

# Improving melody extraction using Probabilistic Latent Component Analysis

Jinyu. Han<sup>1</sup> Ching-Wei. Chen<sup>2</sup>

<sup>1</sup>Interactive Audio Lab  
Northwestern University, USA

<sup>2</sup>Media Technology Lab  
Gracenote, Inc

May 19, 2011

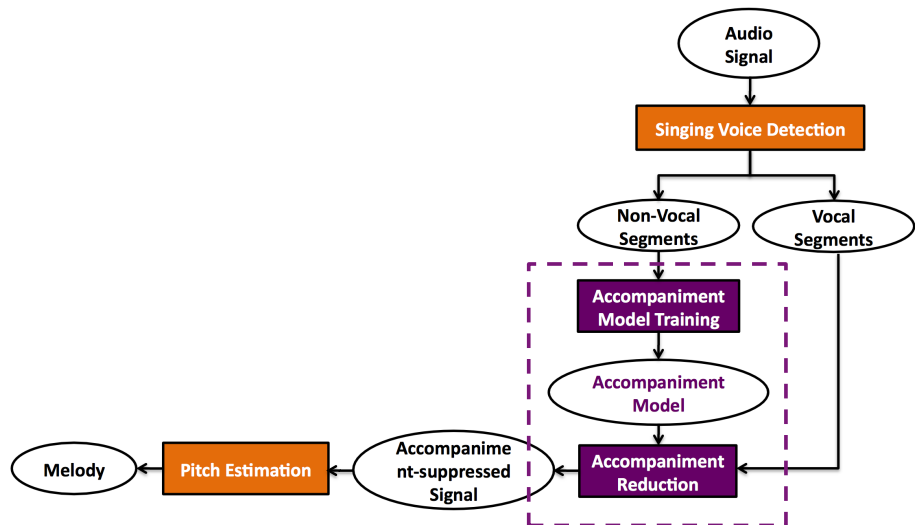
# Agenda

- 1 Introduction
- 2 Modeling the Spectrogram
  - Multinomial Model
  - Probabilistic Latent Component Analysis
- 3 System Description
- 4 Experiment Results
  - Illustration Example
  - System Comparison
- 5 Conclusion

Pick only the singing voice as the Melody

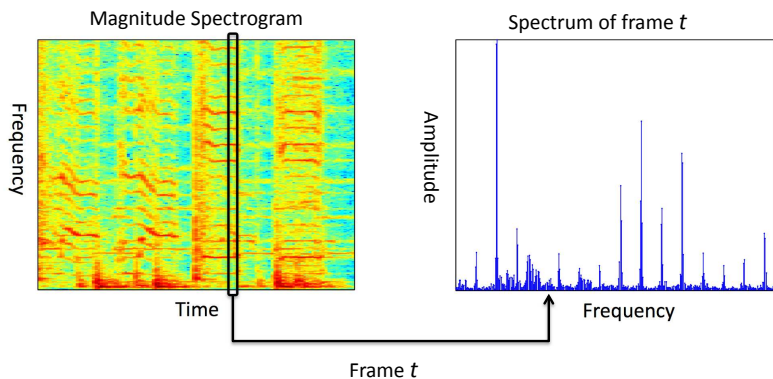


# System Overview



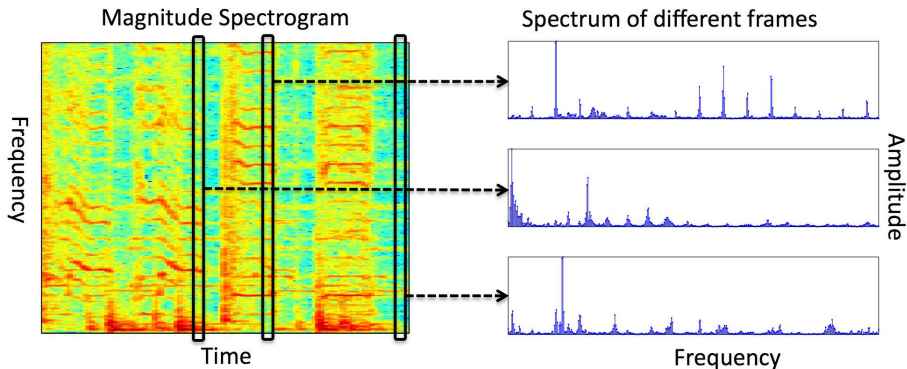
# Multinomial Distribution for Spectrogram

