# REpeating Pattern Extraction Technique (REPET)

## EECS 352: Machine Perception of Music & Audio

# Observation

- **Repetition** is a fundamental element in generating and perceiving structure
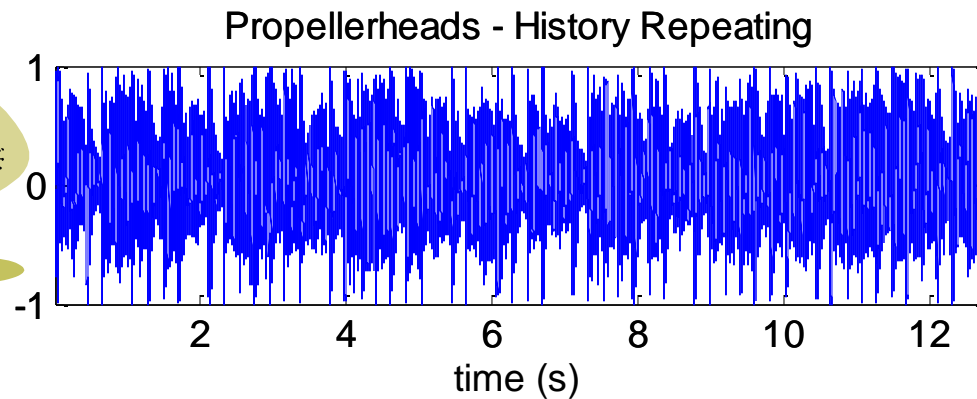


… in nature



… in art
[http://http://en.wikipedia.org/wiki/Campbell's_Soup_Cans]



… in audio

# Observation

- Musical works are often characterized by an **underlying repeating structure** over which varying elements are superimposed



Propellerheads - History Repeating
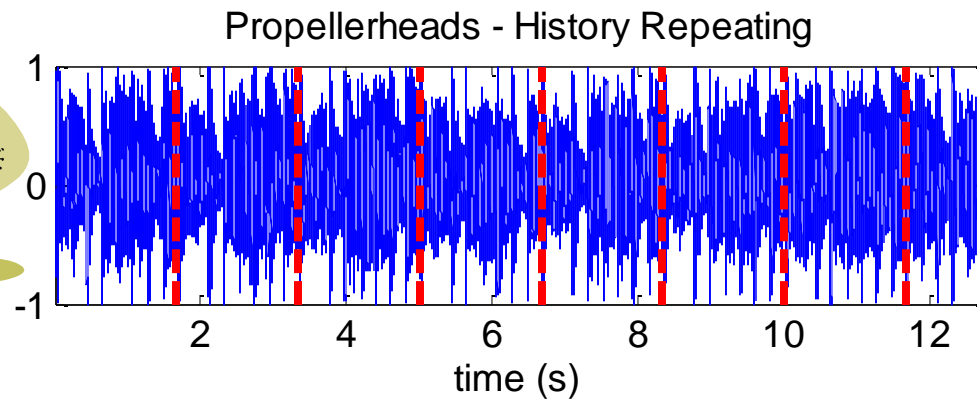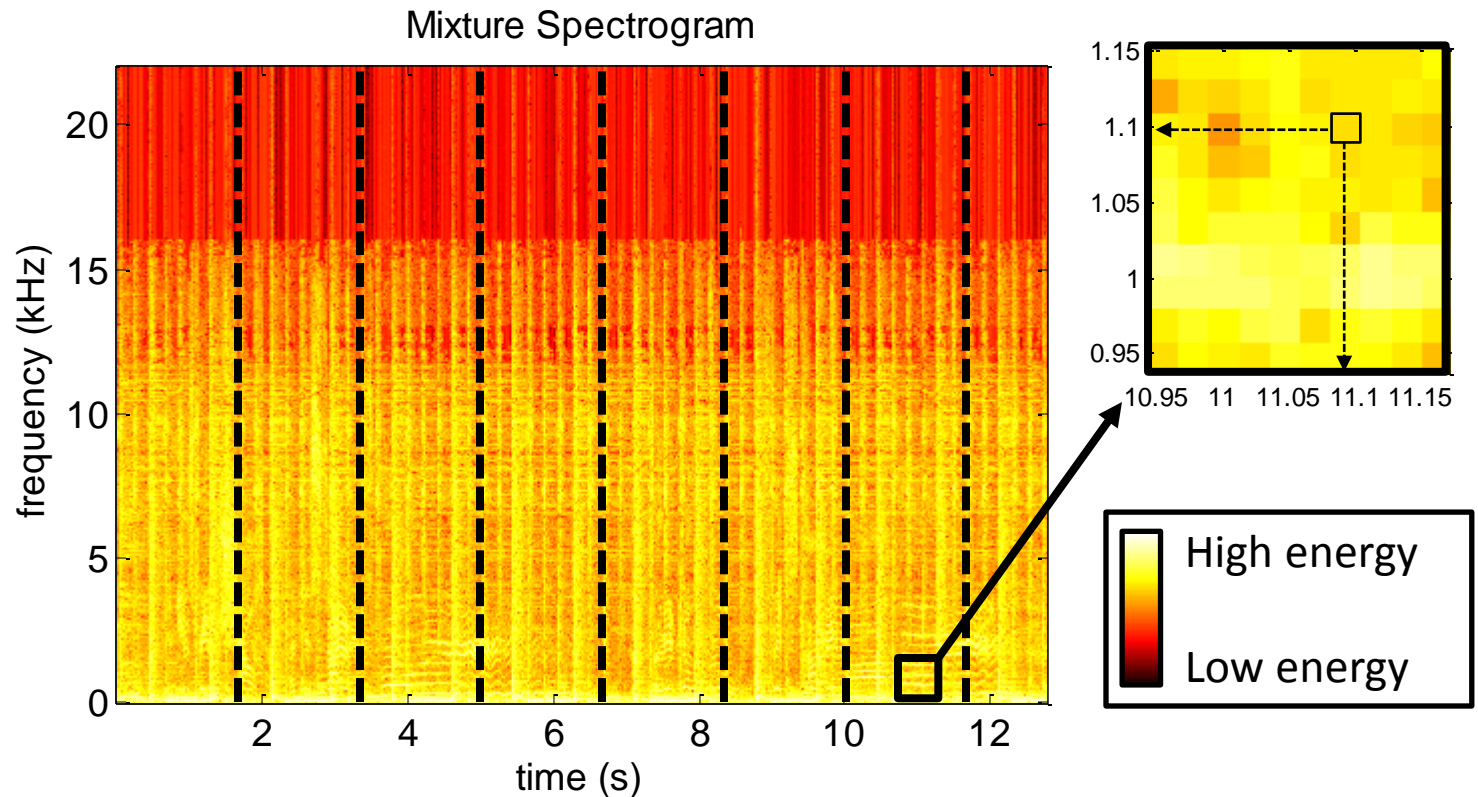
# Observation

- Musical works are often characterized by an **underlying repeating structure** over which varying elements are superimposed
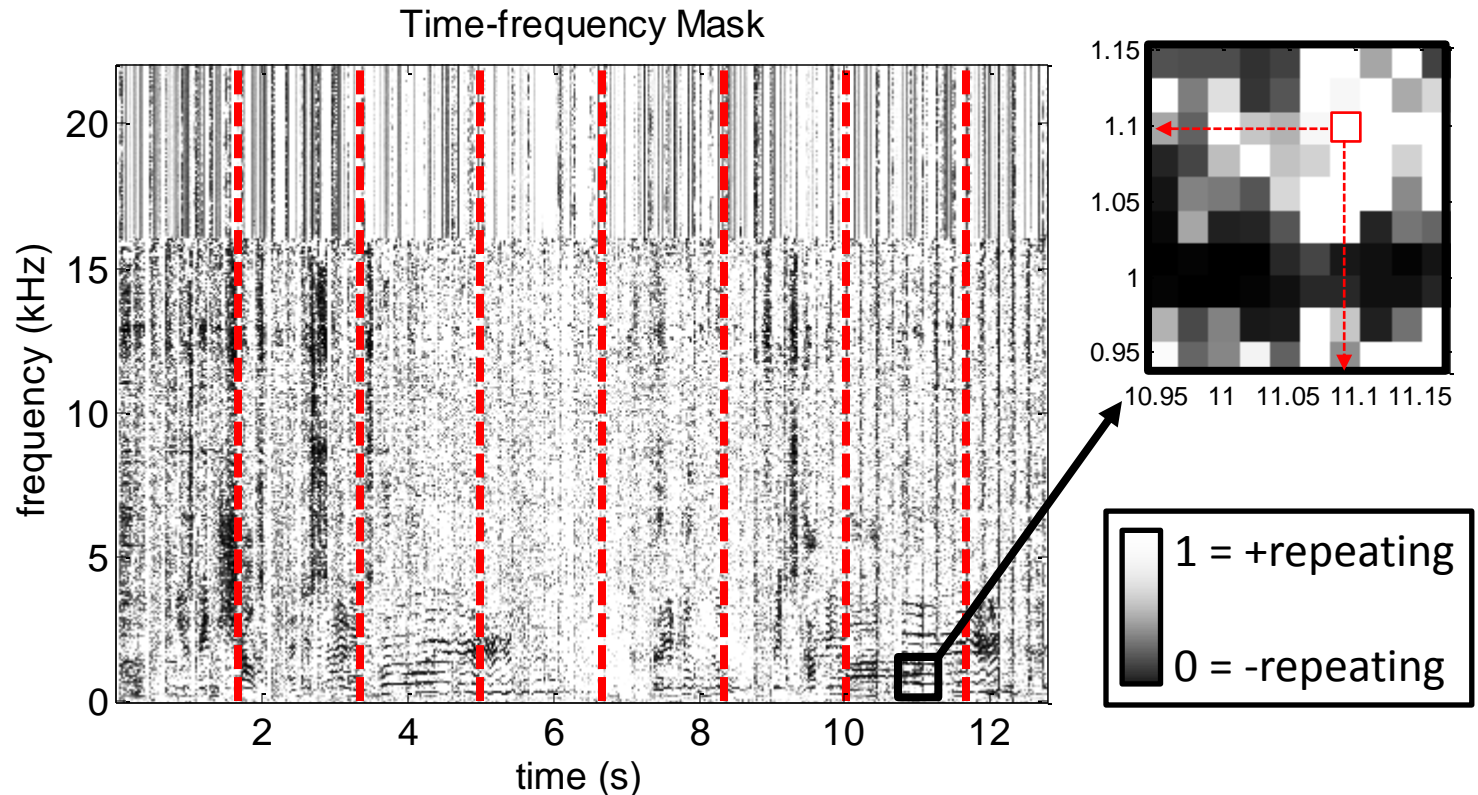
# Assumption

- There should be patterns that are more or less **repeating in time and frequency**

