

ON THE IMPORTANCE OF
TEMPORAL CONTEXT IN PROXIMITY KERNELS:
A VOCAL SEPARATION CASE STUDY



SOURCE 1



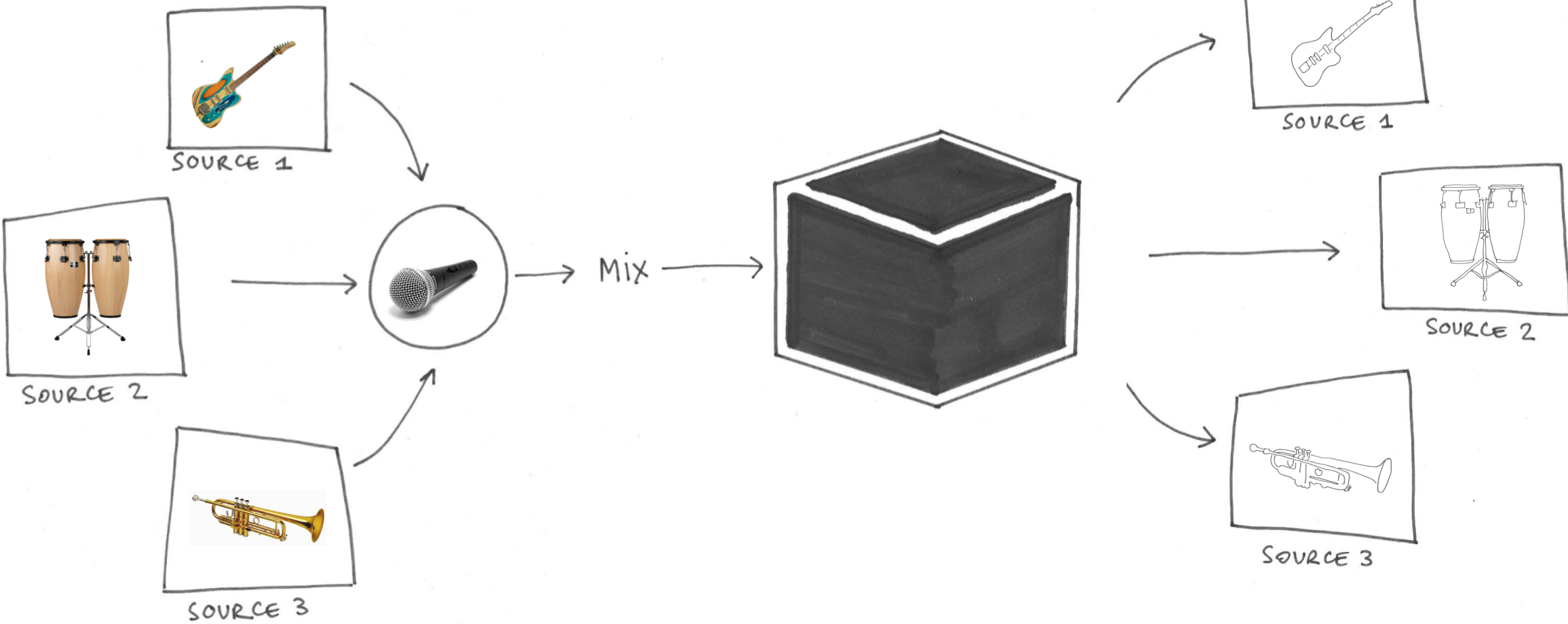
SOURCE 2



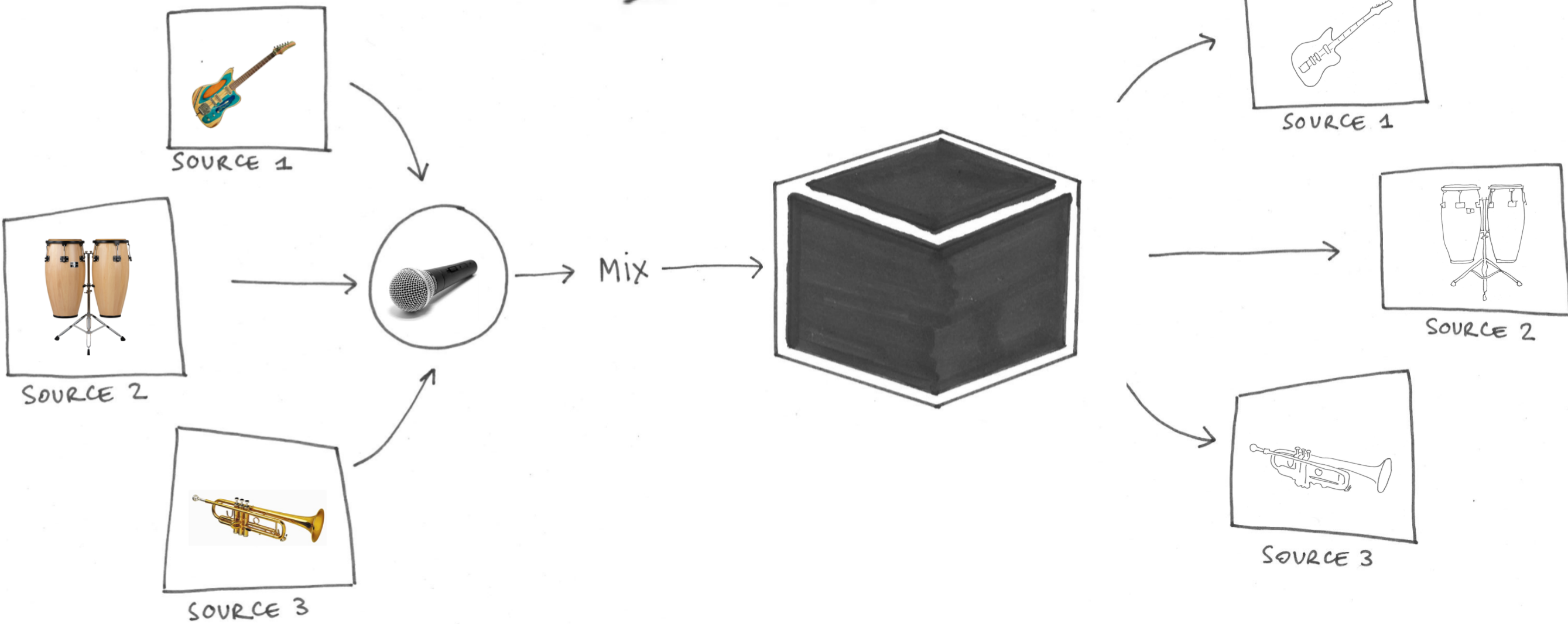
SOURCE 3

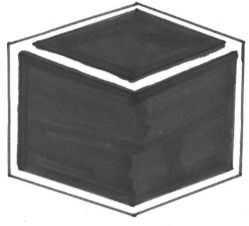


Mix

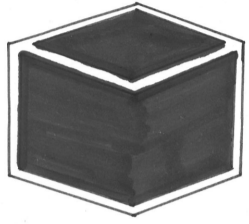


SOURCE SEPARATION



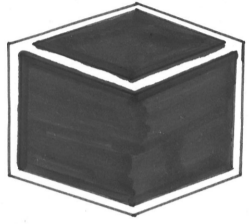


... WHY IS IT USEFUL?



... WHY IS IT USEFUL?

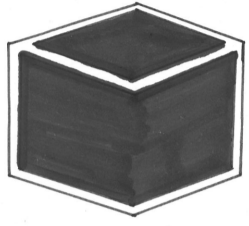




..... WHY IS IT USEFUL?

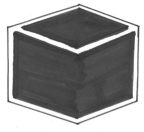
- UP MIXING
- SIGNAL DENOISING
- AUDIO RESTORATION
- KARAOKE
- MUSICAL FEATURE EXTRACTION
- AUTOMATIC TRANSCRIPTION
-





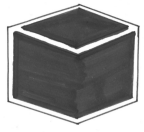
...

HOW DOES IT WORK?

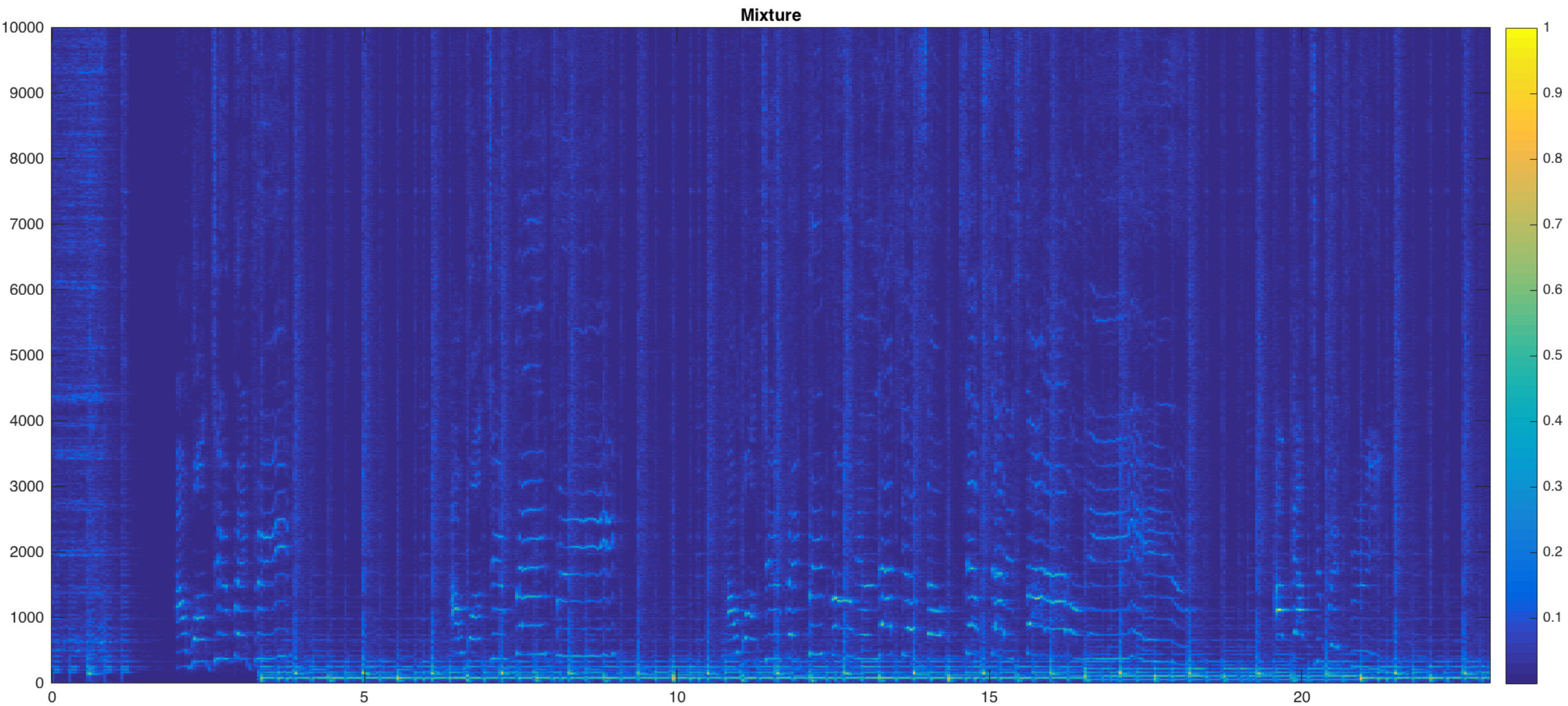


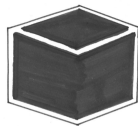
A SIMPLE EXAMPLE : NO OVERLAP

The screenshot shows an audio workstation interface with a track named "09stan by meottis reding.mp3". The track is equipped with several GraphicEQ plugins. The main window displays the audio waveform with a vertical red line indicating the current playback position. A purple highlight is visible on the waveform, and a blue highlight is visible on the EQ window. The EQ window is titled "FX: Track 1 '09stan by meottis reding'" and shows the "AU: Apple: AUGraphicEQ" plugin. The EQ window displays a 31-band equalizer with a 1.6kHz band width and -20 dB gain. The frequency range is from 20Hz to 20kHz, and the gain range is from -20 dB to 20 dB. The EQ window also shows "Flatten EQ" and "31 Bands" options. The CPU usage is 0.5%/0.5% and the sample rate is 0/0 spls.

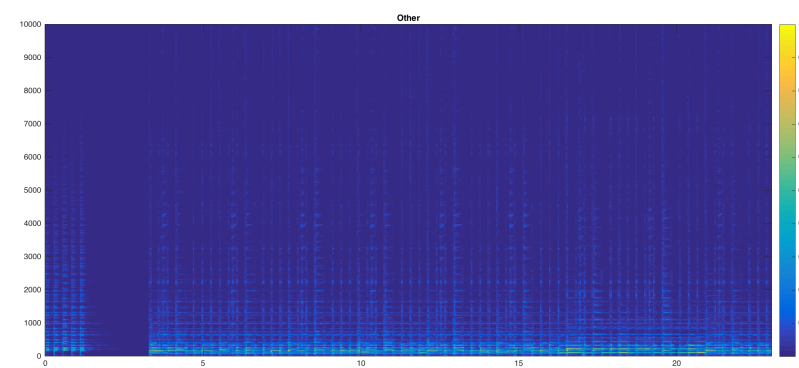
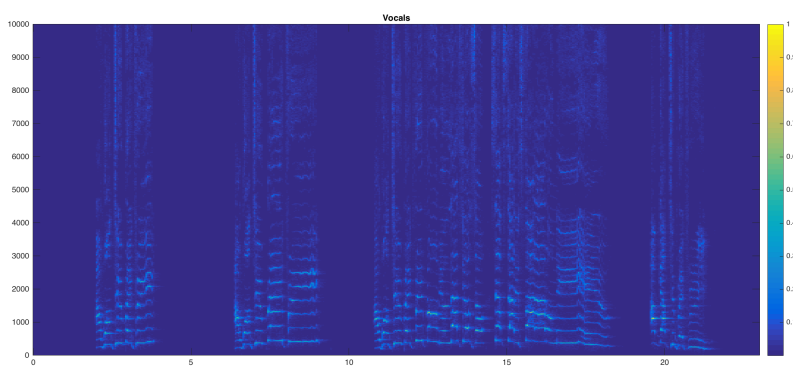
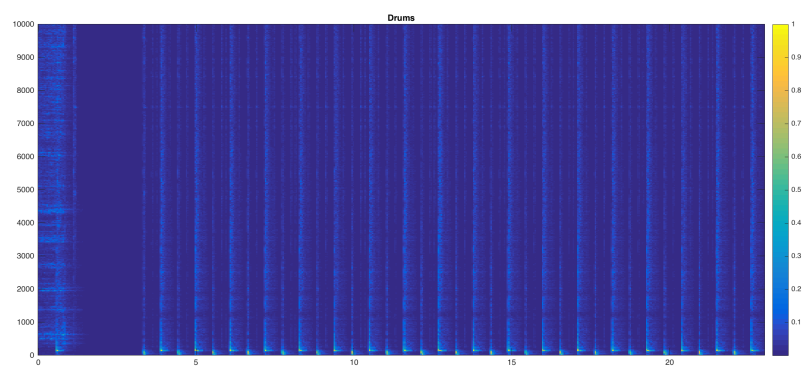
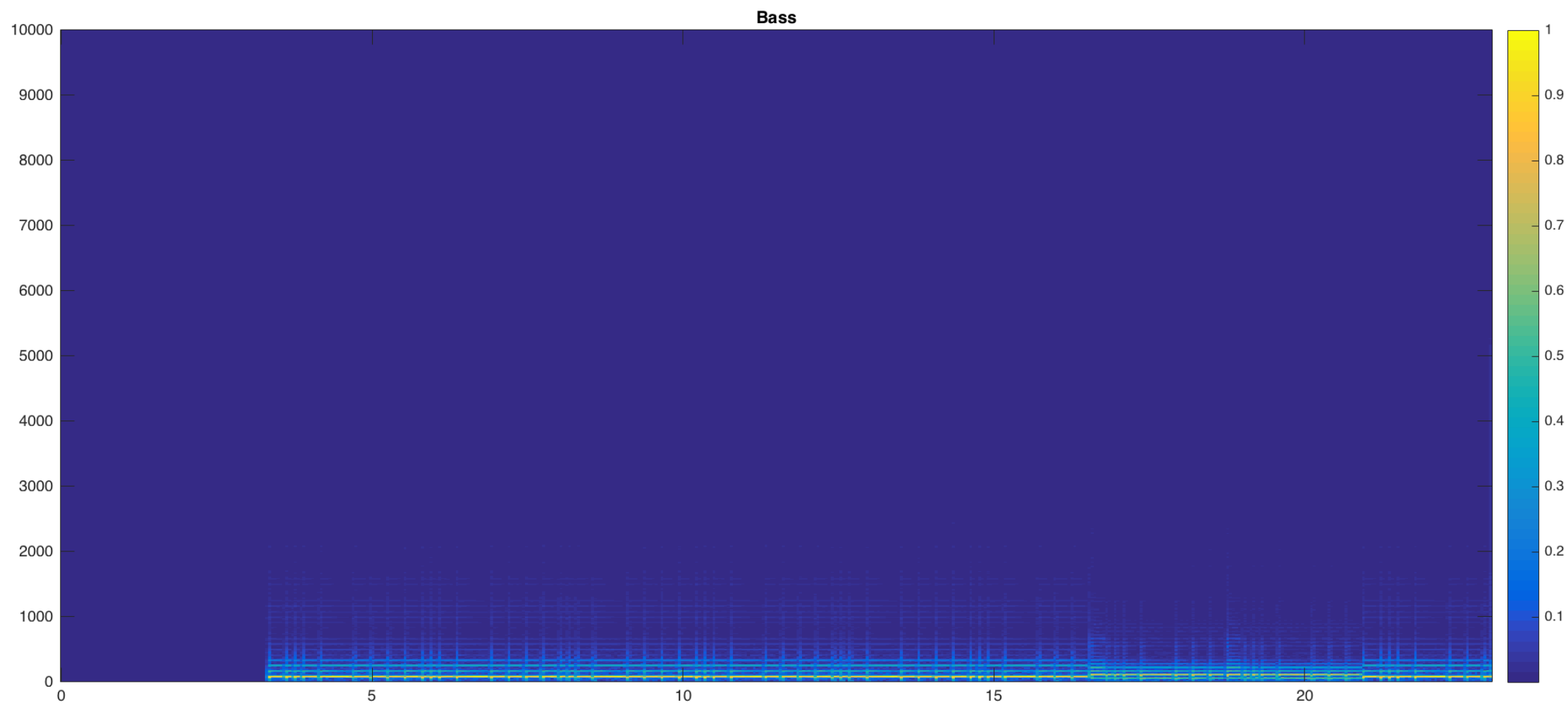


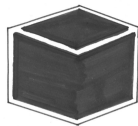
..... MOST POPULAR SONGS: SOURCES DO OVERLAP



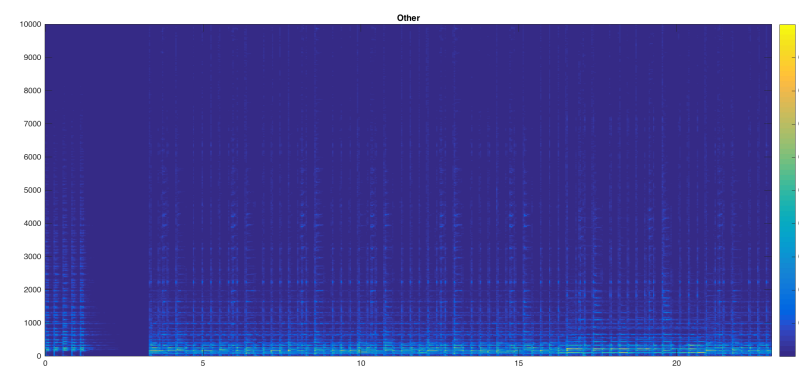
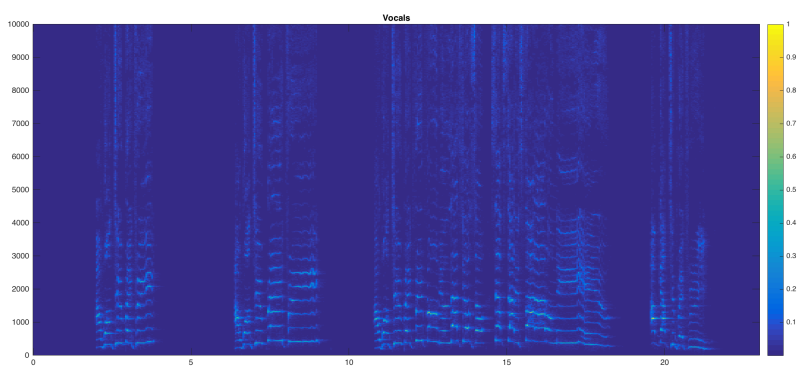
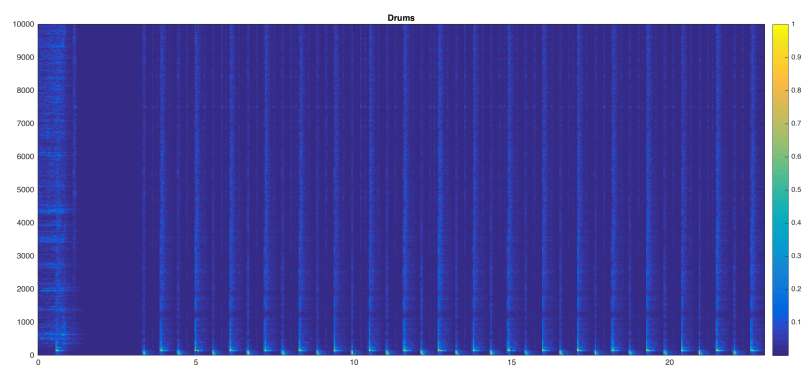
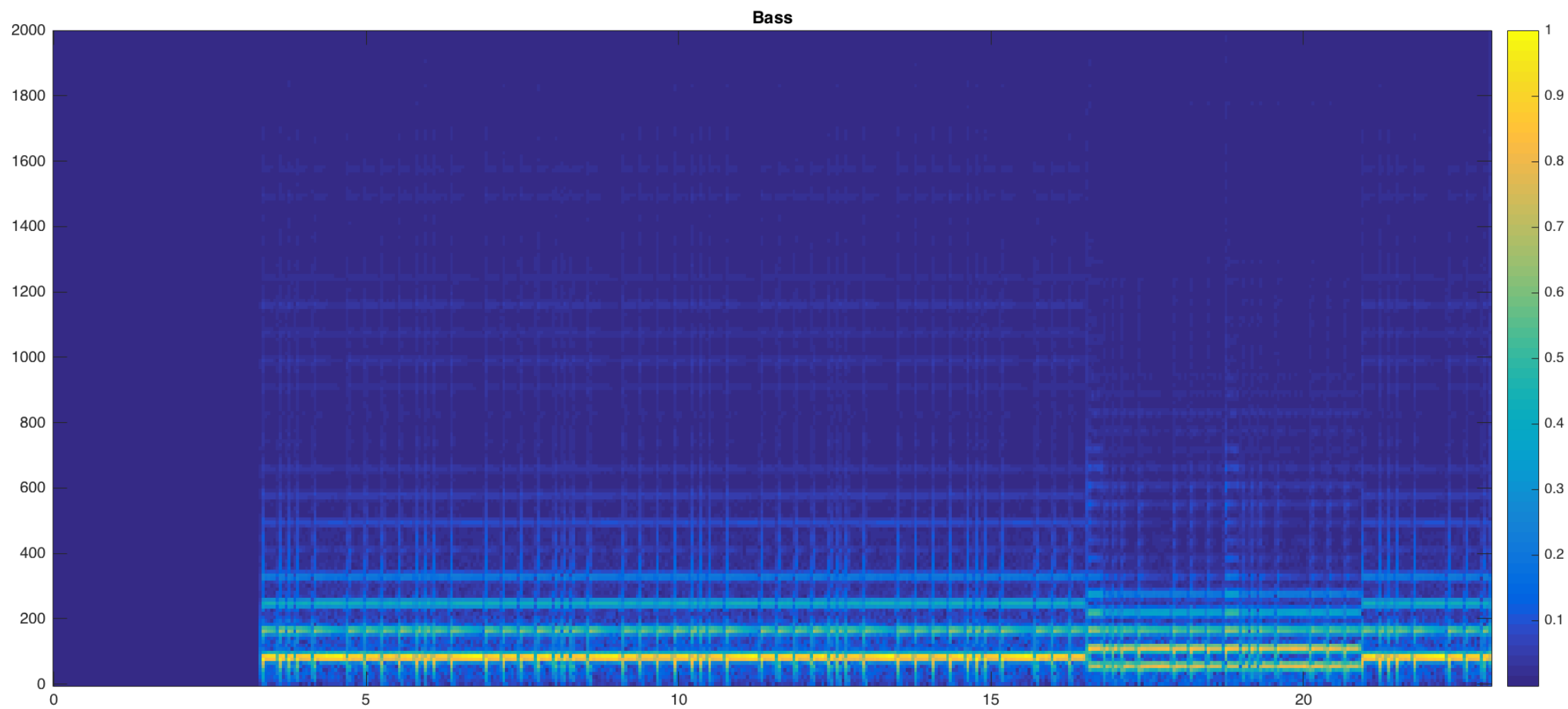


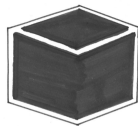
WHAT WE ARE AIMING FOR



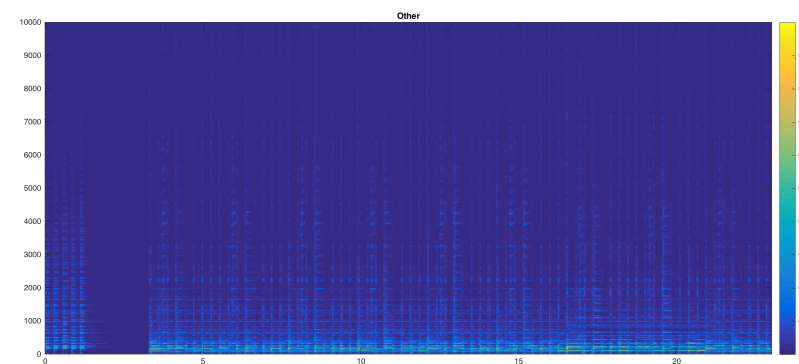
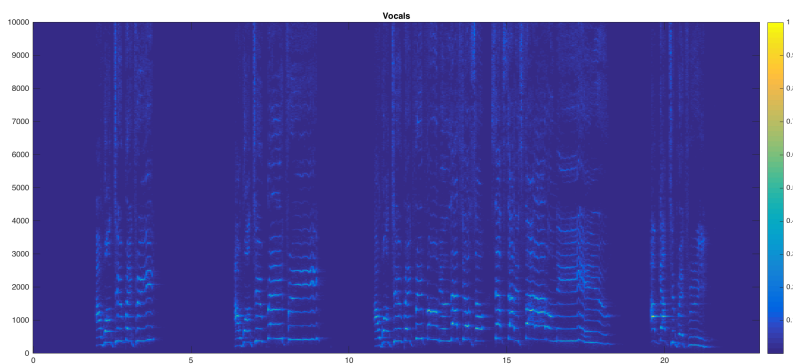
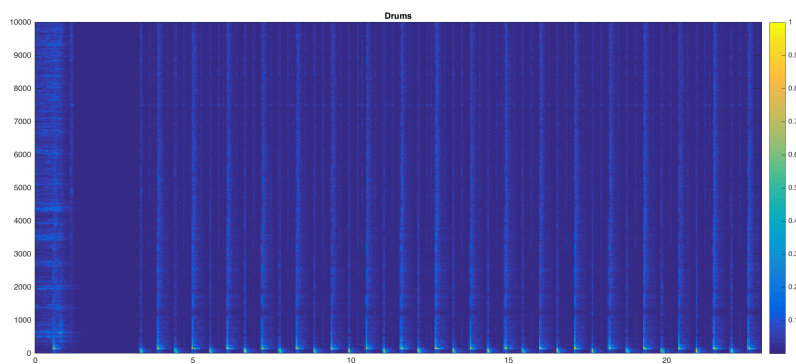
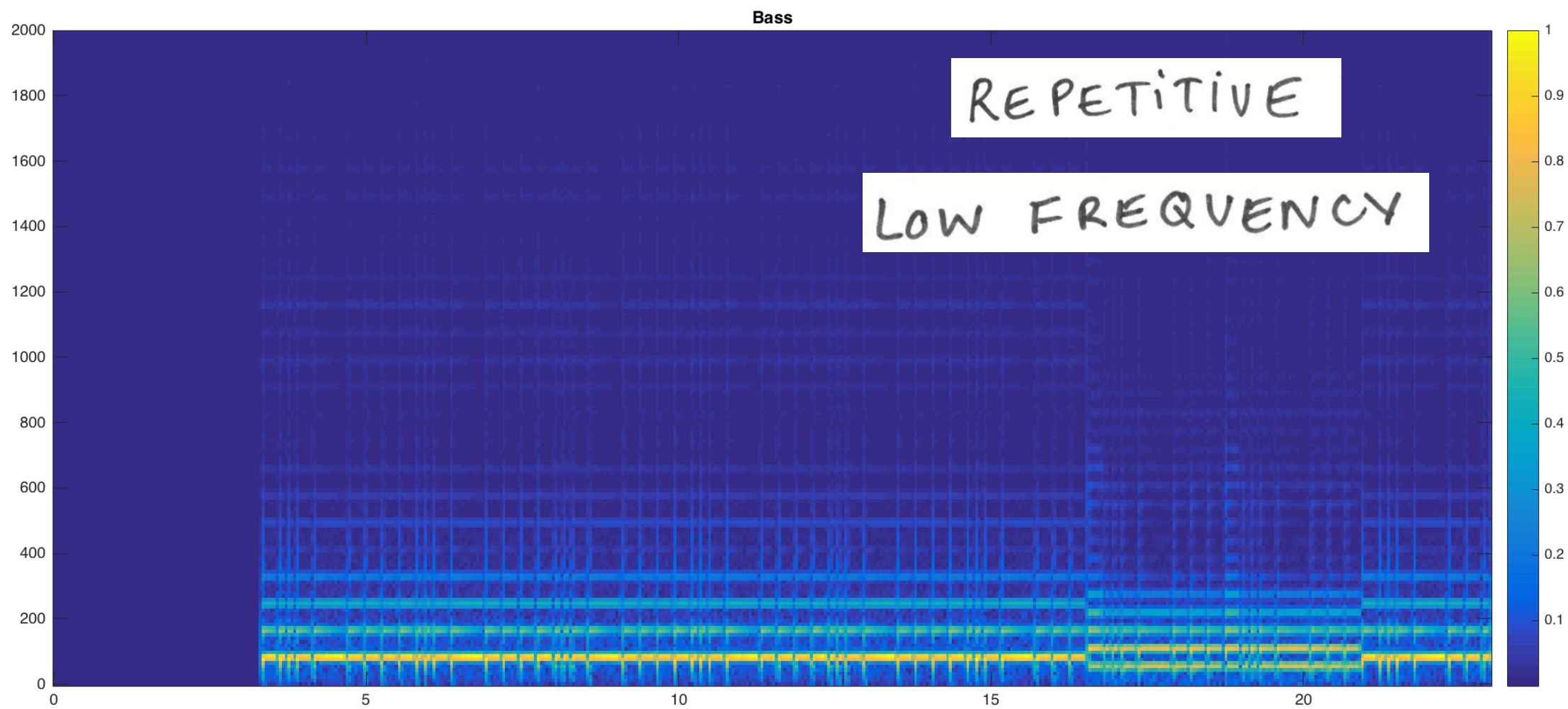


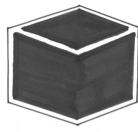
WHAT WE ARE AIMING FOR



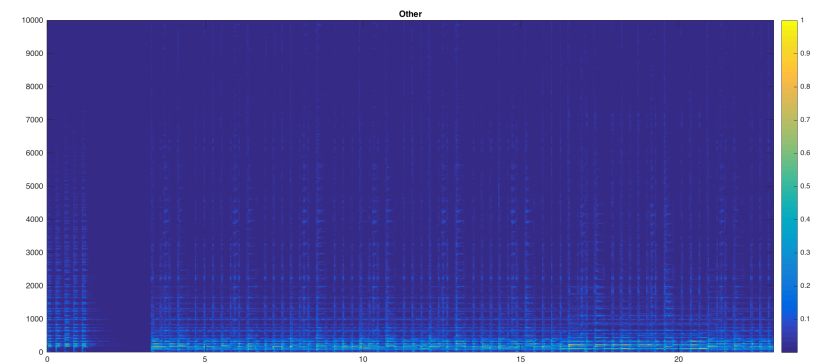
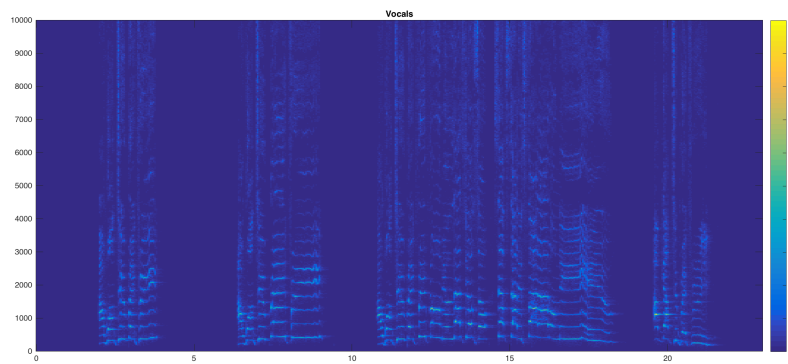
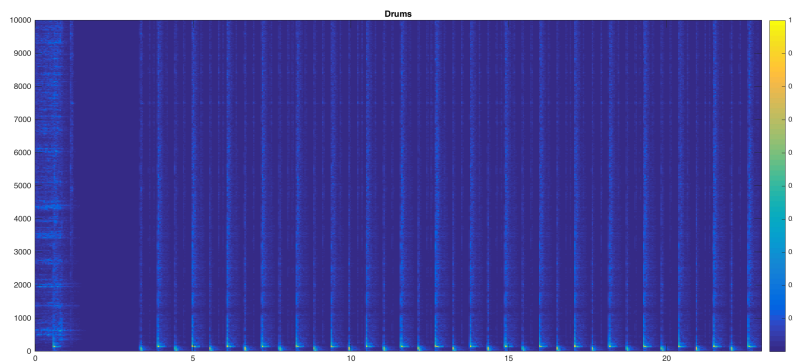
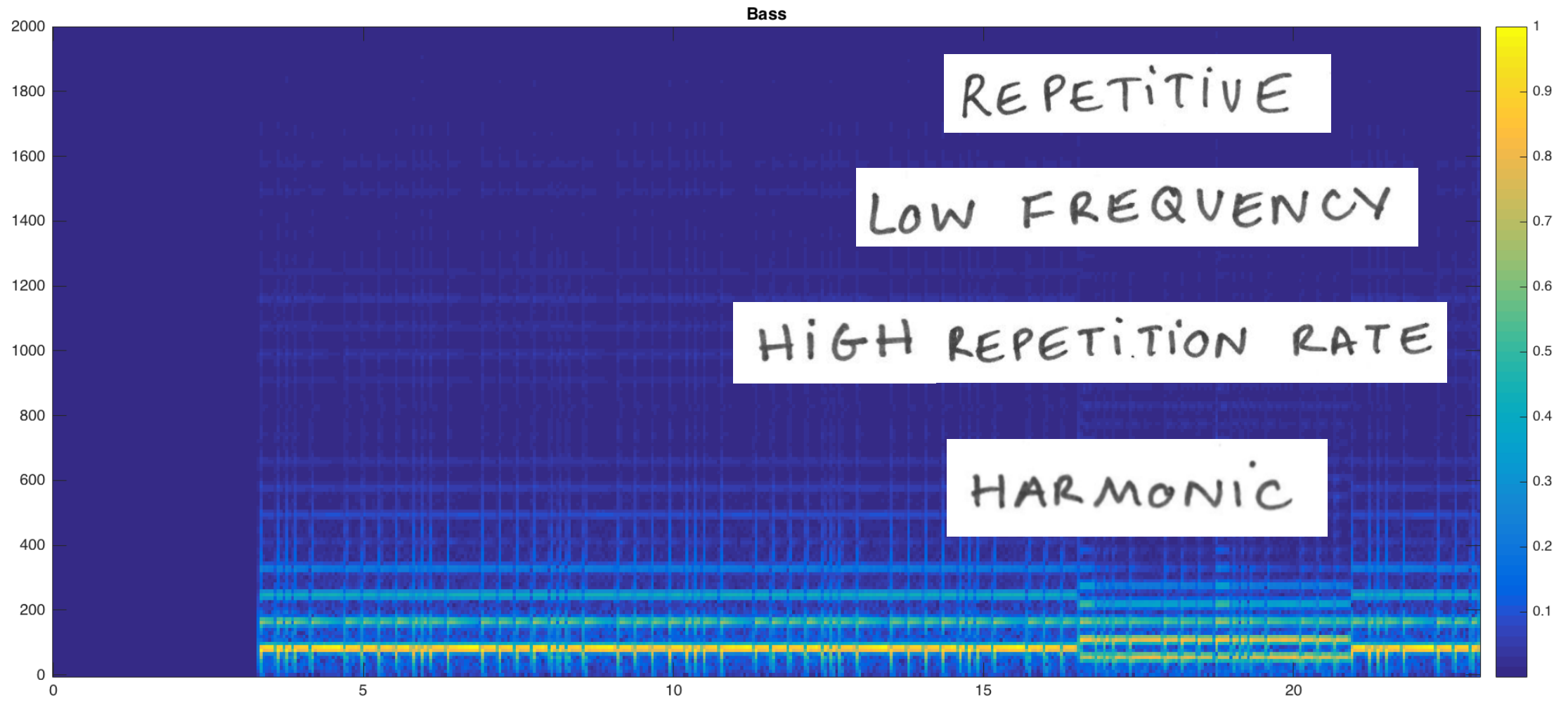


