

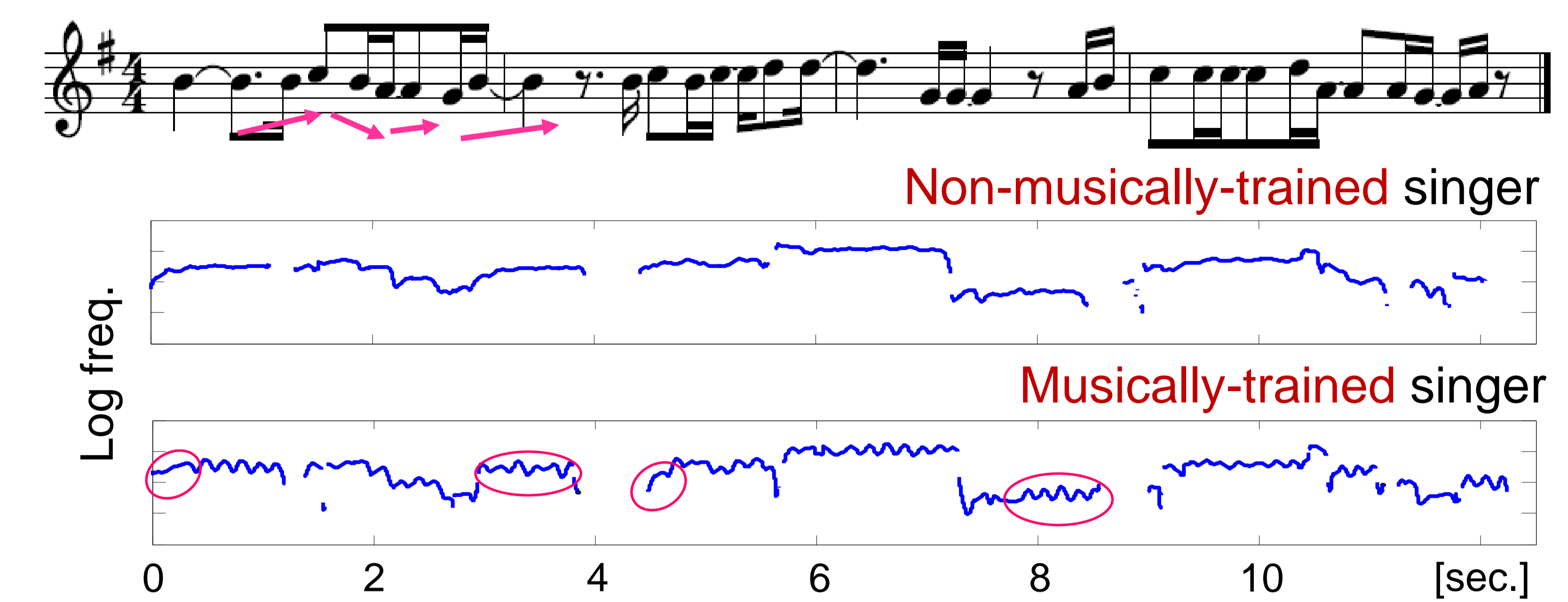
Yasunori Ohishi<sup>1</sup>, Daichi Mochihashi<sup>2</sup>, Hirokazu Kameoka<sup>1</sup>, Kunio Kashino<sup>1</sup>

<sup>1</sup>NTT Communication Science Laboratories, NTT Corporation, Japan, <sup>2</sup>The Institute of Statistical Mathematics

The goal is to build a generative model that can characterize the singing style of a singer in the sung melodic contours, i.e.,  $F_0$  contours, and predict the  $F_0$  contour that reflects the singing style for an arbitrary musical score, using a mixture of Gaussian process experts (MoGPEs). The experimental results showed that our model is promising for predicting the  $F_0$  contour for a given musical score. An application of this work is singing style conversion which automatically converts a poor singing voice into an expressive singing voice.