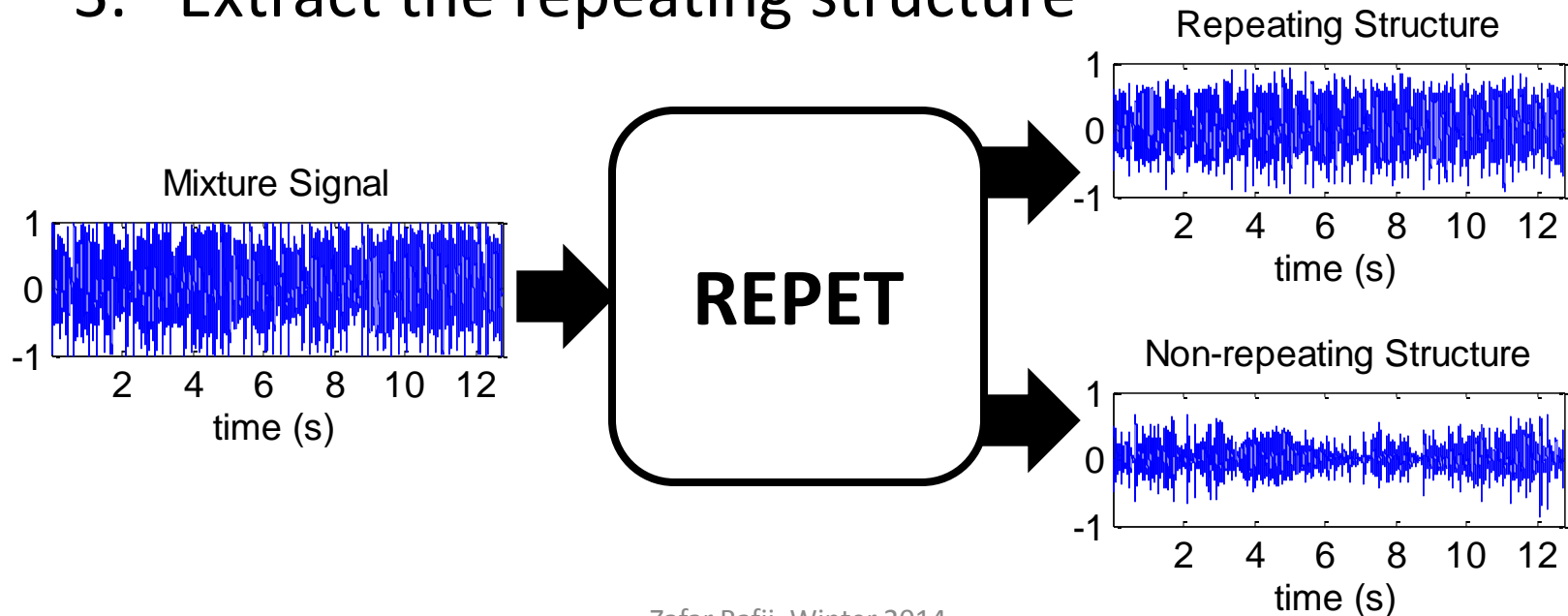Mixture Spectrogram

# Assumption

- The repeating patterns could be identified and extracted using a **time-frequency mask**
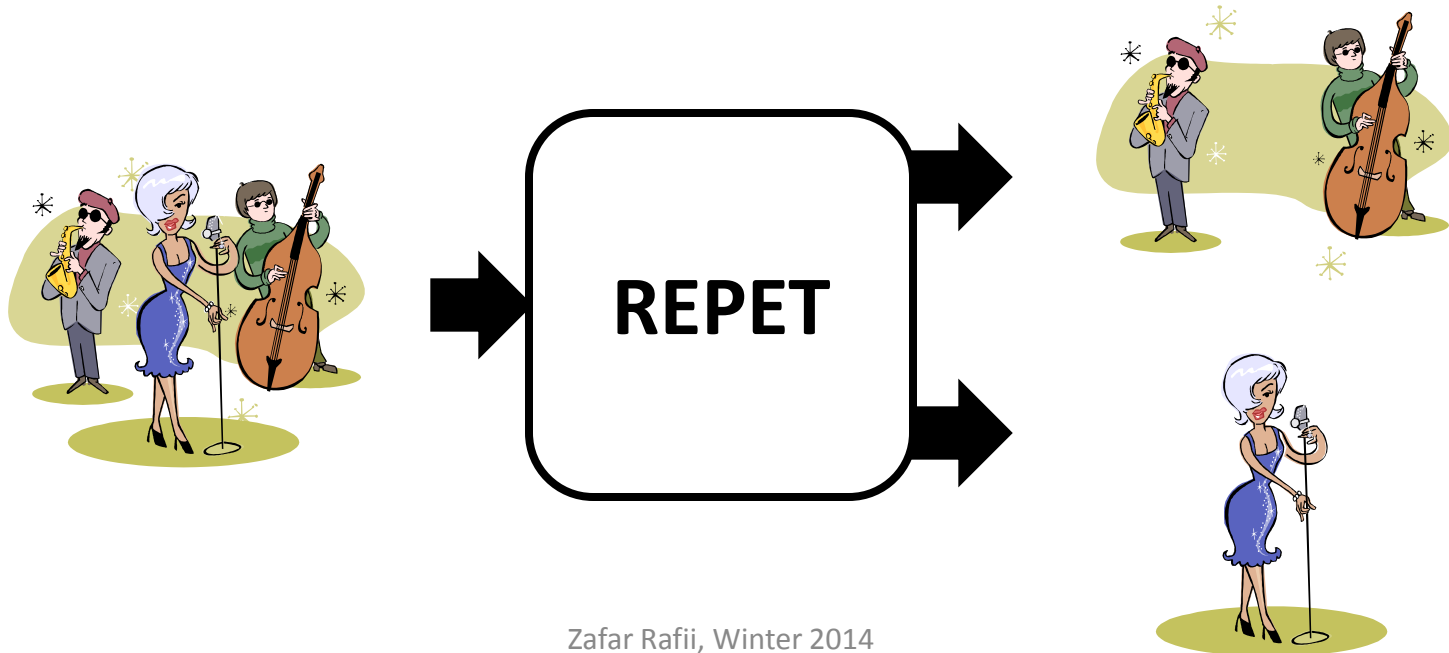


Time-frequency Mask

# Idea

- **REpeating Pattern Extraction Technique!**

  1. Identify the repeating elements
  2. Derive a repeating model
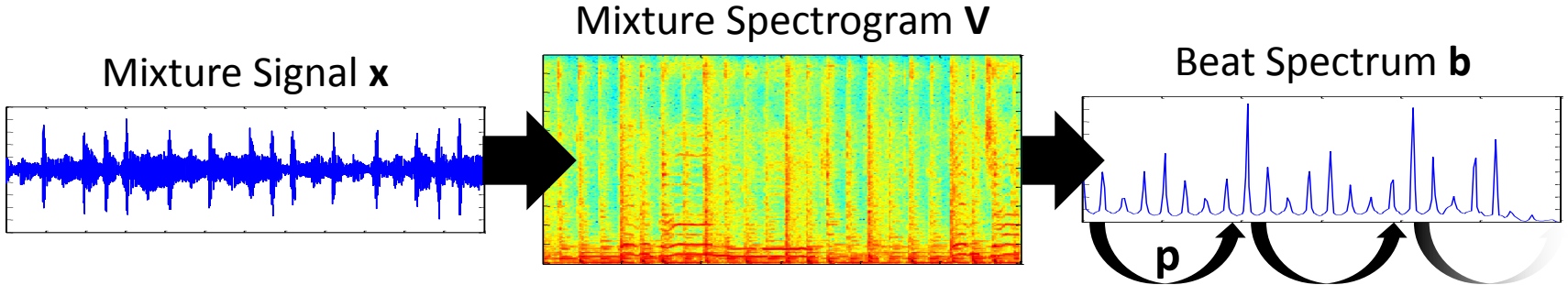  3. Extract the repeating structure

# Idea

- Simple **music/voice separation** method!
  - Repeating structure = background music
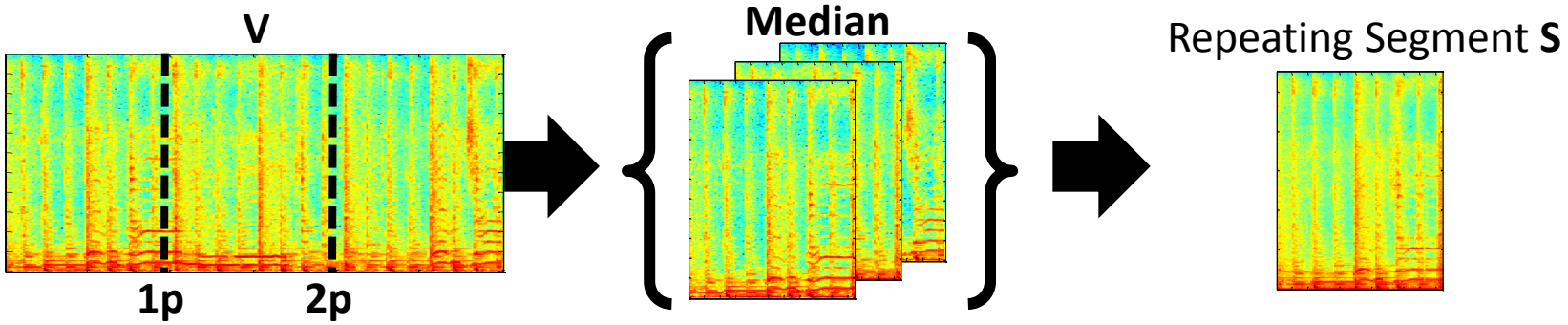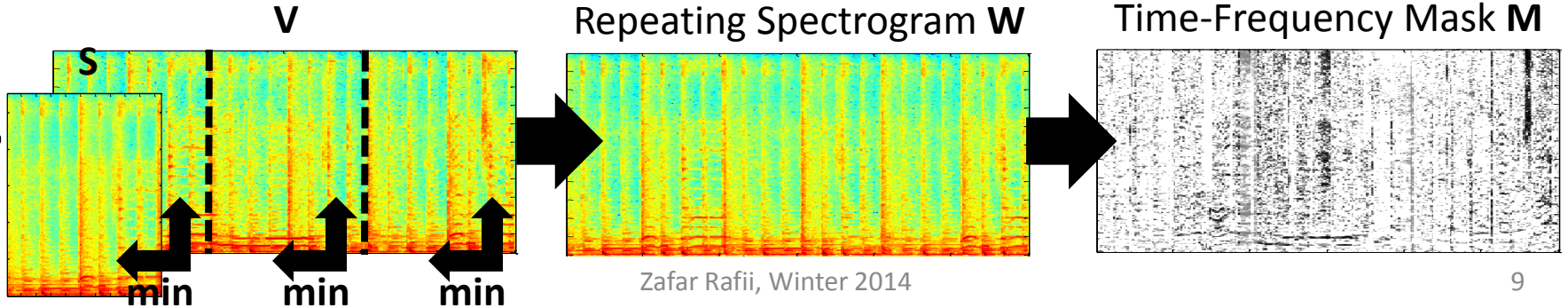  - Non-repeating structure = foreground voice



**REPET**

# REPET



**Step 1**
Mixture Signal **x** → Mixture Spectrogram **V** → Beat Spectrum **b**

**p**

**Step 2**
**V** → Median → Repeating Segment **S**

1p  2p

**Step 3**
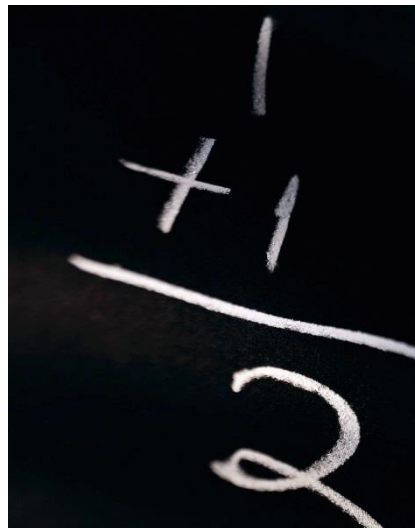**V** → Repeating Spectrogram **W** → Time-Frequency Mask **M**

**S**

min  min  min

# Practical Advantages

- Does not depend on special parametrizations
- Does not rely on complex frameworks
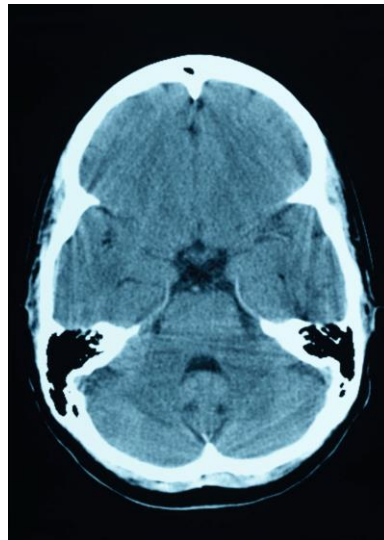- Does not require external information

# Practical Interests

- Karaoke gaming (need the music)
- Query-by-humming (need the voice)
- Audio remixing (need both components)
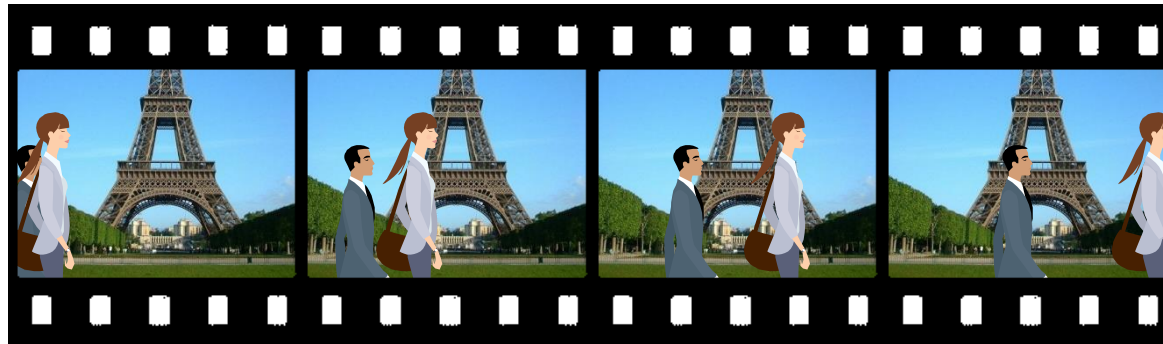
# Intellectual Interests