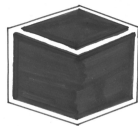
WHAT WE ARE AIMING FOR



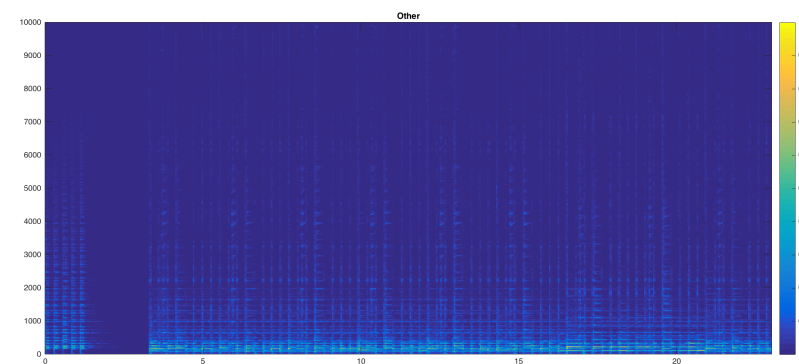
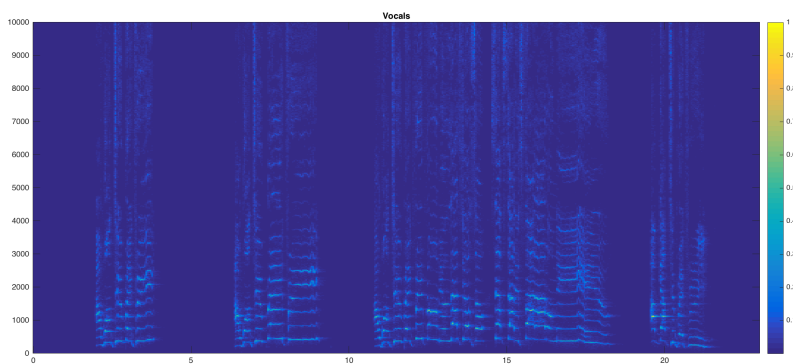
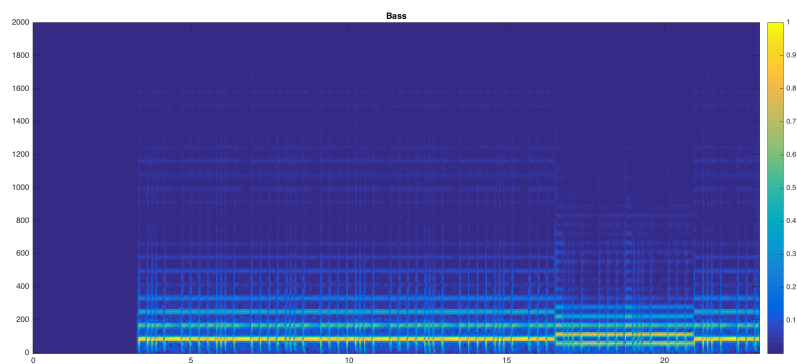
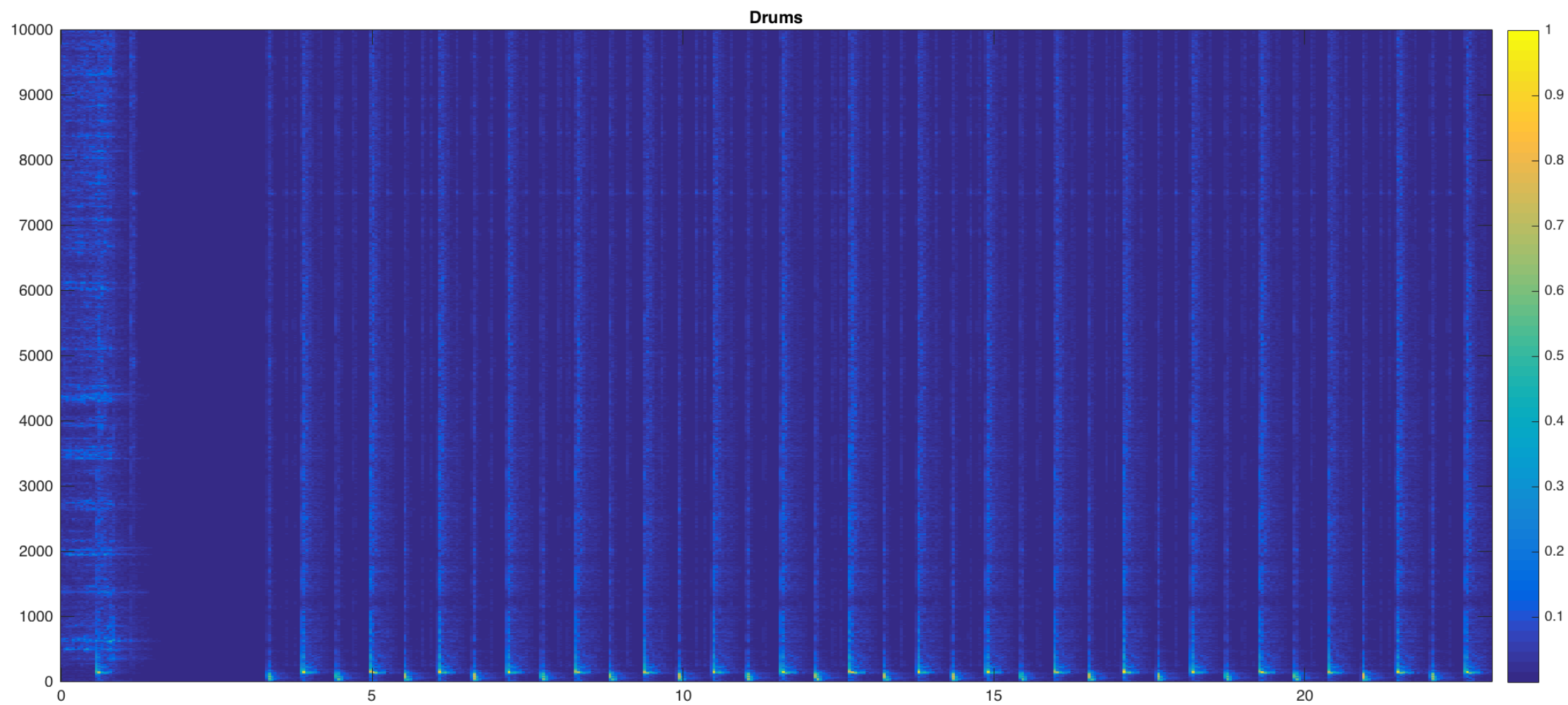


WHAT WE ARE AIMING FOR

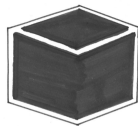




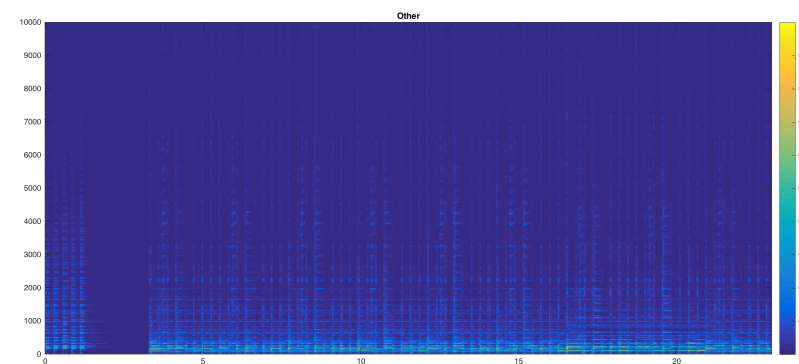
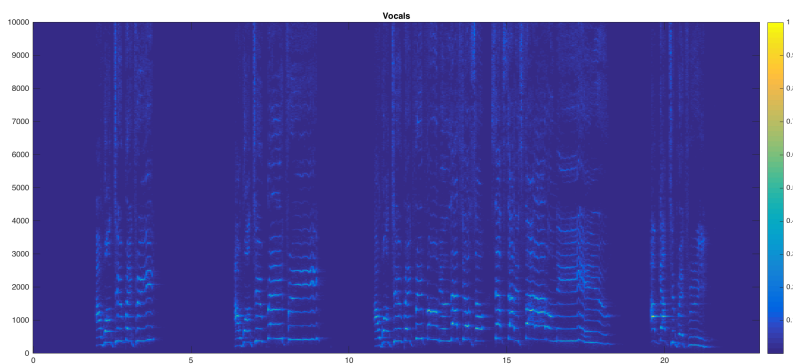
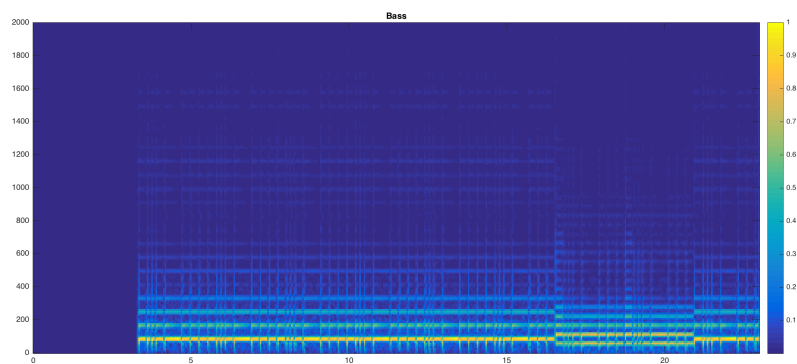
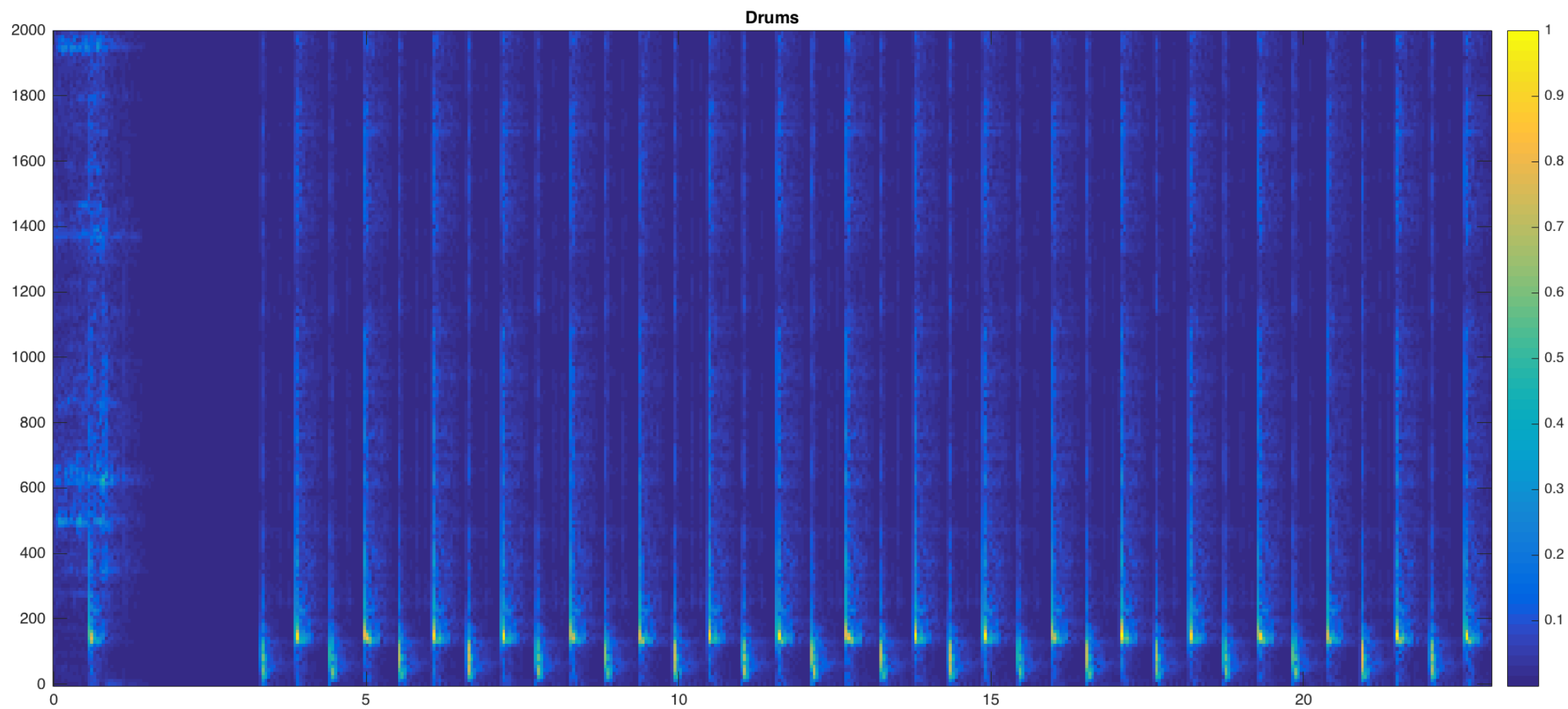
WHAT WE ARE AIMING FOR



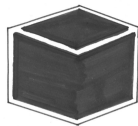
REPETITIVE
LOW FREQUENCY
REPETITION RATE
HARMONIC



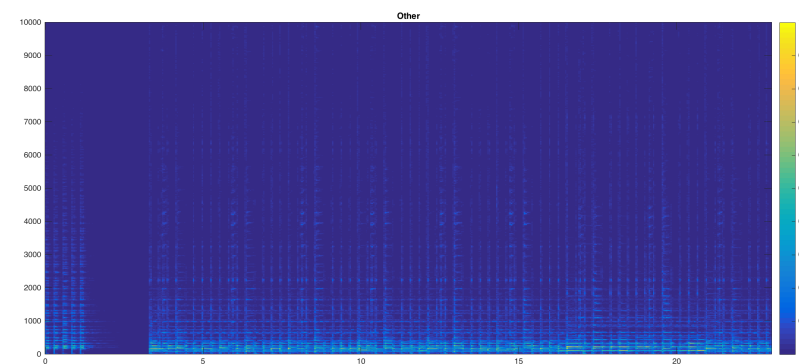
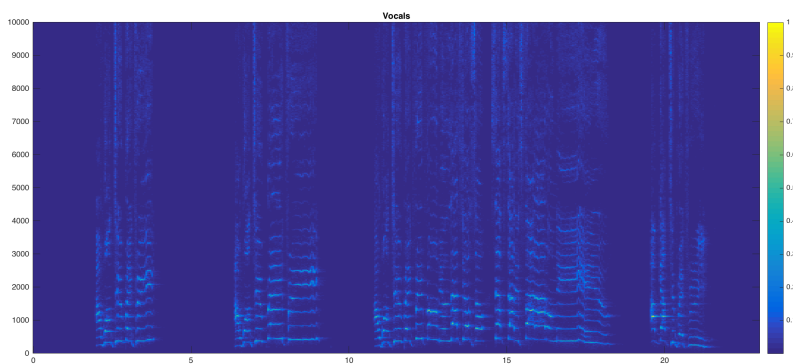
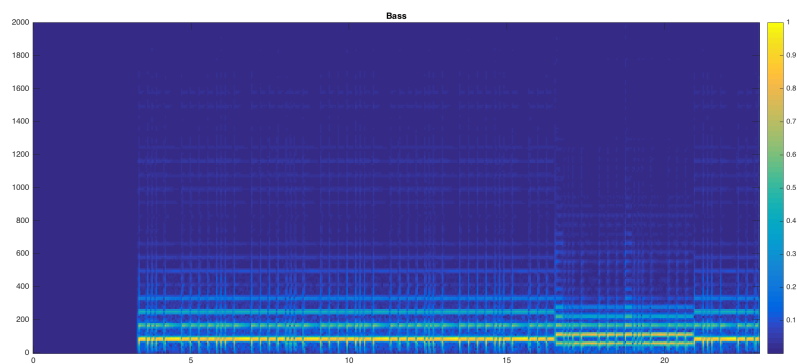
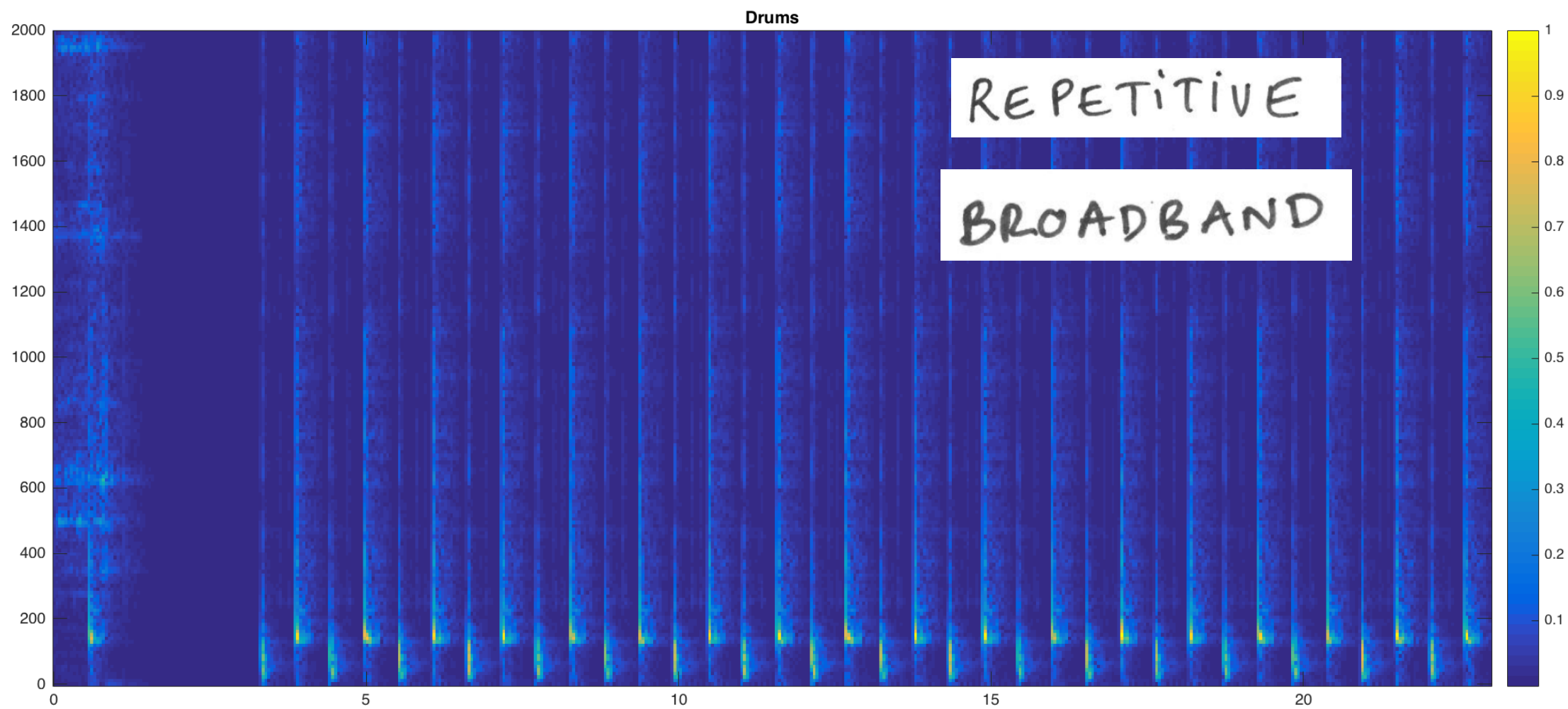
WHAT WE ARE AIMING FOR



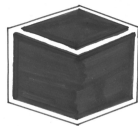
REPETITIVE
LOW FREQUENCY
REPETITION RATE
HARMONIC



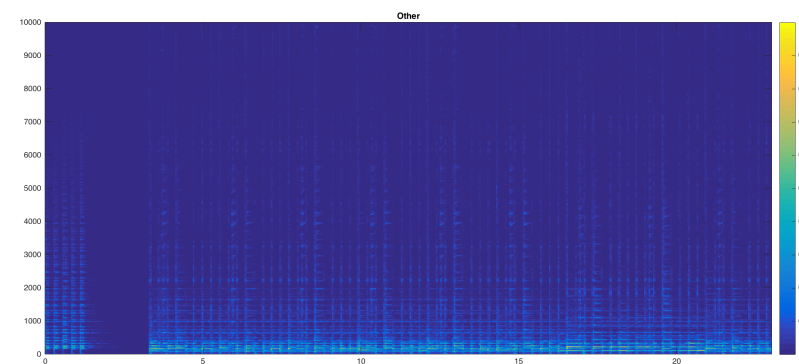
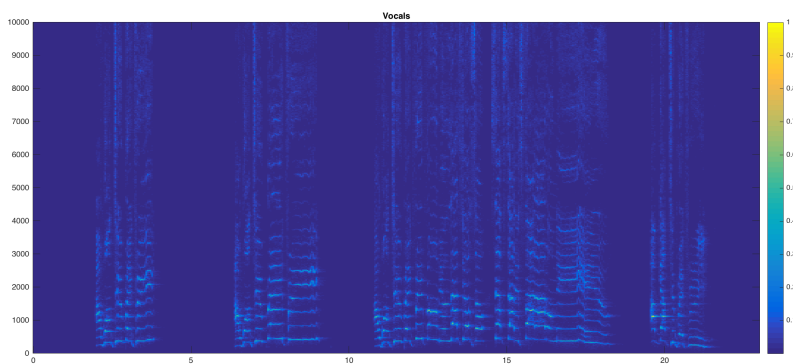
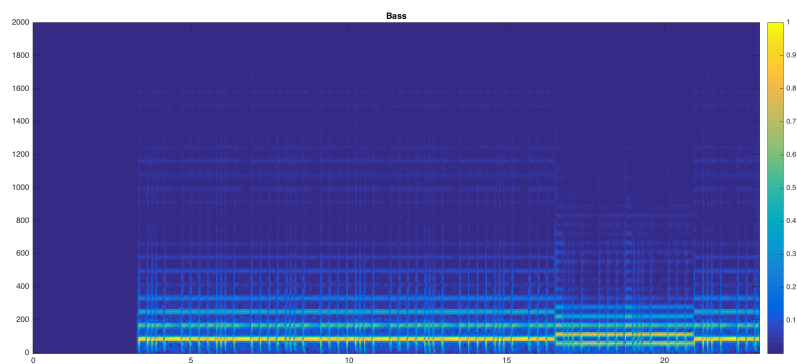
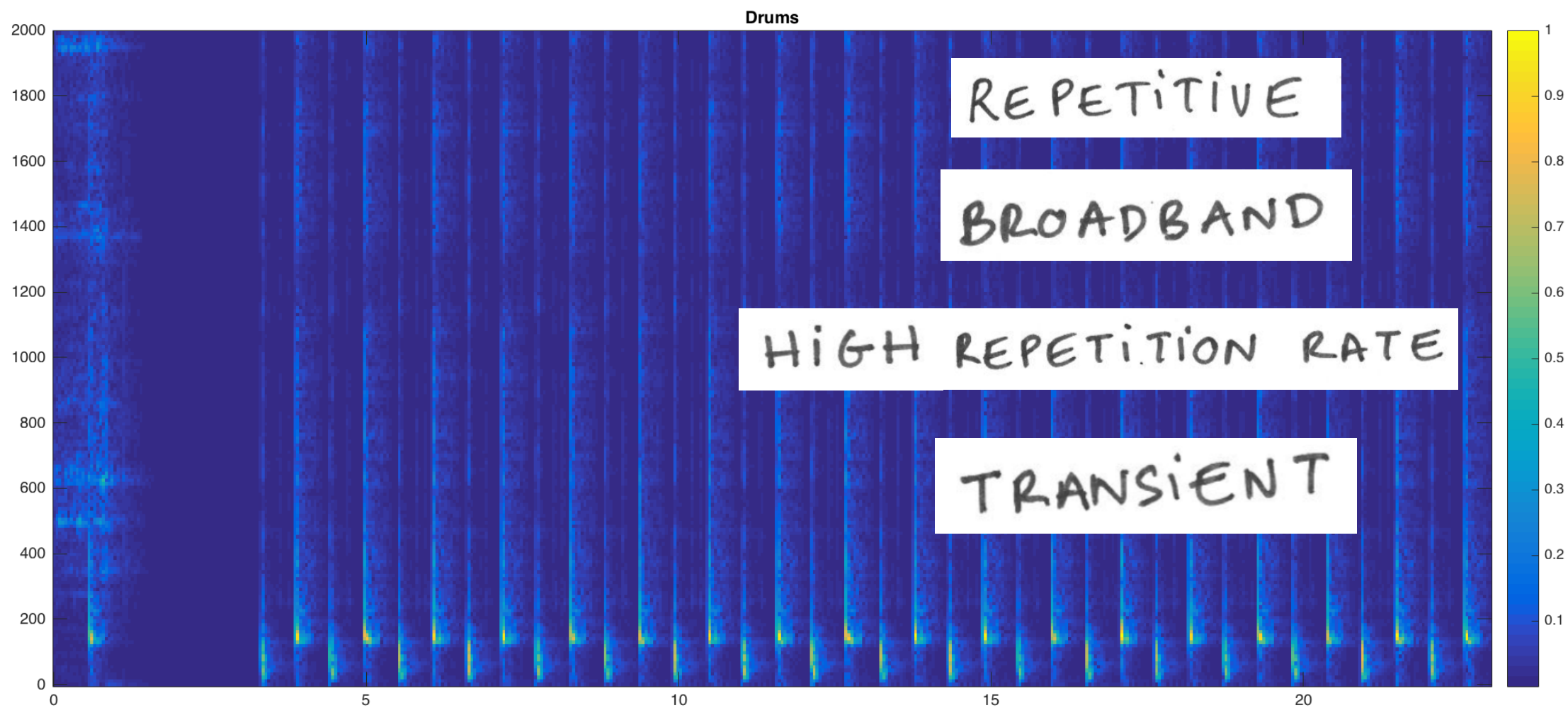
WHAT WE ARE AIMING FOR



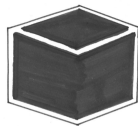
REPETITIVE
LOW FREQUENCY
REPETITION RATE
HARMONIC



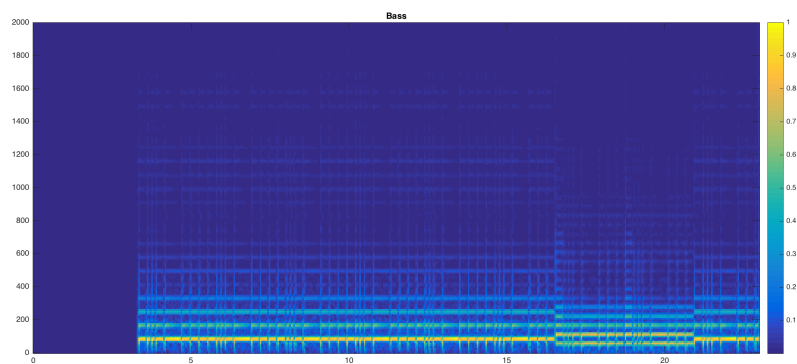
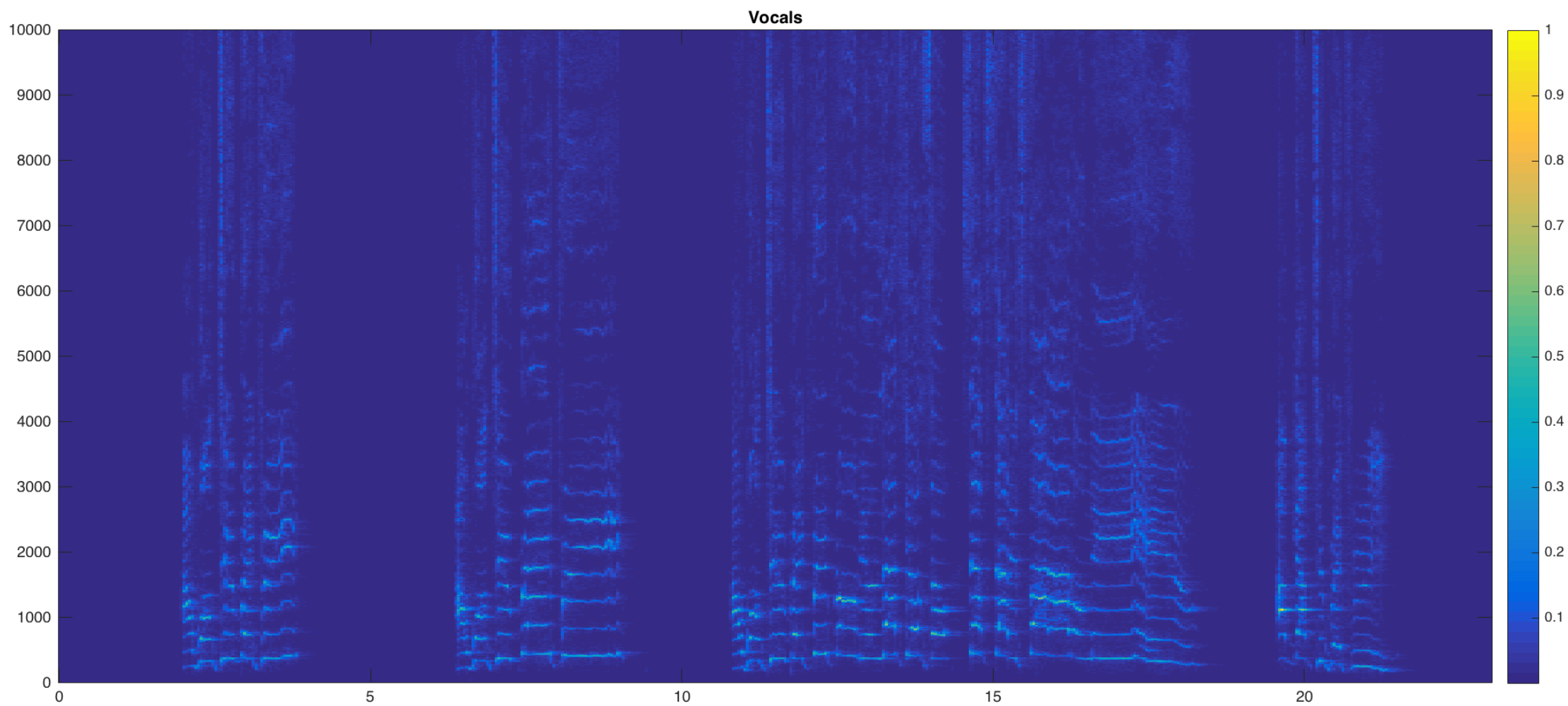
WHAT WE ARE AIMING FOR



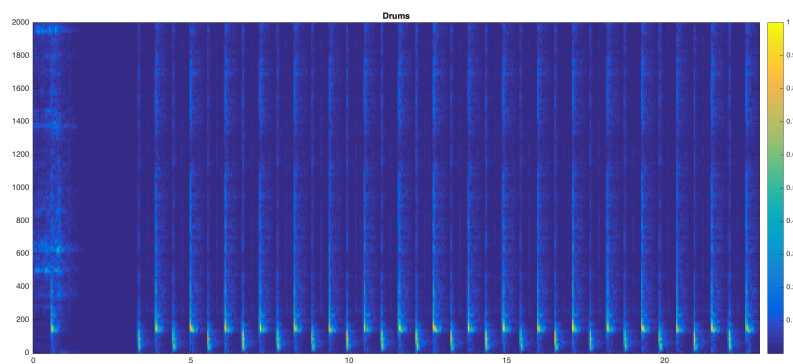
REPETITIVE
LOW FREQUENCY
REpetition RATE
HARMONIC



