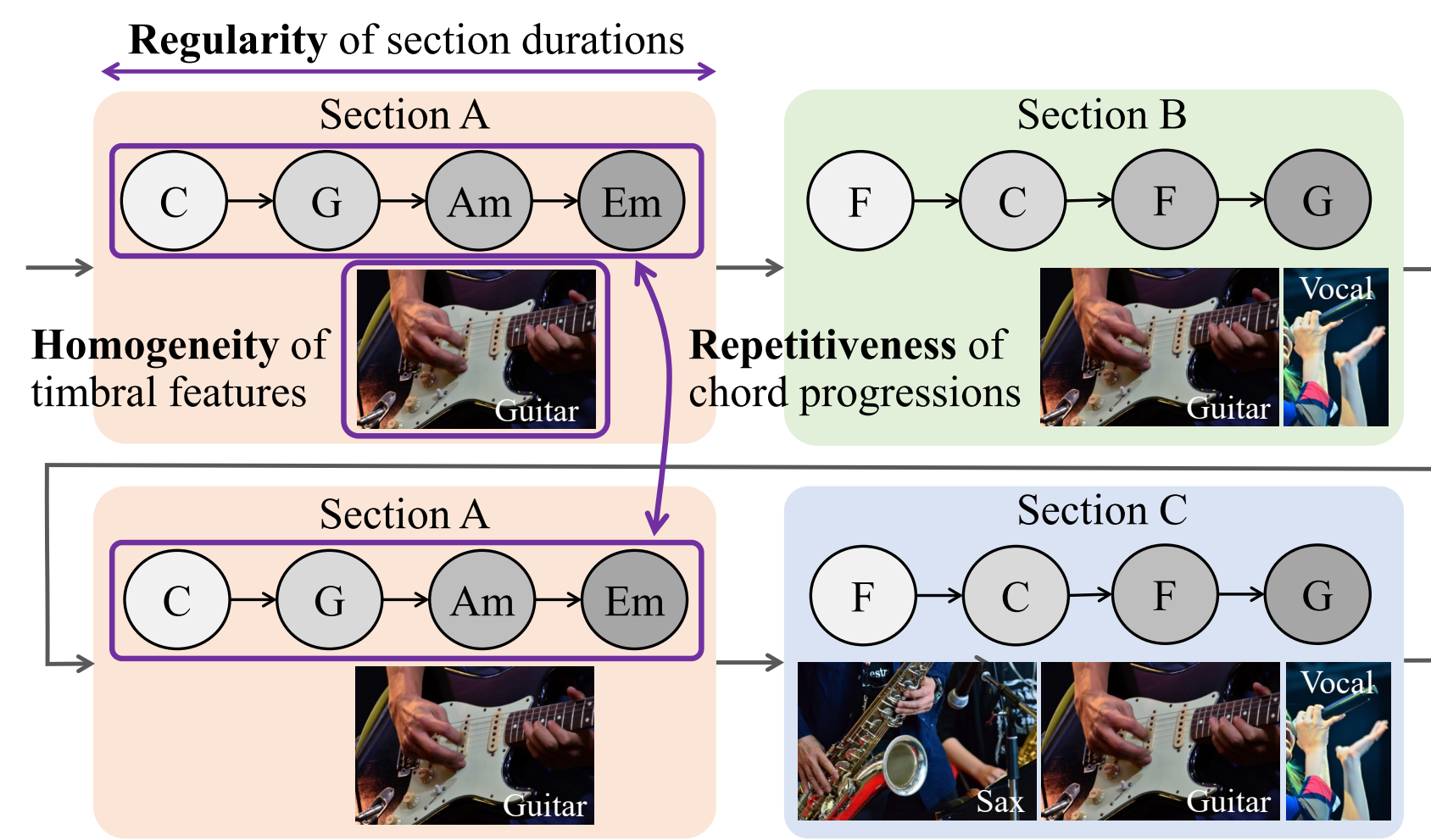
# Music Structure Analysis Based on an LSTM-HSMM Hybrid Model