## Motivation

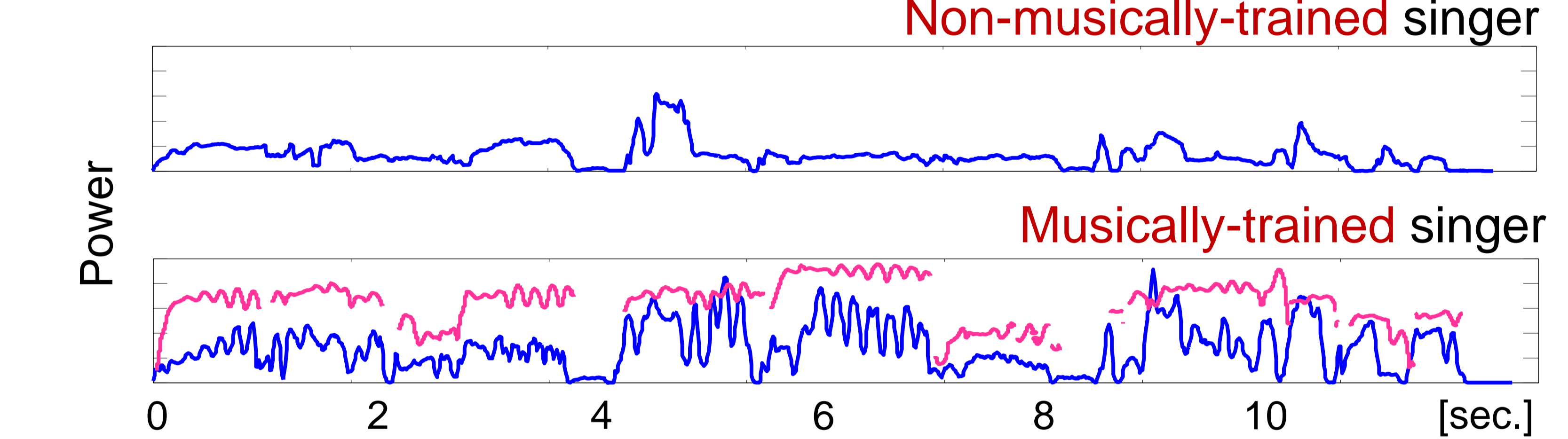
Dynamic fluctuations in a sung melodic contour ( $F_0$  contour)



- Vibrato*: quasi-periodic variation
- Portamento*: gradual sliding of pitch

Related to singing skill, style, and expression [1-3]  
Musically-trained singer can skillfully control the  $F_0$  fluctuations and expressively sing a melody based on the personal style

Dynamic fluctuations in a power contour



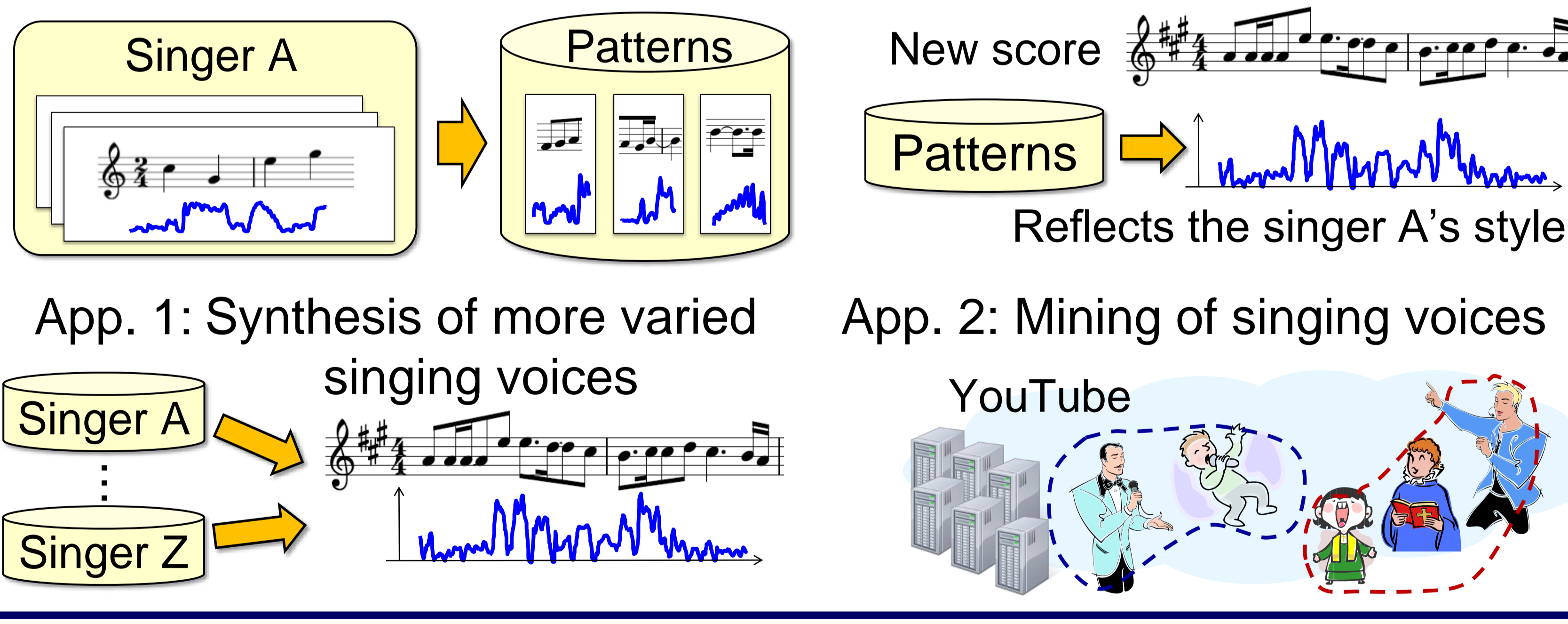
- Rise and fall of power according to the pitch of musical notes
- Fluctuations correlated with *vibrato* in the  $F_0$  contour