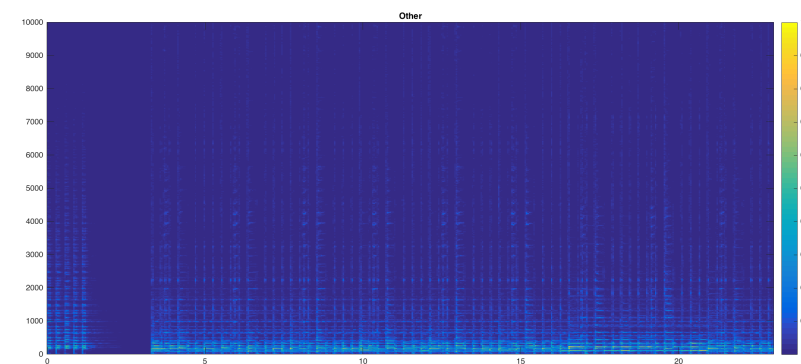
WHAT WE ARE AIMING FOR

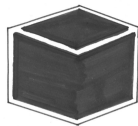


REPETITIVE
LOW FREQUENCY
REPETITION RATE
HARMONIC

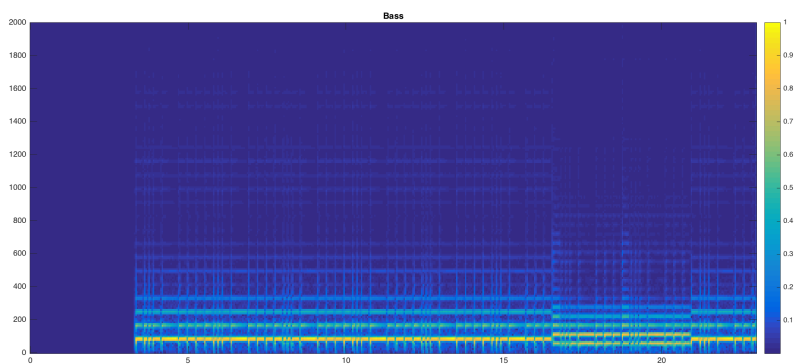
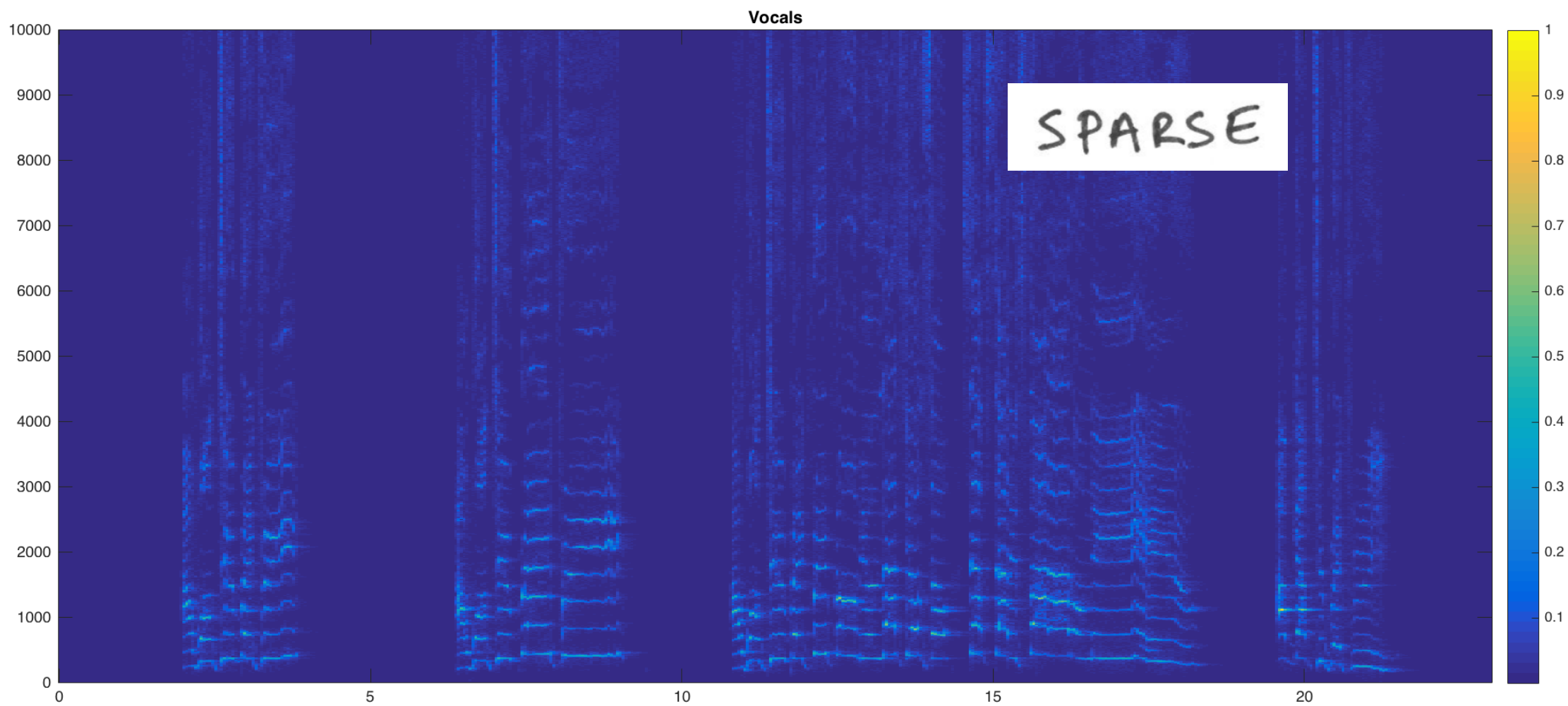


REPETITIVE
BROADBAND
REPETITION RATE
TRANSIENT

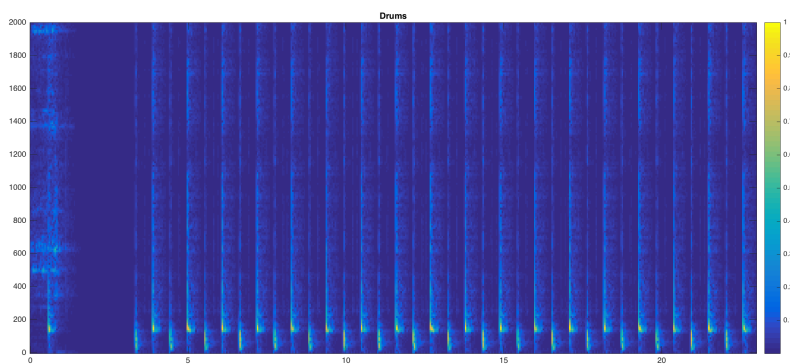




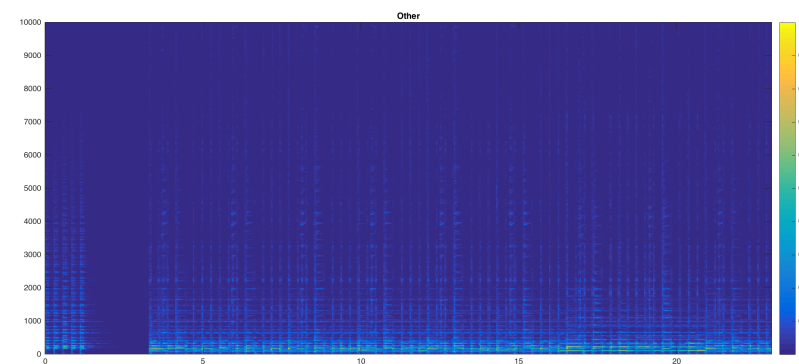
WHAT WE ARE AIMING FOR

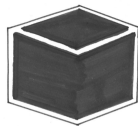


REPETITIVE
LOW FREQUENCY
REPETITION RATE
HARMONIC

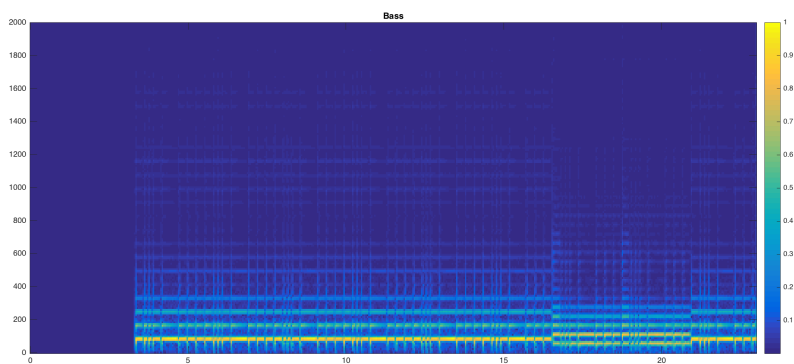
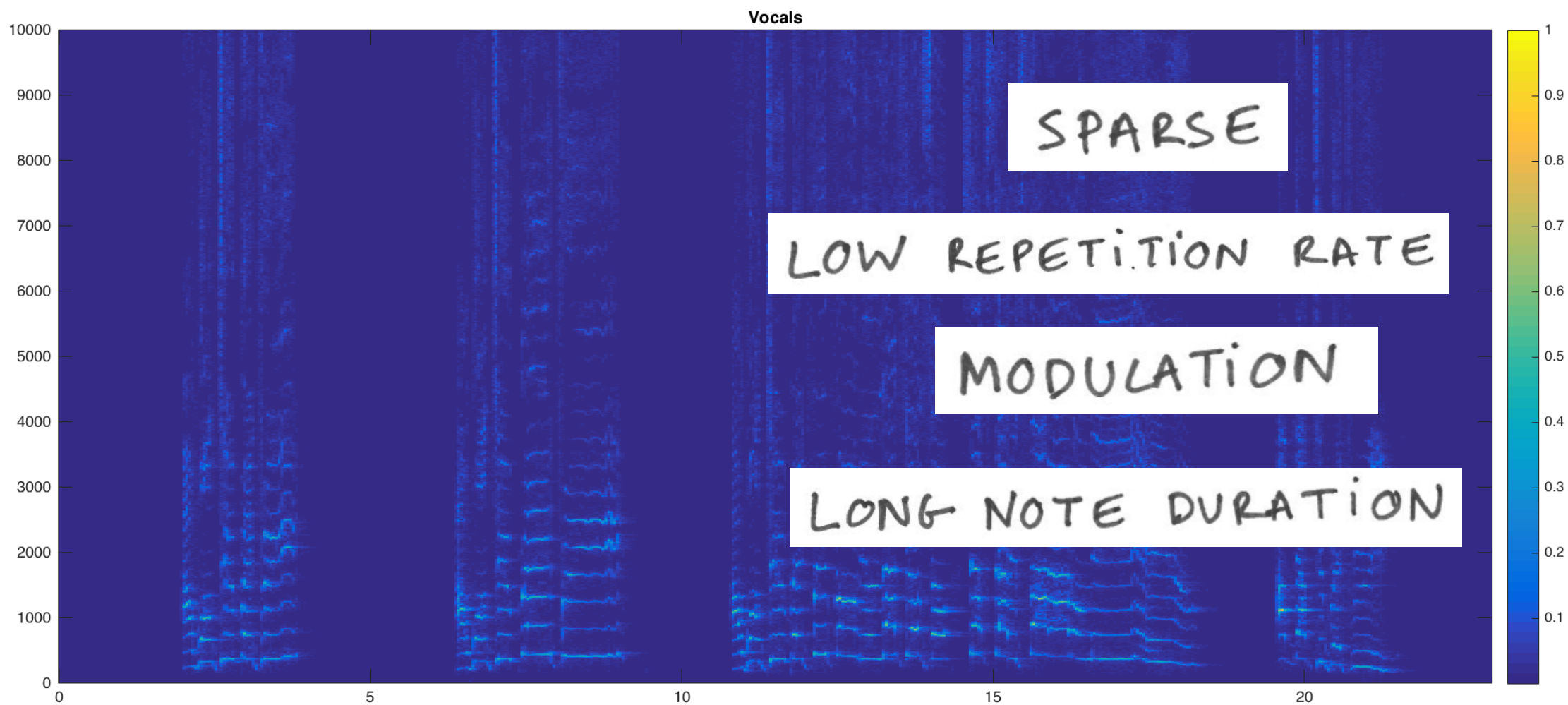


REPETITIVE
BROADBAND
REPETITION RATE
TRANSIENT

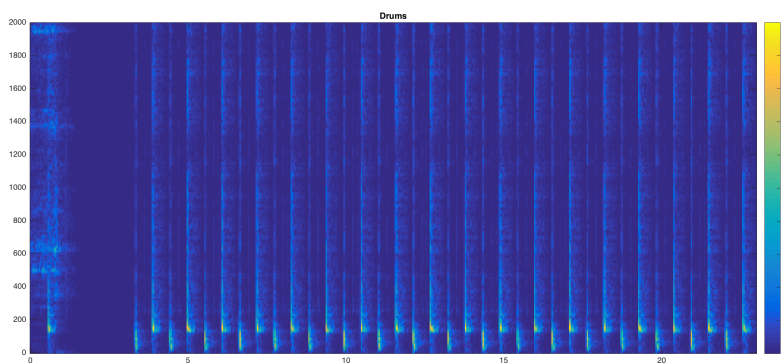




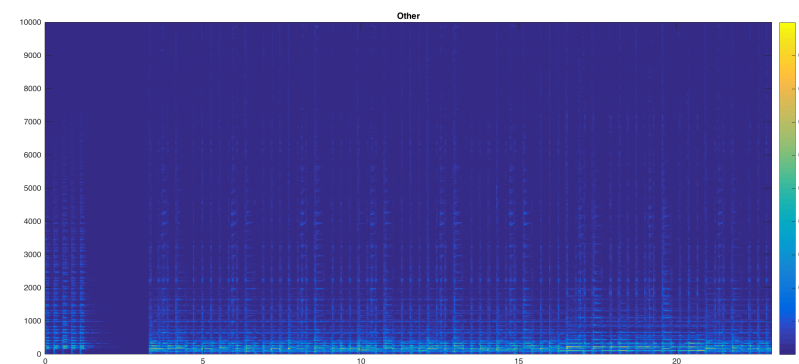
WHAT WE ARE AIMING FOR

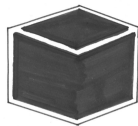


REPETITIVE
LOW FREQUENCY
REPETITION RATE
HARMONIC

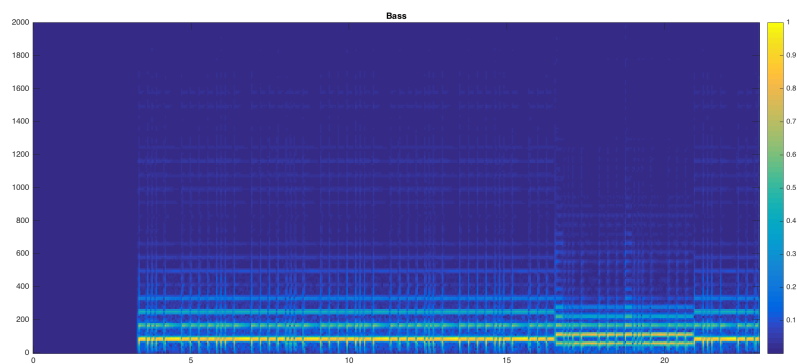
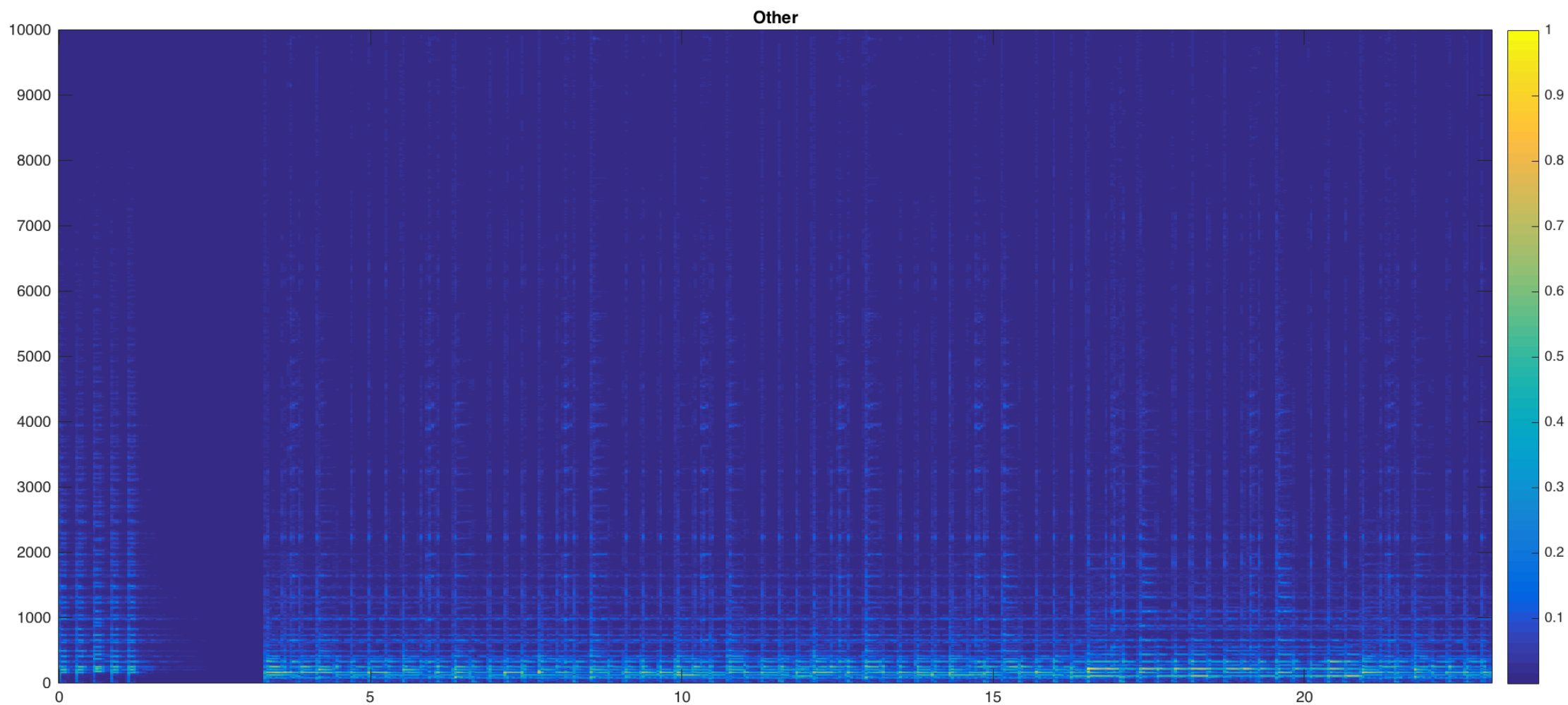


REPETITIVE
BROADBAND
REPETITION RATE
TRANSIENT

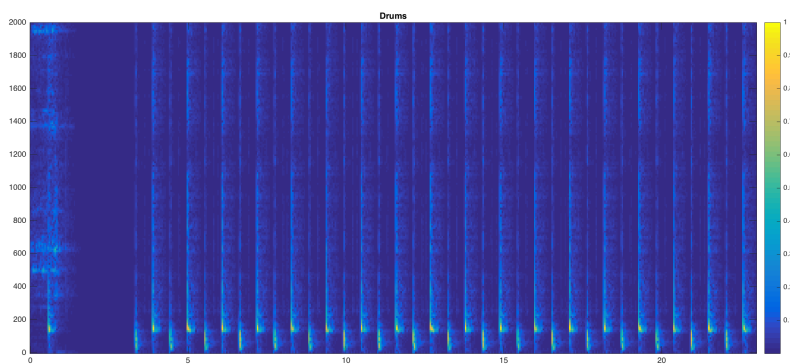




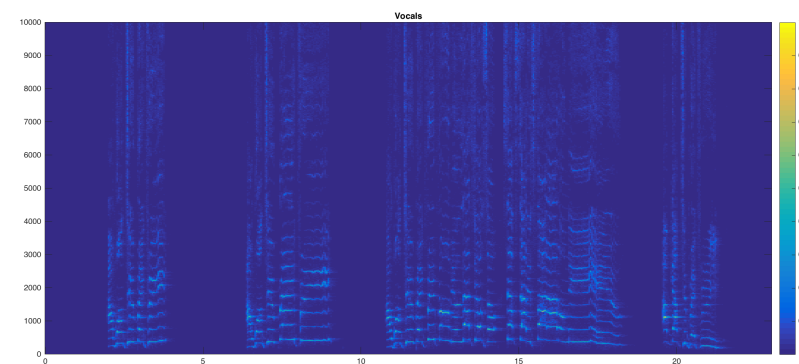
WHAT WE ARE AIMING FOR



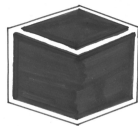
REPETITIVE
LOW FREQUENCY
REPETITION RATE
HARMONIC



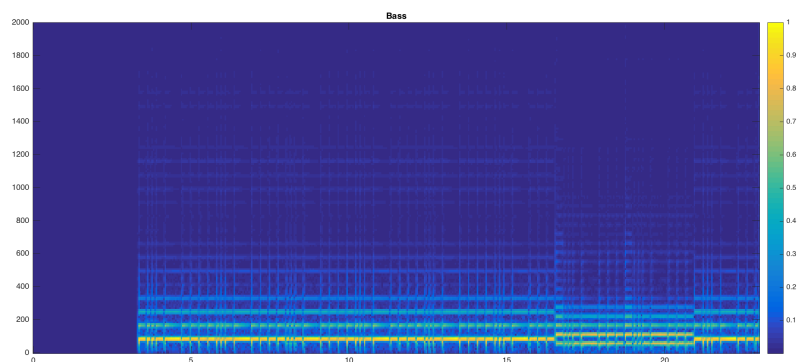
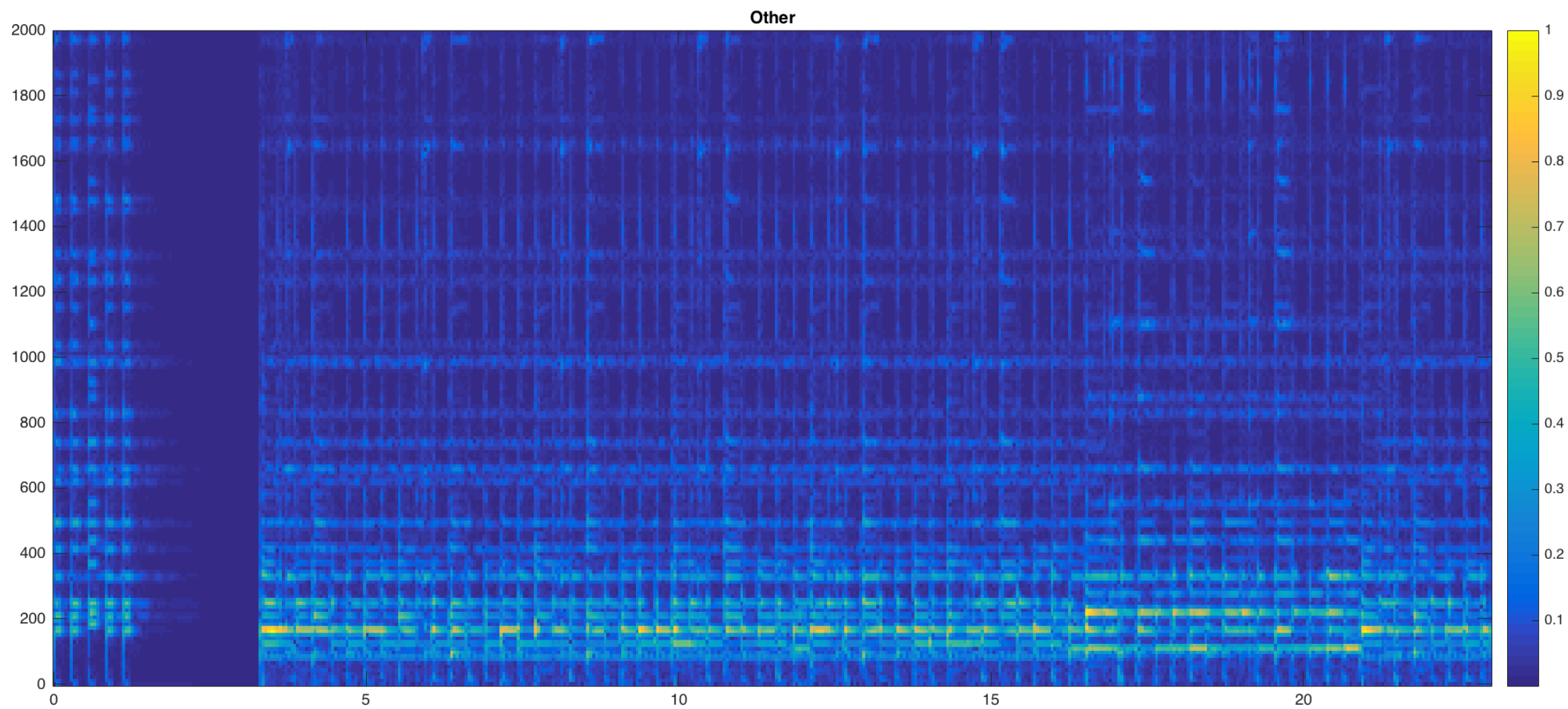
REPETITIVE
BROADBAND
REPETITION RATE
TRANSIENT



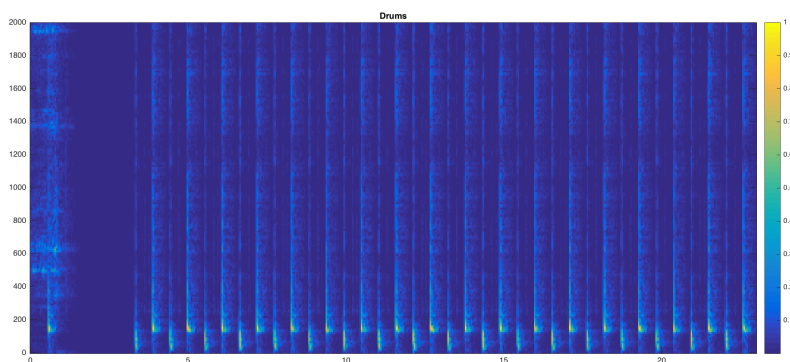
REPETITIVE
SPARSE
LOW REPETITION RATE
MODULATION
LONG NOTE DURATION



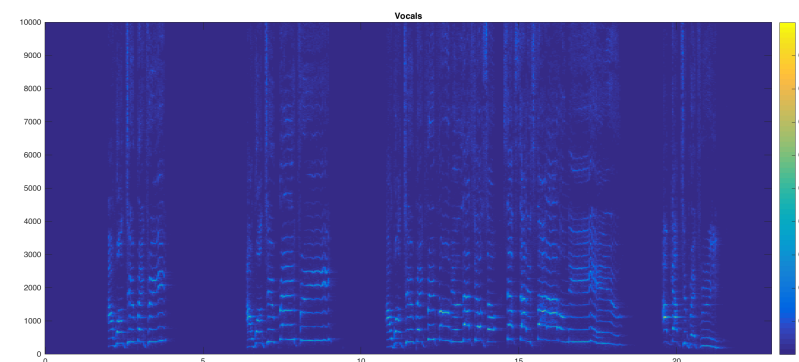
WHAT WE ARE AIMING FOR



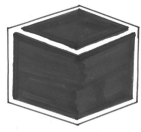
REPETITIVE
LOW FREQUENCY
REPETITION RATE
HARMONIC



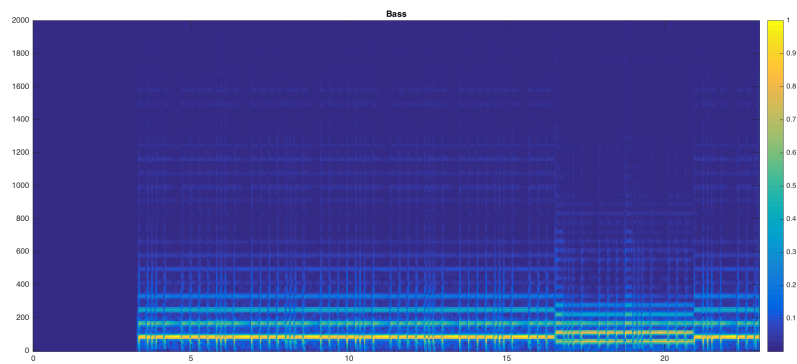
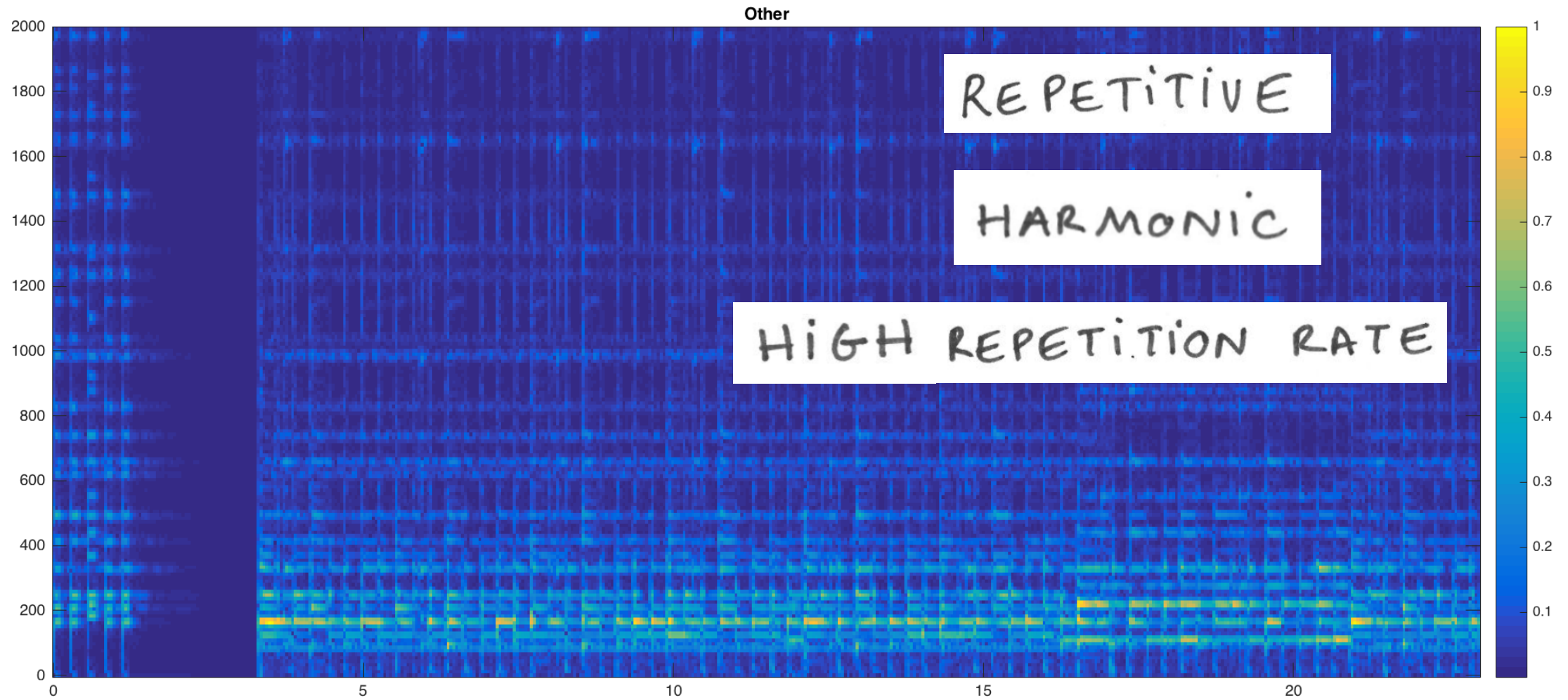
REPETITIVE
BROADBAND
REPETITION RATE
TRANSIENT