- Music understanding
- Music perception
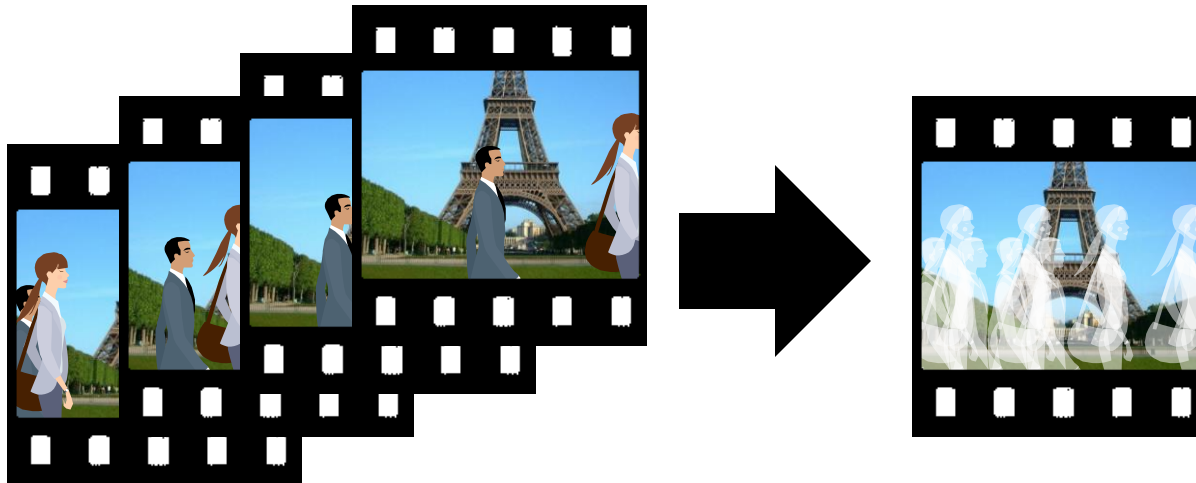- Simply based on repetition!

# Parallels

- **Background subtraction** in computer vision
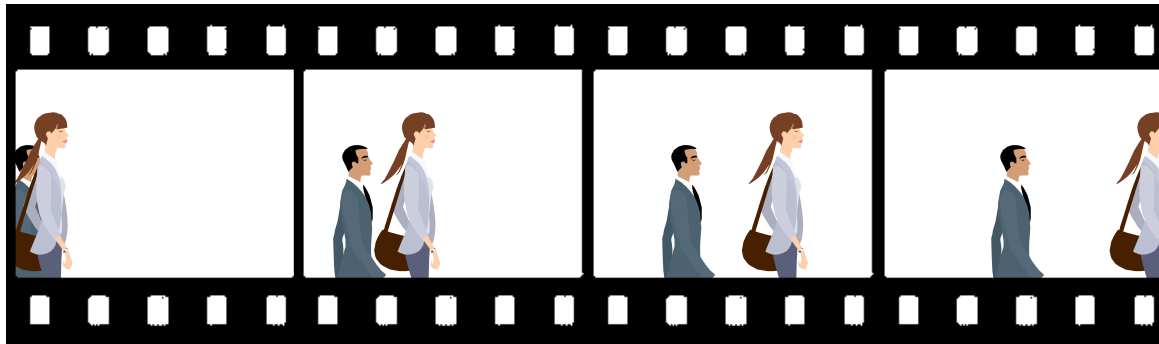
Sequence of video frames



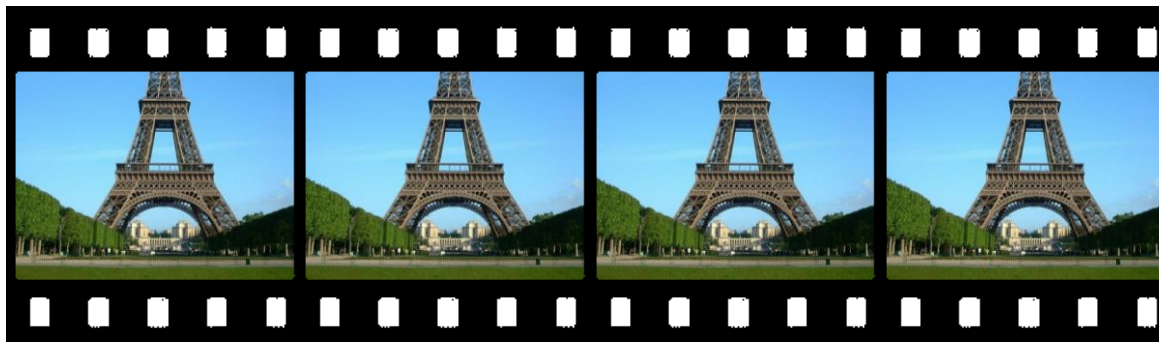Compare frames to estimate a background model

# Parallels

- **Background subtraction** in computer vision
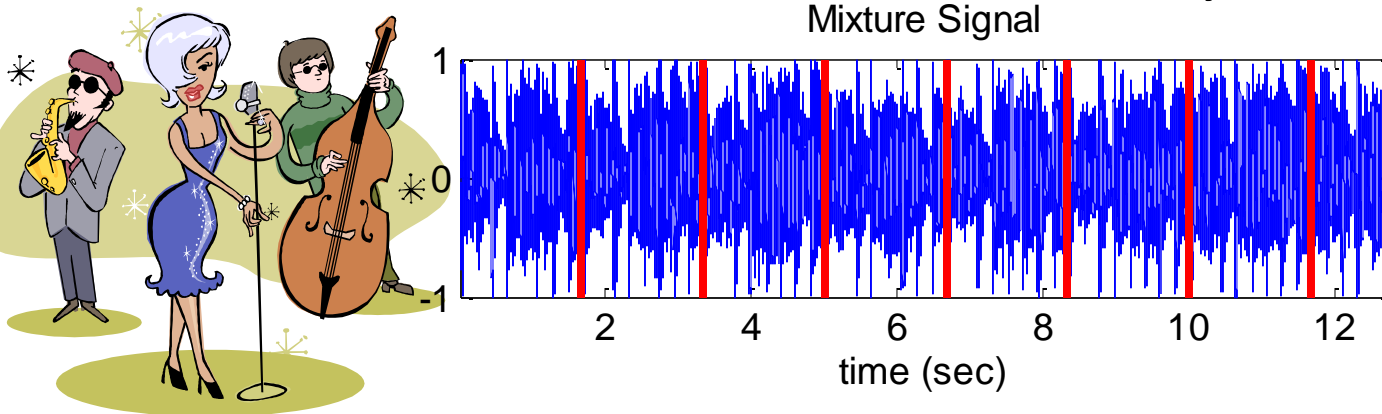
Extracted varying foreground scene



Extracted fixed background scene

# Parallels

- **Background subtraction** in computer vision
  - In audio, we also need to identify the repetitions!
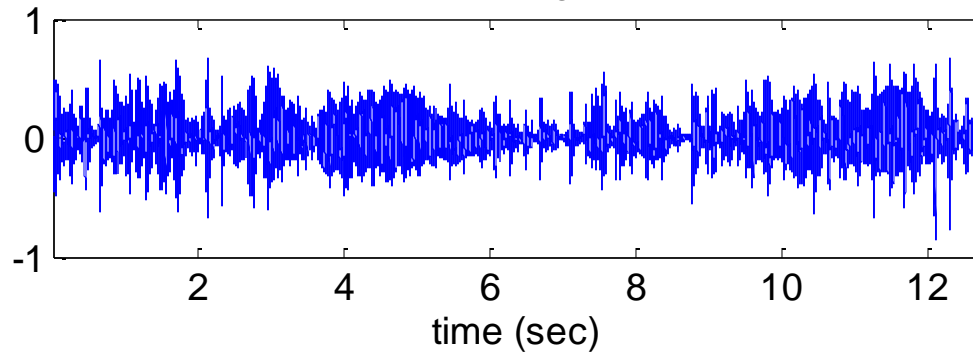


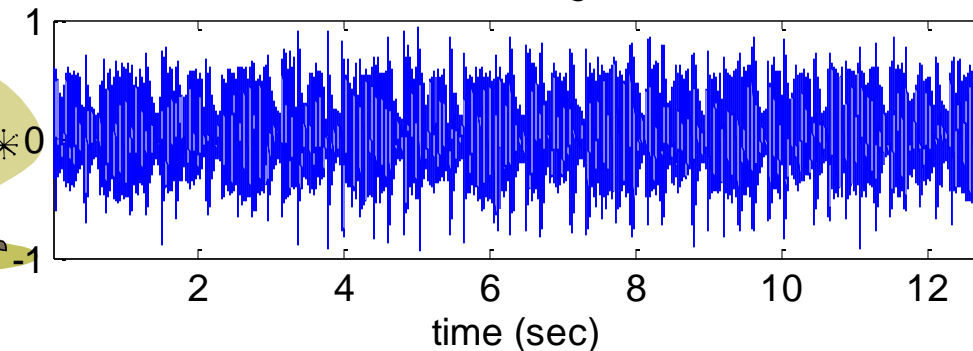Mixture Signal

# Parallels

- **Background subtraction** in computer vision
  - In audio, we also need to identify the repetitions!



Vocal Foreground

Musical Background

# Parallels

- **Auditory segregation** in human listeners



Target identified as
the repeating object

Handsome professor

Unknown audio mixtures
with the same target
and different distractors

# Parallels

- **Auditory segregation** in human listeners



1 mixture:

2 mixtures:

3 mixtures:

5 mixtures:

10 mixtures:

red/black = target/probe,
other colors = distractors

As the number of mixtures increases,
the target becomes more apparent...
[courtesy of Josh McDermott]

# REPET



Zafar Rafii, Winter 2014

19

# 1. Repeating Period

- We compute the **autocorrelations** of the frequency rows of the mixture spectrogram



Mixture Spectrogram

Autocorrelation Plots
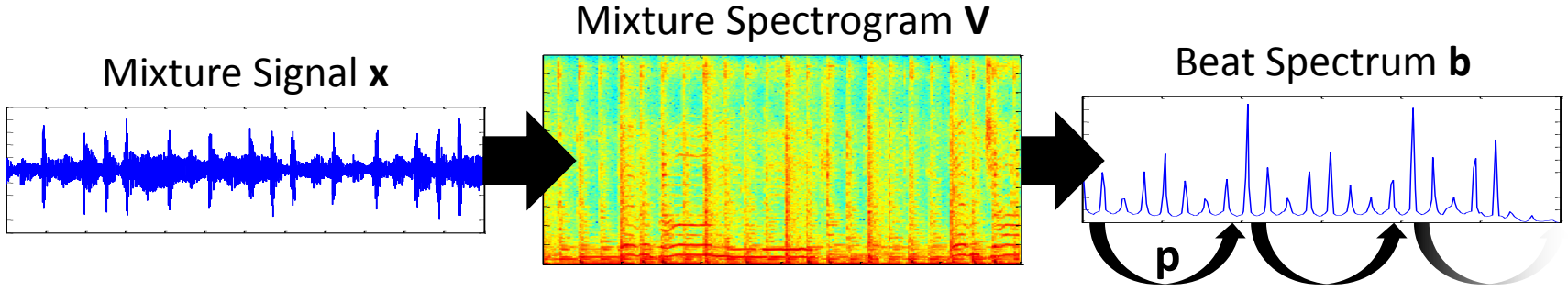
Spectrum at 10 kHz

Autocorrelation at 10 kHz

# 1. Repeating Period

- We take the mean of the autocorrelation rows and obtain the **beat spectrum**



Mixture Spectrogram → xcorr → Autocorrelation Plots → mean → Beat Spectrum

# 1. Repeating Period

- The beat spectrum reveals the **repeating period p** of the underlying repeating structure



Mixture Signal

Beat Spectrum

# REPET

**Step 1**

Mixture Signal **x**

Mixture Spectrogram **V**

Beat Spectrum **b**

**p**

**Step 2**

**V**

1p    2p

**Median**

Repeating Segment **S**

**Step 3**

**V**

S

min    min    min

Repeating Spectrogram **W**

Time-Frequency Mask **M**

# 2. Repeating Segment

- We then use the repeating period to **segment** the mixture spectrogram at period rate