Musically-trained singer can also control the power of singing voice

## Goal

Prediction of the  $F_0$  contour with the expressive dynamic fluctuations for an arbitrary musical note sequence

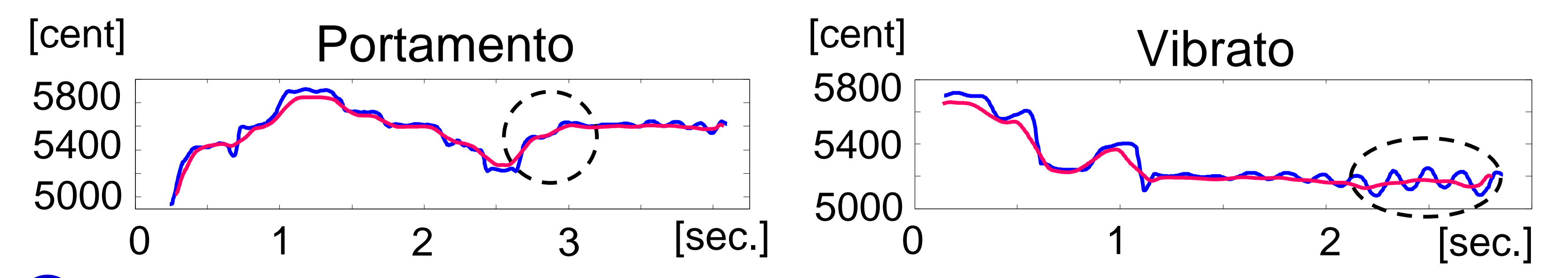
- Learning of the  $F_0$  patterns for musical note sequences
- Prediction of the  $F_0$  contour for a new musical note sequence



## Previous works

Hidden Markov model (HMM) [4-6]

- Difficult to characterize the continuous and rapid changes in the fluctuations, because its hidden-state space is discrete
- Context-dependent decision-tree clustering causes an over-smoothing effect of the  $F_0$  contours