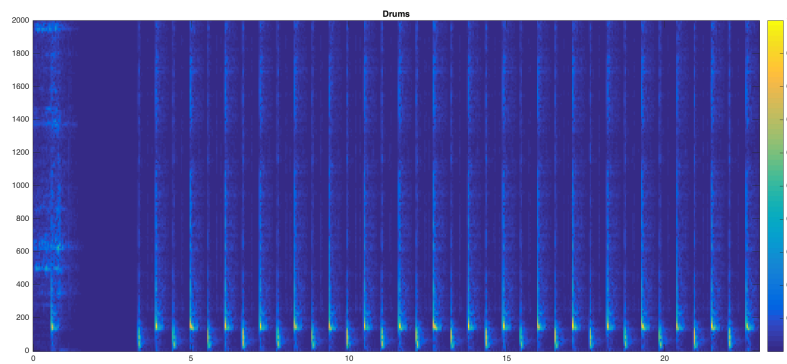
REPETITIVE
SPARSE
LOW REPETITION RATE
MODULATION
LONG NOTE DURATION



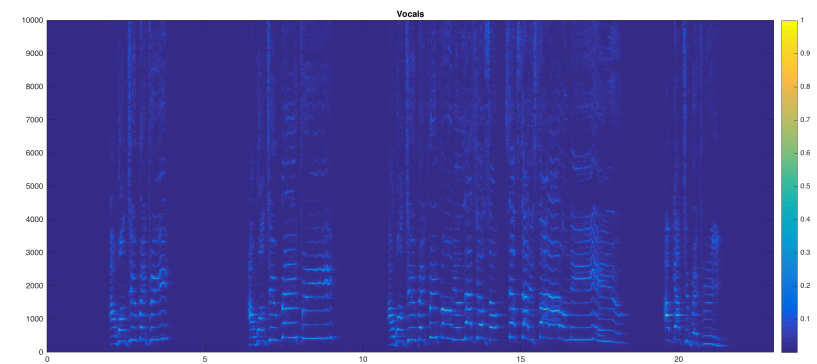
WHAT WE ARE AIMING FOR



REPETITIVE
LOW FREQUENCY
REPETITION RATE
HARMONIC



REPETITIVE
BROADBAND
REPETITION RATE
TRANSIENT



REPETITIVE
SPARSE
LOW REPETITION RATE
MODULATION
LONG NOTE DURATION

BASS

REPETITIVE

LOW FREQUENCY

HIGH REPETITION RATE

HARMONIC

DRUMS

REPETITIVE

BROADBAND

HIGH REPETITION RATE

TRANSIENT

VOCALS

SPARSE

LOW REPETITION RATE

MODULATION

LONG NOTE DURATION

OTHER

REPETITIVE

HARMONIC

HIGH REPETITION RATE

BASS

REPETITIVE

LOW FREQUENCY

HIGH REPETITION RATE

HARMONIC

DRUMS

REPETITIVE

BROADBAND

HIGH REPETITION RATE

TRANSIENT

VOCALS

SPARSE

LOW REPETITION RATE

MODULATION

LONG NOTE DURATION

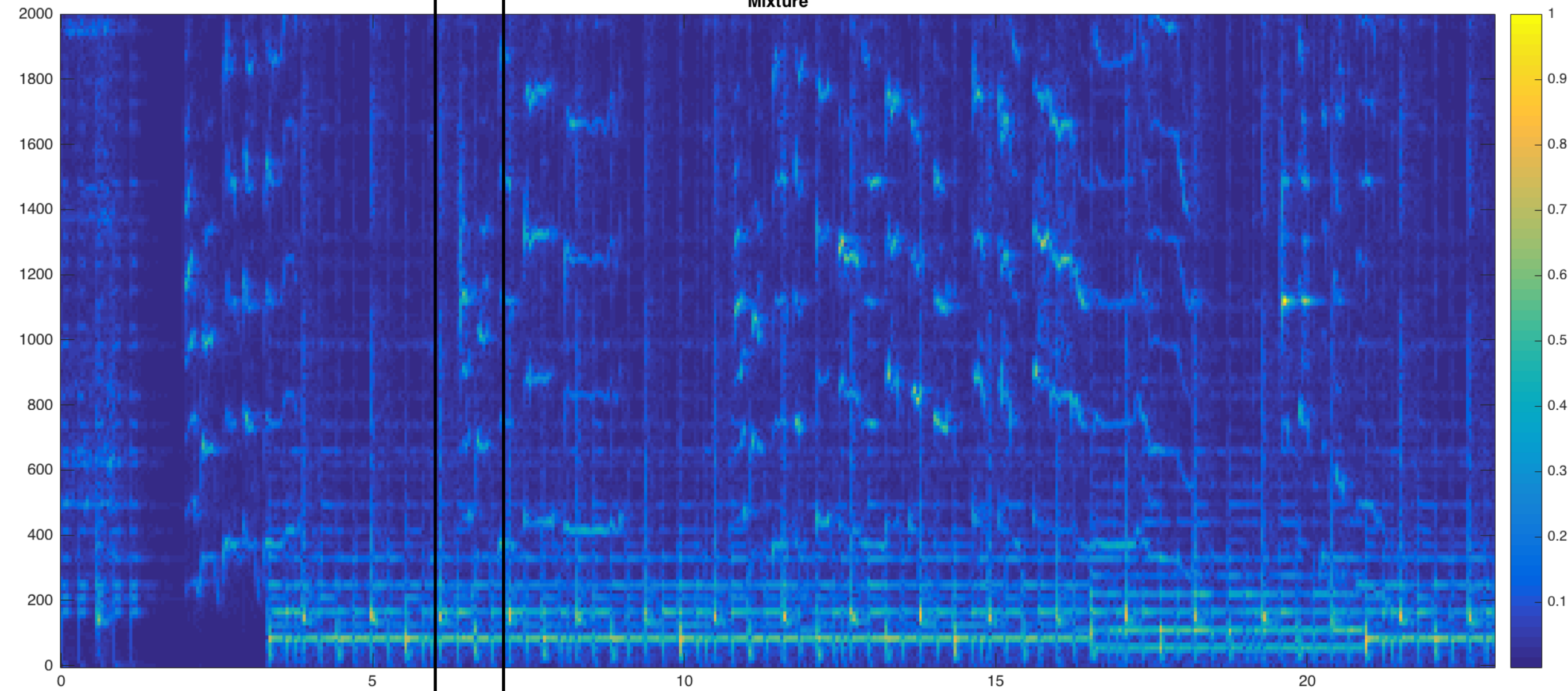
OTHER

REPETITIVE

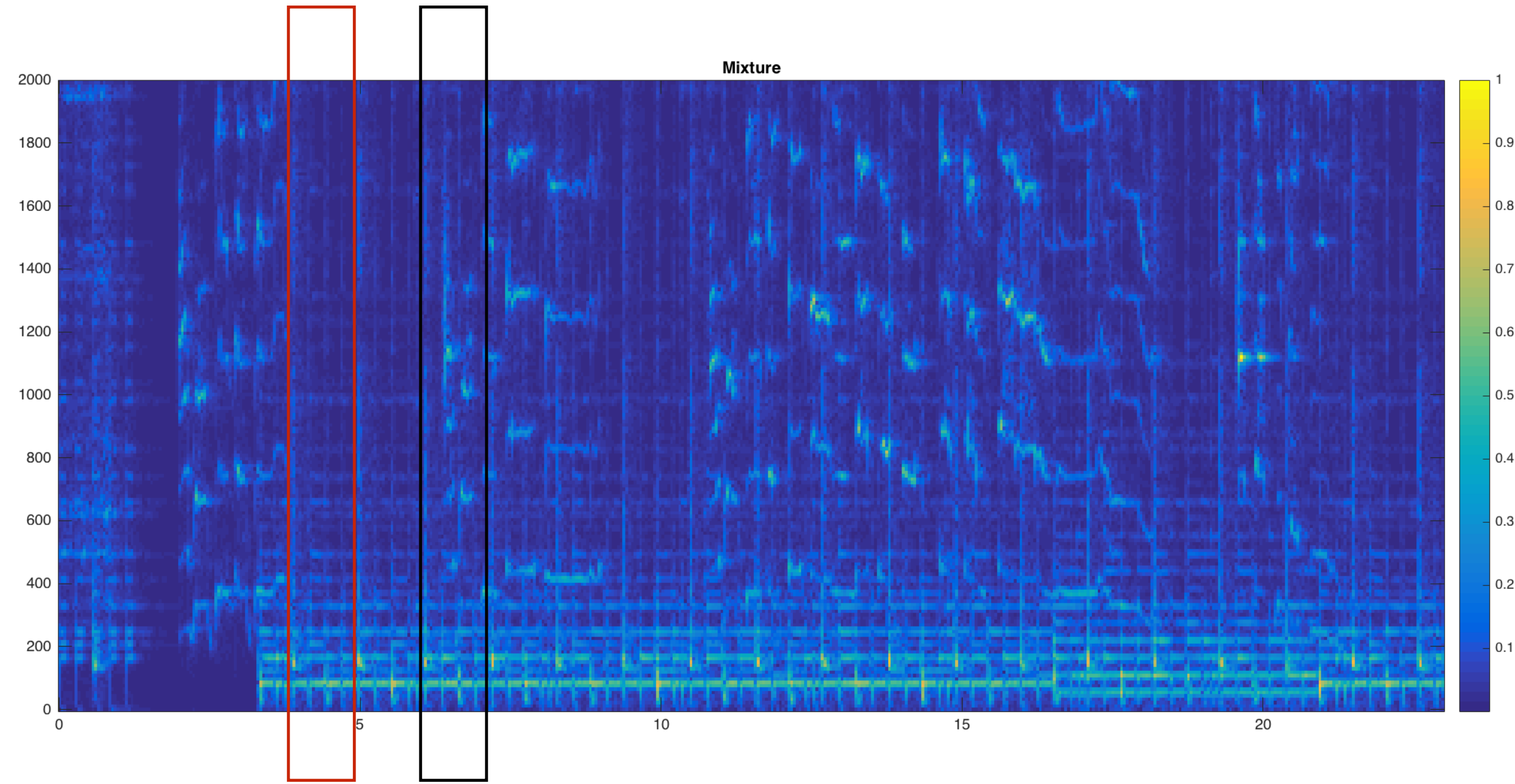
HARMONIC

HIGH REPETITION RATE

Mixture



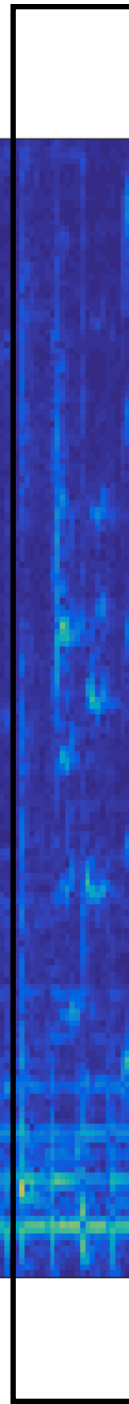
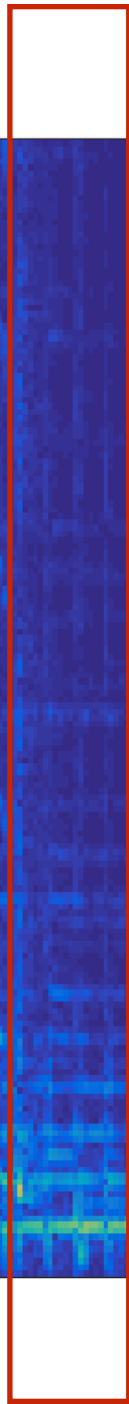
Mixture



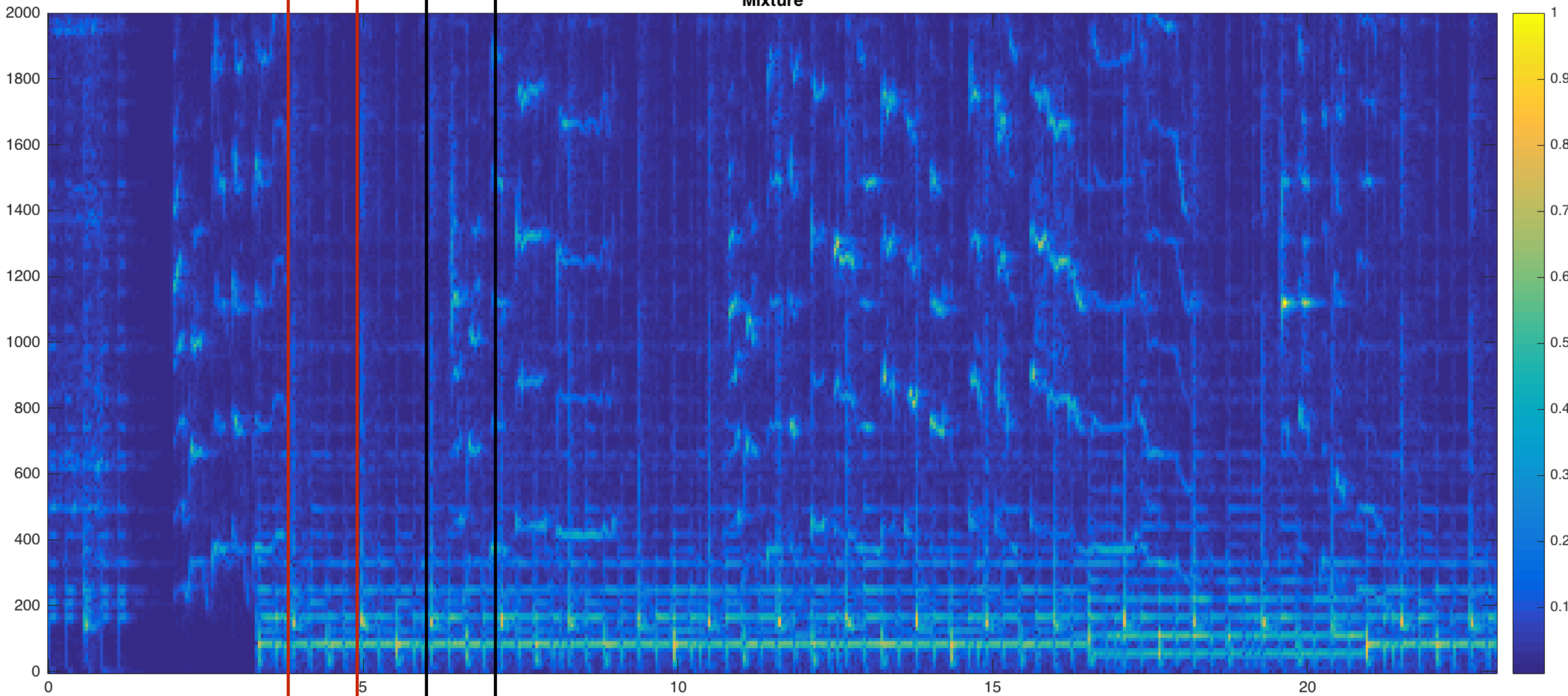
KERNEL ADDITIVE MODELLING

K.A.M.