Mixture Spectrogram

Segmented Spectrogram

Beat Spectrum

# 2. Repeating Segment

- We derive a **repeating segment model** by taking the element-wise median of segments


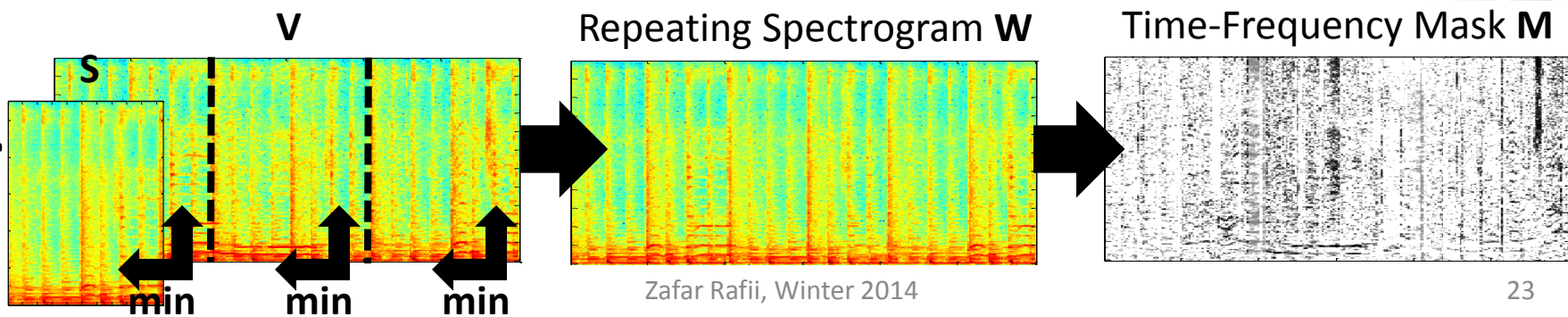
Mixture Spectrogram

Segmented Spectrogram

Repeating Segment

median

# 2. Repeating Segment

- The **median** helps to derive a clean repeating segment, removing the non-repeating outliers



Mixture Spectrogram

Repeating Segment

median

+ energy

- energy

# REPET



**Step 1**

Mixture Signal **x**

Mixture Spectrogram **V**

Beat Spectrum **b**

**p**

**Step 2**

**V**

1p    2p

**Median**

Repeating Segment **S**

**Step 3**

**V**

S

min    min    min

Repeating Spectrogram **W**

Time-Frequency Mask **M**

# 3. Repeating Structure

- We take the element-wise **min** between the repeating segment model and the segments

# 3. Repeating Structure

- We obtain a **repeating spectrogram model** for the repeating background

# 3. Repeating Structure

- The repeating spectrogram **should not have values higher than** the mixture spectrogram

Mixture Spectrogram $\geq 0$ **=** Repeating Spectrogram $\geq 0$ **+** Non-repeating Spectrogram $\geq 0$

# 3. Repeating Structure

- We then **divide**, element-wise, the repeating spectrogram by the mixture spectrogram

Mixture Spectrogram

Repeating Spectrogram

**divides**

# 3. Repeating Structure

- We obtain a **soft time-frequency mask** (with values between 0 and 1)



Mixture Spectrogram

Repeating Spectrogram

Time-frequency Mask

divides

# 3. Repeating Structure

- In the soft t-f mask, the **more/less** a t-f bin is repeating, the more it is weighted toward **1/0**

# 3. Repeating Structure

- We could further derive a **binary t-f mask** by fixing a threshold between 0 and 1

# 3. Repeating Structure

- We **multiply**, element-wise, the t-f mask with the mixture STFT to get the background STFT



Mixture Spectrogram

Background Spectrogram

Time-frequency Mask

.x

# 3. Repeating Structure

- We obtain the **repeating background** signal by inverting its STFT into the time domain



Mixture Spectrogram → iSTFT → Background Signal

# 3. Repeating Structure

- We obtain the **non-repeating foreground** signal by subtracting background from mixture

# Summary

- Repeating background ≈ **music component**
- Non-repeating foreground ≈ **voice component**

# Music/Voice Separation

- A variety of techniques has been proposed to separate **music** and **voice** from a mixture
    - Accompaniment modeling, Pitch-based inference, Non-negative Matrix Factorization (NMF), etc.

# Music/Voice Separation

- **Accompaniment modeling**
    - Modeling of the musical accompaniment from the non-vocal segments in the mixture

Mixture spectrogram    Vocal/non-vocal segmentation    Music spectrogram

→ Need an accurate vocal/non-vocal segmentation!
→ Need a sufficient amount of non-vocal segments!

# Music/Voice Separation

- **Pitch-based inference**
  - Separation of the vocals using the predominant pitch contour extracted from the vocal segments

Mixture spectrogram                Predominant pitch detection                Voice spectrogram

→ Need an accurate predominant pitch detection!

→ Cannot extract unvoiced vocals!

# Music/Voice Separation

- **Non-negative Matrix Factorization (NMF)**
  - Iterative factorization of the mixture spectrogram into non-negative additive basic components

Mixture spectrogram    Bases    Activations    Music & voice spectrograms

$\rightarrow$ Need to know the number of components!

$\rightarrow$ Need a proper initialization!

# Evaluation

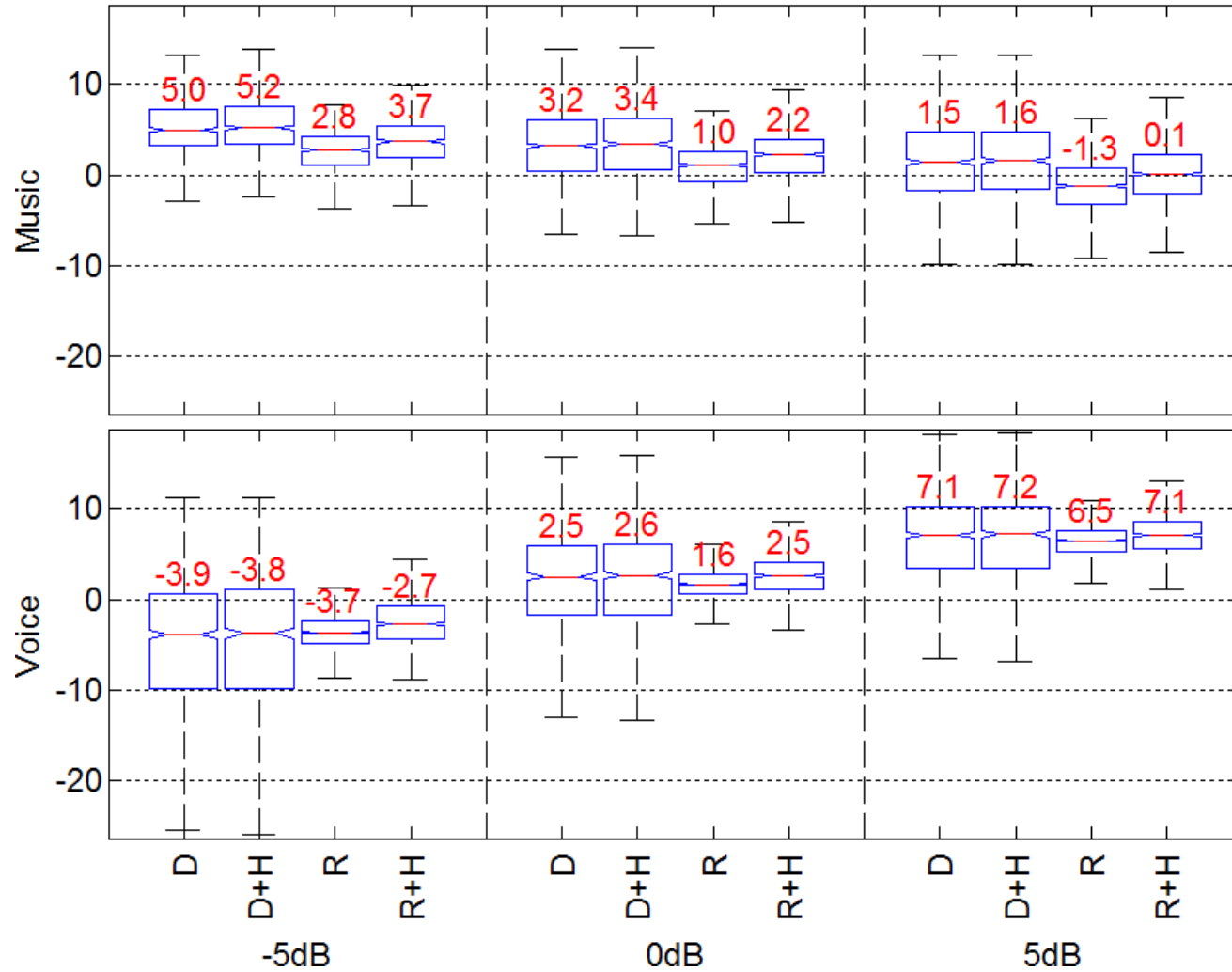- **REPET** [Rafii et al., 2013]
  - Automatic period finder
  - Soft time-frequency masking

- **Competitive method** [Durrieu et al., 2011]
  - Source-filter modeling with NMF framework
  - Unvoiced vocals estimation

- **Data set** [Hsu et al., 2010]
  - 1,000 song clips (from karaoke Chinese pop songs)
  - 3 voice-to-music mixing ratios (-5, 0, and 5 dB)
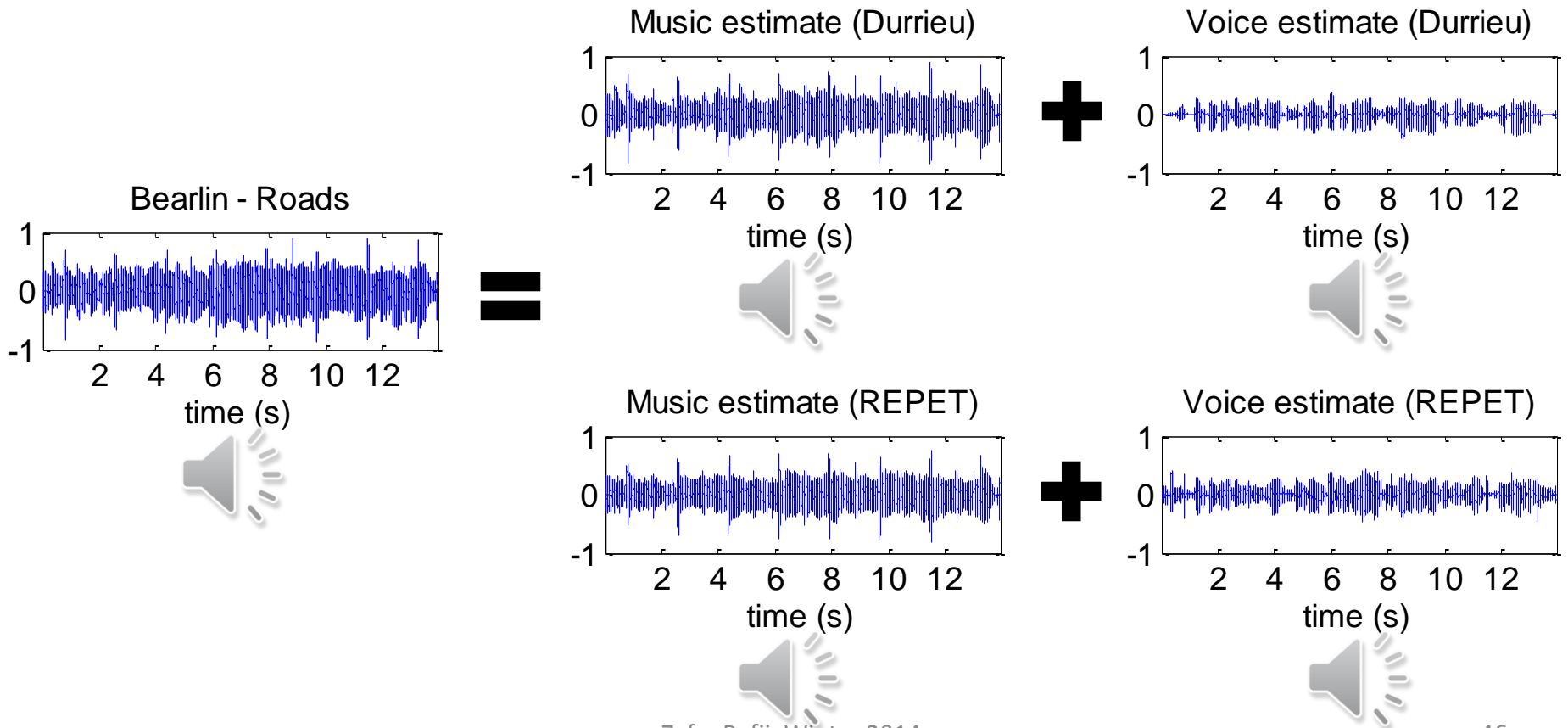
# Evaluation