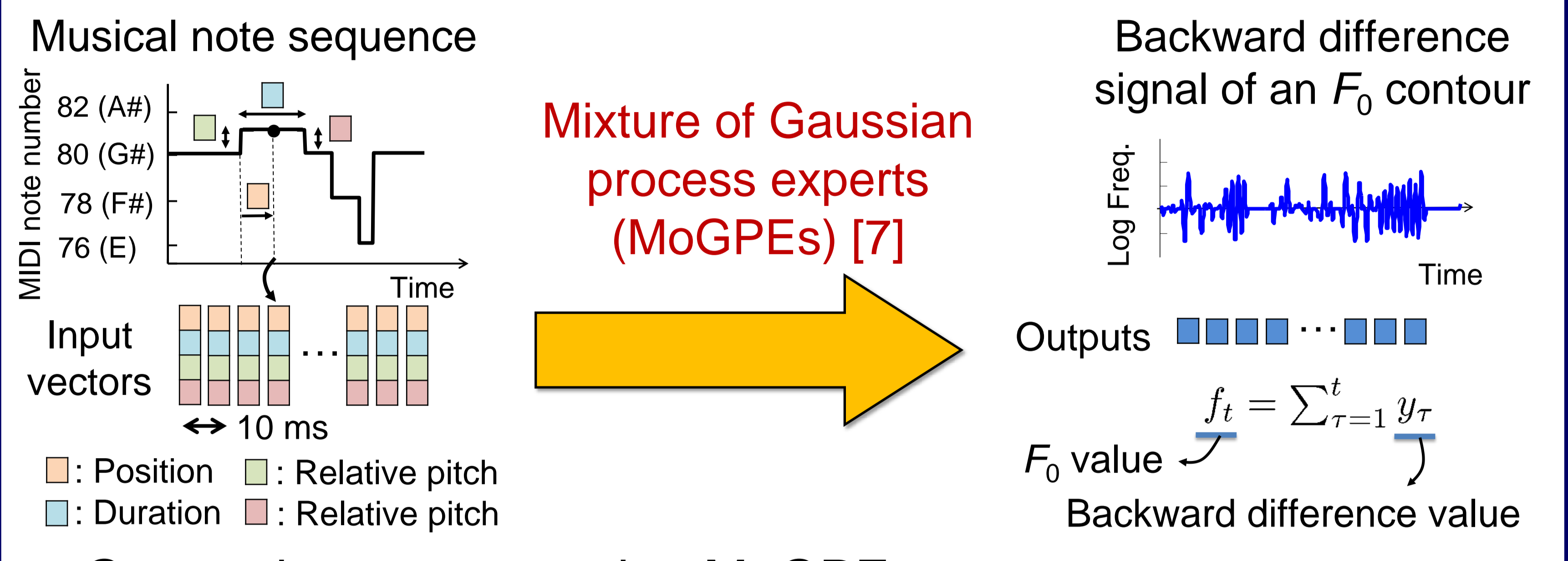


The  $F_0$  fluctuations disappear because of the over-smoothing effect

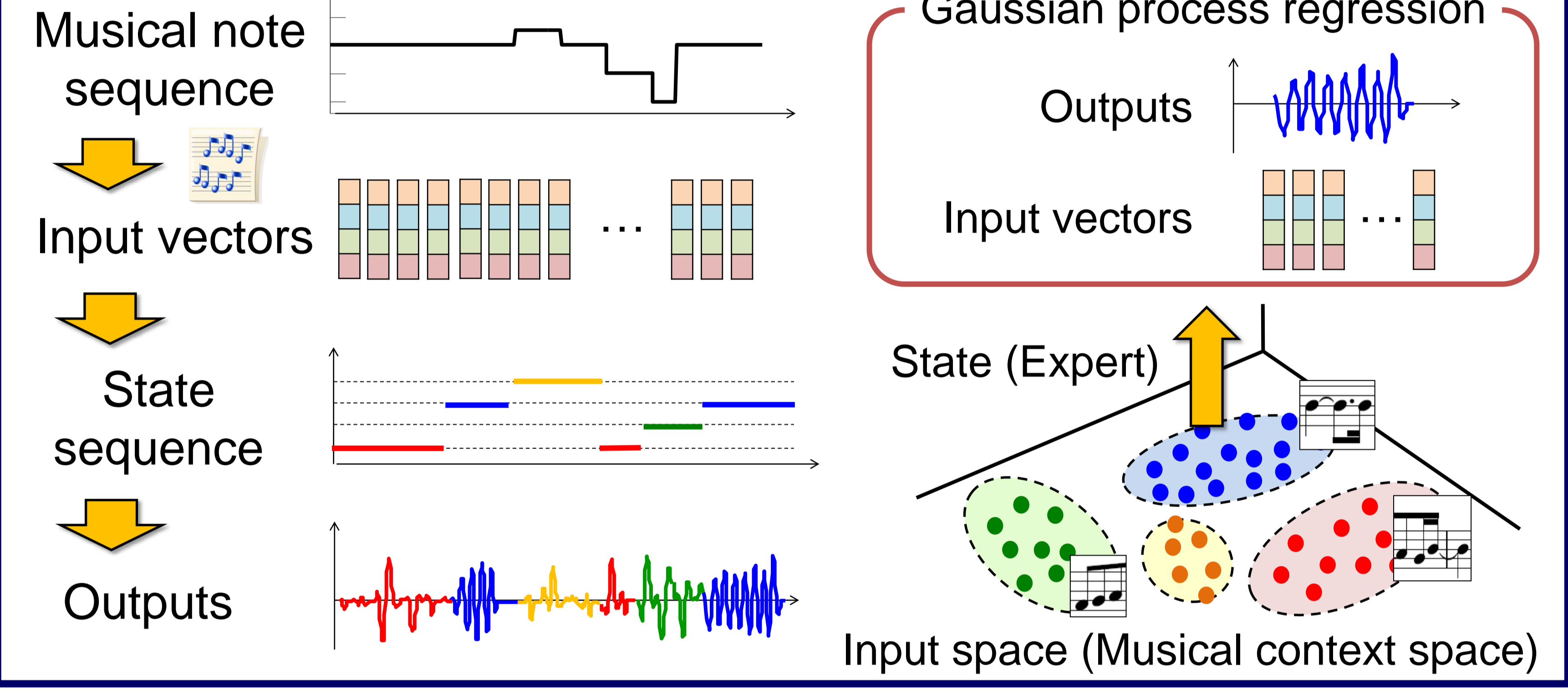
Characterize the  $F_0$  fluctuations using Gaussian process regression which is a sophisticated machine learning algorithm