SIMILARITY KERNEL

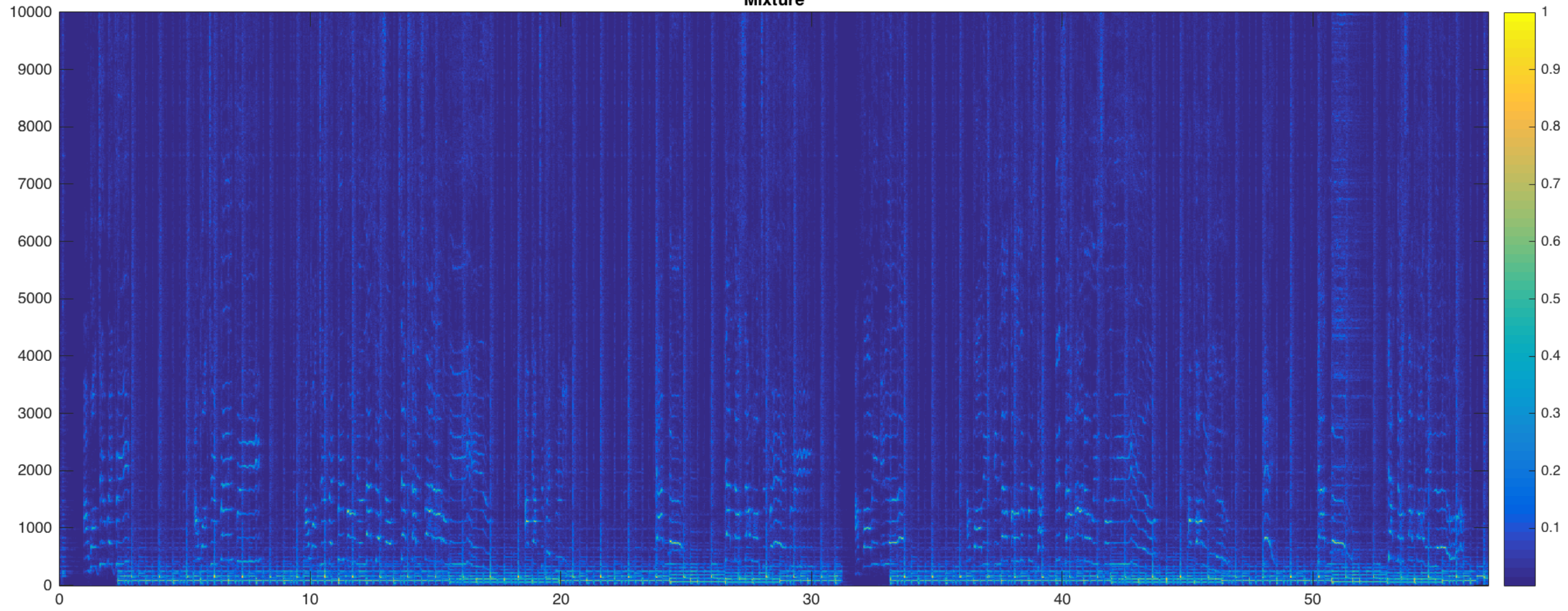


Mixture



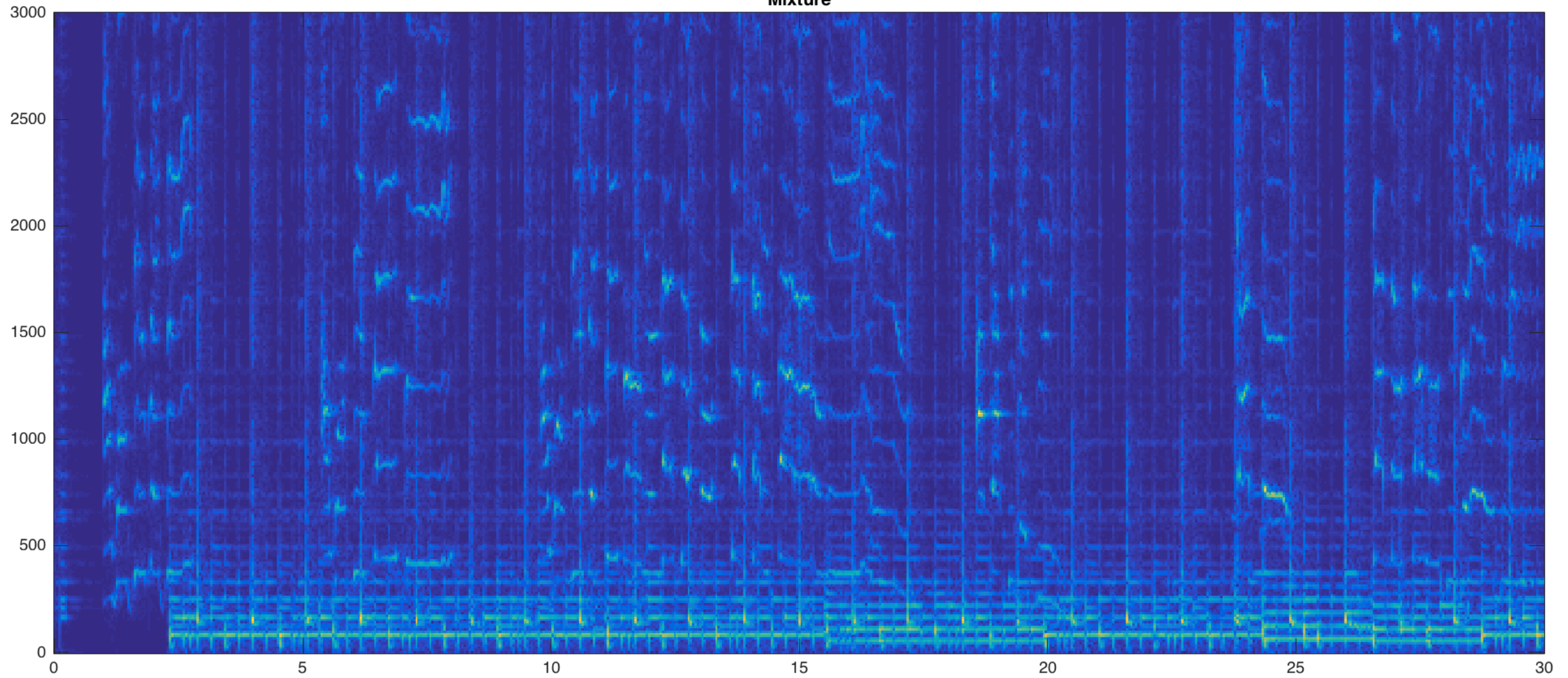
VOCAL SEPARATION WITH K.A.M

Mixture

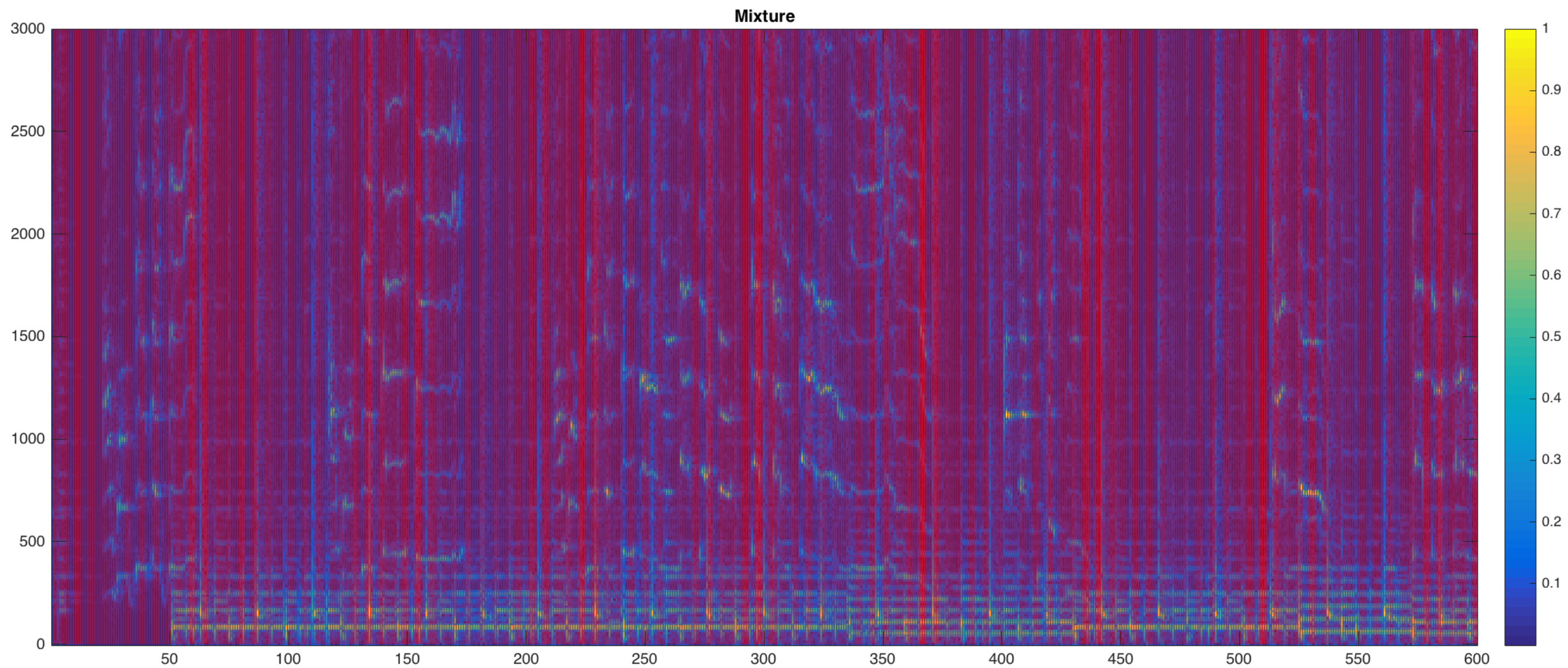


VOCAL SEPARATION WITH K.A.M

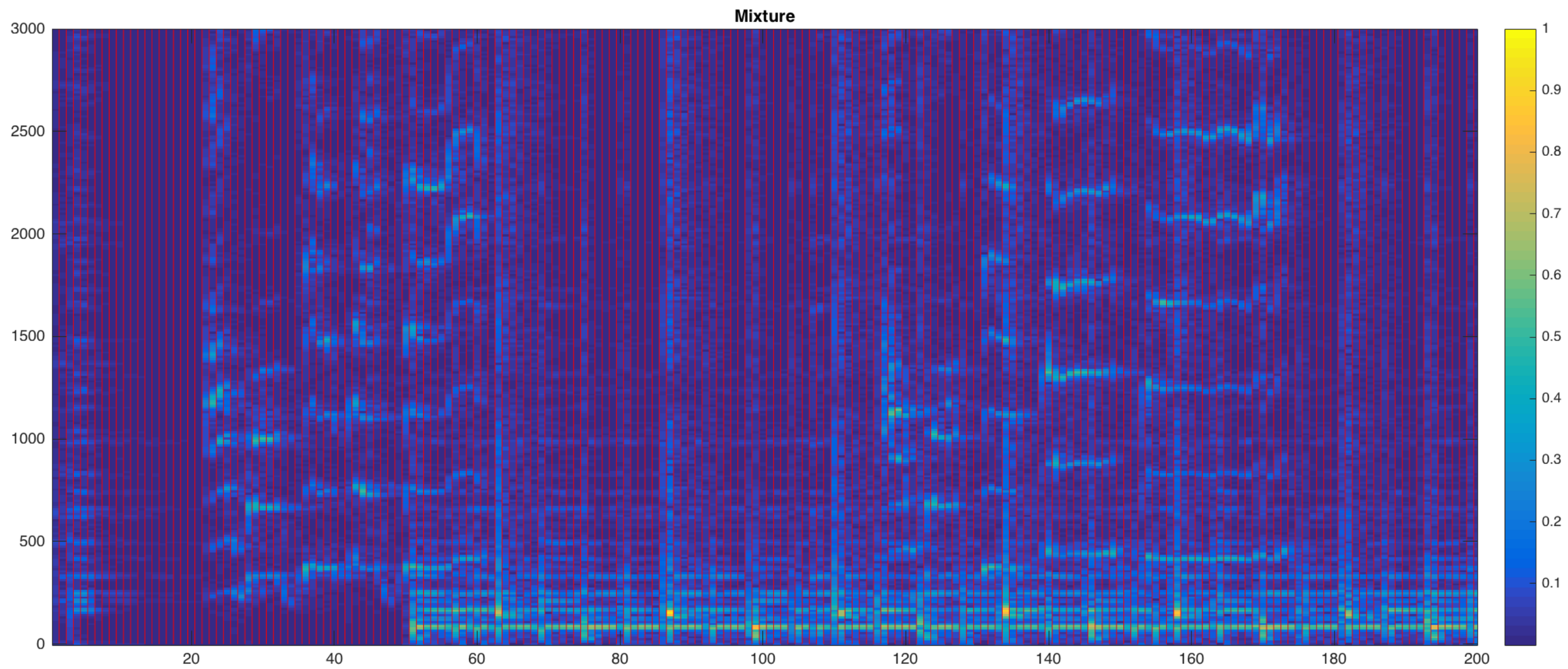
Mixture



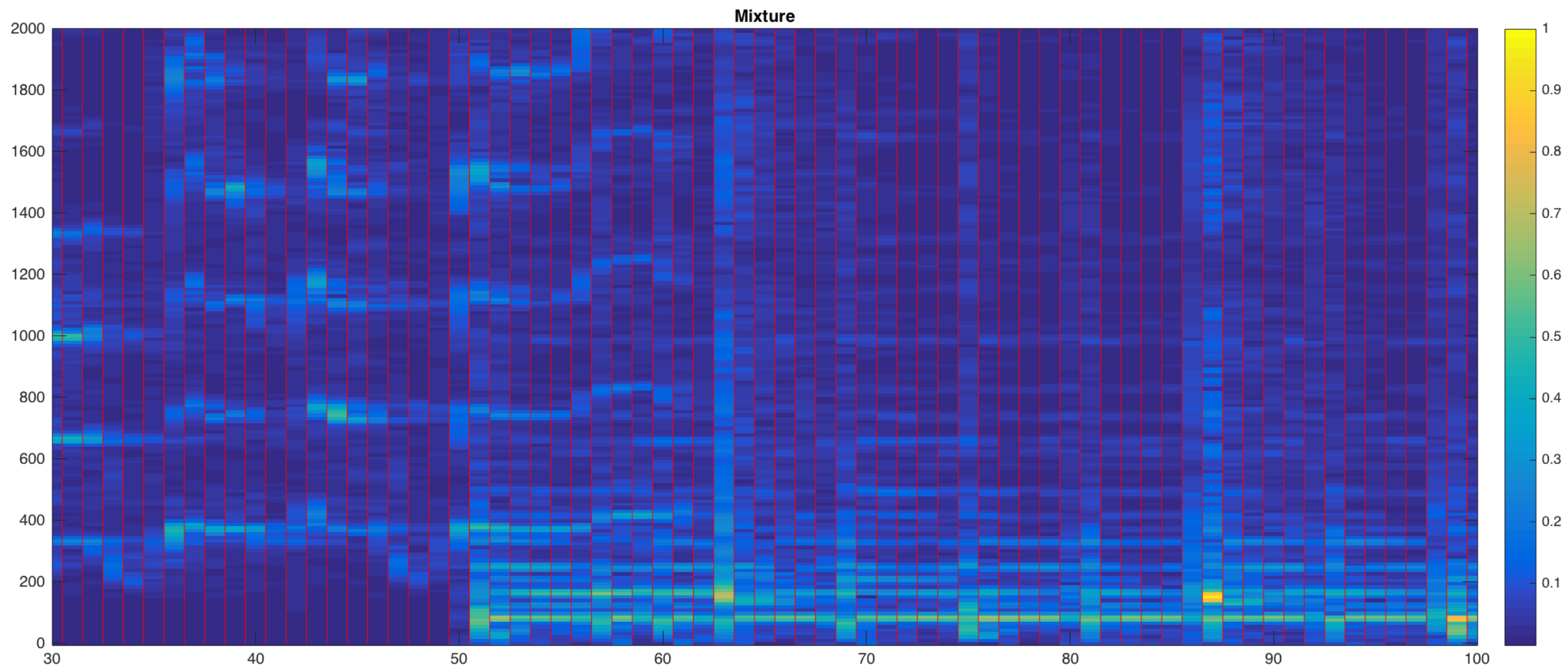
1. DIVIDE IN TIME FRAMES



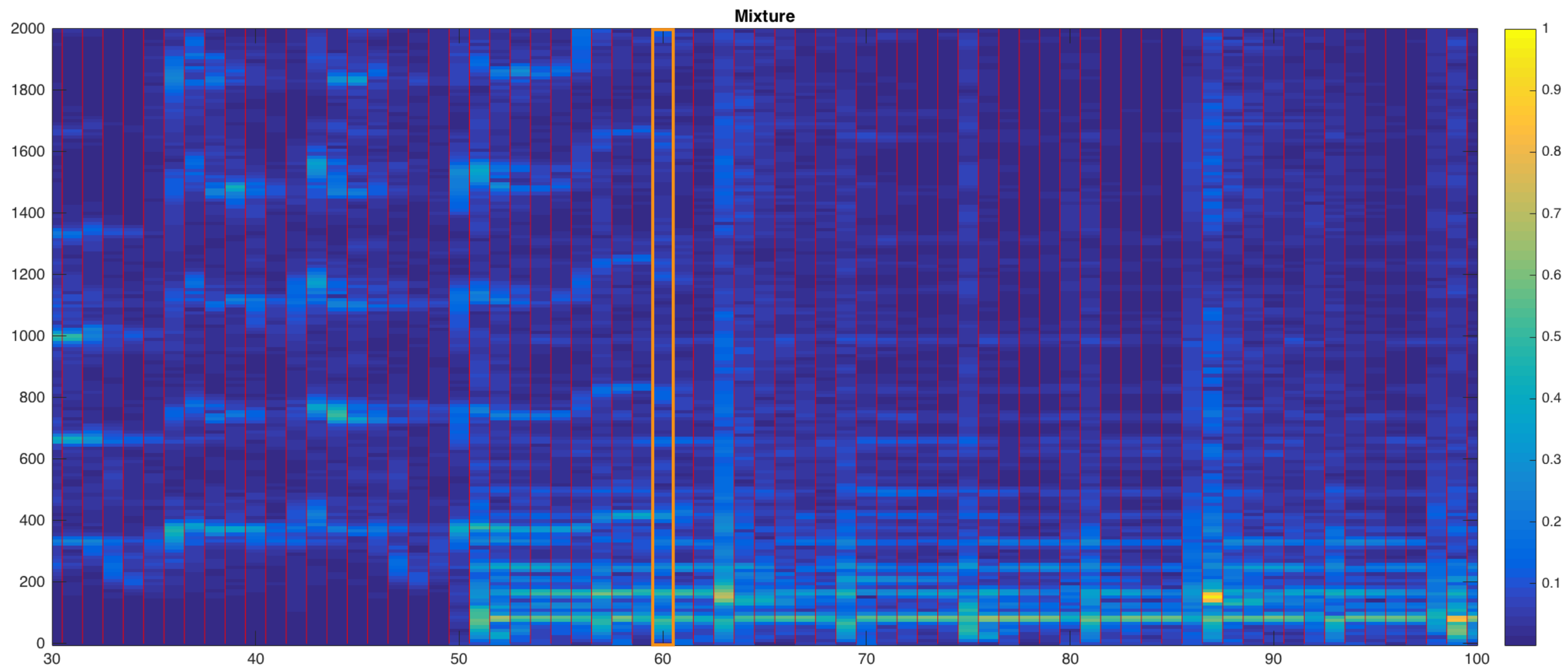
1. DIVIDE IN TIME FRAMES



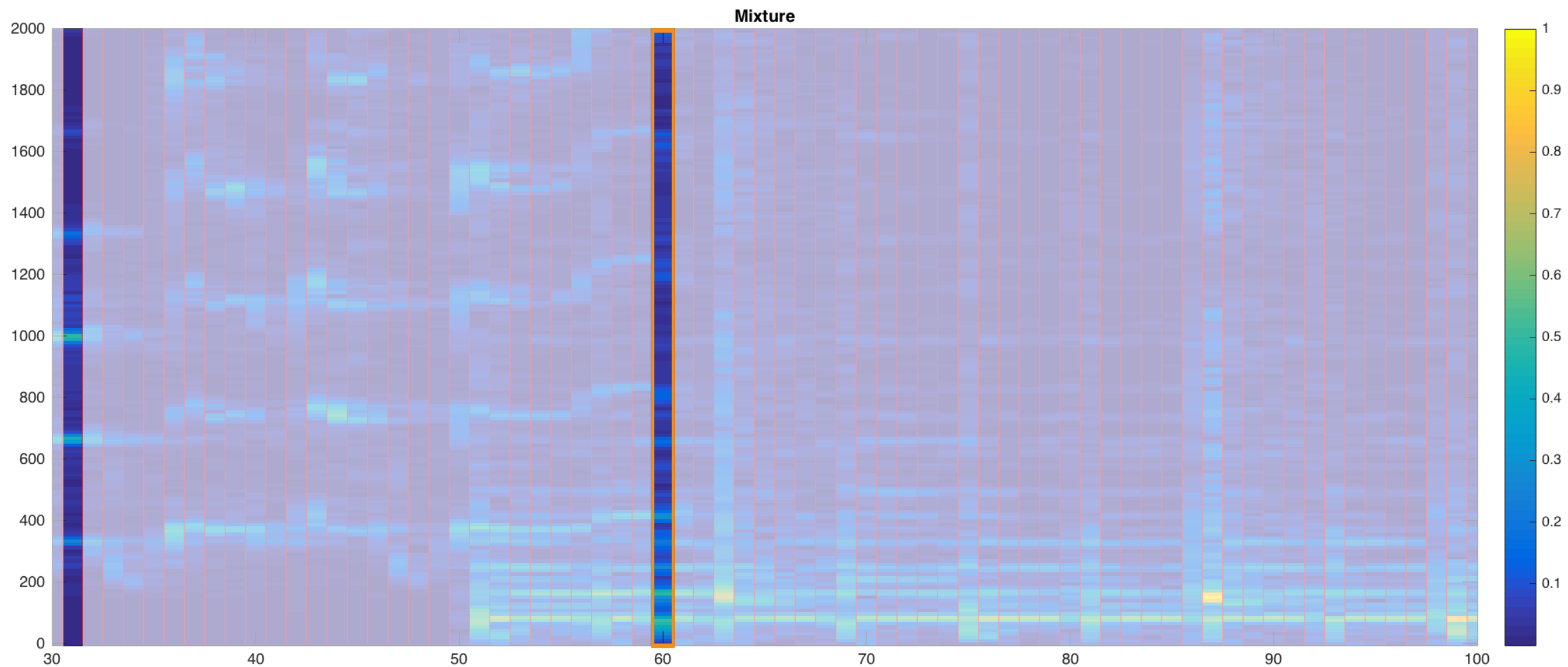
1. DIVIDE IN TIME FRAMES



1. DIVIDE IN TIME FRAMES

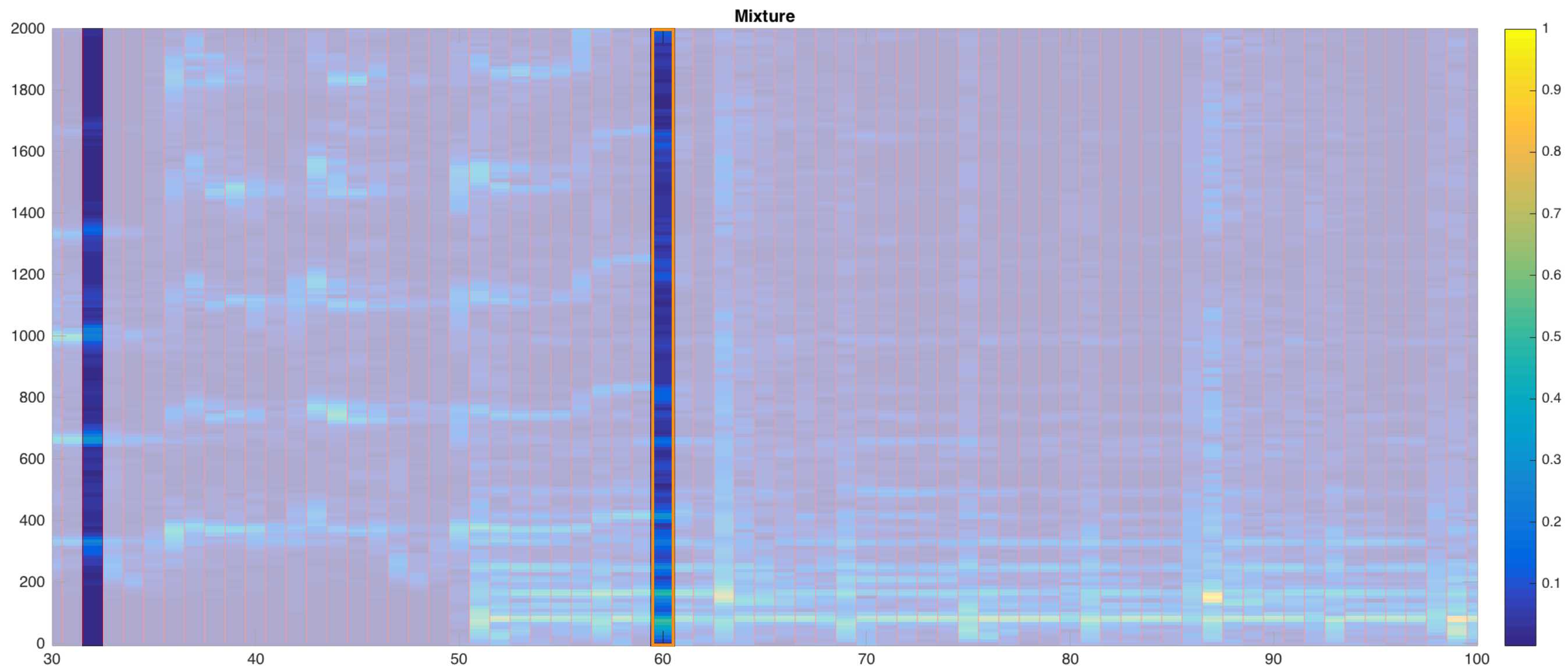


2. LOOK FOR SIMILAR TIME FRAMES



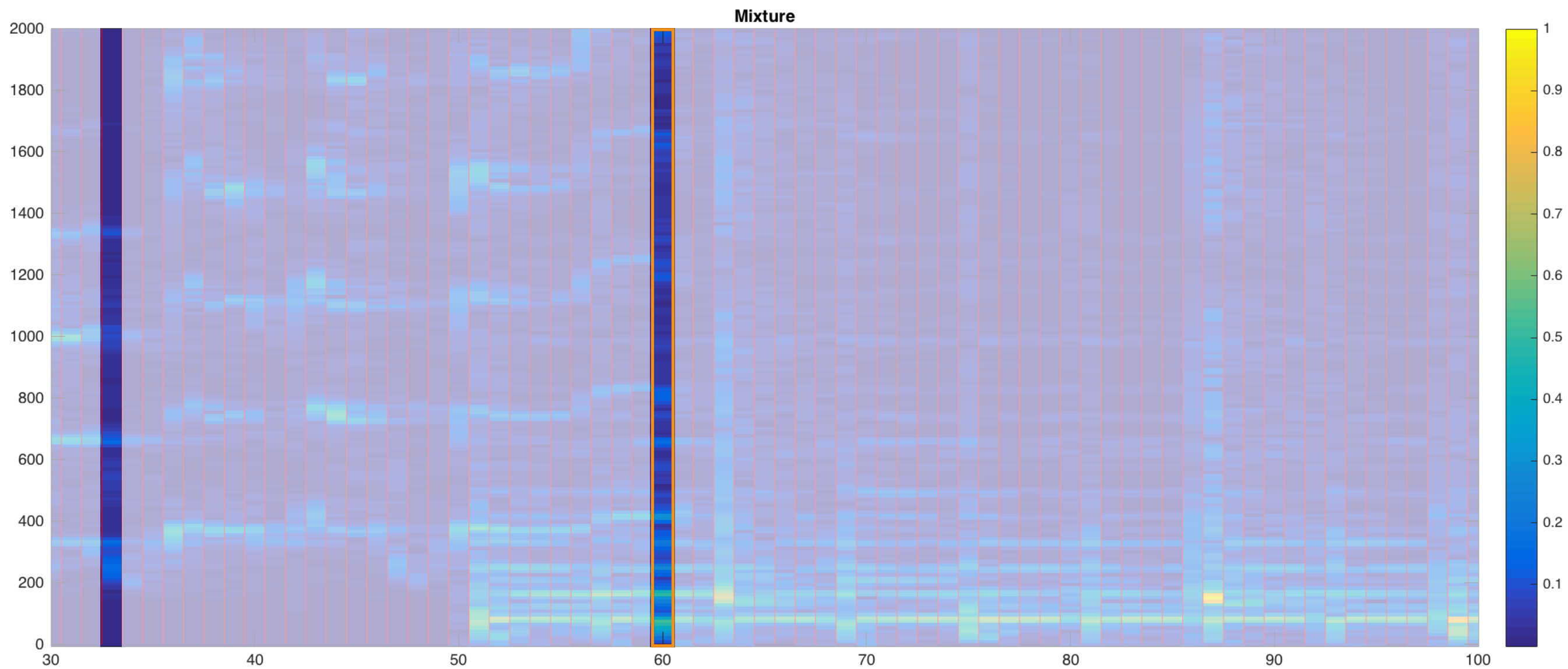
DISTANCE BETWEEN FRAMES

2. LOOK FOR SIMILAR TIME FRAMES



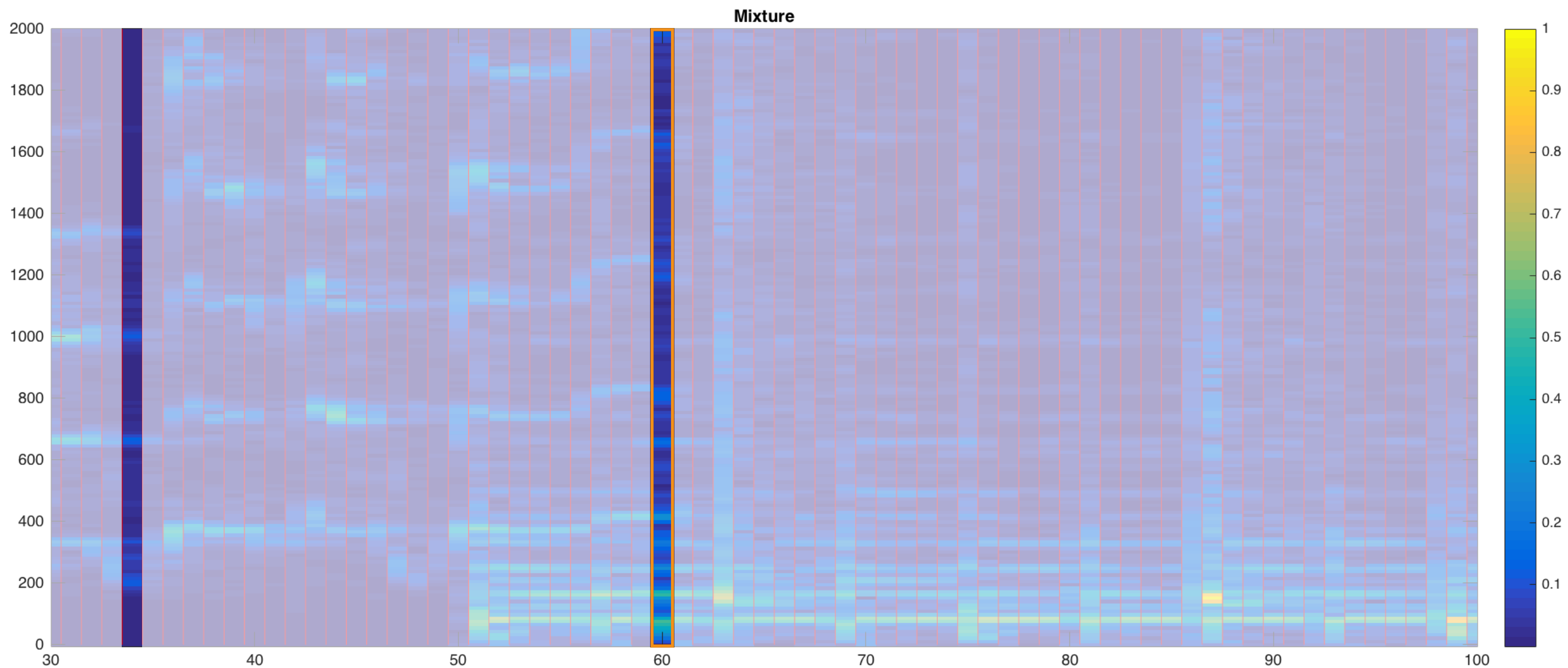
DISTANCE BETWEEN FRAMES

2. LOOK FOR SIMILAR TIME FRAMES



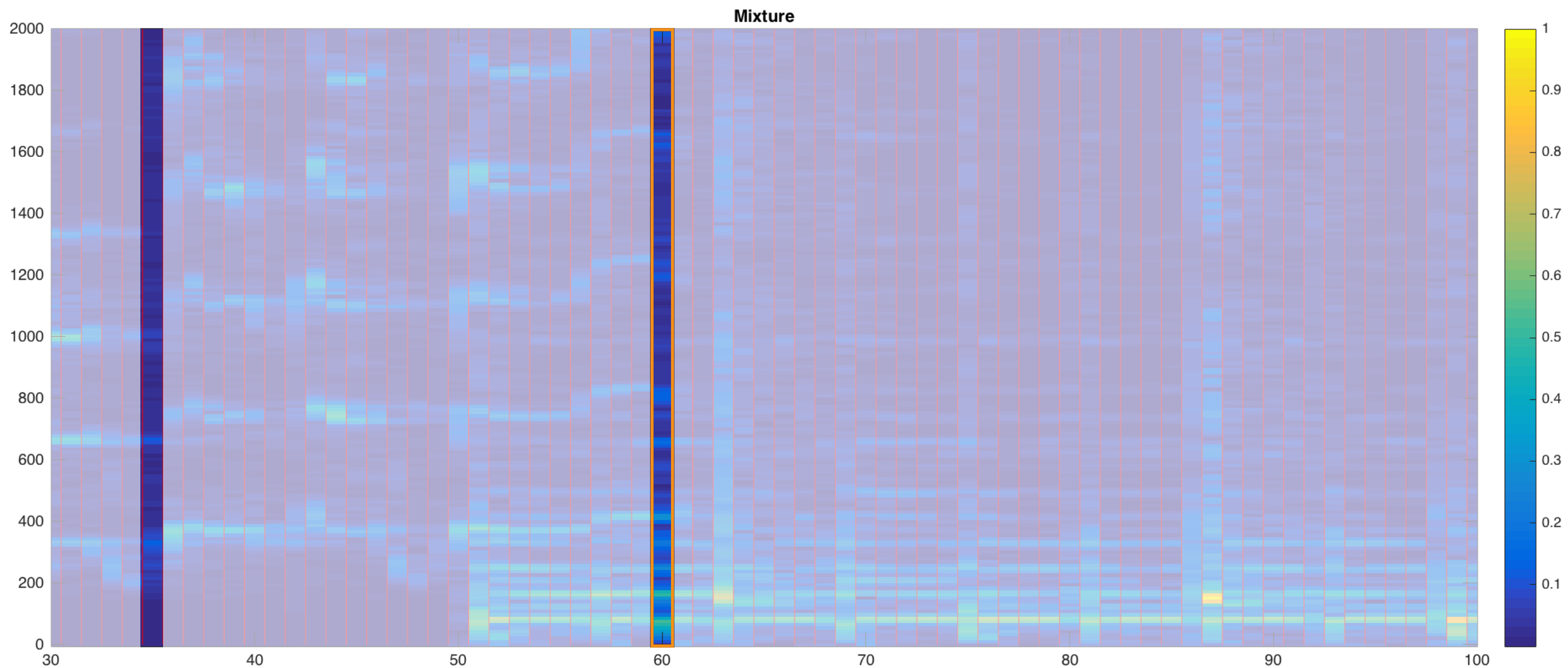
DISTANCE BETWEEN FRAMES

2. LOOK FOR SIMILAR TIME FRAMES



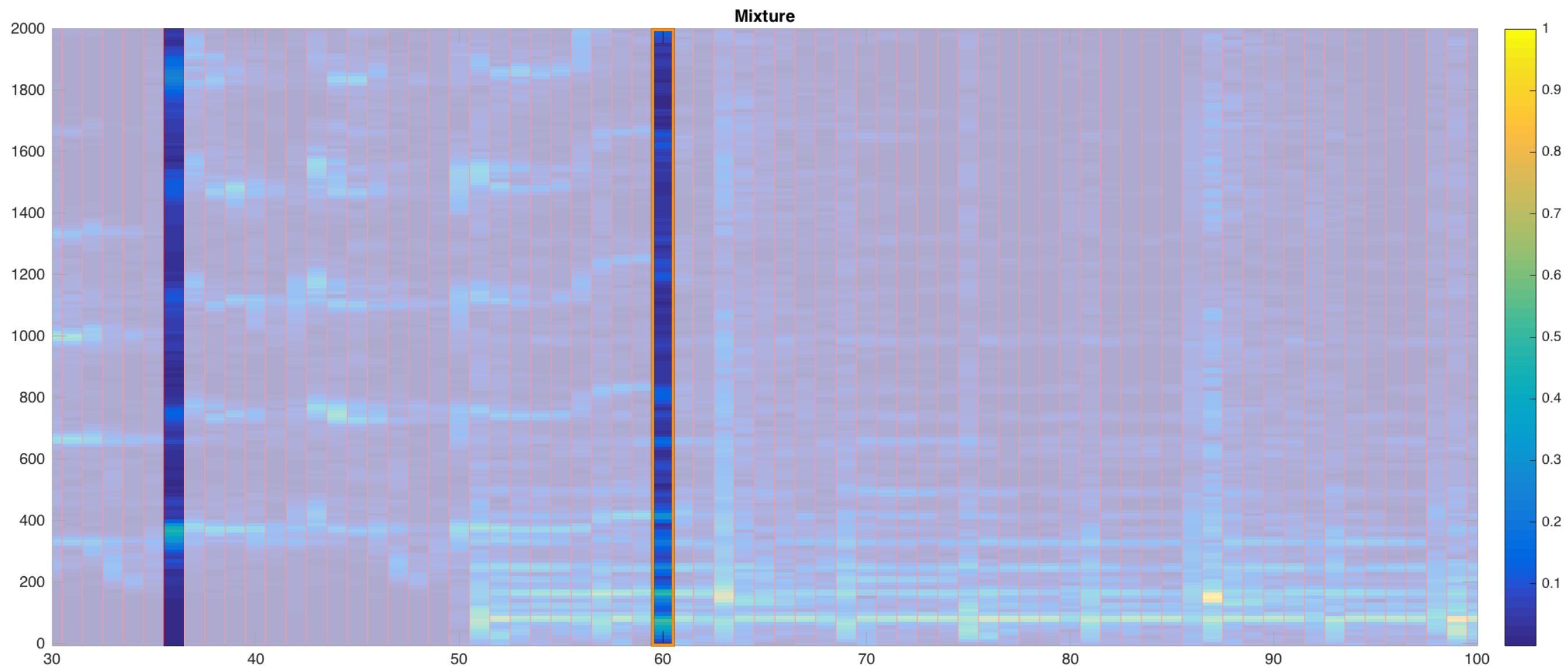
DISTANCE BETWEEN FRAMES

2. LOOK FOR SIMILAR TIME FRAMES



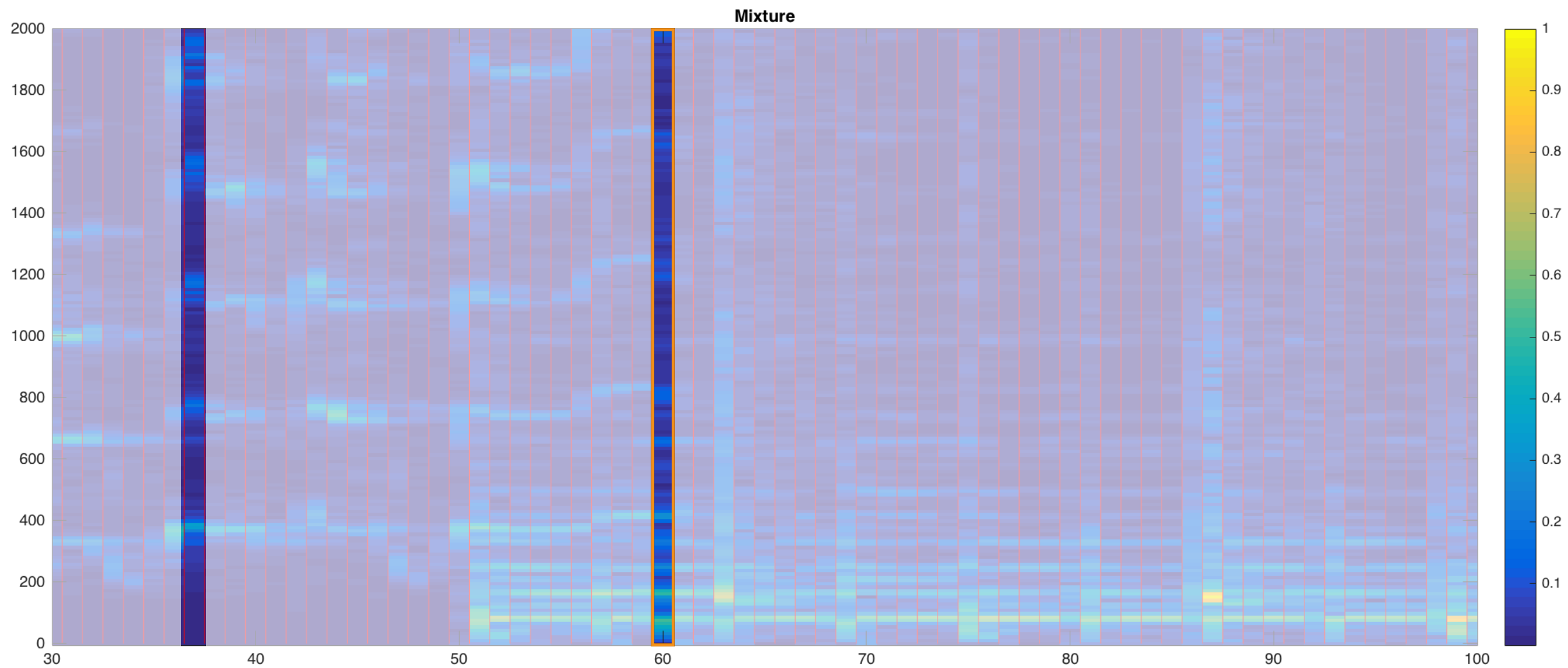