# Evaluation

- **Conclusions**

  - REPET can compete with state-of-the-art (and more complex) music/voice separation methods

  - There is room for improvement (+ high-pass, + optimal period, + vocal frames)

  - Average computation time: 0.016 second for 1 second of mixture! (vs. 3.863 seconds for Durrieu)

# Examples

- REPET vs. Durrieu (source-filter + NMF)



Music estimate (Durrieu)

Voice estimate (Durrieu)

Bearlin - Roads

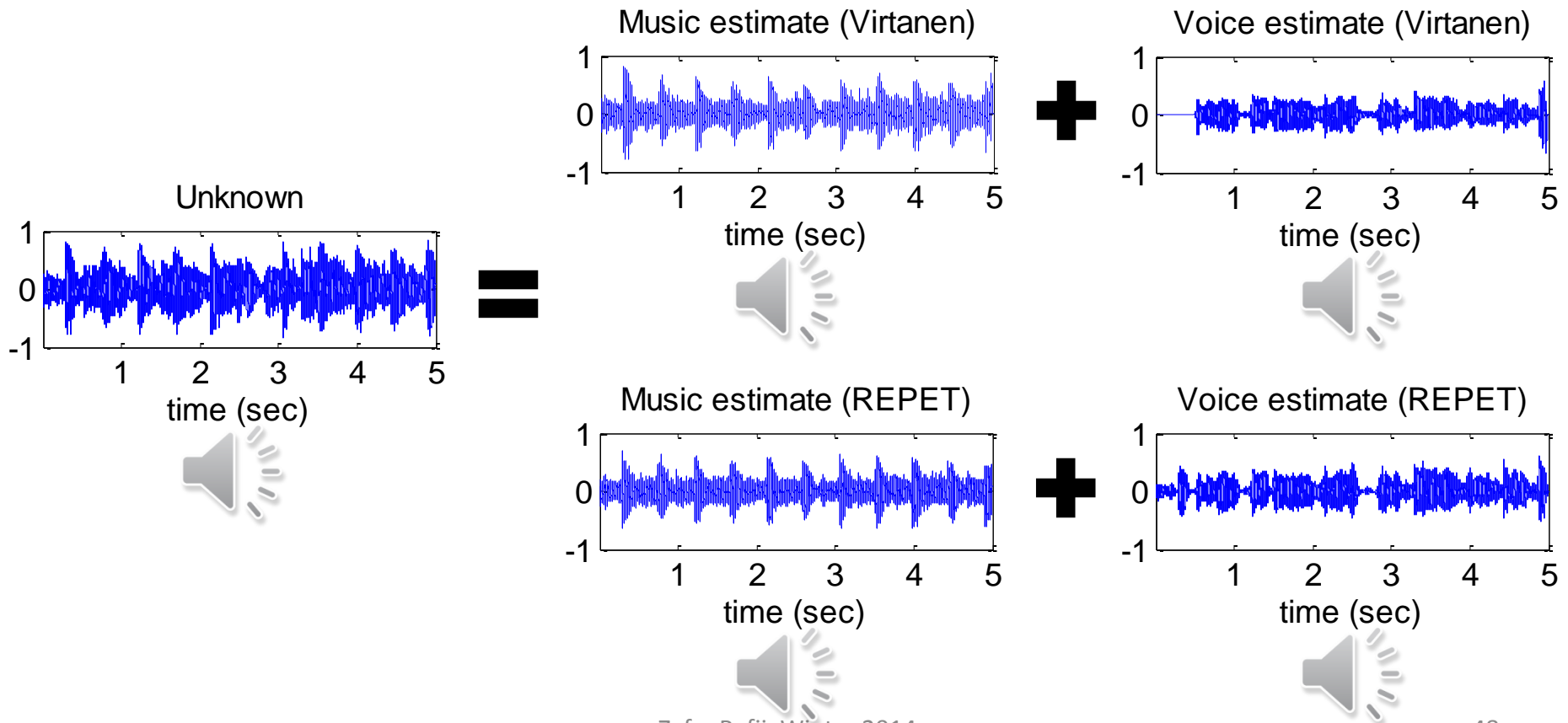Music estimate (REPET)

Voice estimate (REPET)

# Examples

- REPET vs. Ozerov (accompaniment modeling)
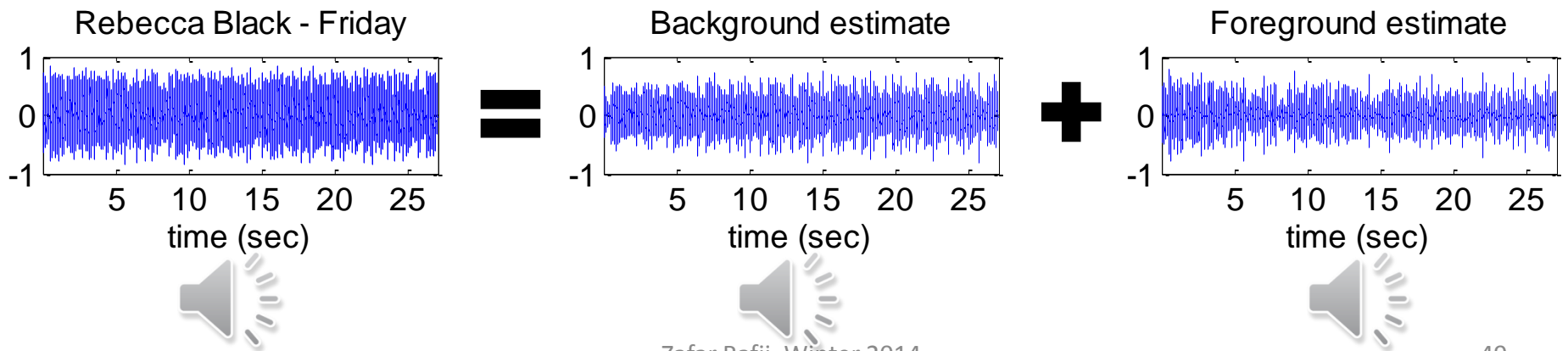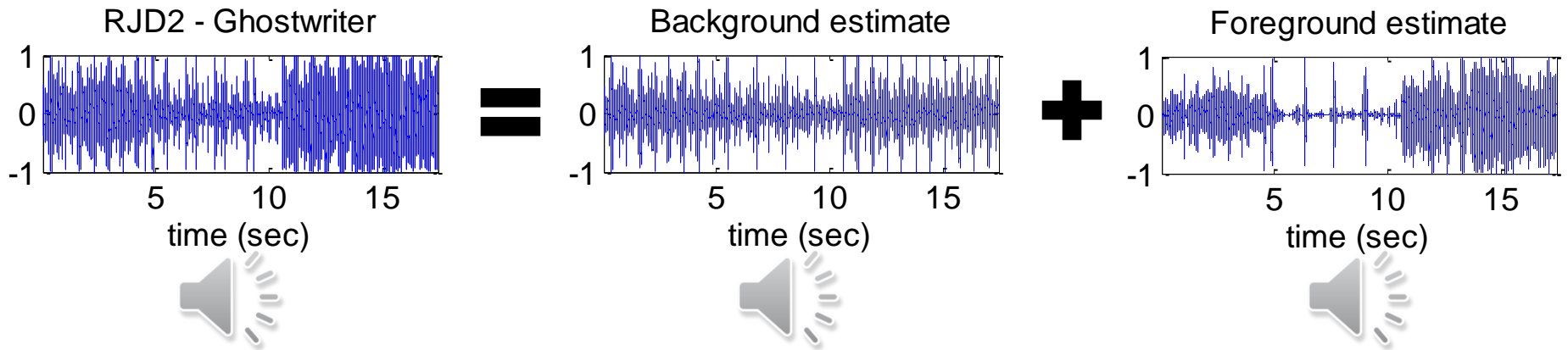
# Examples

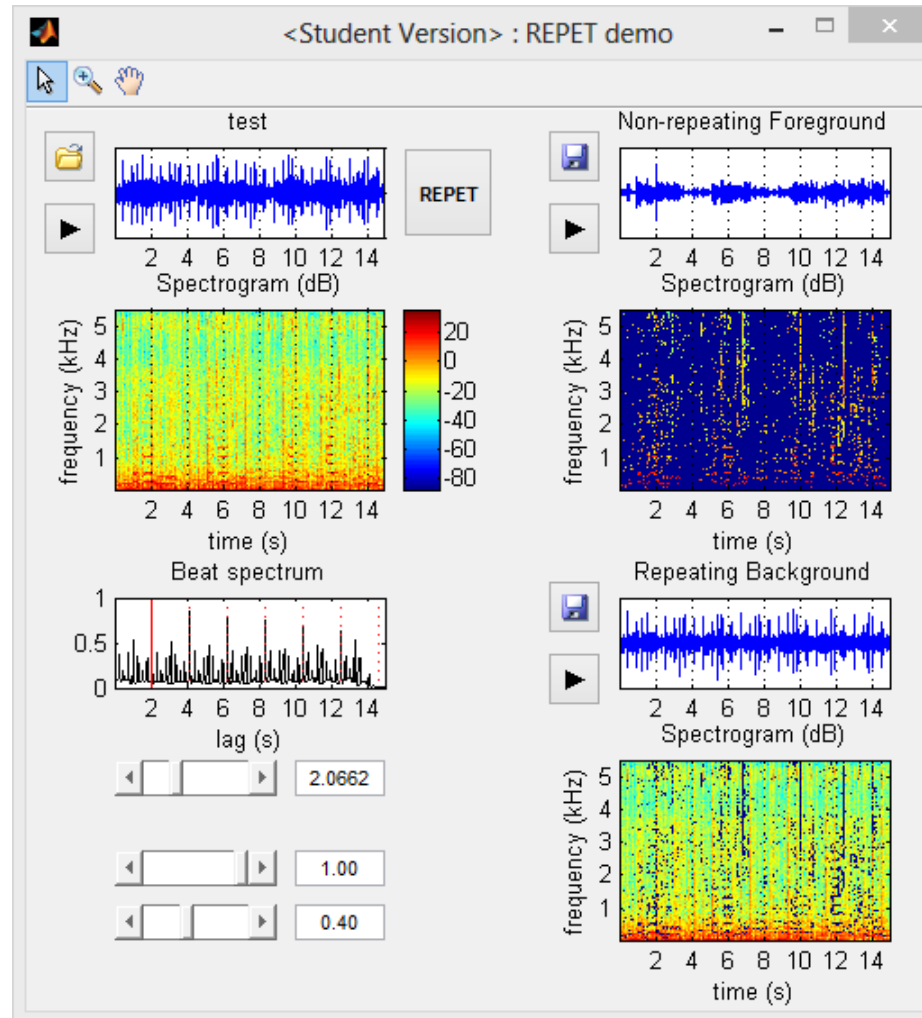- REPET vs. Virtanen (NMF + pitch-based)

# Examples

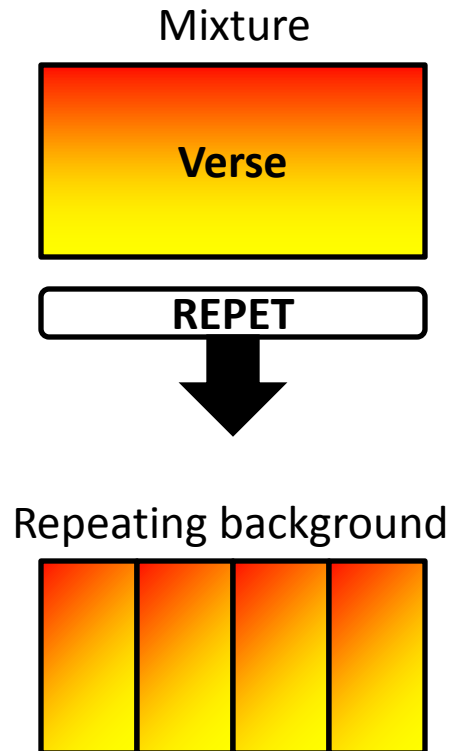- REPET (more examples…)

# Demo

# Thank you!

# References

- J.-L. Durrieu, B. David, and G. Richard, "A Musically Motivated Mid-level Representation for Pitch Estimation and Musical Audio Source Separation," *IEEE Journal on Selected Topics on Signal Processing*, vol. 5, no. 6, pp. 1180-1191, October 2011.