## Statistical modeling of $F_0$ contour

Use the singing voice signals synchronized with melodies

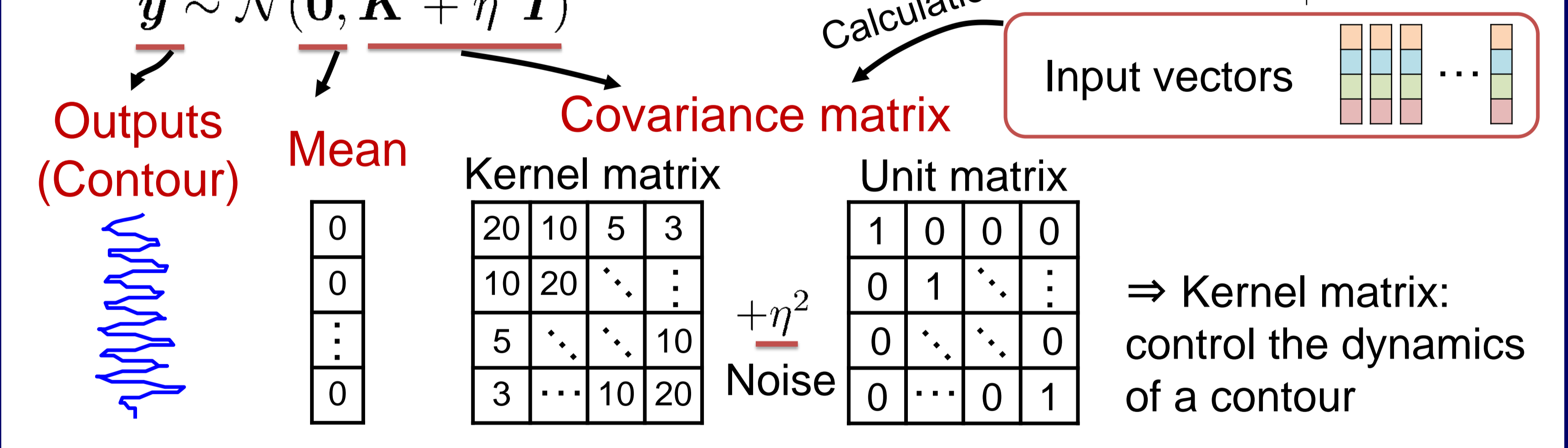


Generative process using MoGPEs

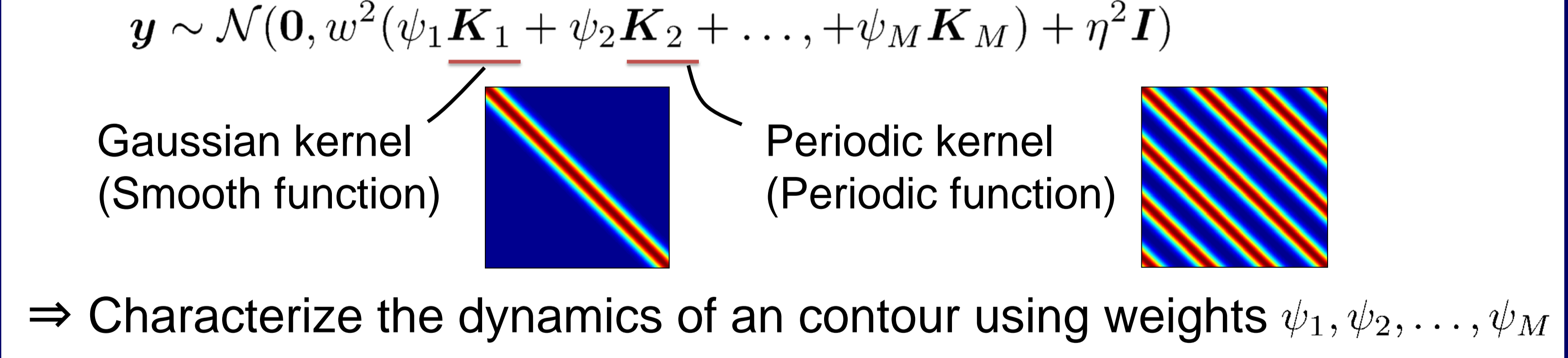


## Gaussian process regression (GPR) [8]

Probability distribution of a contour



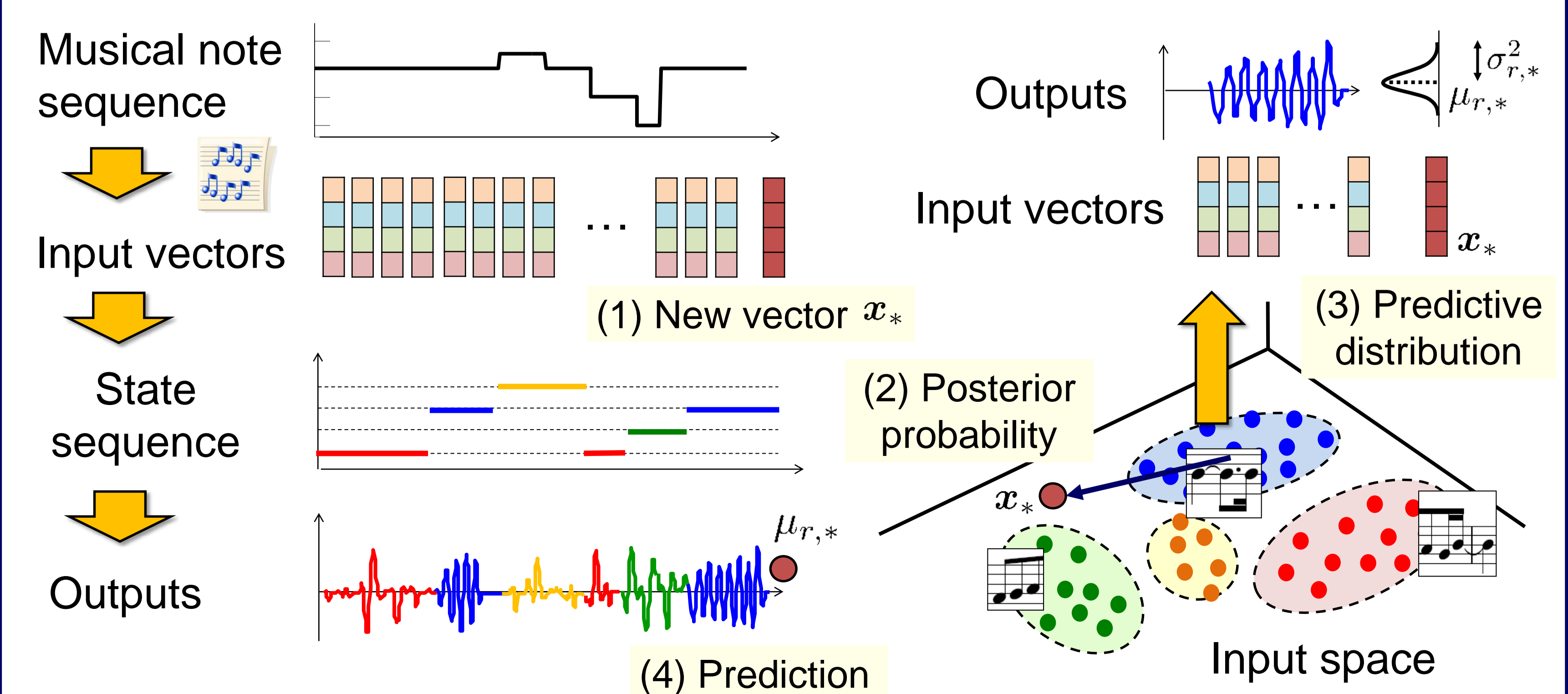