DISTANCE BETWEEN FRAMES

2. LOOK FOR SIMILAR TIME FRAMES



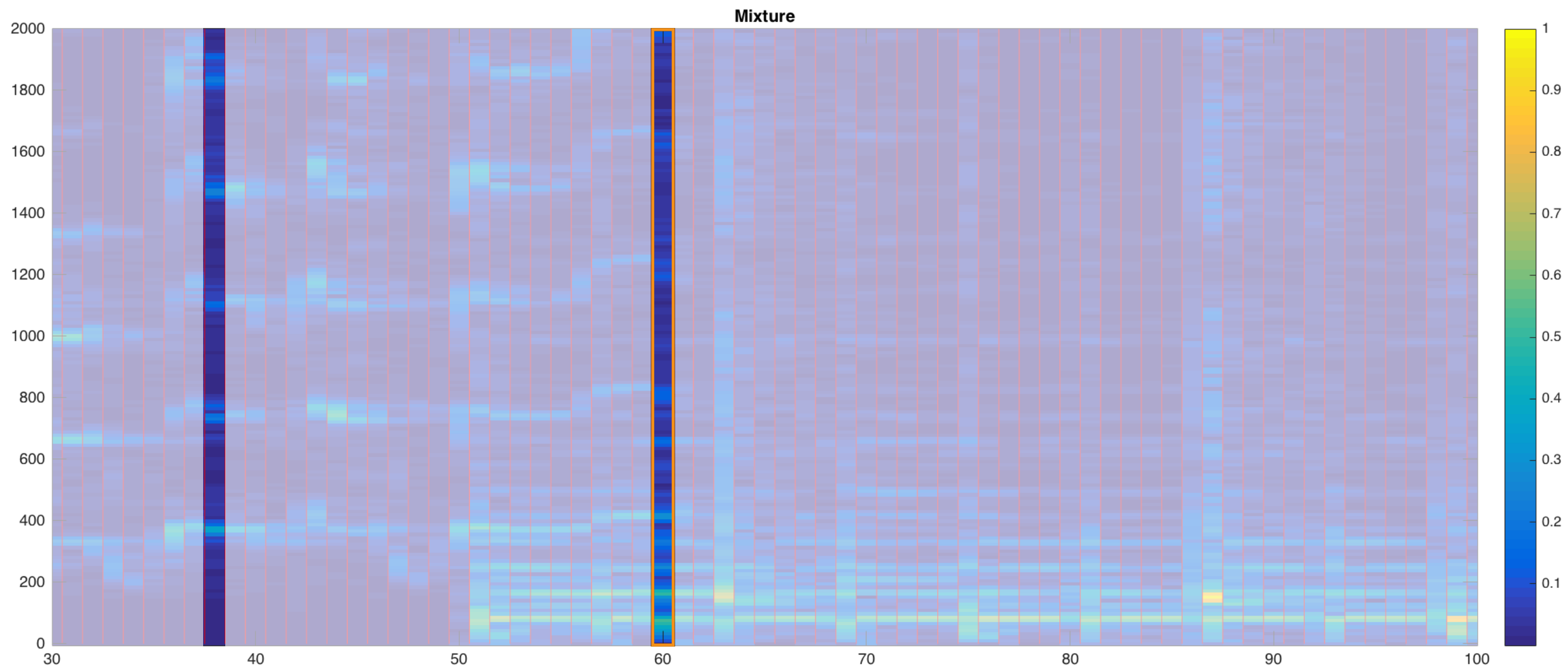
DISTANCE BETWEEN FRAMES

2. LOOK FOR SIMILAR TIME FRAMES



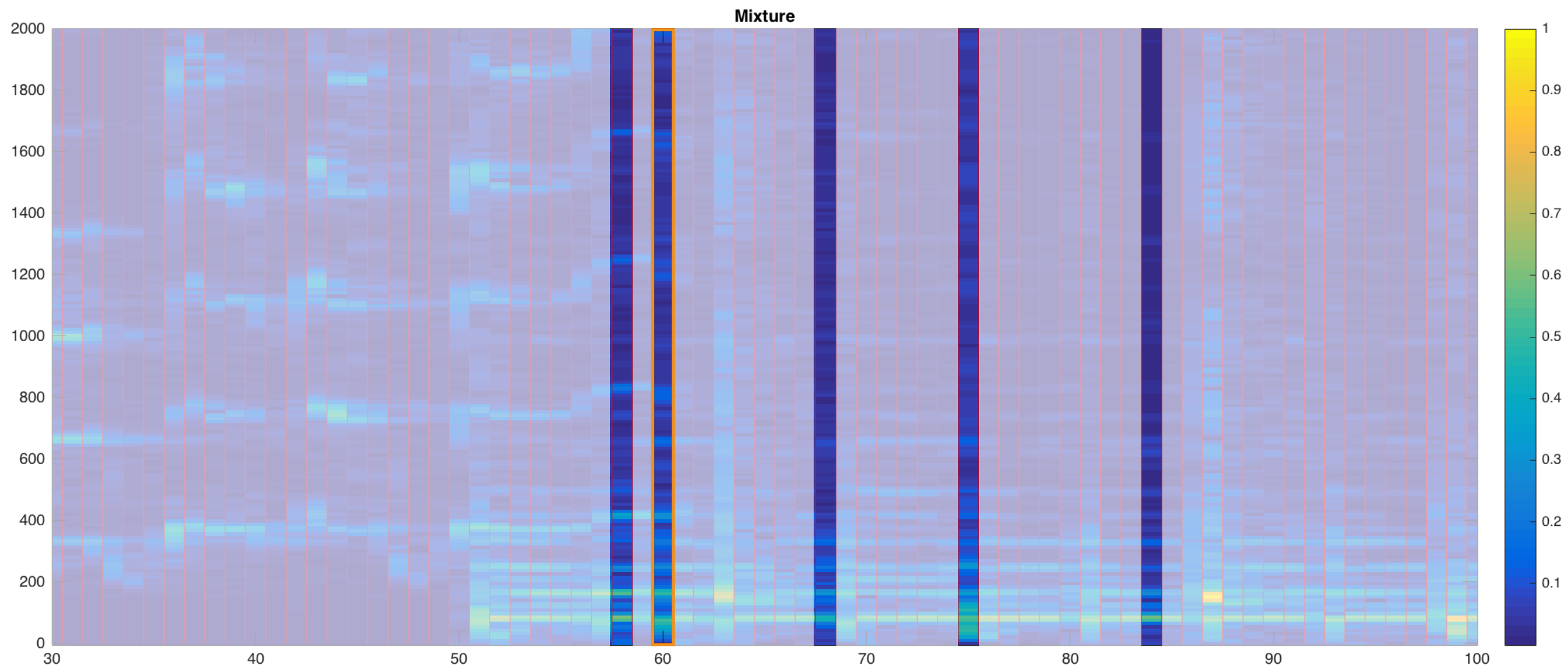
DISTANCE BETWEEN FRAMES

2. LOOK FOR SIMILAR TIME FRAMES

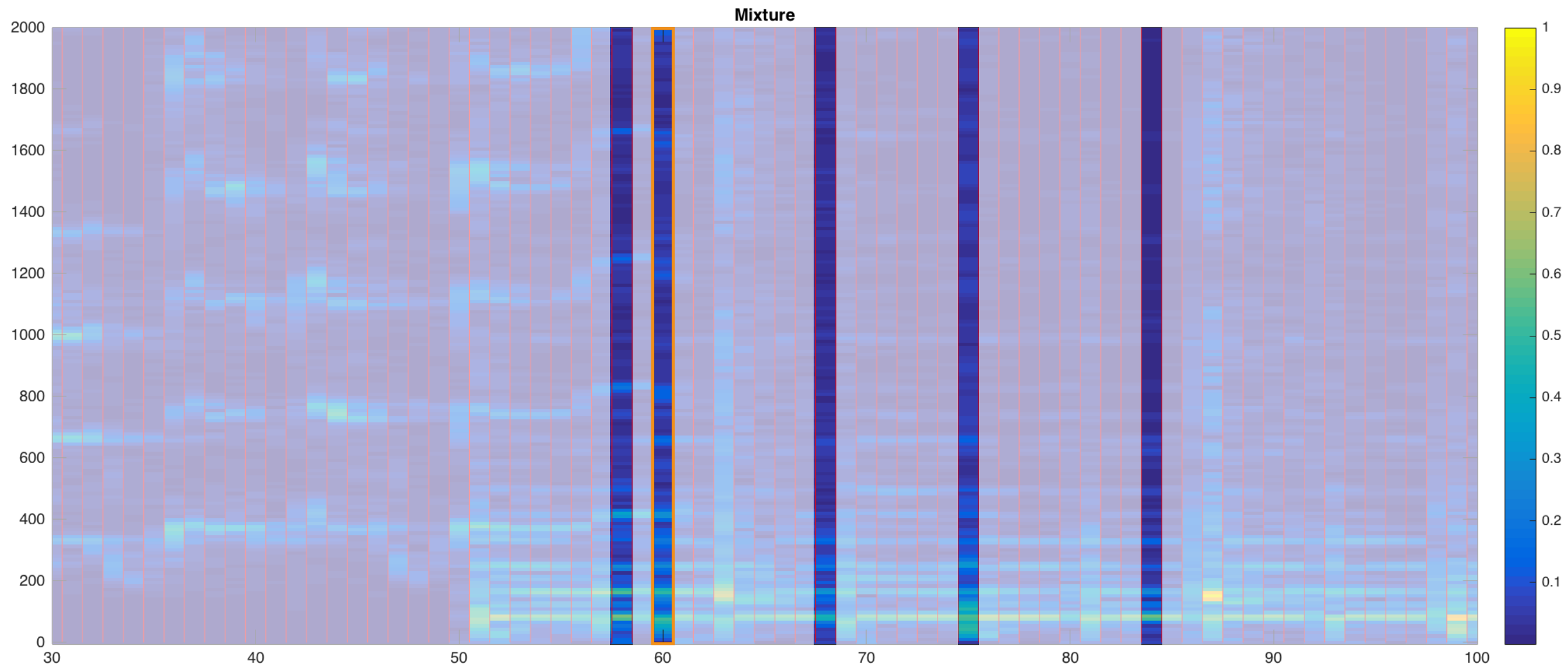


DISTANCE BETWEEN FRAMES

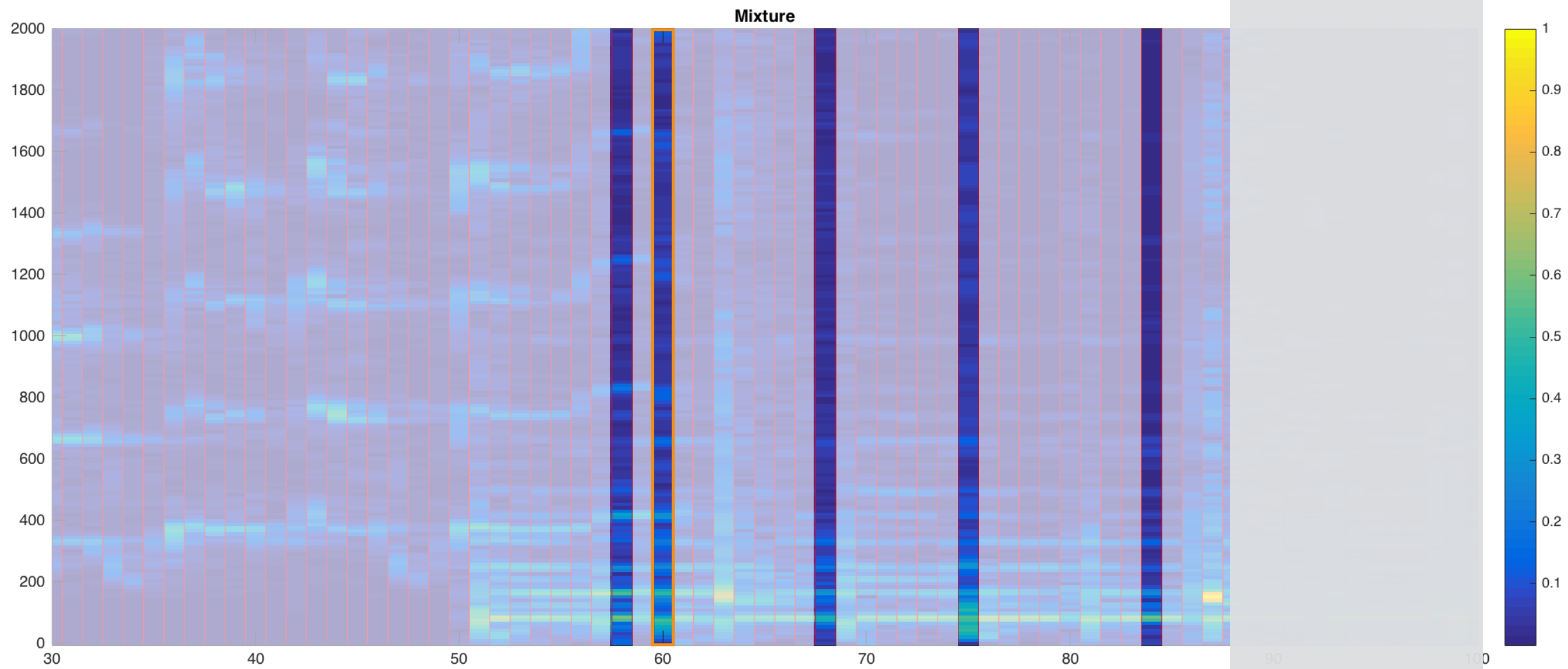
3. STORE THE N-CLOSEST TIME FRAMES



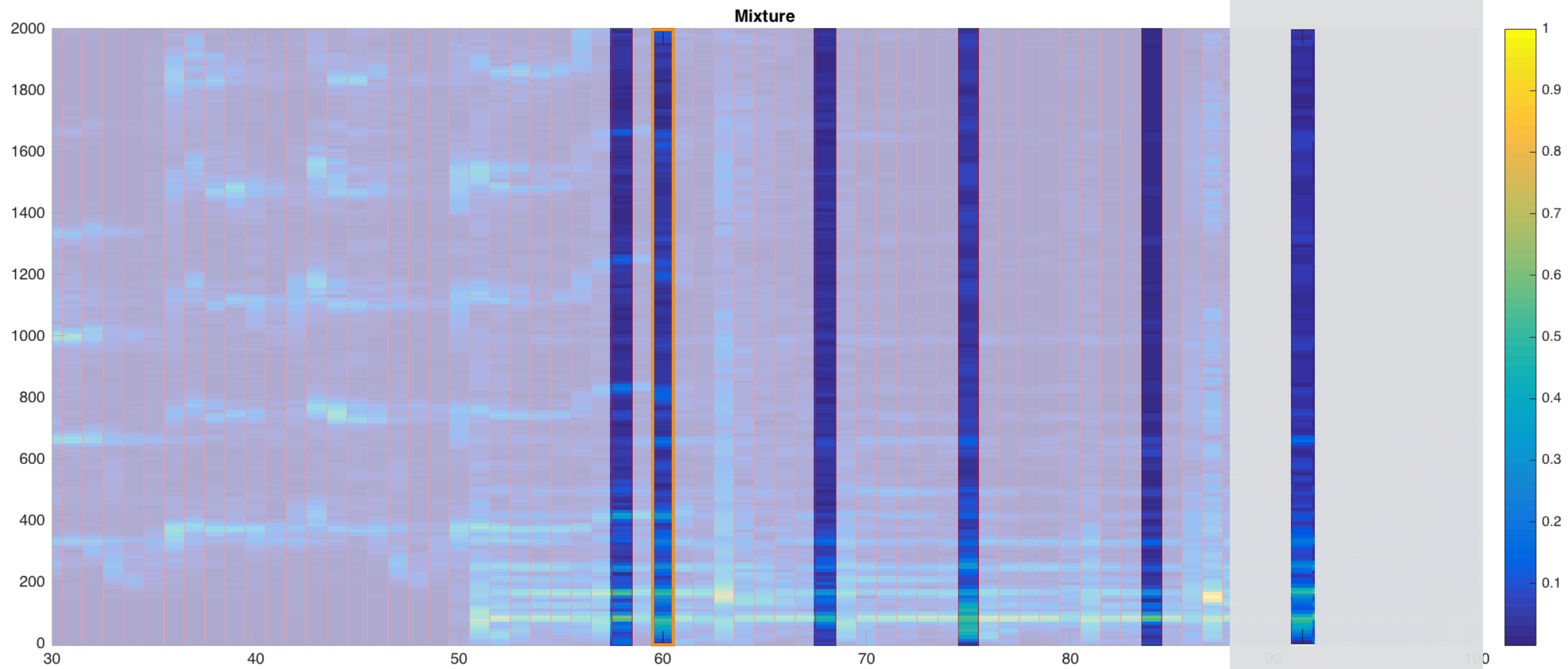
4. MEDIAN FILTERING



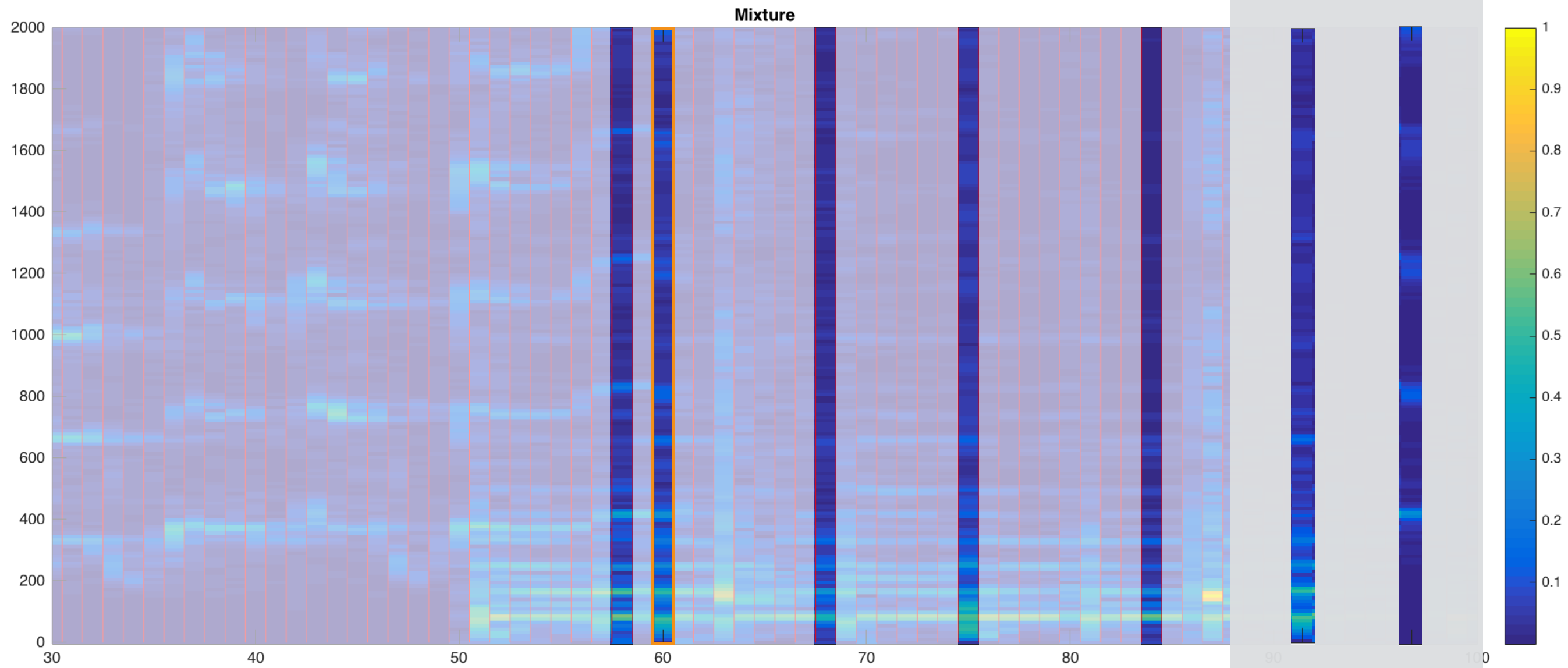
4. MEDIAN FILTERING



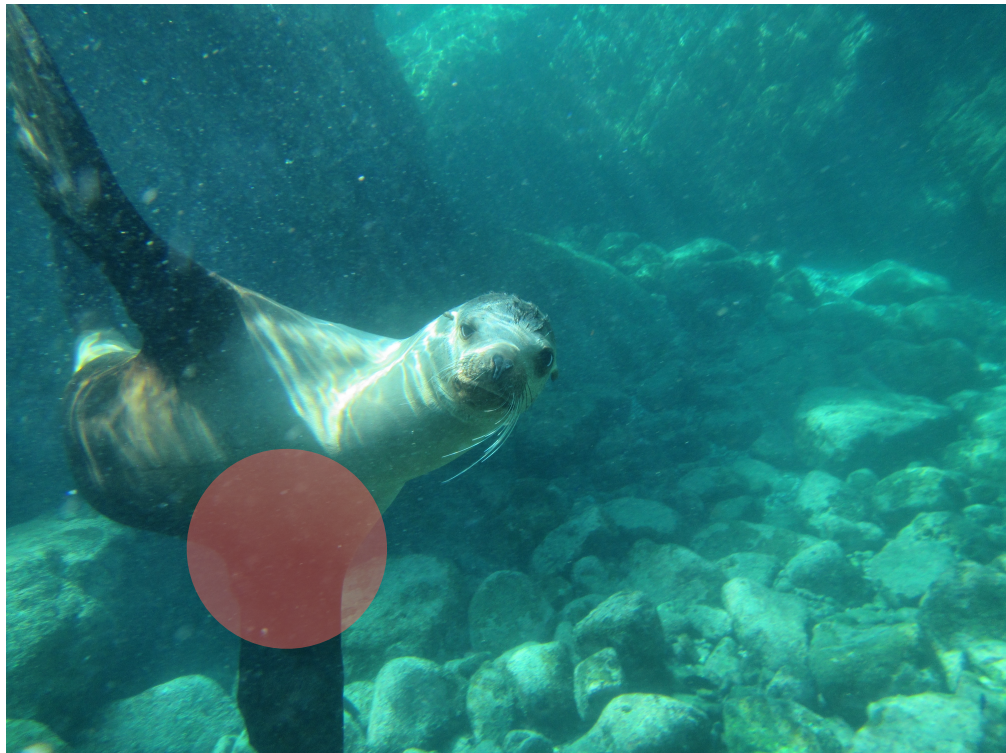
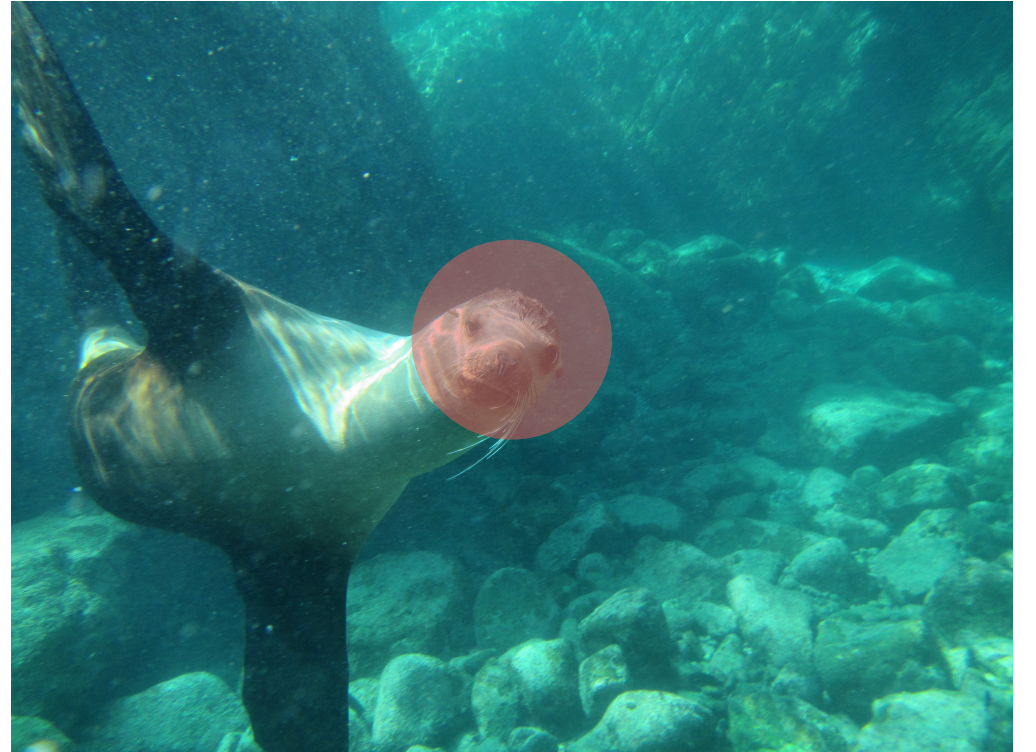
4. MEDIAN FILTERING

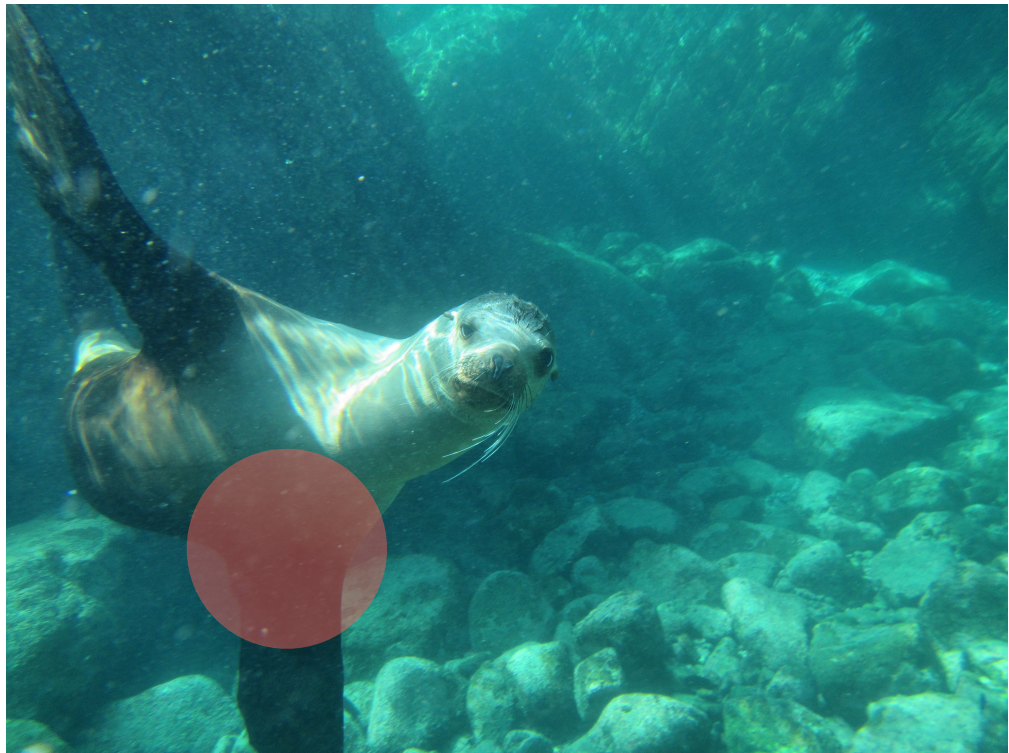
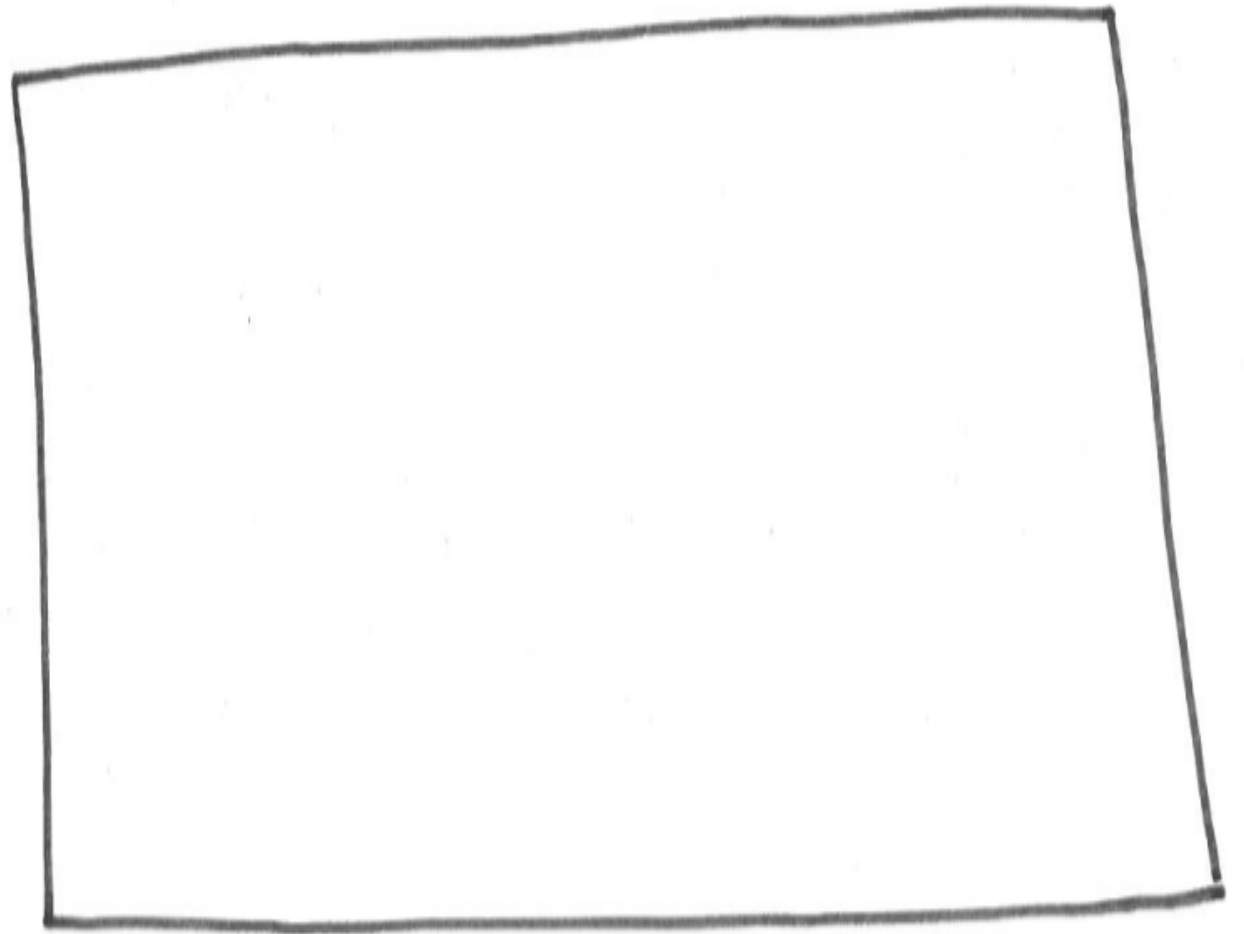


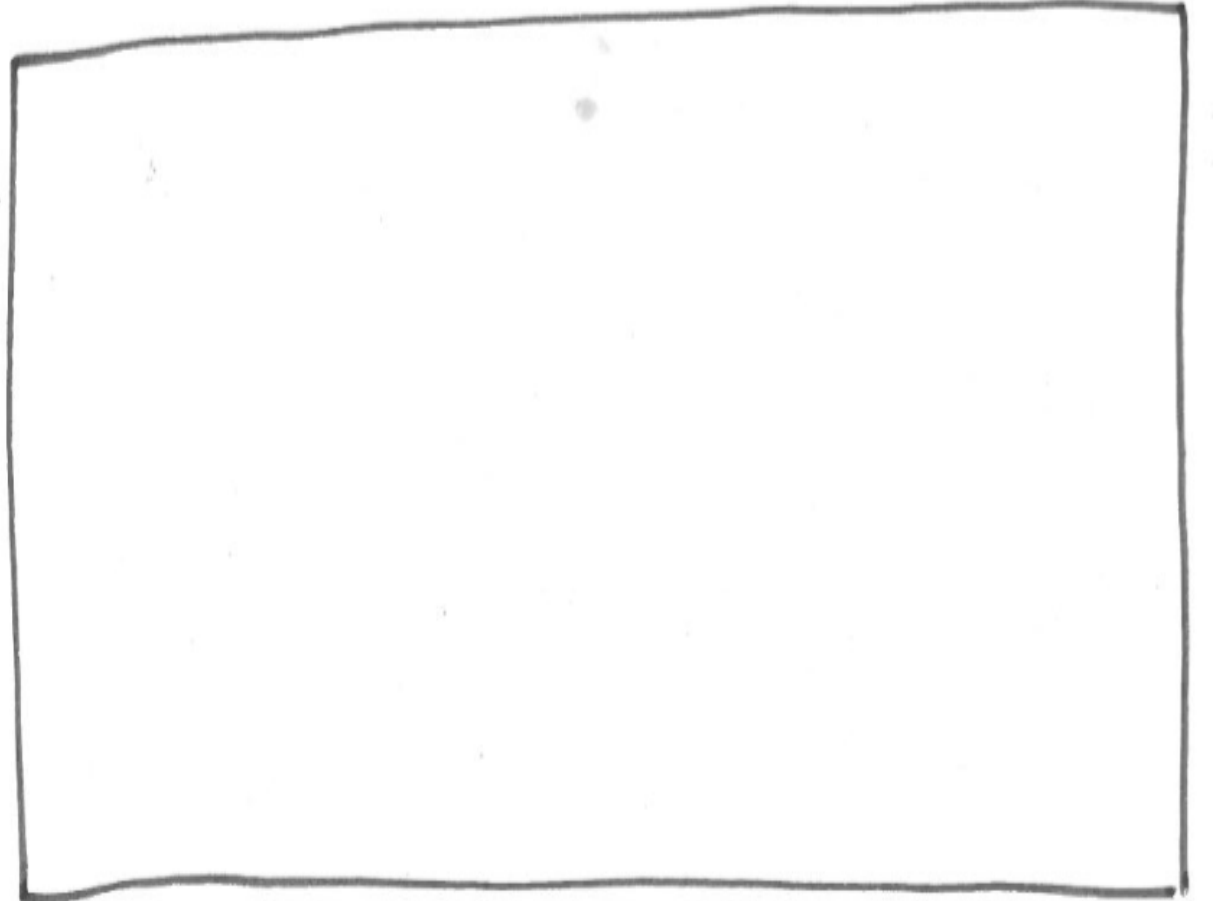
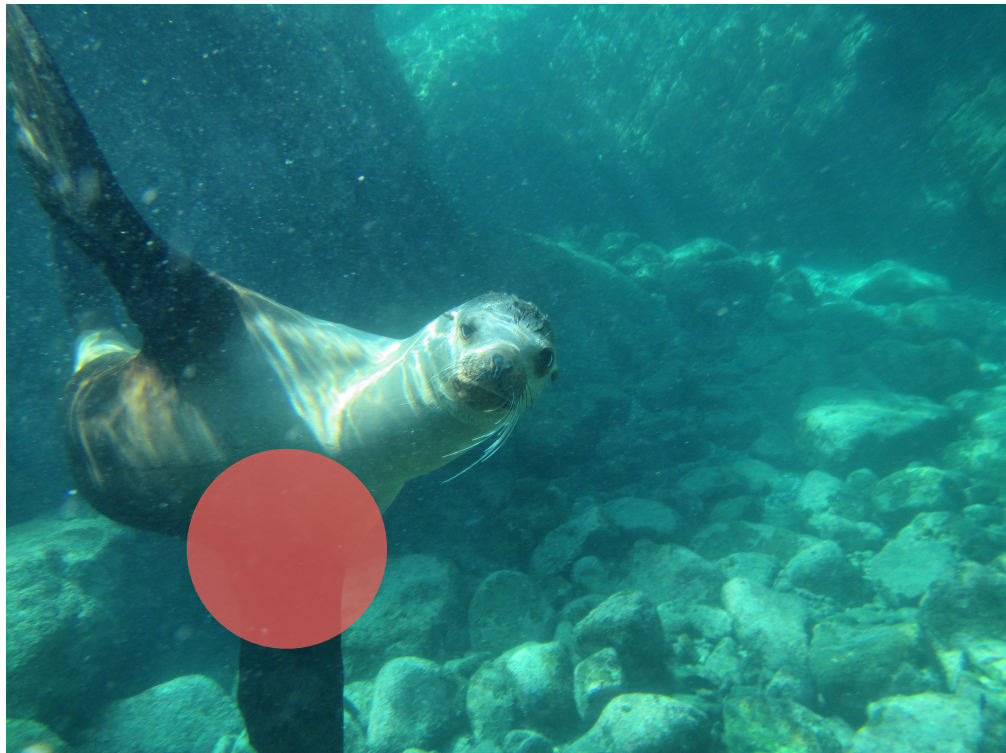
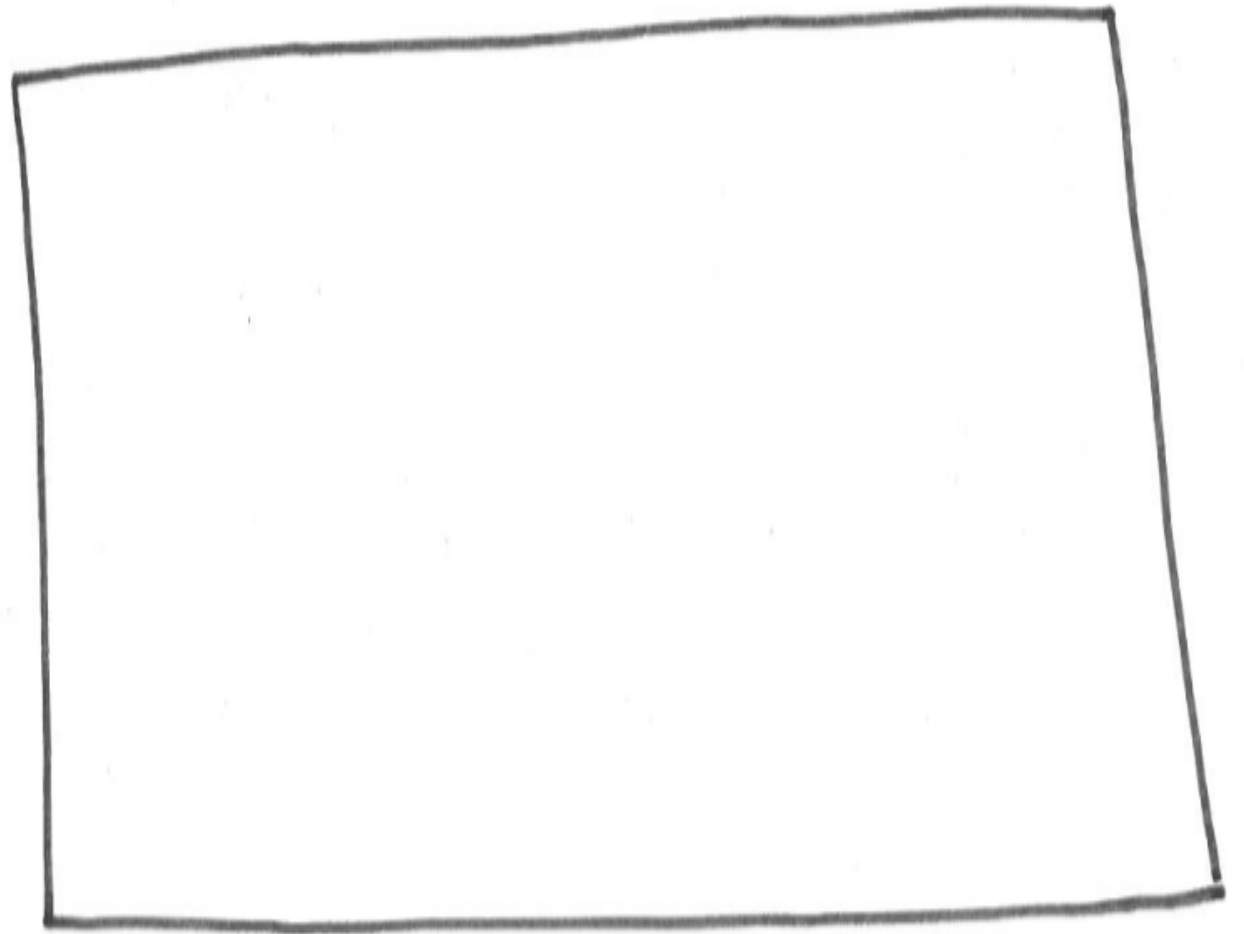
4. MEDIAN FILTERING

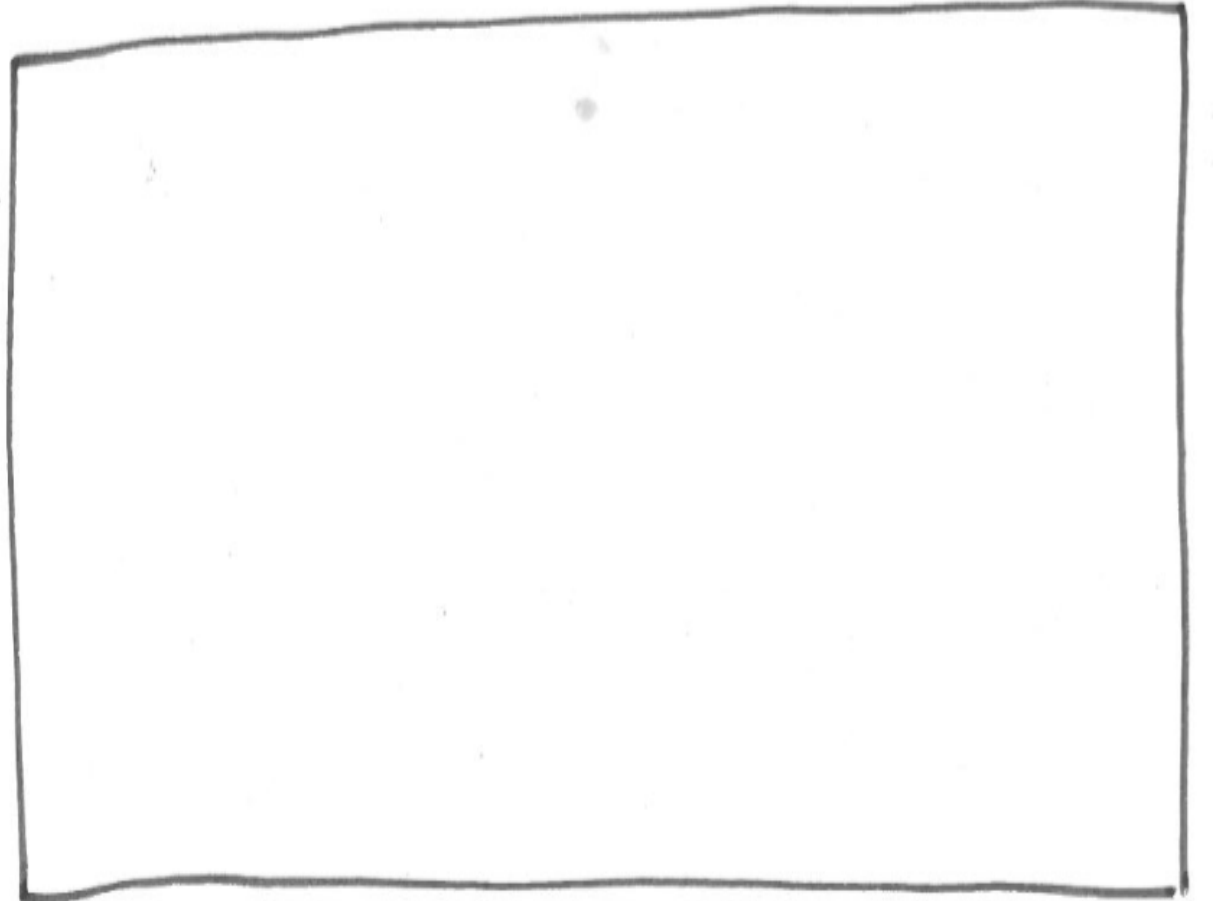
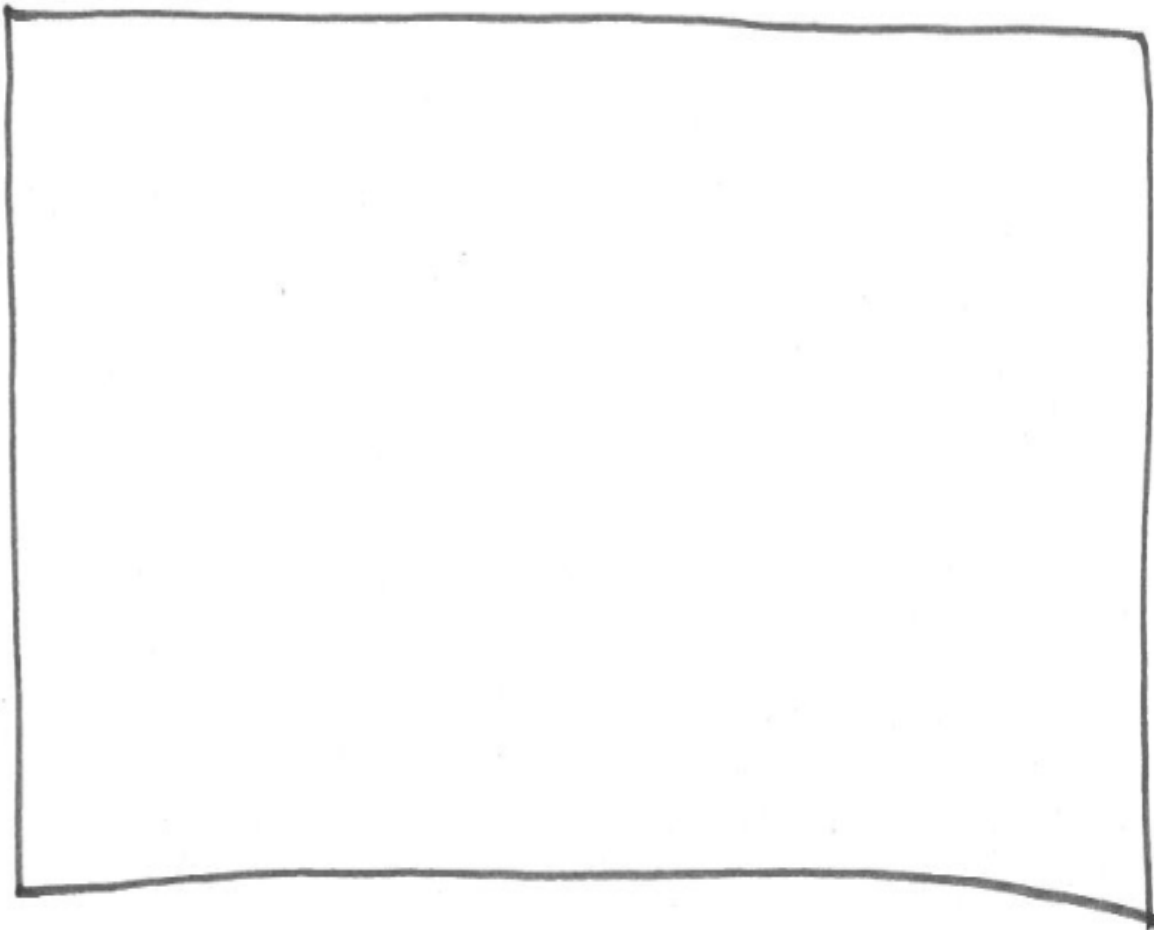
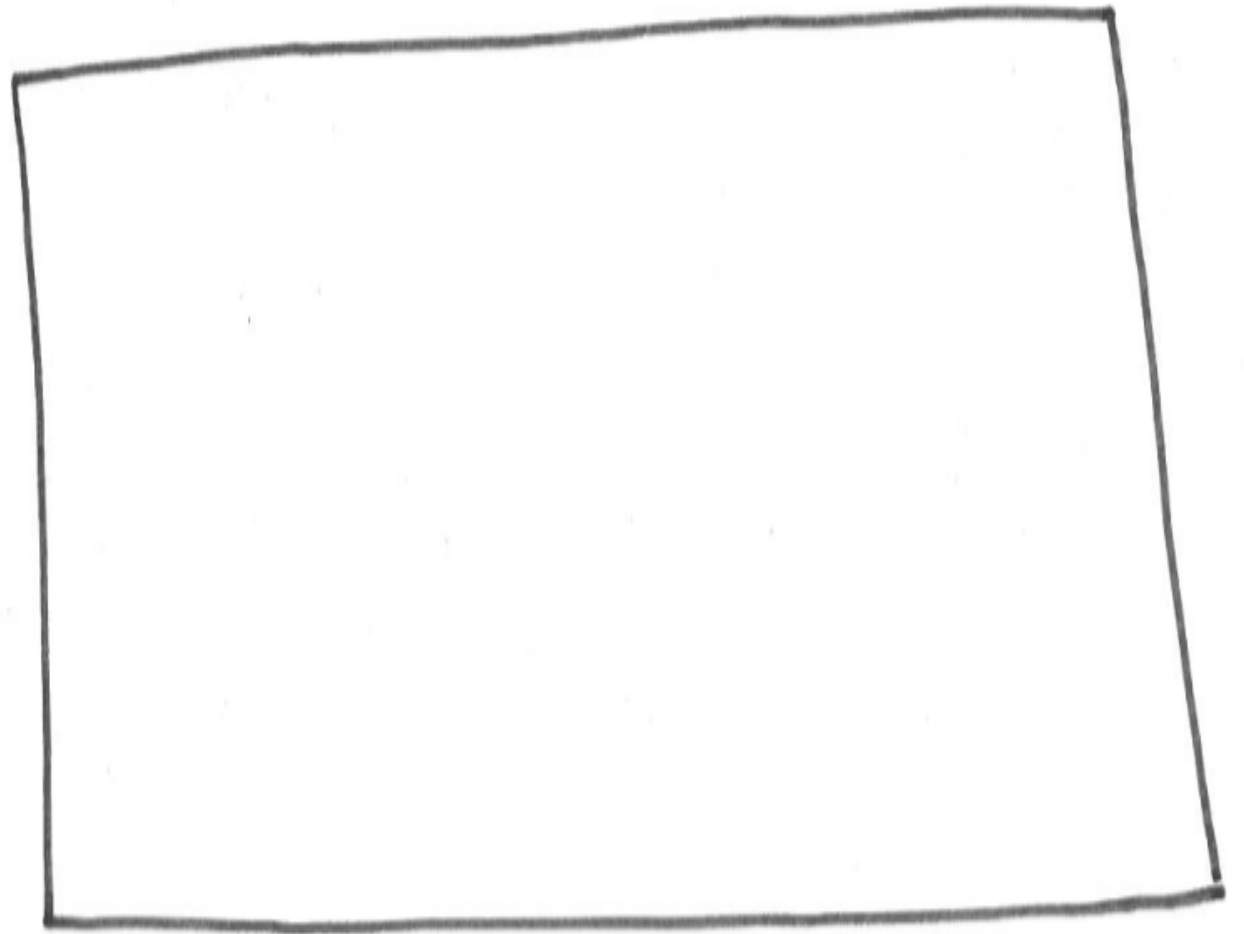