Go Shibata<sup>1</sup> Ryo Nishikimi<sup>1</sup> Kazuyoshi Yoshii<sup>1</sup>

<sup>1</sup> Graduate School of Informatics, Kyoto University, Japan

## Abstract

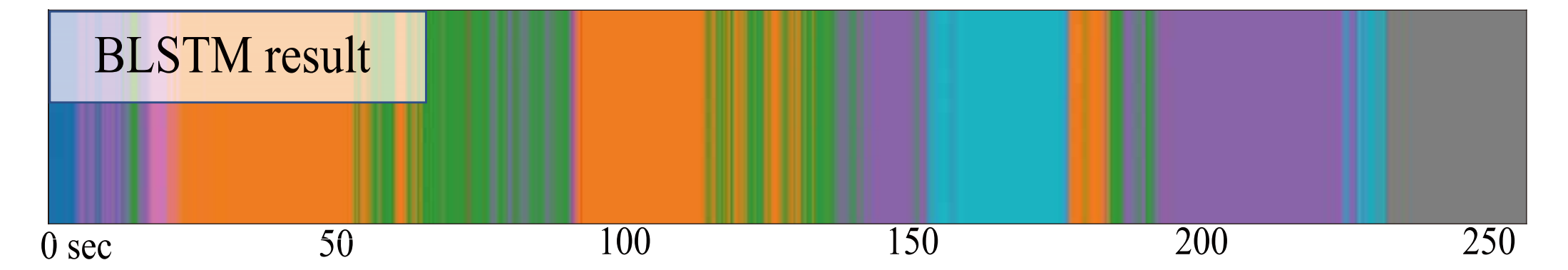
- Recognize meaningful musical sections
- Probabilistic formulation based on **music knowledge** about musical sections: **homogeneity, repetitiveness, and regularity**
- Emission probabilities of mel spectra computed by **bidirectional LSTM**
- Unsupervised learning based on Bayesian inference



## Background

Supervised learning of a DNN **does not work well**

- Gives unnaturally-frequent label switching because of the lack of annotated training data
- Thus we need to make effective use of **music knowledge** about musical sections



## Model Formulation and Inference

### LSTM-HSMM Hybrid Model

The proposed method deals with segmentation and labeling simultaneously

### Model Formulation

#### Two-Level Hierarchical Markov Chains

- Upper-level: ergodic semi-Markov model

Generate a sequence of sections  $Z$  and their durations  $D$

$$\begin{aligned} p(z_1, d_1) &= \rho_{z_1} \psi_{d_1} \\ p(z_n, d_n | z_{n-1}, d_{n-1}) &= \pi_{z_{n-1}z_n} \psi_{d_n} \\ p(z_N, d_N | z_{N-1}, d_{N-1}) &= \pi_{z_{N-1}z_N} \psi_{d_N} v_{z_N} \end{aligned}$$

$\rho_z$  initial probability of sections  
 $\psi_d$  duration probability of sections  
 $\pi_{zz'}$  transition probability of sections  
 $v_z$  terminal probability of sections

- Lower-level: left-to-right Markov model

Generate a sequence of chords  $S$

$$p(s_{n,\tau} | z_n, s_{n,\tau-1}) = \phi_{s_{n,\tau-1}s_{n,\tau}}^{(z_n)}$$

$\phi_{ss'}^{(z)}$  transition probability of chords  
 $s_{n,1} = 1$   
 $\tau_1 < \tau_2 \Rightarrow s_{n,\tau_1} \leq s_{n,\tau_2}$

### Acoustic Model

Outputs chroma vectors, MFCCs, and mel spectra

- Chroma vectors

Describe the **repetitiveness** of chord sequences

Generated from Gaussian distributions corresponding to sections and chords

- MFCCs (Mel-Frequency Cepstrum Coefficients)