Multiple kernel learning [9]



## Parameter estimation and prediction

MCMC-EM algorithm [10]: estimate a state sequence, parameters of GPRs, and parameters of Gaussian distributions in the input space

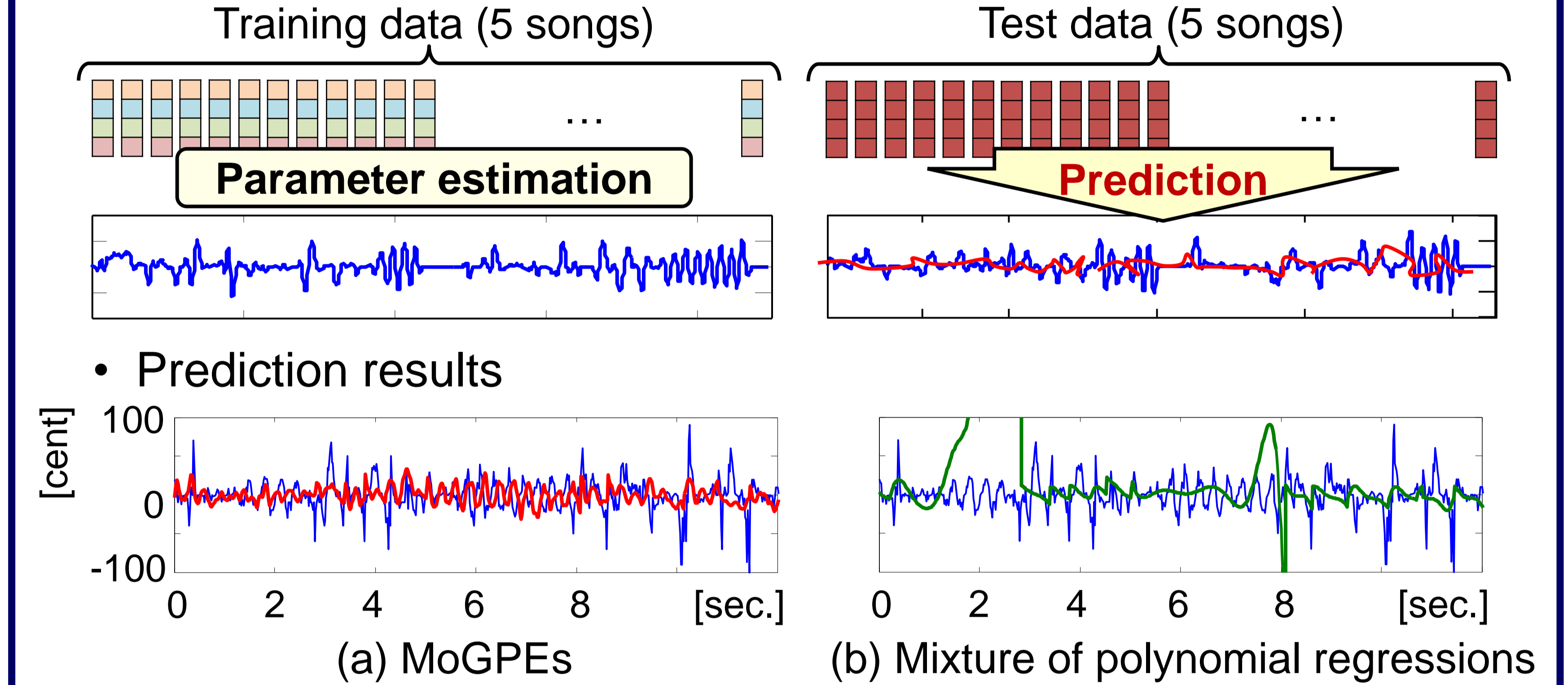
Prediction: use the mean value of the predictive distribution