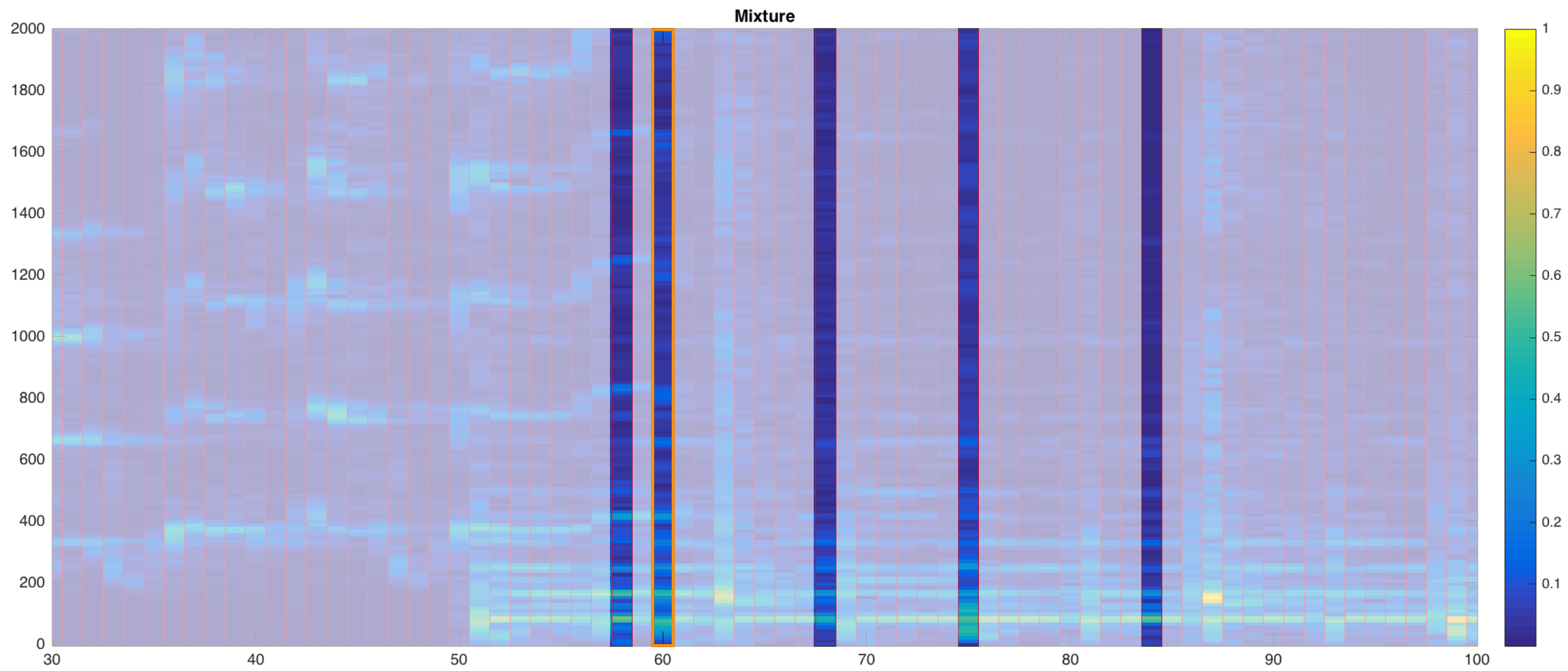


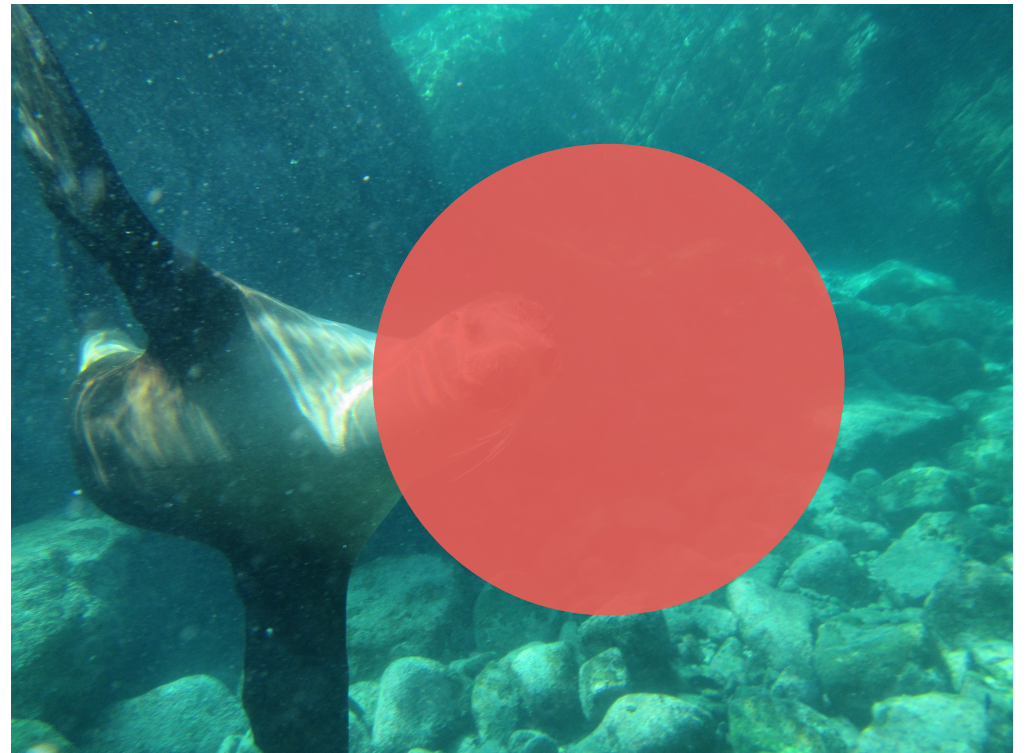
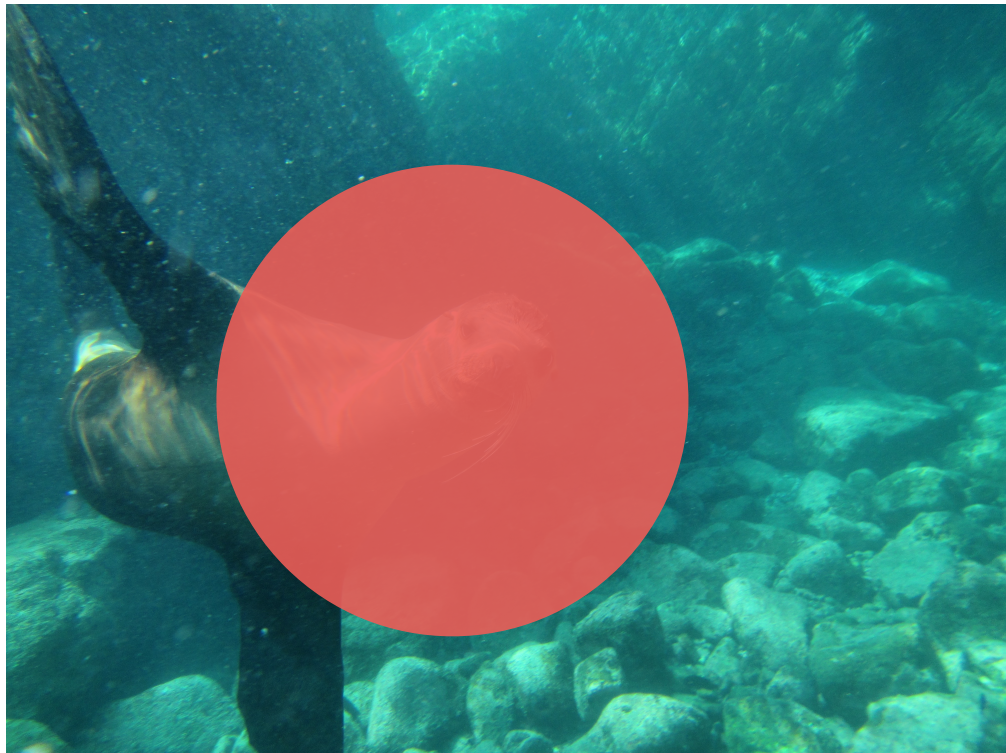
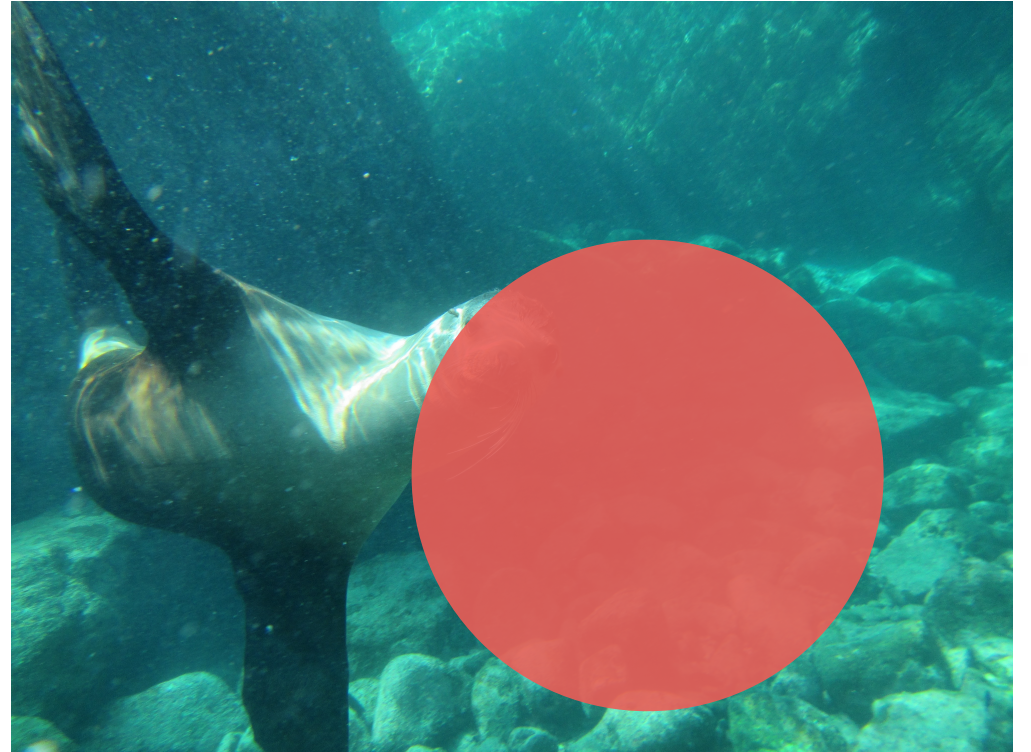
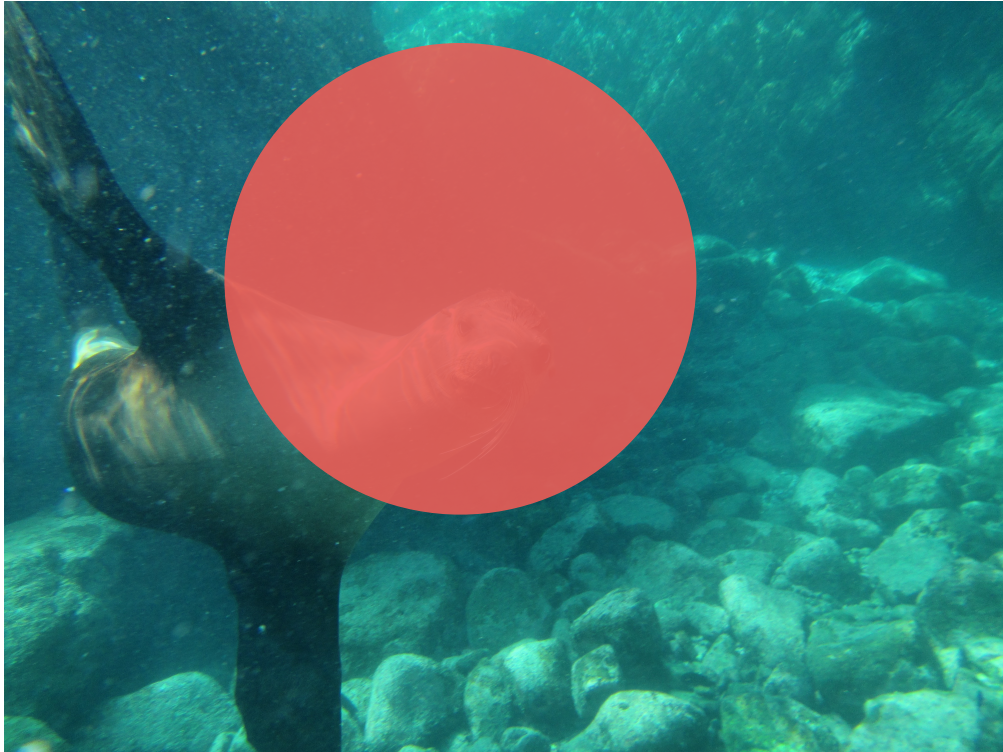


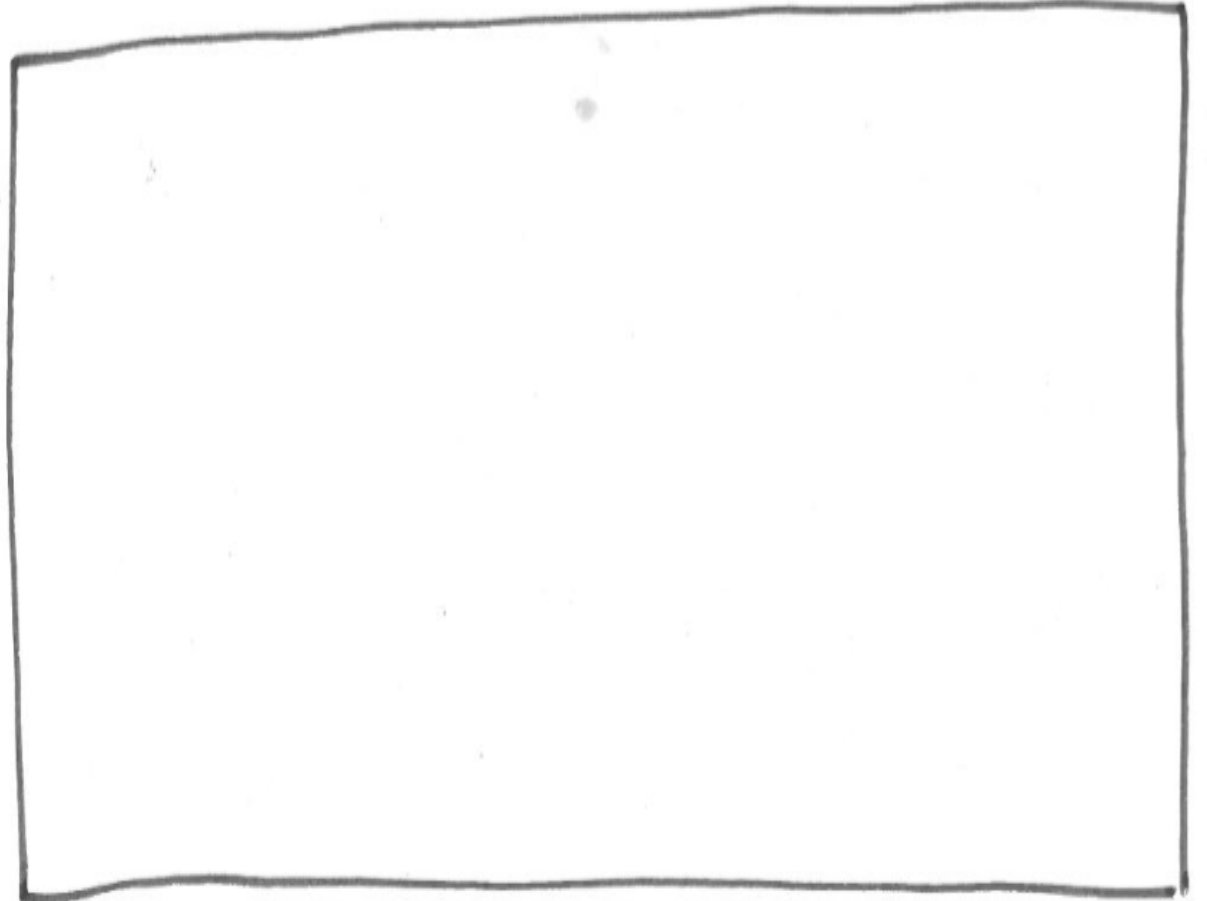
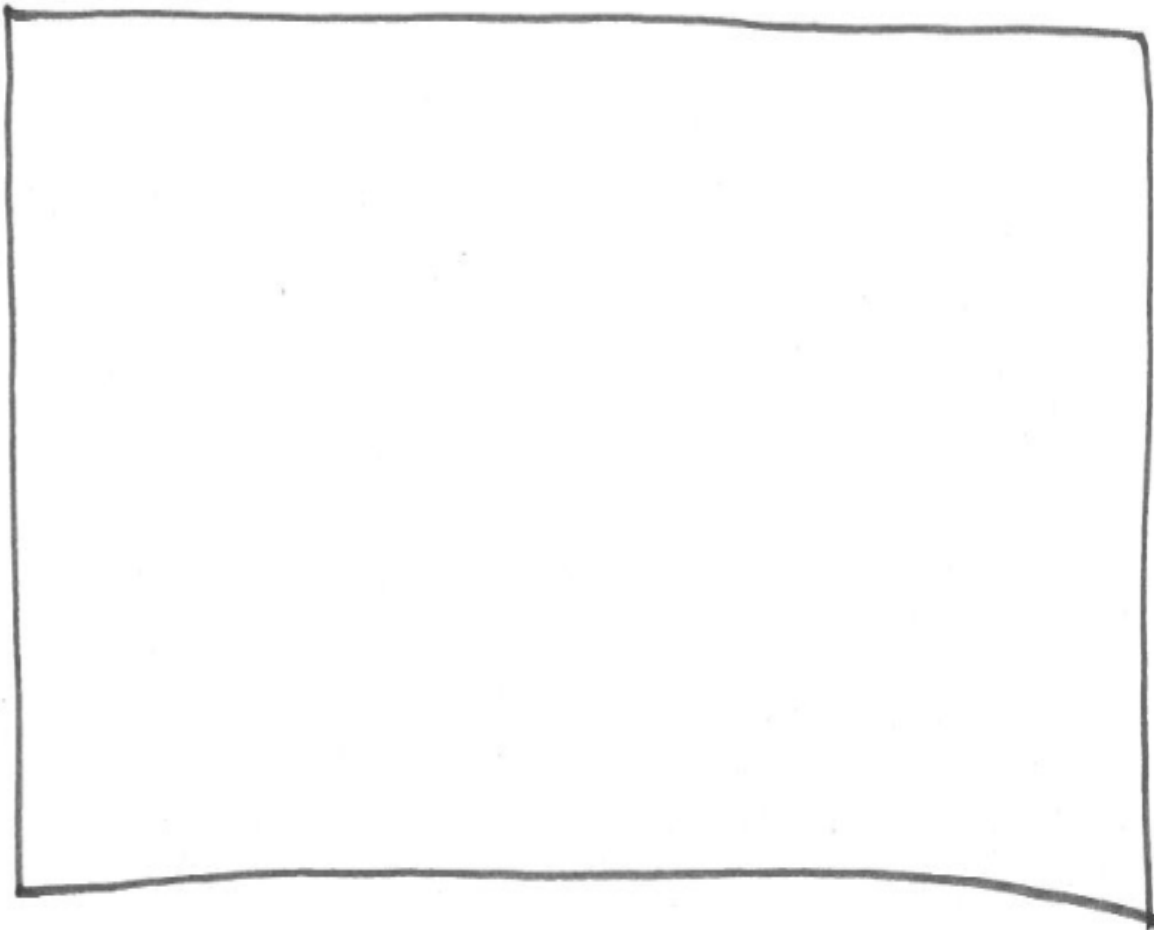
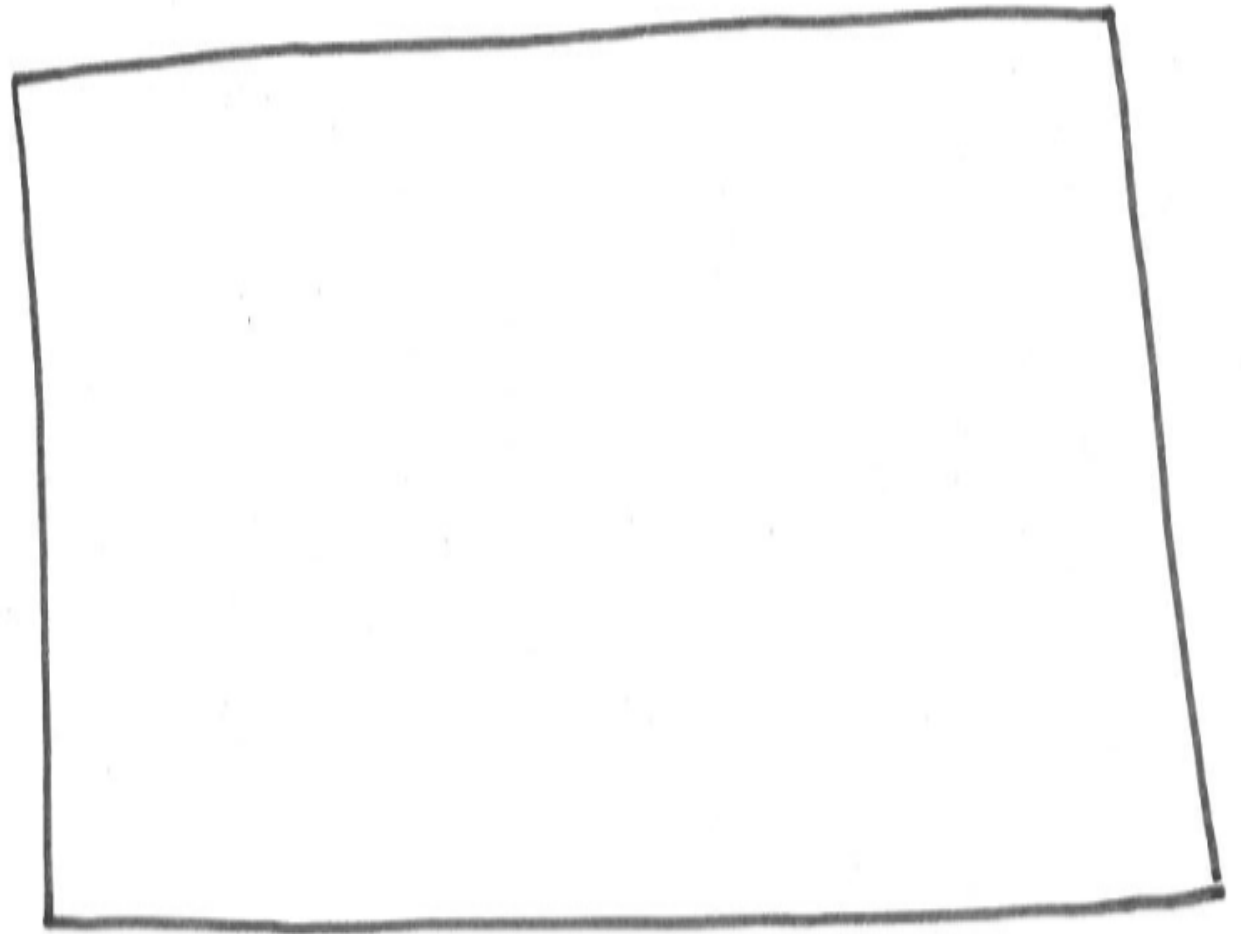
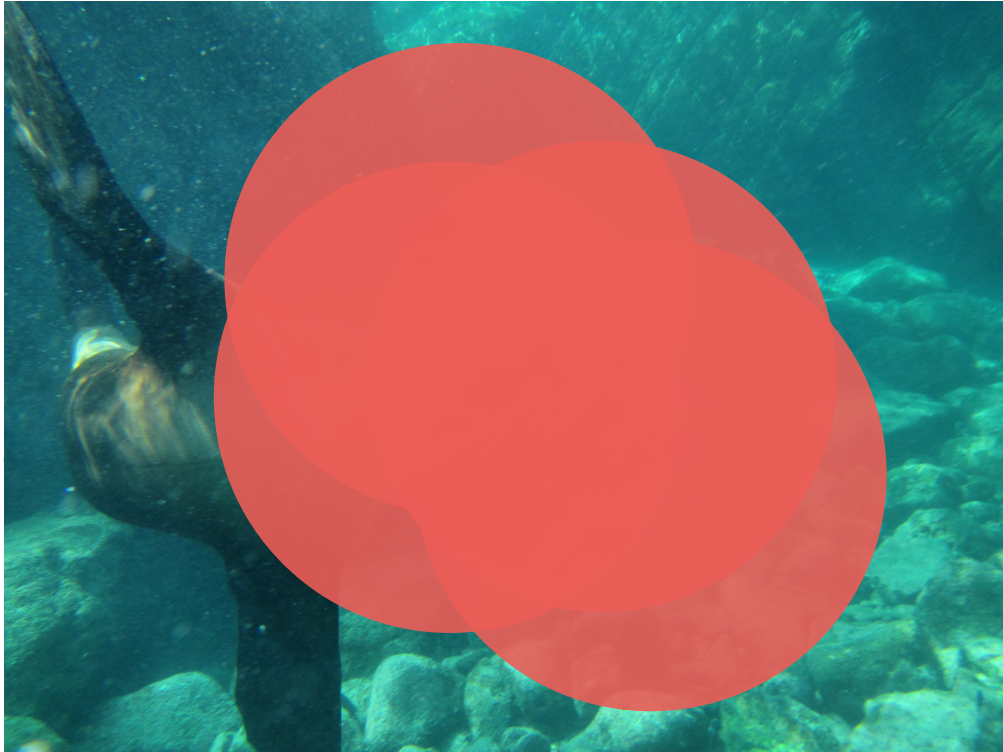