- C.-L. Hsu and J.S. R. Jang, "On the Improvement of Singing Voice Separation for Monaural Recordings Using the MIR-1K Dataset," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 2, pp. 310-319, February 2010.

- A. Liutkus, Z. Rafii, R. Badeau, B. Pardo, and G. Richard, "Adaptive Filtering for Music/Voice Separation exploiting the Repeating Musical Structure," in *37th International Conference on Acoustics, Speech and Signal Processing*, Kyoto, Japan, March 25-30, 2012.

- J. H. McDermott, D. Wrobleski, and A. J. Oxenham, "Recovering Sound Sources from Embedded Repetition," in *National Academy of Sciences*, vol. 108, pp. 1188-1193, 2011.

- A. Ozerov, P. Philippe, F. Bimbot, and R. Gribonval, "Adaptation of Bayesian Models for Single-Channel Source Separation and its Application to Voice/Music Separation in Popular Songs," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 5, pp. 1564-1578, July 2007.

- M. Piccardi, "Background Subtraction Techniques: a Review," *IEEE International Conference on Systems, Man and Cybernetics*, The Hague, Netherlands, October 10-13, 2004.

- Z. Rafii and B. Pardo, "A Simple Music/Voice Separation Method based on the Extraction of the Repeating Musical Structure," *36th International Conference on Acoustics, Speech and Signal Processing*, Prague, Czech Republic, May 22-27, 2011.

- Z. Rafii and B. Pardo, "Music/Voice Separation using the Similarity Matrix," in *13th International Society for Music Information Retrieval*, Porto, Portugal, October 8-12, 2012.

- Z. Rafii and B. Pardo, "REpeating Pattern Extraction Technique (REPET): A Simple Method for Music/Voice Separation," in *IEEE Transactions on Audio, Speech, and Language Processing*, Vol. 21, no. 1, pp. 22-27, January, 2013.

- T. Virtanen, A. Mesaros, and M. Ryynänen, "Combining Pitch-based Inference and Non-Negative Spectrogram Factorization in Separating Vocals from Polyphonic Music," *ISCA Tutorial and Research Workshop on Statistical and Perceptual Audition*, Brisbane, Australia, pp. 17-20, September 21, 2008.
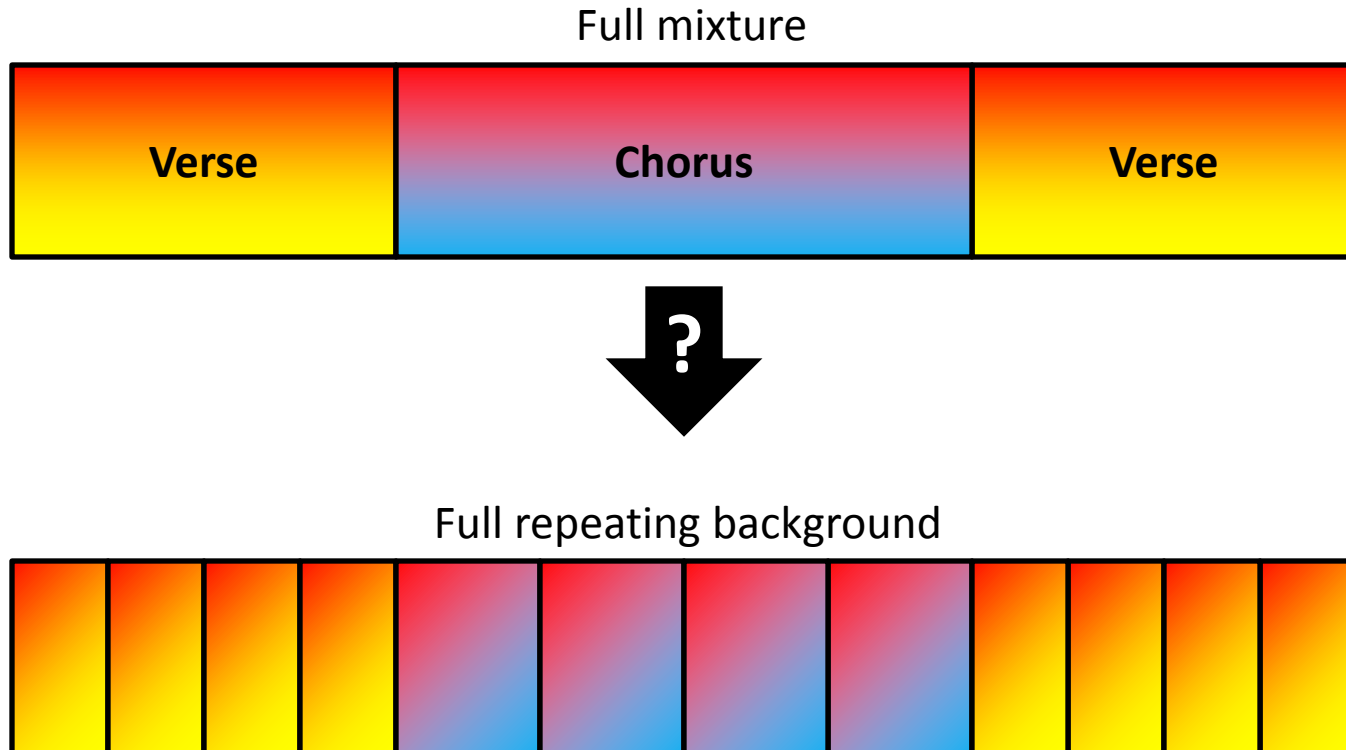
# Extensions

- REPET works well on excerpts with a relatively **stable repeating background** (e.g., 10 s verse)
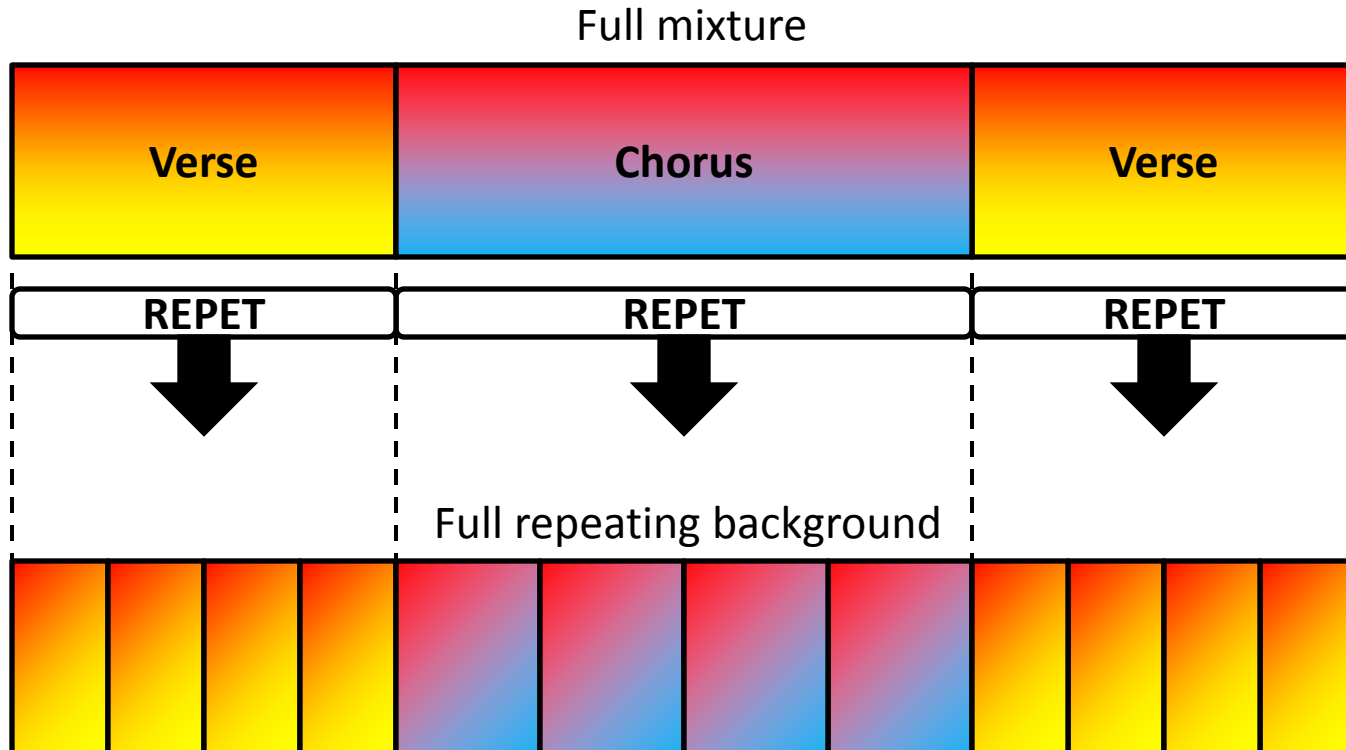
Mixture

Verse

REPET

Repeating background

# Extensions

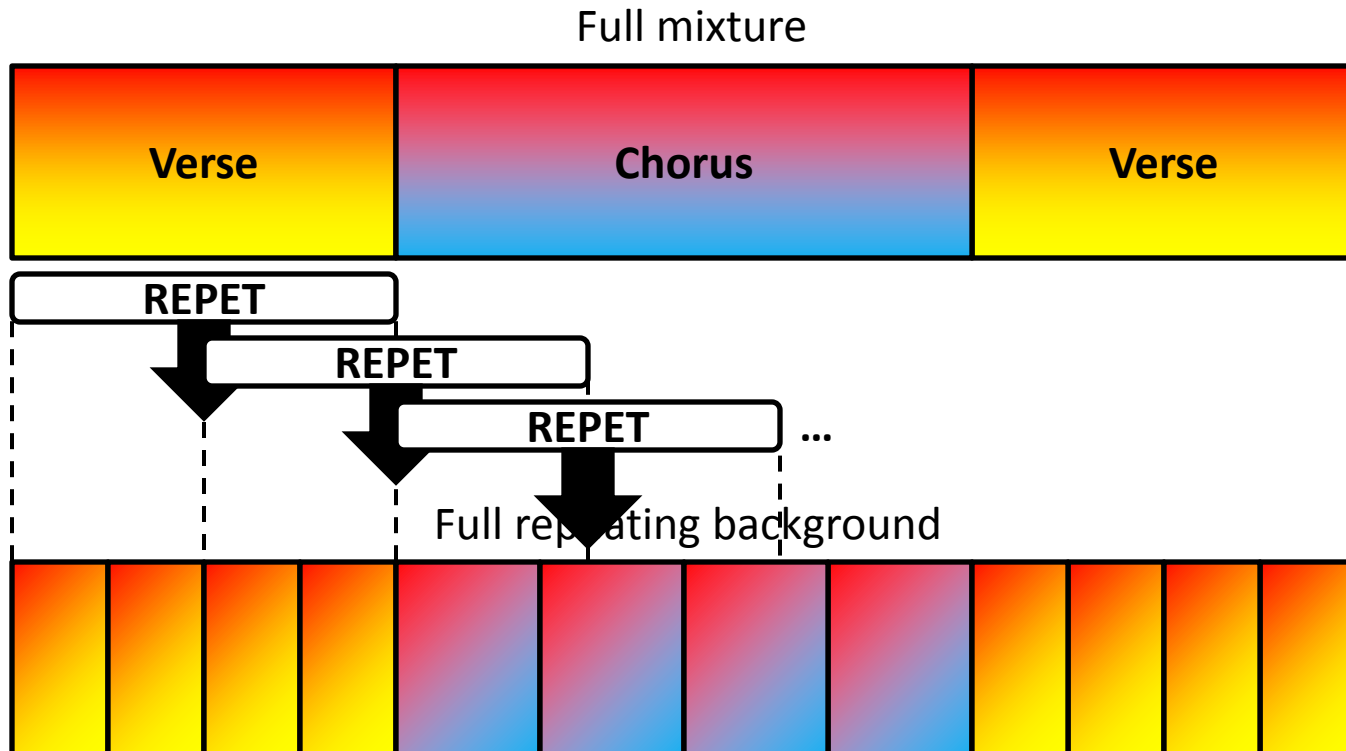- For full-track songs, the repeating background is likely to **vary over time** (e.g., verse/chorus)

Full mixture

| Verse | Chorus | Verse |

?