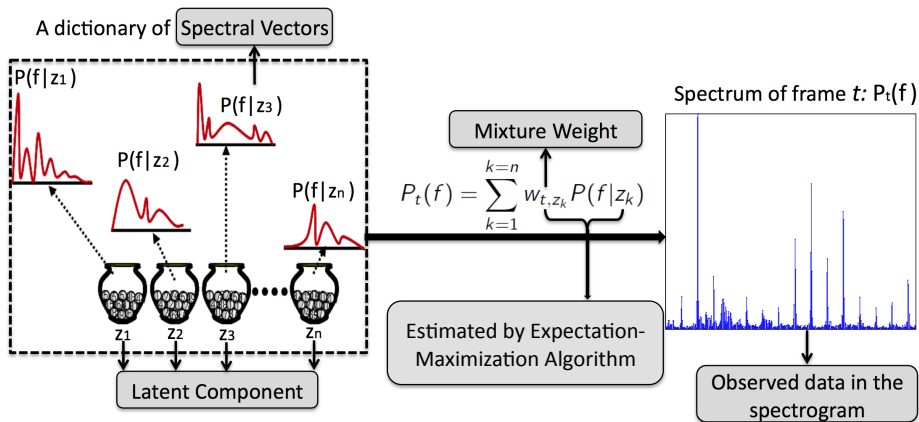
Figure: Probability distribution underlying the  $t$ -th spectrum



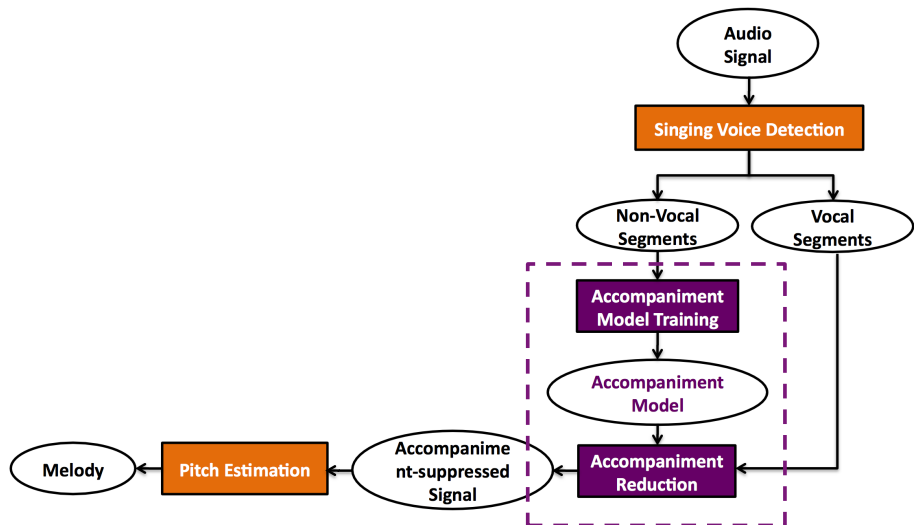
- Treat the spectrum in each time slice as a histogram
- Treat the histogram as a probability distribution