## Evaluation

Predictive performance of the output for a new input vector

- Used the  $F_0$  contours and the MIDI signals of melodies annotated manually in the AIST annotation [11]
- Training data: Song No. 38, 39, 42, 44, and 45 (649.3 sec.)
- Test data: Song No. 46, 64, 72, 74, and 76 (625.6 sec.)



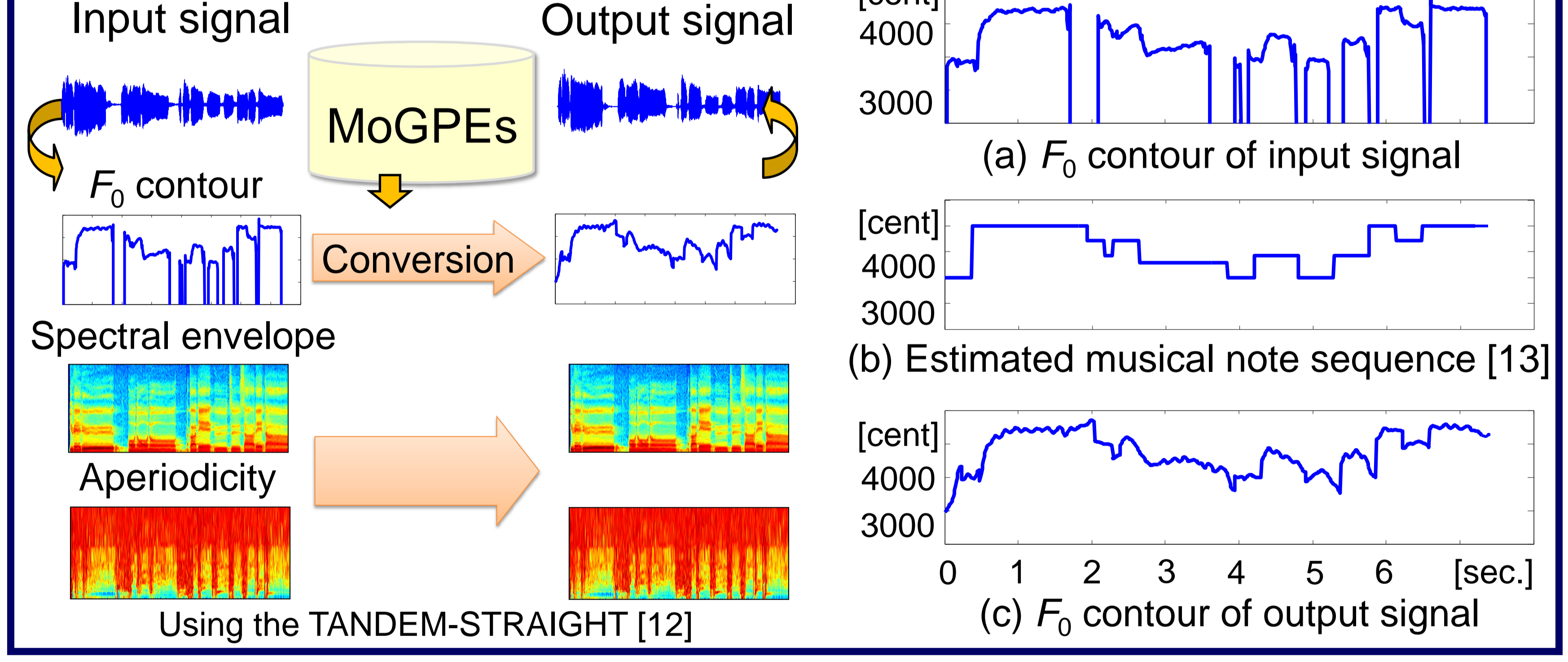
Evaluation measure: root mean square error (RMSE)