Full repeating background

# Prior Segmentation

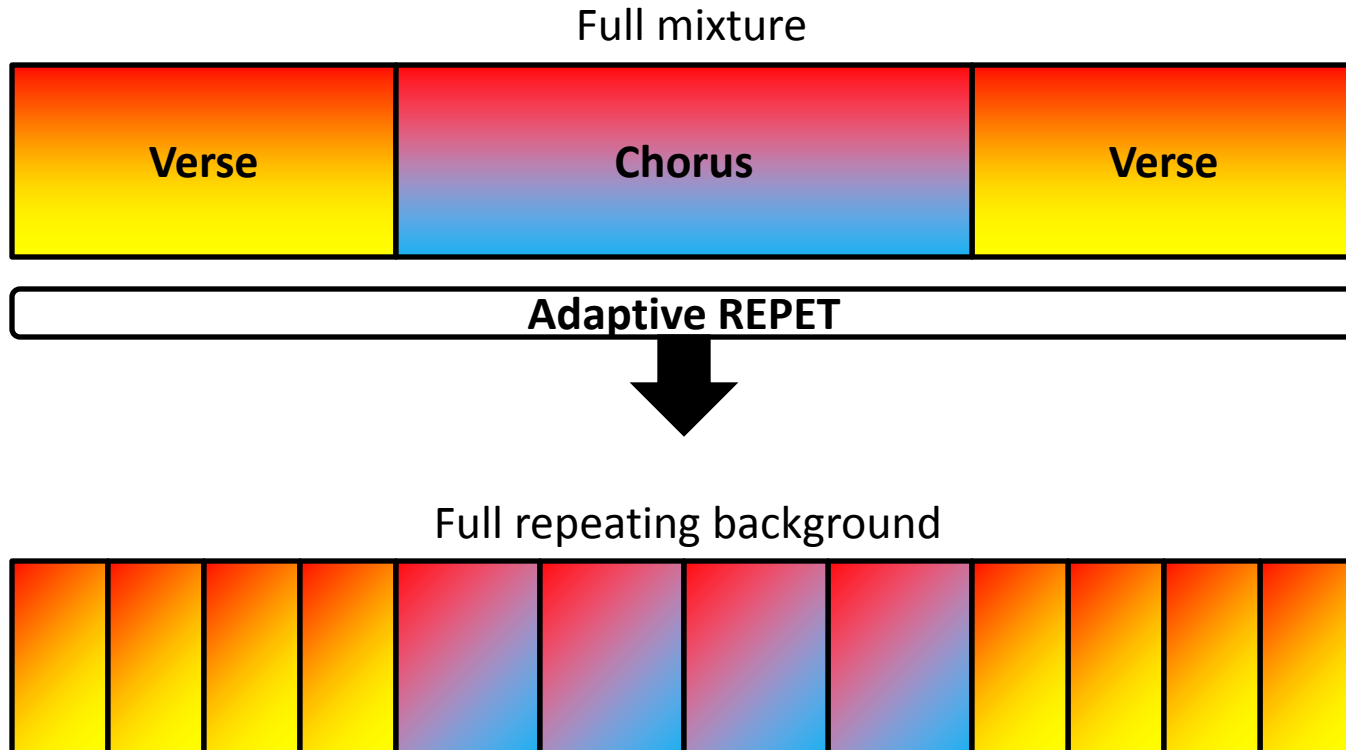- We could do a **prior segmentation** of the song and apply REPET to the individual sections

Full mixture

| Verse | Chorus | Verse |
|---|---|---|
| REPET | REPET | REPET |

Full repeating background

# Sliding Window

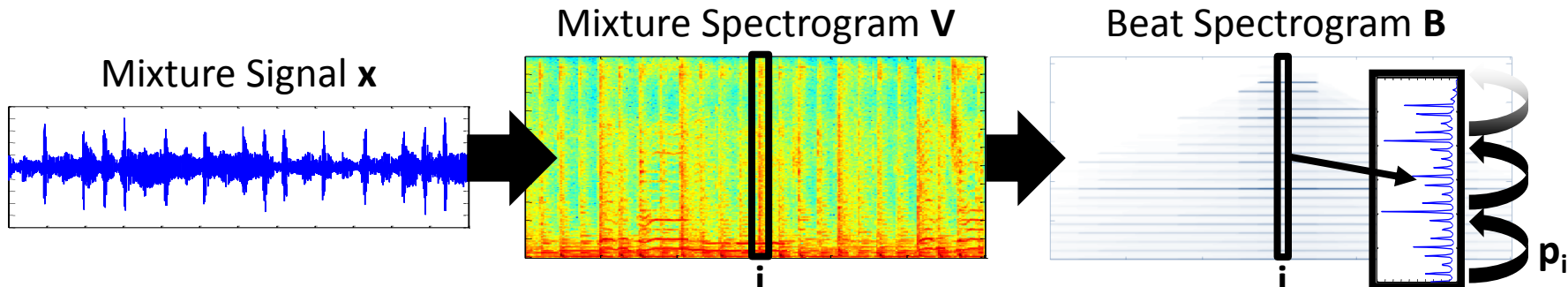- We could apply REPET to local sections of the song over time via a fixed **sliding window**

Full mixture

# Adaptive REPET

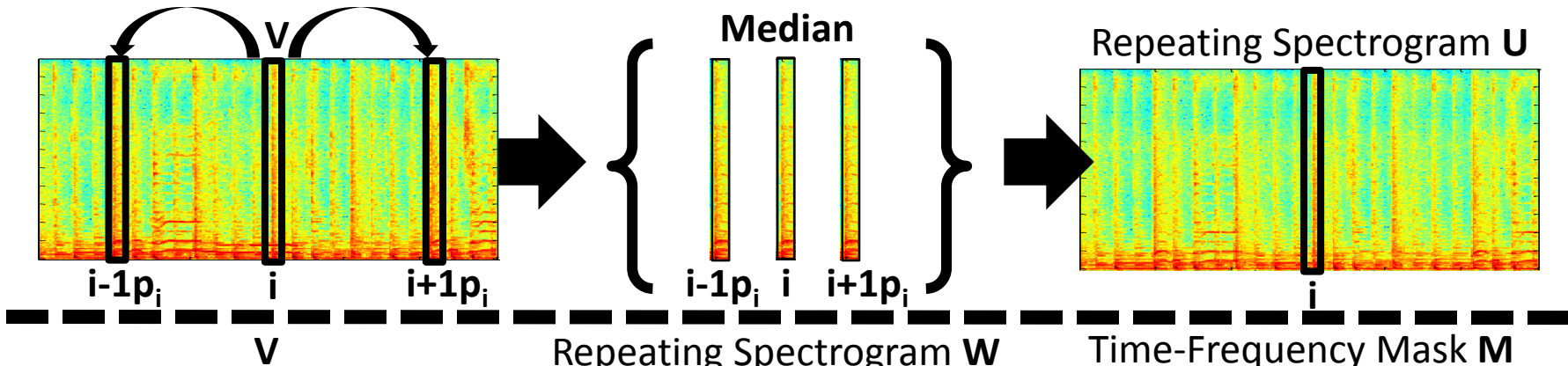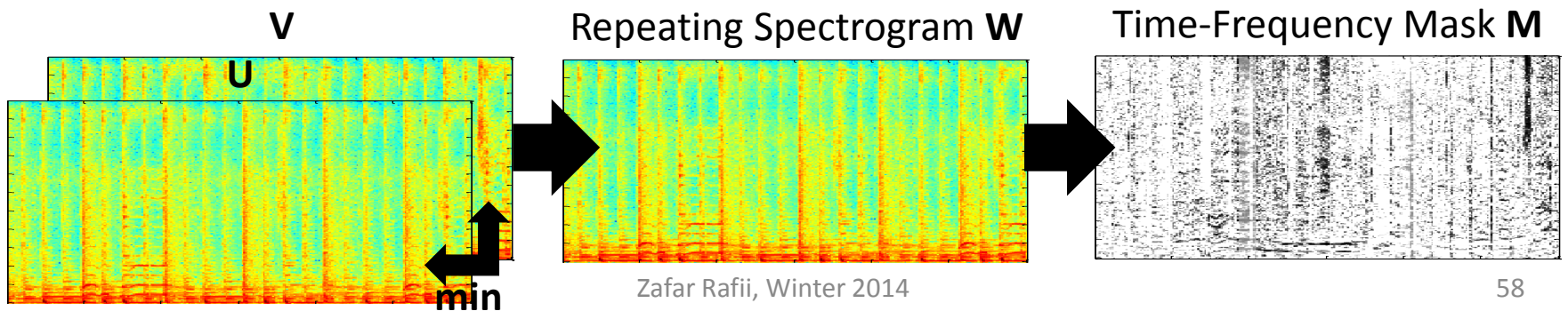- We could directly **adapt REPET** along time by locally modeling the repeating background

Full mixture

| Verse | Chorus | Verse |

**Adaptive REPET**

Full repeating background

# Adaptive REPET

# Original REPET



**Step 1**

Mixture Signal **x** → Mixture Spectrogram **V** → Beat Spectrum **b**

**p**

**Step 2**

**V** → **Median** → Repeating Segment **S**

1p    2p

**Step 3**

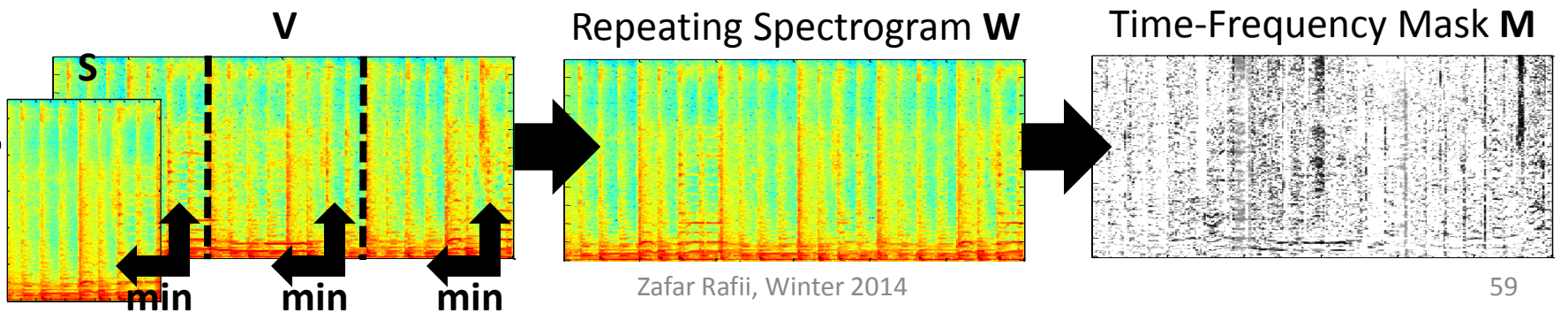**V** (with **S**) → Repeating Spectrogram **W** → Time-Frequency Mask **M**