Describe the **homogeneity** of timbral features

Generated from Gaussian distributions corresponding to sections

- Mel spectra

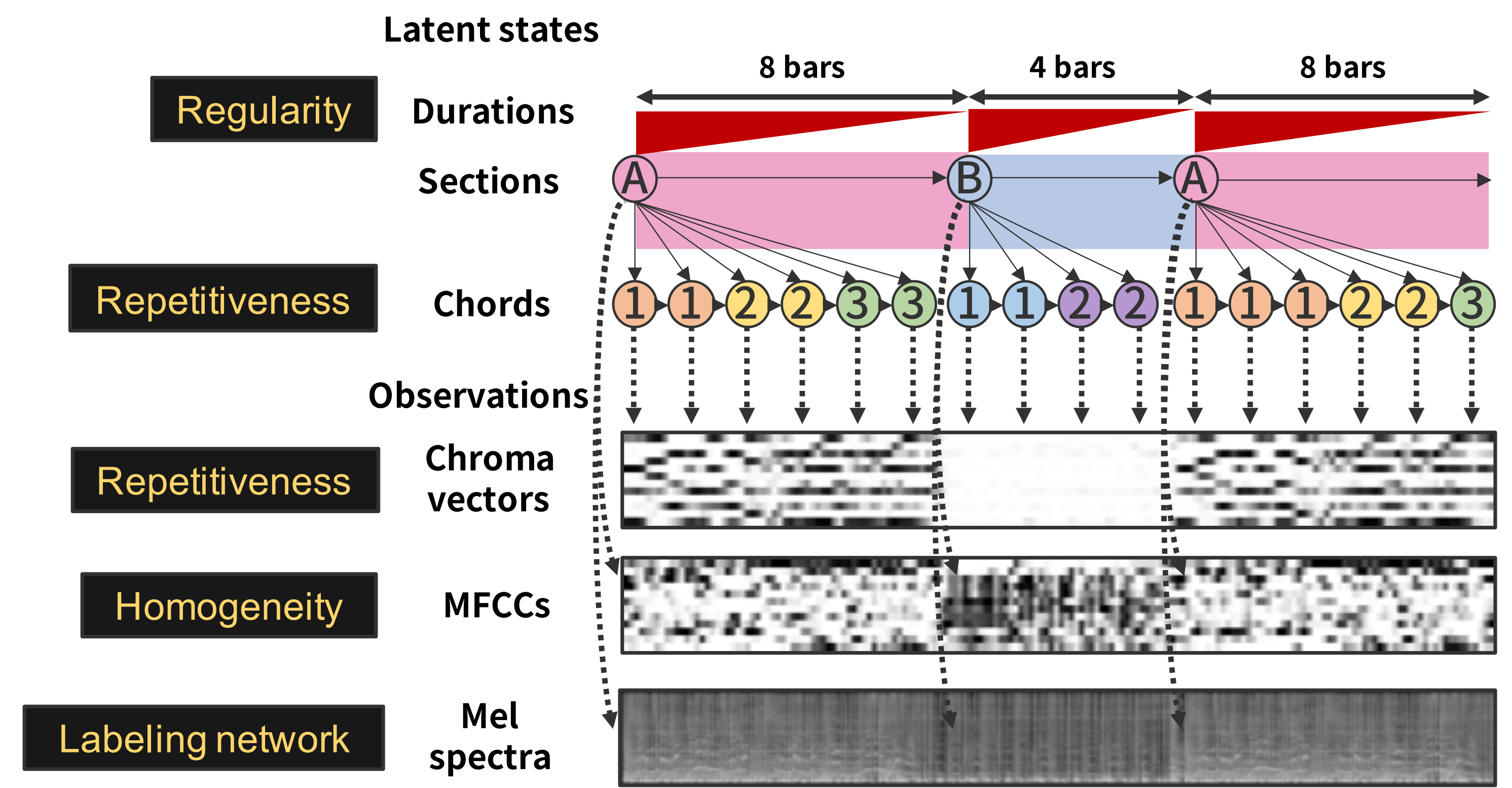
Associate sections with labels

Generated using probabilities based on a **bidirectional LSTM**

### Prior Distributions

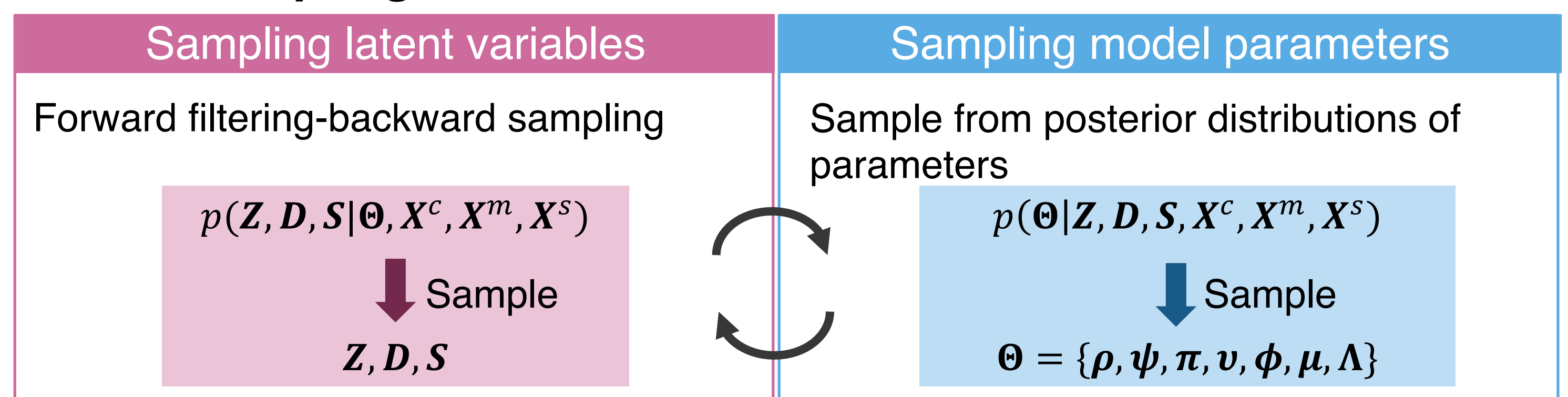
We put conjugate prior distributions for parameters of the model

- Degenerate unnecessary sections during the Bayesian sparse learning
- Incorporate the **regularity** of section durations as prior knowledge
- Use empirical distributions  $\alpha_{\text{emp}}^p, \alpha_{\text{emp}}^{\pi_z}, \alpha_{\text{emp}}^{\psi}$ , and  $\alpha_{\text{emp}}^v$  for the prior distributions