$$RMSE = \sqrt{\sum_{t=1}^{T_t} (y_t - \mu_{*,t})^2 / T_t}$$

$y_t$ : Actual output,  $T_t$ : Signal length,  $\mu_{*,t}$ : Predictive mean

Num. of experts	5	10	20	30	40	50
MoPRs	65.6	71.8	59.7	73.1	79.6	56.6
MoGPEs	<b>26.4</b>	<b>25.0</b>	<b>24.0</b>	<b>23.9</b>	<b>23.1</b>	<b>22.3</b>

## Singing style conversion



## Appendix

Joint distribution

$$p(\{x_t, y_t\}_{t=1}^T, \{z_t\}_{t=1}^T, \{\theta_r^{GP}\}_{r=1}^R, \{\theta_r^{\text{GP}}\}_{r=1}^R | \Omega) = \prod_{r=1}^R [p(\theta_r^{\text{GP}} | \Omega) p(X_r | \theta_r^{\text{GP}}) p(y_r | X_r, \theta_r^{\text{GP}}, \Omega)] \times p(\{z_t\}_{t=1}^T | \Omega)$$

$$\begin{cases} p(\theta_r^{\text{GP}} | \Omega) = \mathcal{N}(\mu_r; \mathbf{m}_0, \Sigma_r / \beta_0) \mathcal{W}(\Sigma_r^{-1}; \mathbf{W}_0, \nu_0), & p(X_r | \theta_r^{\text{GP}}) = \mathcal{N}(X_r; \mu_r, \Sigma_r) \\ p(y_r | X_r, \theta_r^{\text{GP}}, \Omega) = \mathcal{GP}(y_r; \mathbf{0}, \mathbf{K}_r + \eta_r^2 \mathbf{I}_r), & p(\{z_t\}_{t=1}^T | \Omega) = \frac{\Gamma(\alpha)}{\Gamma(\alpha + T)} \prod_{r=1}^R \frac{\Gamma(T_r + \alpha / R)}{\Gamma(\alpha / R)} \end{cases}$$

Predictive distribution

$$p(y_* | \mathcal{D}, x_*, \theta, \Omega) = \sum_{r=1}^R p(y_* | y_r, X_r, x_*, z_* = r, \theta_r^{\text{GP}}) p(z_* = r | x_*, \theta_r^{\text{GP}}) = \mathcal{N}(y_*; \mu_*, \sigma_*^2)$$

$$\mu_* = \sum_{r=1}^R c_r k_{r,*}^T (\mathbf{K}_r + \eta_r^2 \mathbf{I}_r)^{-1} y_r, \quad \sigma_*^2 = \sum_{r=1}^R c_r^2 (k_r(x_*, x_*) - k_{r,*}^T (\mathbf{K}_r + \eta_r^2 \mathbf{I}_r)^{-1} k_{r,*}), \quad c_r = p(z_* = r | x_*, \theta_r^{\text{GP}})$$

## References

- J. Sundberg, Northern Illinois University Press, 1987.
- T. Saitou et al., *Speech Communication*, vol. 45, no. 3-4, pp. 405-417, 2005.
- L. Regnier, Ph.D. thesis, IRCAM / UPMC in Paris, France, 2013.
- T. Yoshimura et al., in *Proc. EUROSPEECH* 1999.
- H. Zen et al., *Speech Communication*, vol. 51, no. 11, pp. 1039-1064, 2009.
- T. Nose et al., in *Proc. INTERSPEECH* 2013, pp. 378-382.
- E. Meeds et al., in *Proc. NIPS* 2006.
- C. E. Rasmussen et al., MIT Press, Cambridge, Mass, USA, 2006.
- K. Yoshii et al., in *Proc. ICASSP* 2013, pp. 463-467.
- C. Andrieu et al., *Machine Learning*, vol. 50, no. 1-2, pp. 5-43, 2003.
- M. Goto, in *Proc. ISMIR* 2006.
- H. Kawahara et al., in *Proc. ICASSP* 2008, pp. 3933-3936.
- Y. Ohishi et al., in *Proc. INTERSPEECH* 2012.