min    min    min

# Generalization

- REPET (and its extension) assumes **periodically repeating patterns**
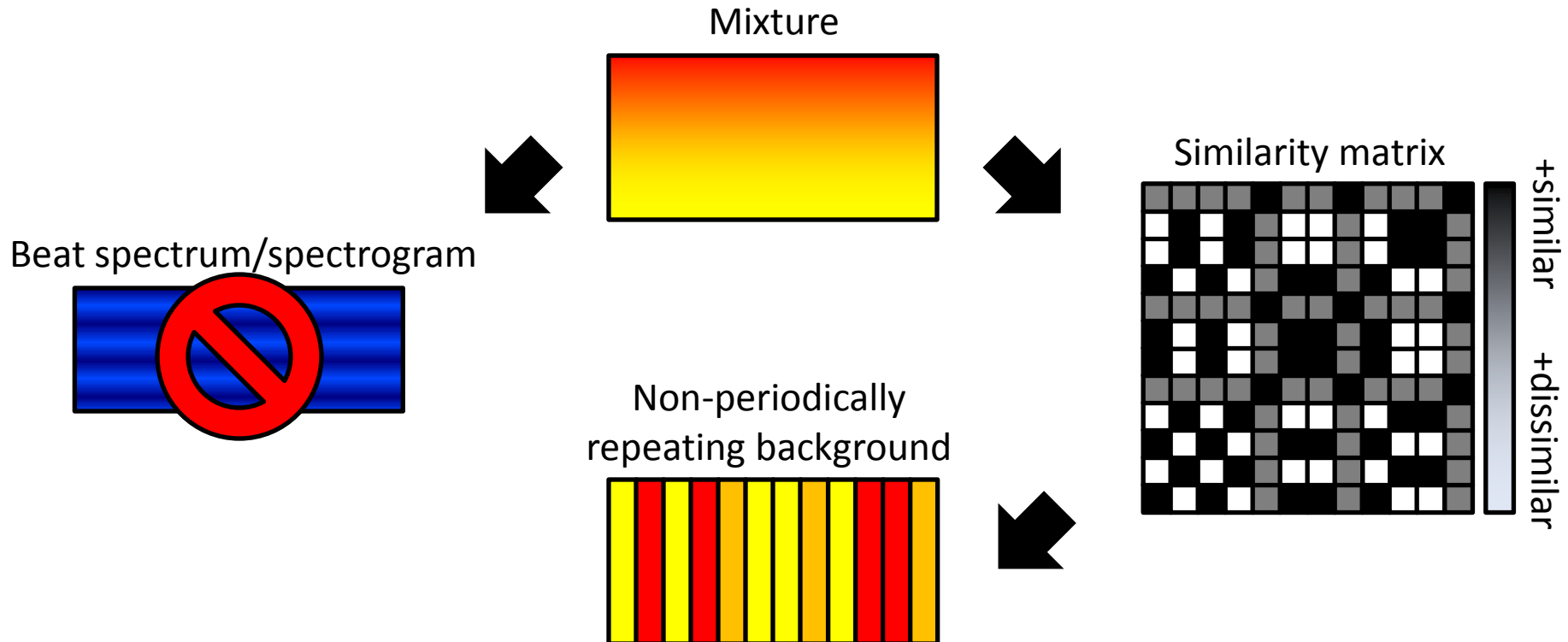
Mixture

REPET

Periodically
repeating background

# Generalization

- Repetitions can also happen **intermittently** or **without a global (or local) period**

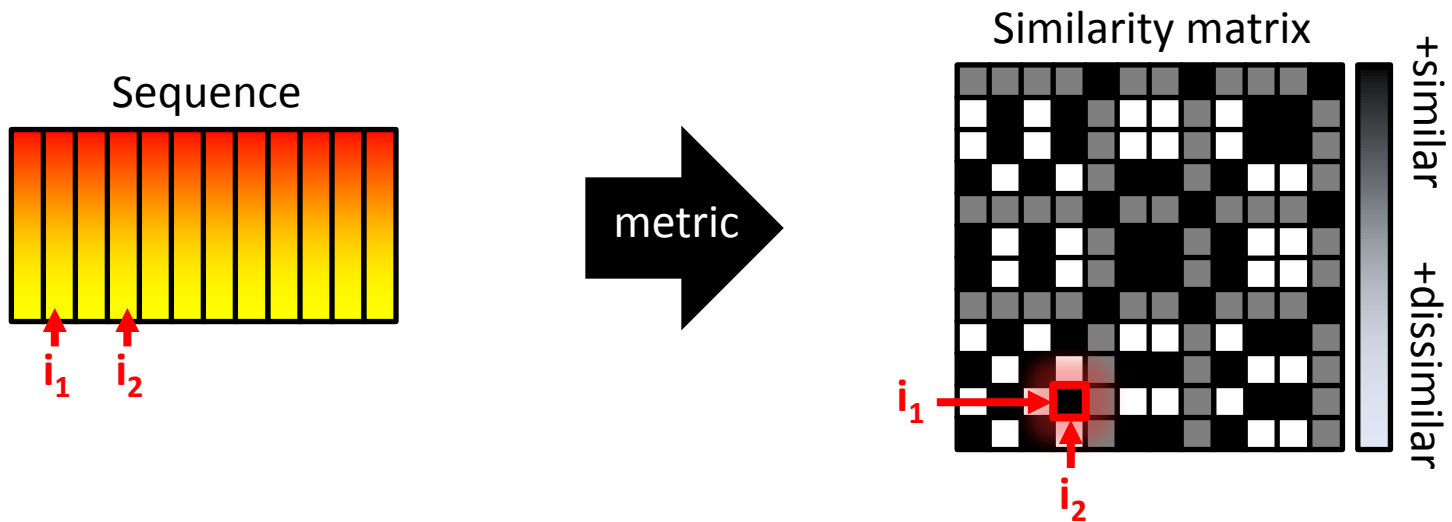Mixture

?

Non-periodically
repeating background

# Generalization

- Instead of looking for periodicities, we can look for similarities, using a **similarity matrix**

Mixture

Beat spectrum/spectrogram

Similarity matrix

+similar

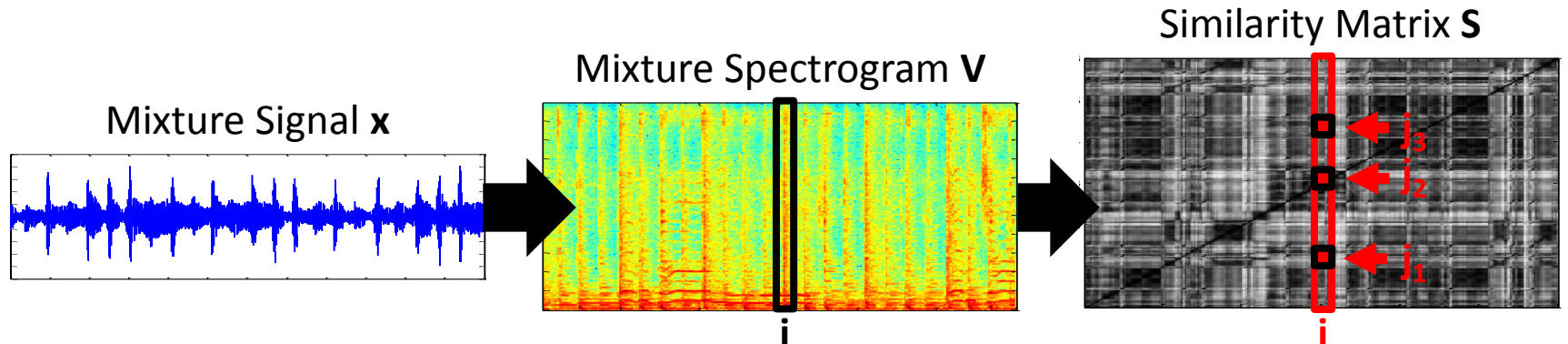+dissimilar

Non-periodically
repeating background

# Generalization

- The **similarity matrix** is a matrix where each bin measures the (dis)similarity between any two elements of a sequence given a metric
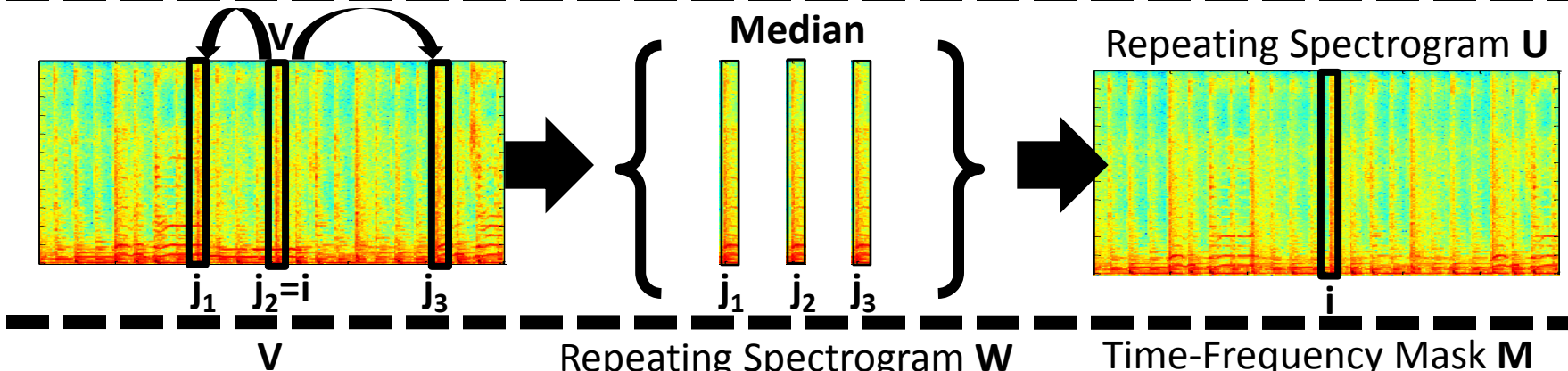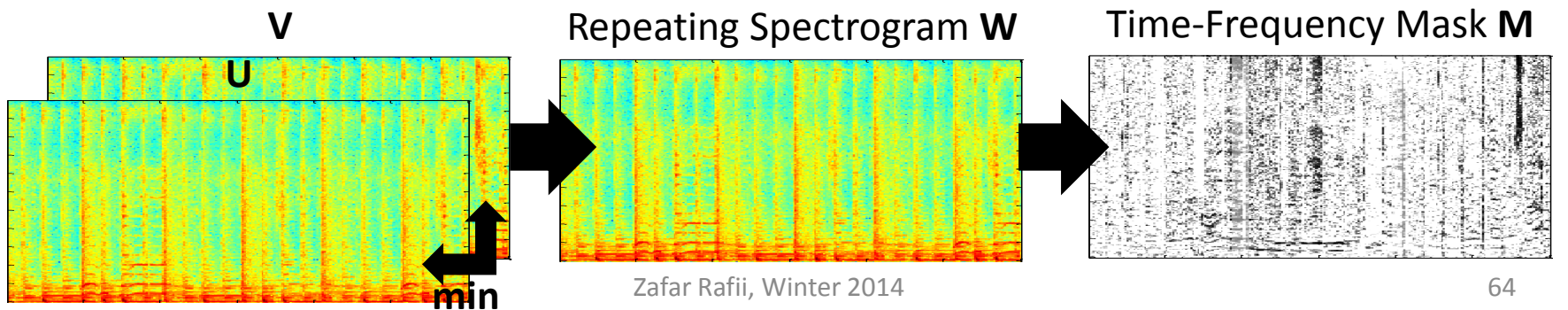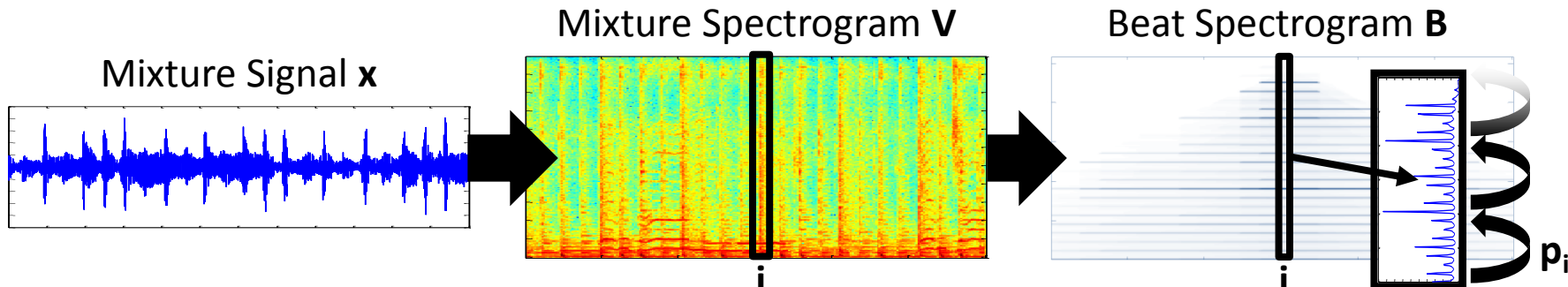


Sequence

metric

Similarity matrix

+similar

+dissimilar

# REPET-SIM

# Adaptive REPET

# REPET + Pitch

- ## REPET models the **background rhythm**

Mixture spectrogram     "Period"     Rhythmic mask     Background rhythm

*frequency* ↑  *time* →

*repetitions (in time)*

- ## Pitch-based methods model the **lead melody**

Mixture spectrogram     "Pitch"     Melodic mask     Lead melody

*frequency* ↑  *time* →

*harmonics (in frequency)*

# REPET + Pitch

- **Auditory processing** in human listeners