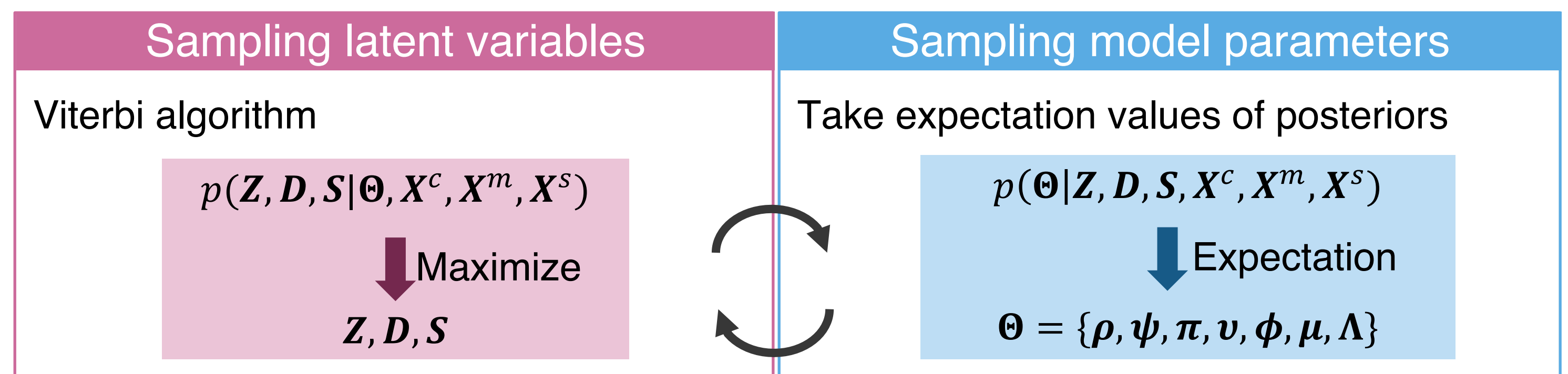


### Inference

#### Gibbs Sampling



#### Viterbi Training

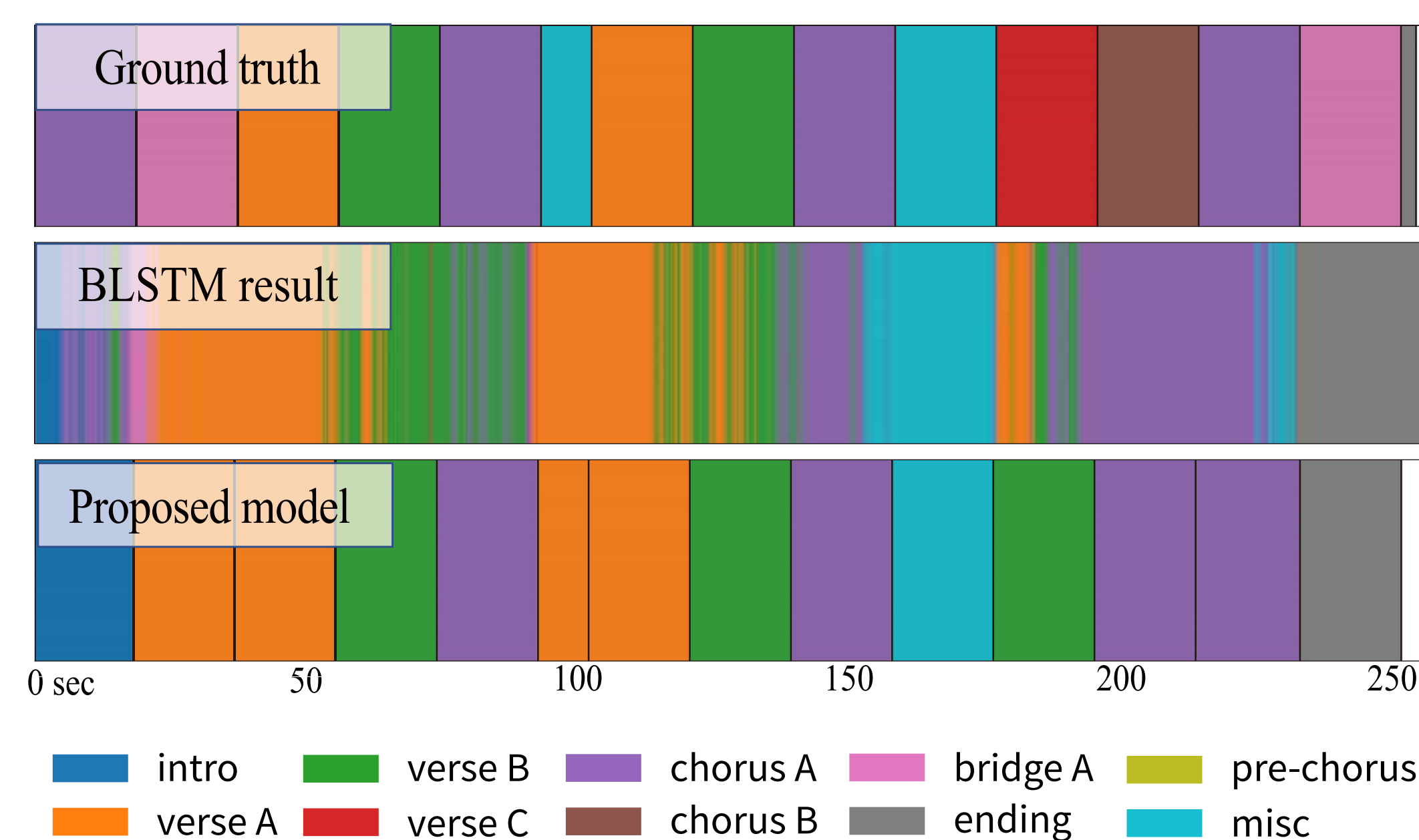


## Evaluation

### Comparison with Conventional Methods

Method	Segmentation		Clustering $F_{\text{pair}}$ (%)	Labeling accuracy(%)
	$F_{0.5}$ (%)	$F_{3.0}$ (%)		
GS3 [Grill+, '15]	<b>52.3</b>	73.5	54.2	n/a
SUG2 [Schlüter+, '14]	25.8	<b>73.7</b>	37.3	n/a
FK2 [Kaiser+, '13]	30.0	65.7	63.4	n/a
[Paulus, '09]	n/a	63.0	<b>63.7</b>	34.4
Proposed	43.3	66.5	54.6	<b>45.3</b>

### Example of Analysis Result



- The bidirectional LSTM was confused between “verse A” and “verse B” or between “verse B” and “chorus A”, while the proposed model correctly recognized these sections
- We need to improve the proposed model to avoid errors such as the confusion between “intro” and “chorus A”

## Future Work

- The proposed method worked best in terms of **labeling accuracy**
- There is much room for improvement except for labeling accuracy

- Refine the model to incorporate the novelty aspect
- Deal with more hierarchies because music has a hierarchical structure, from motive and phrase to section and section group