# Multinomial Distribution for Spectrogram



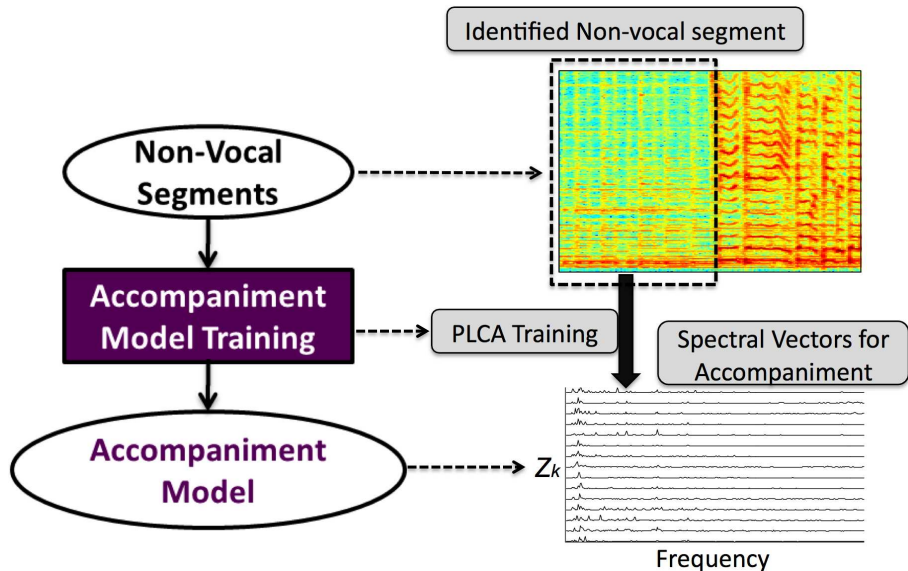


# System Overview

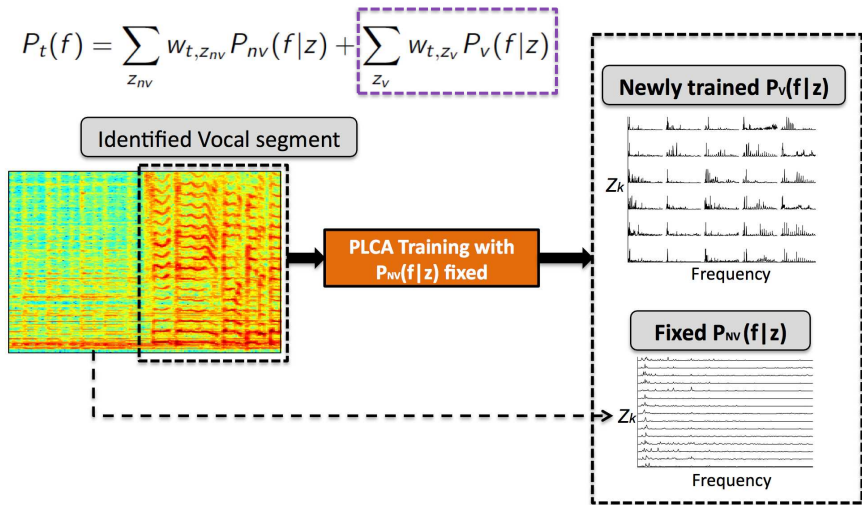




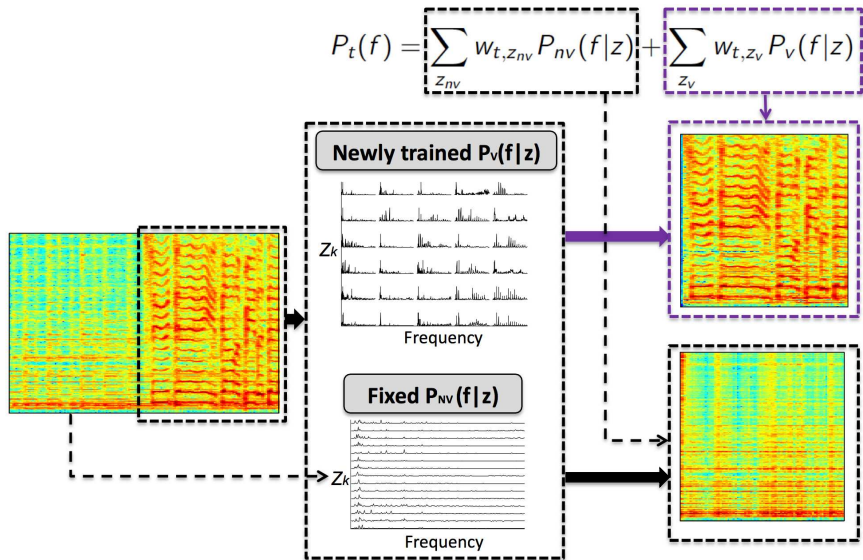
# Train $P_{nv}(f|z)$ from the non-vocal segment



# Extract singing voice in the mixture



# Extract singing voice in the mixture

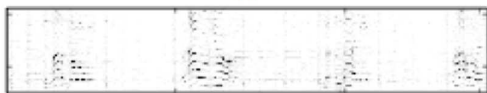


A 14-s clip of "Simple Man" by Lynyrd Skynyrd



Mixture

Extracted singing voice



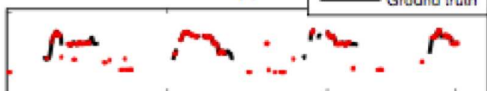
Extracted Voice

Original singing voice



Clean Voice

Melody Line



Time

## Compare our system to DHP[1] and LW[2]

	Precision	Recall	F-measure	Accuracy
DHP	<b>0.52</b>	0.48	0.50	0.48
LW	0.09	0.086	0.09	0.19
Proposed	0.43	<b>0.80</b>	<b>0.55</b>	<b>0.61</b>

Parts of MIREX 2005 dataset: 9 recordings, totalling about 270 seconds of audio.



Z. Duan, J. Han, and B. Pardo, "Harmonically informed pitch tracking", in Proc. ISMIR, 2009.

Y. Li and D. Wang,

"Separation of singing voice from music accompaniment for monaural recordings", IEEE Trans. Audio, Speech, and Language

# Conclusion

- The Probabilistic Latent Variable Model is introduced to model the accompaniment and lead vocal adaptively
- Experimental results show that the melody of the singing voice in mixture audio is successfully extracted to some extent.
- Future directions include improving the vocal/nonvocal segmentation module and the pitch estimation algorithm.

## Acknowledgement

- The first author performed this work with Ching-Wei Chen while at the Gracenote Media Technology Lab. We thank Markus Cremer, Bob Coover, Phillip Popp, Trista Chen, and Peter Dunker for enlightening discussions.
- The authors would like to thank the reviewers for their comments that help improve the paper.
- We also want to thank Bryan Pardo, David Little, Zhiyao Duan, Zafar Rafii, and Mark Cartwright for their suggestions that improve the presentation.