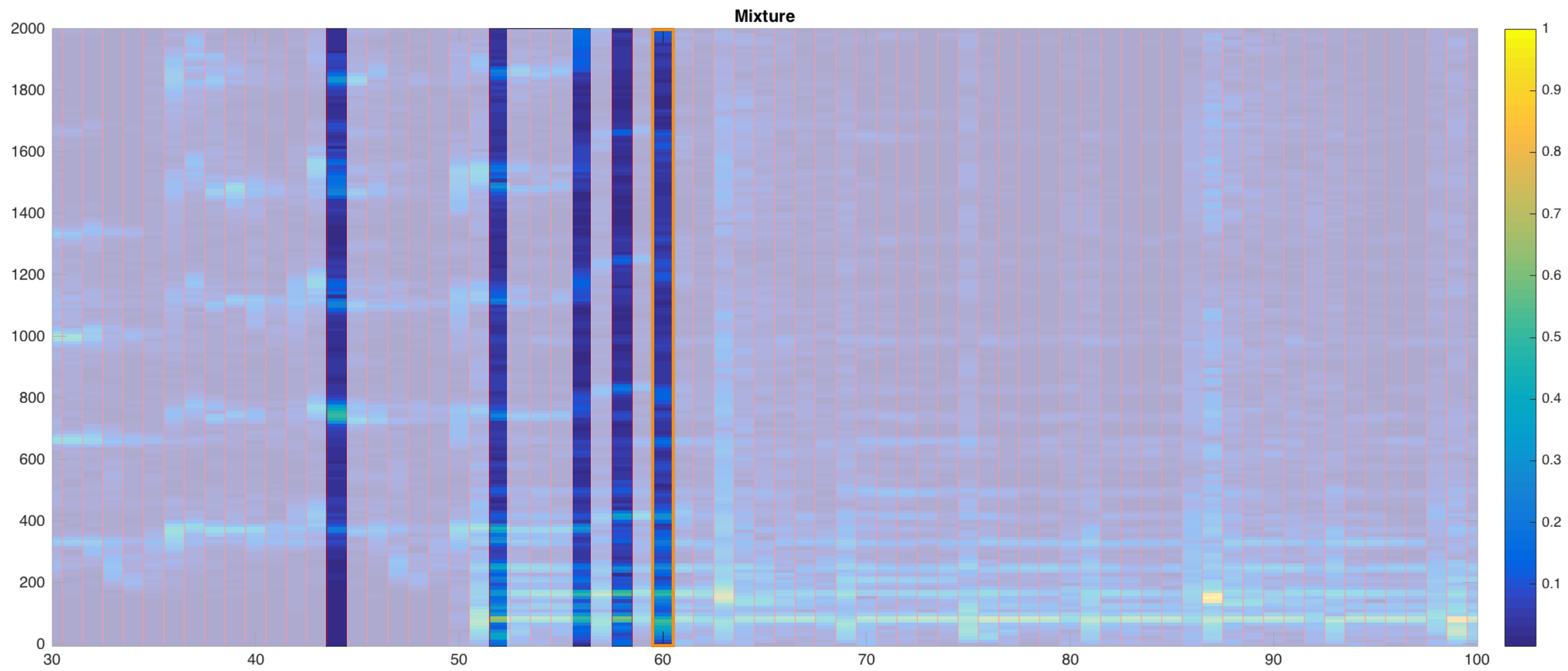
N-CLOSEST TIME FRAMES



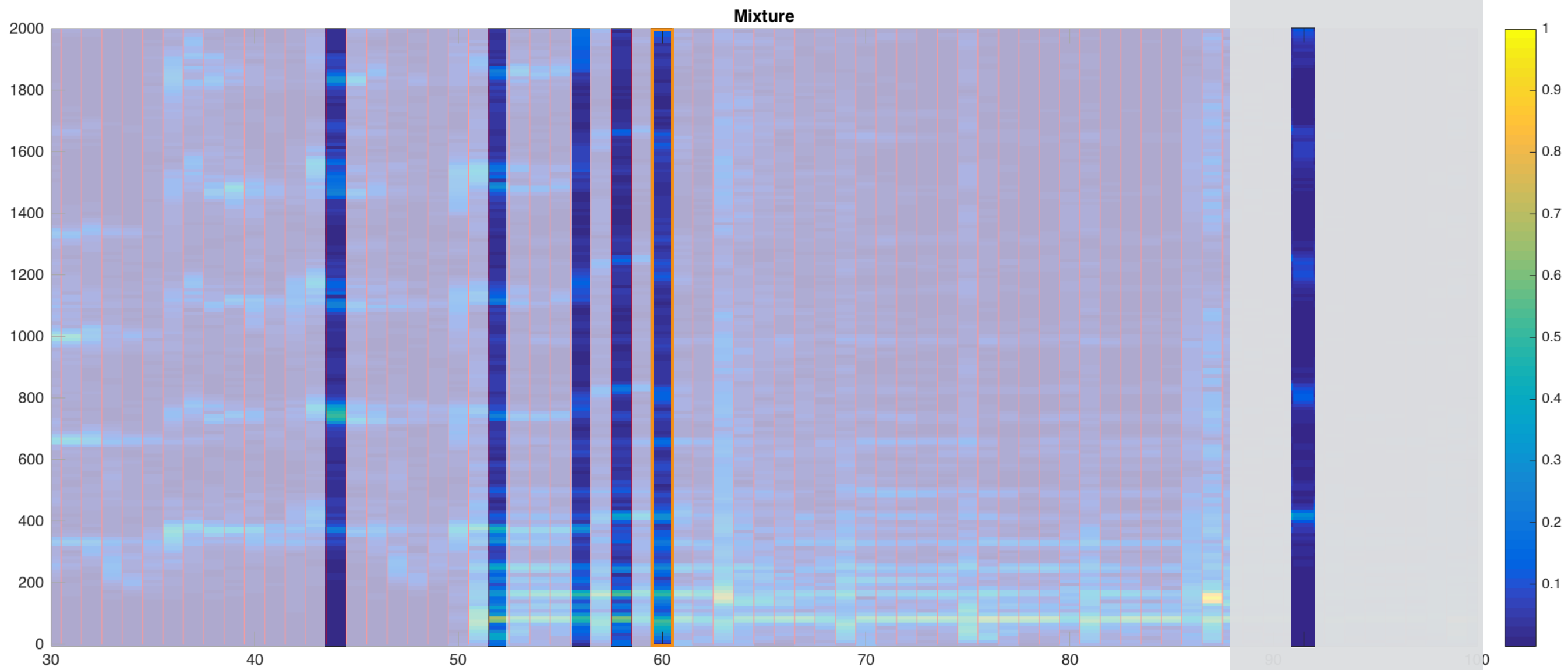




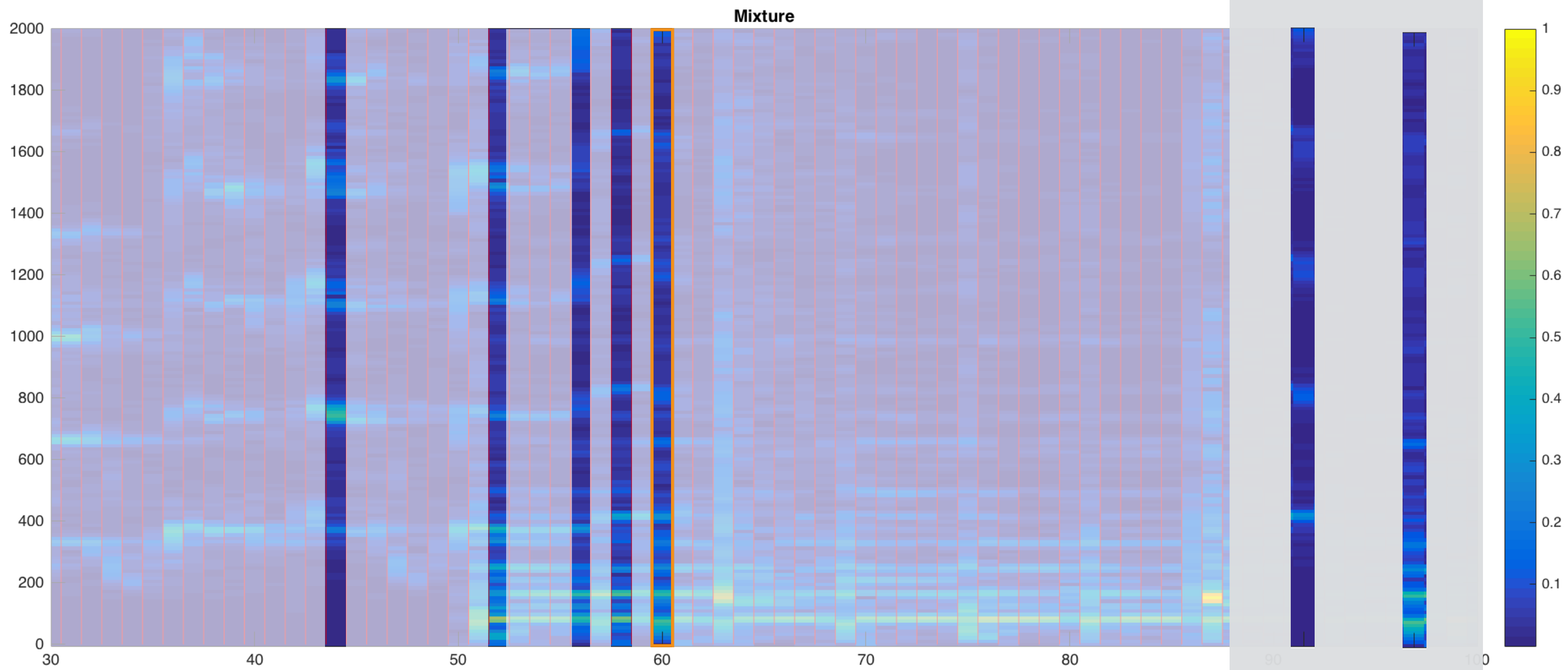
N-CLOSEST TIME FRAMES



N-CLOSEST TIME FRAMES

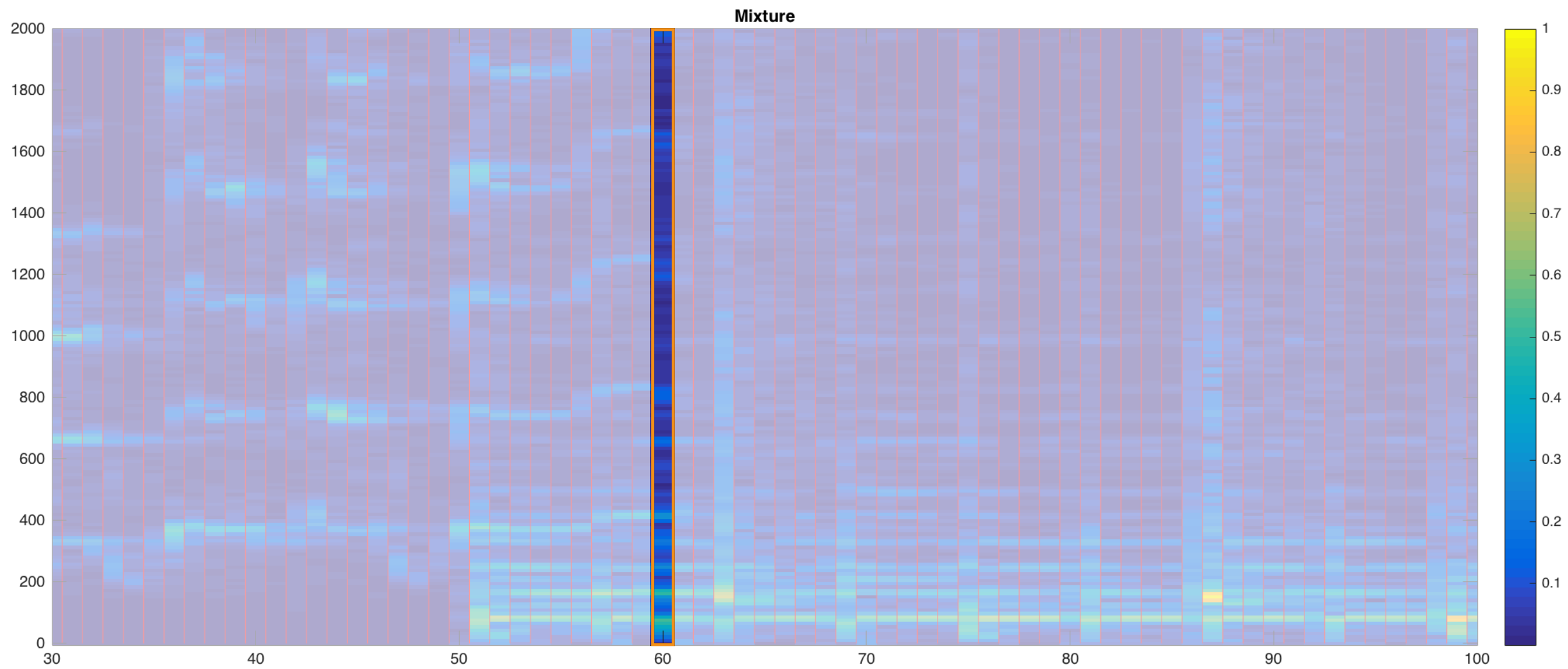


N-CLOSEST TIME FRAMES



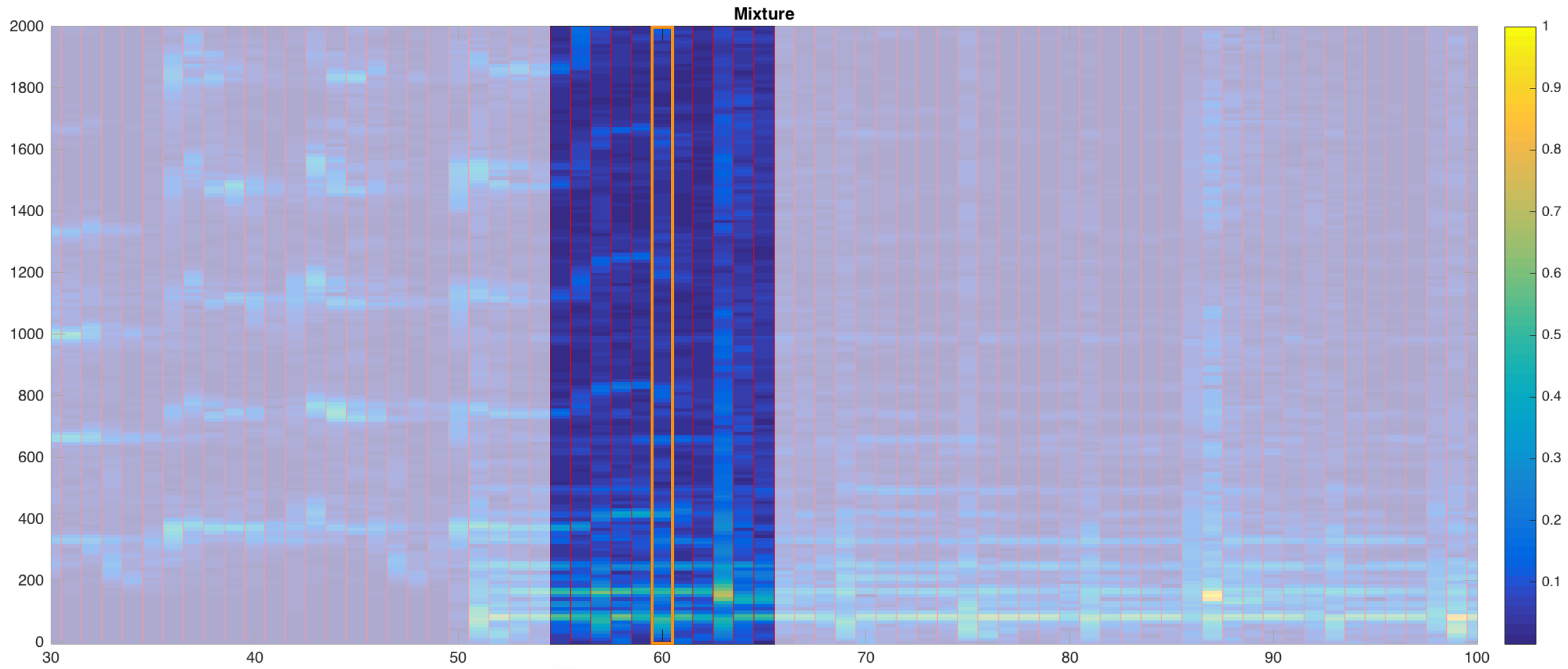
PROPOSED

EXTENSION

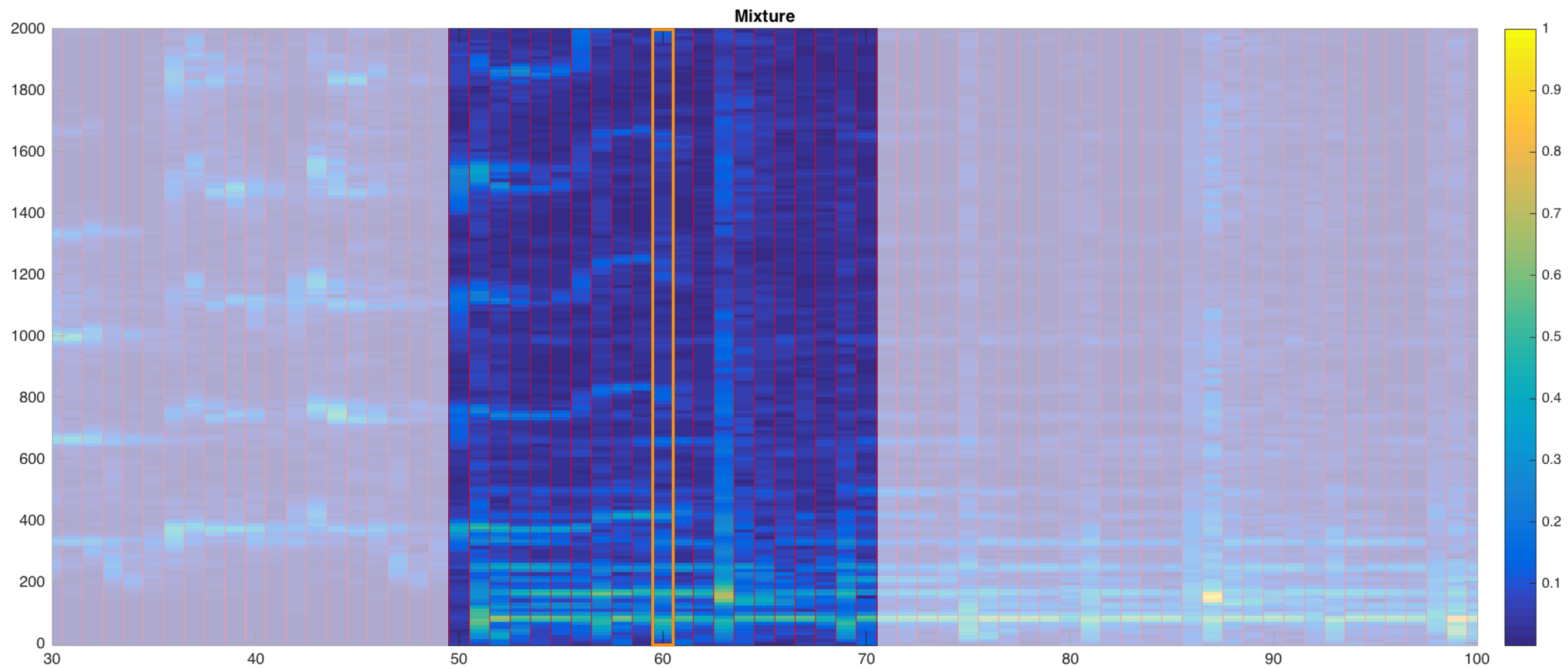


PROPOSED

EXTENSION

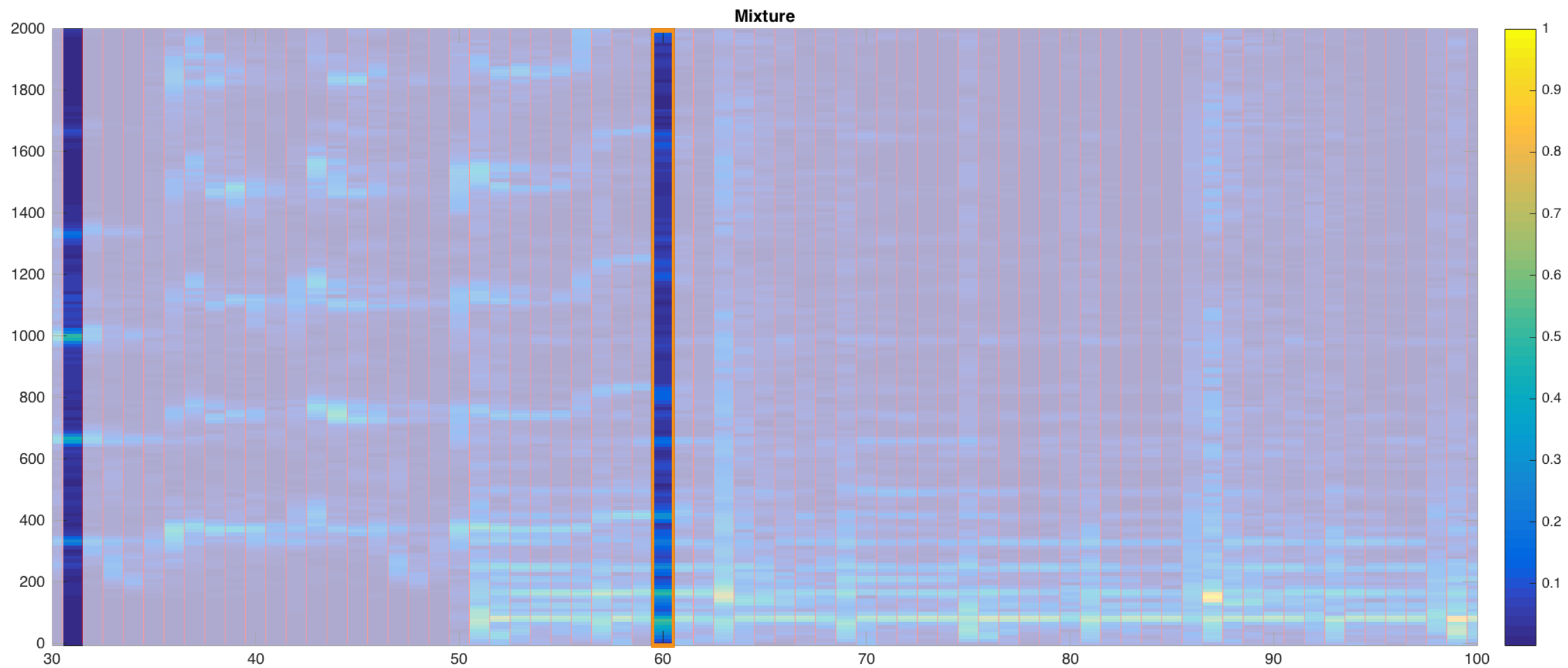


PROPOSED EXTENSION



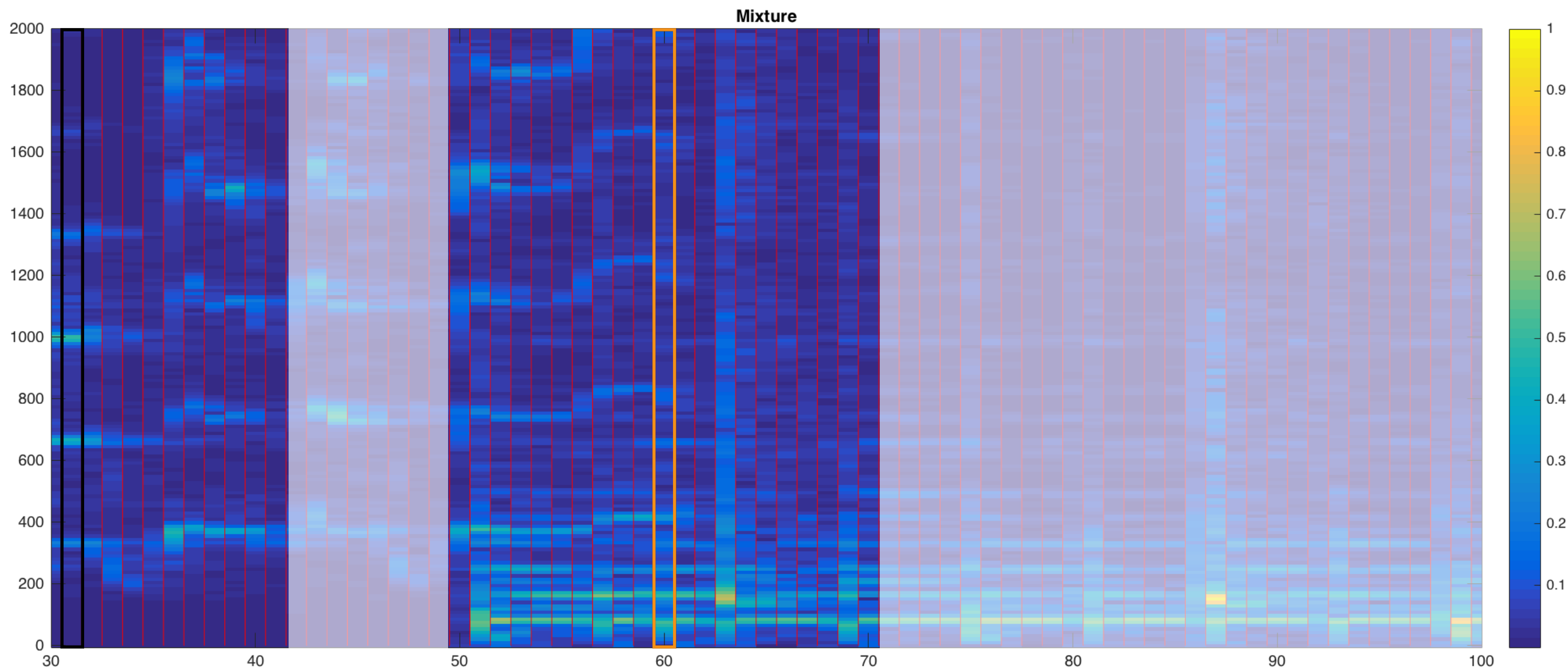
← →
TEMPORAL CONTEXT

2. LOOK FOR SIMILAR TIME FRAMES



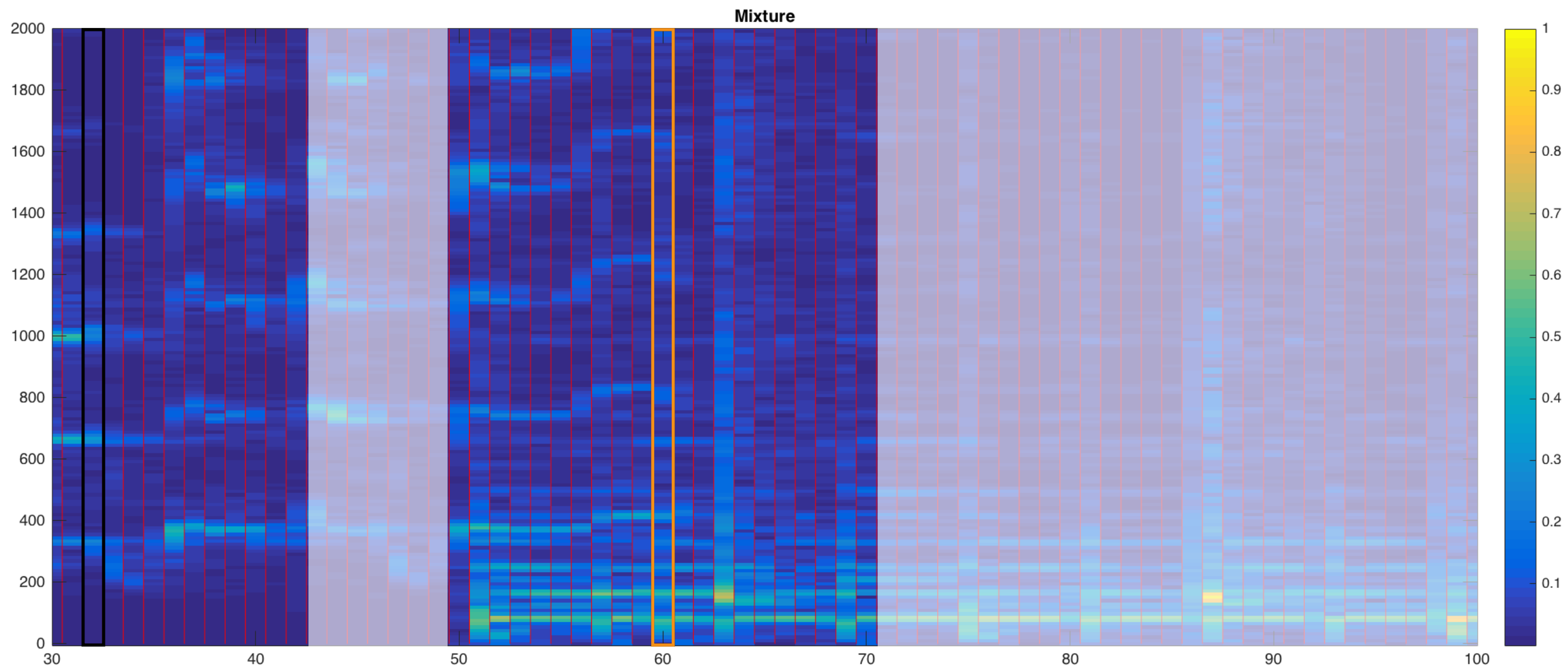
DISTANCE BETWEEN FRAMES

2. LOOK FOR SIMILAR TIME GROUP OF FRAMES



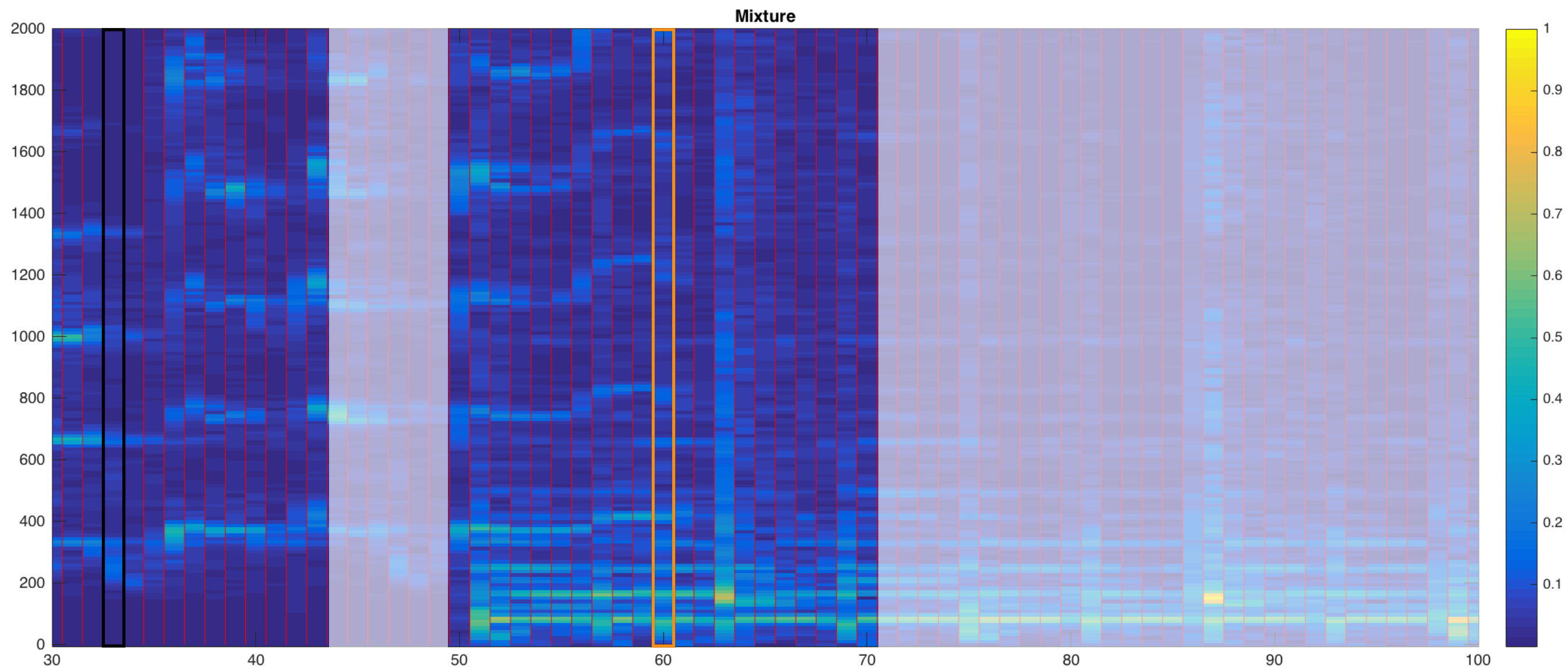
DISTANCE BETWEEN GROUP OF FRAMES

2. LOOK FOR SIMILAR TIME FRAMES



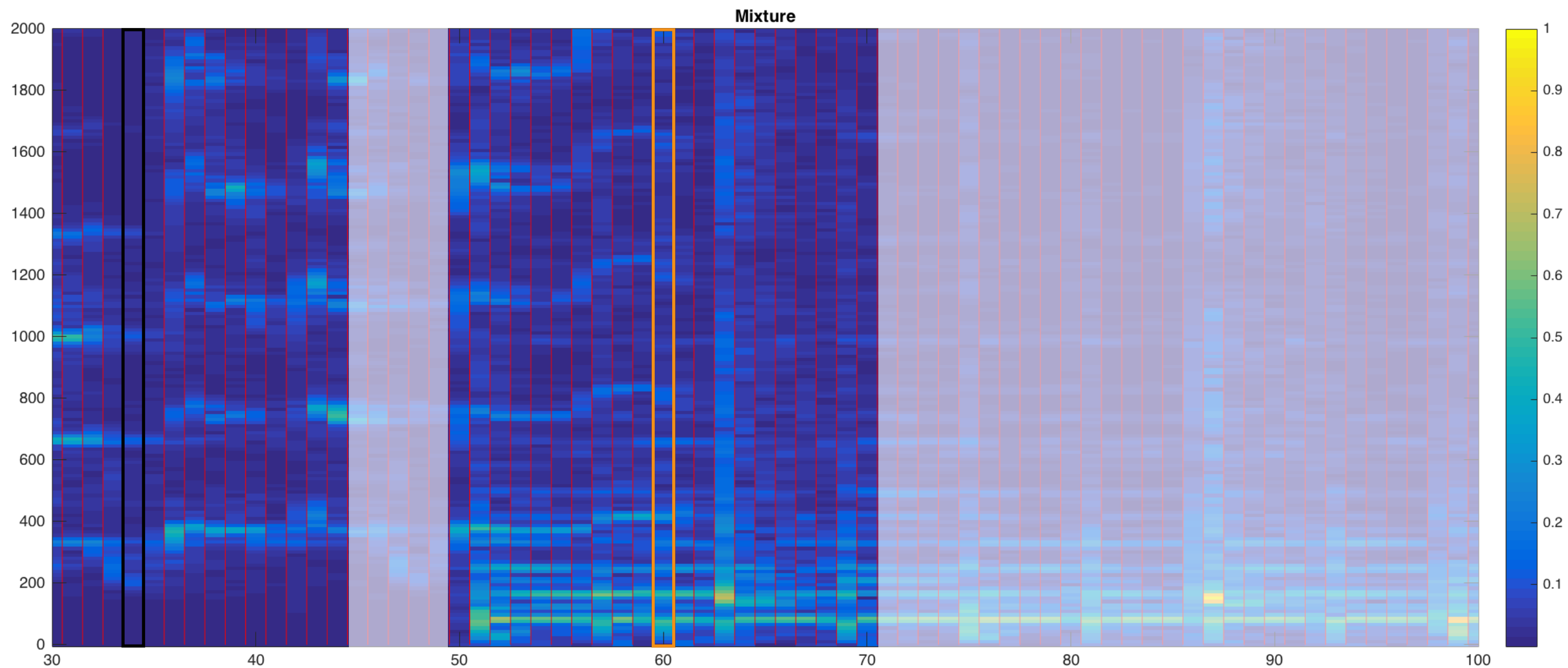
DISTANCE BETWEEN GROUP OF FRAMES

2. LOOK FOR SIMILAR TIME FRAMES



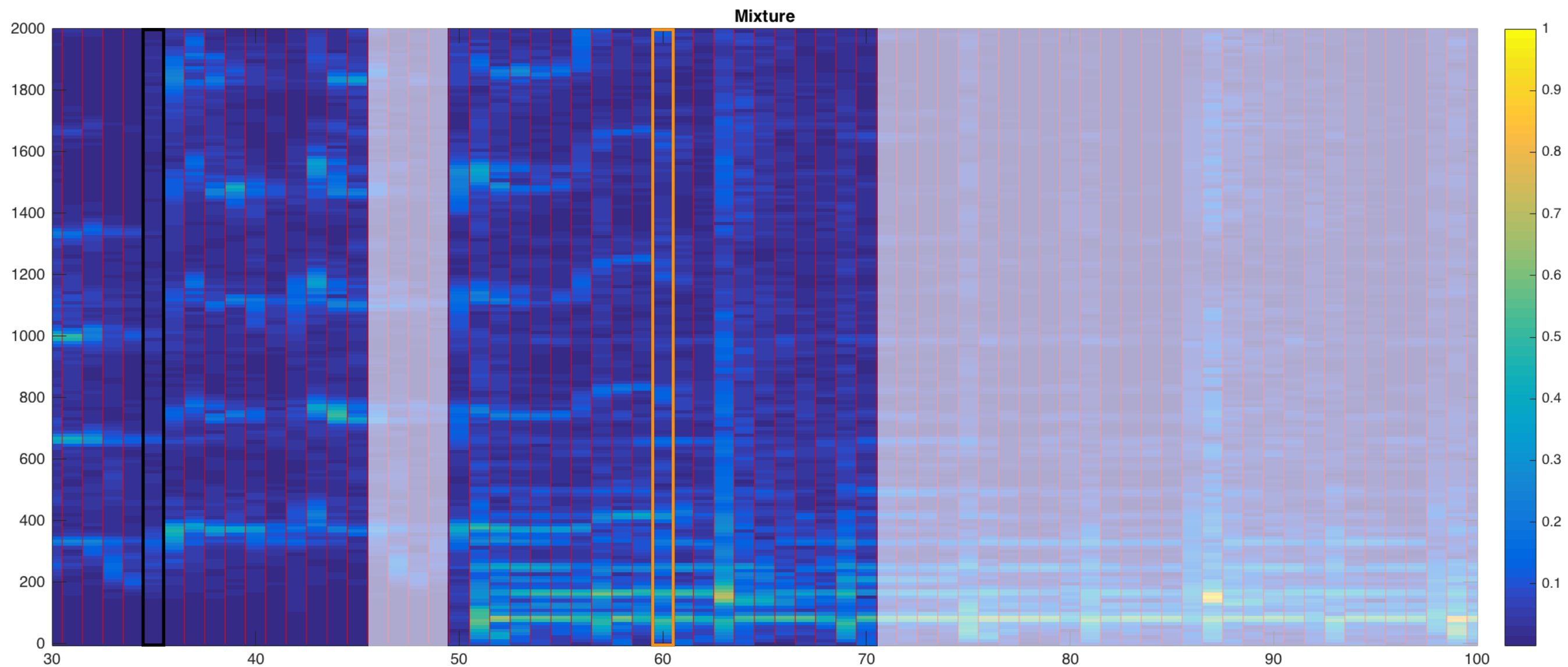
DISTANCE BETWEEN GROUP OF FRAMES

2. LOOK FOR SIMILAR TIME FRAMES



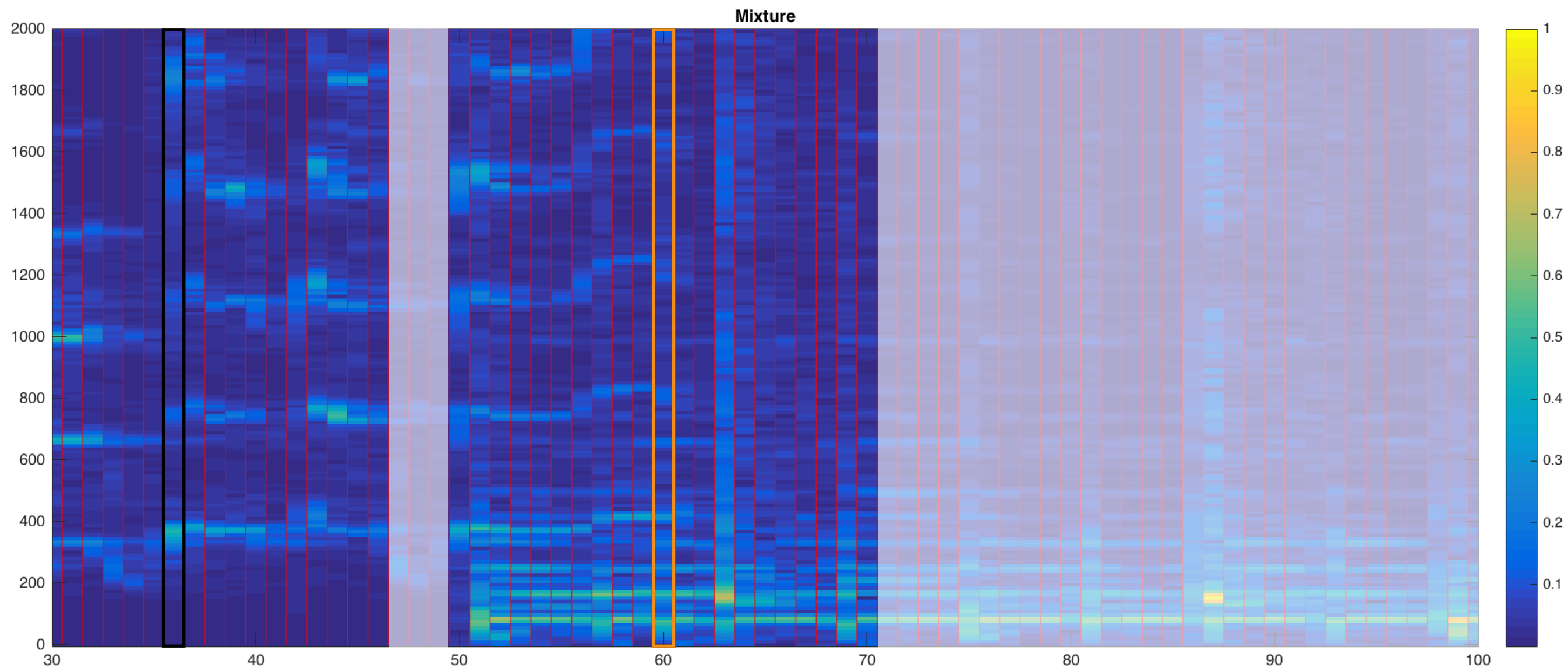
DISTANCE BETWEEN GROUP OF FRAMES

2. LOOK FOR SIMILAR TIME FRAMES

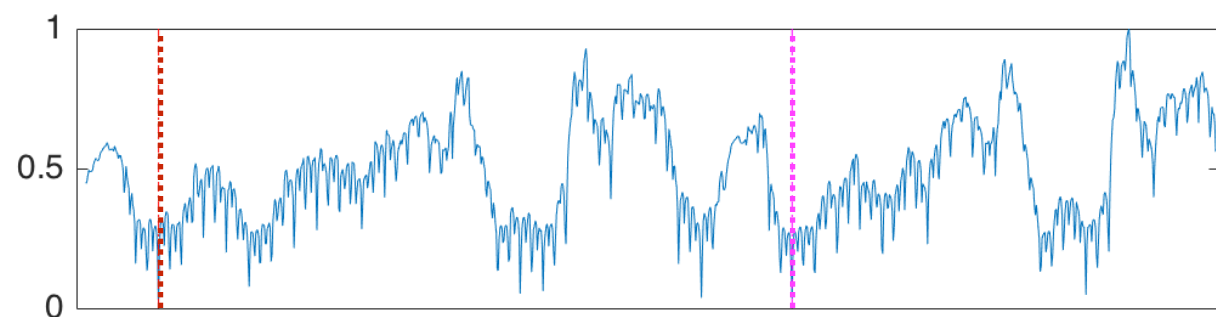
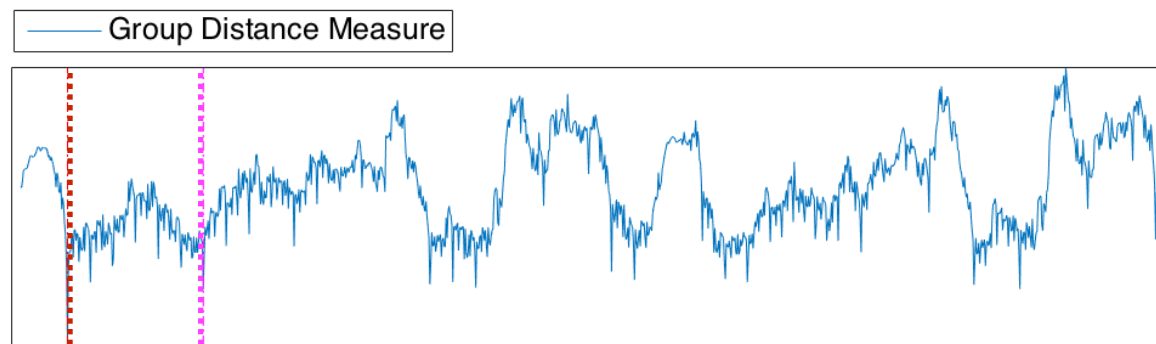
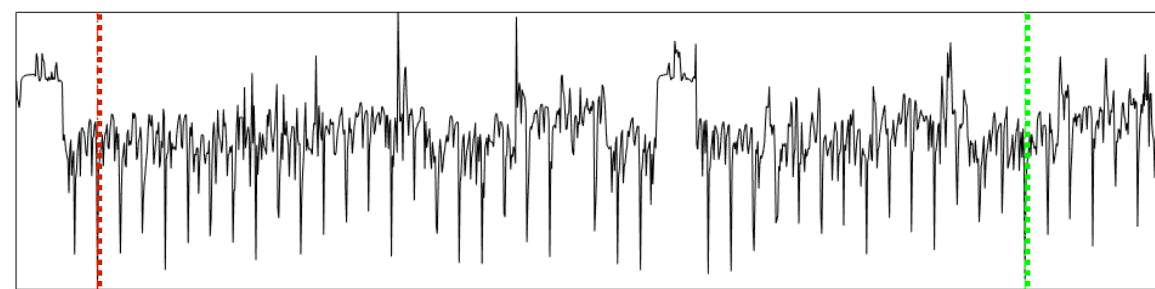
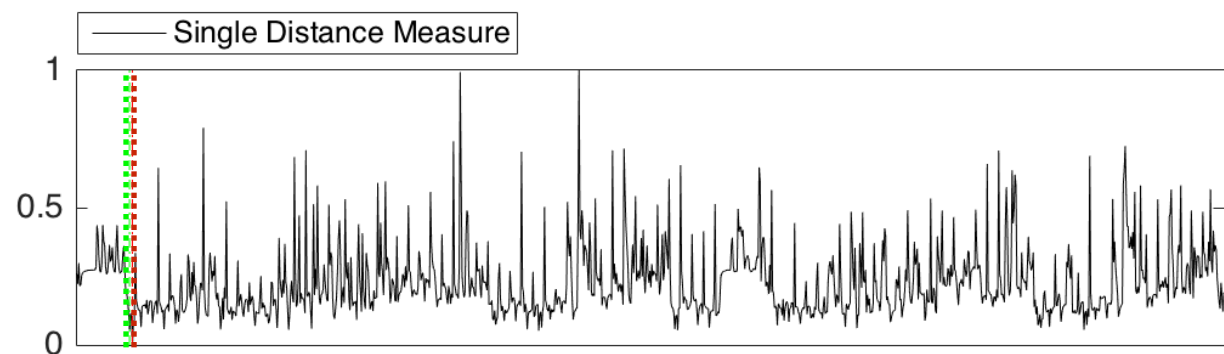


DISTANCE BETWEEN GROUP OF FRAMES

2. LOOK FOR SIMILAR TIME FRAMES

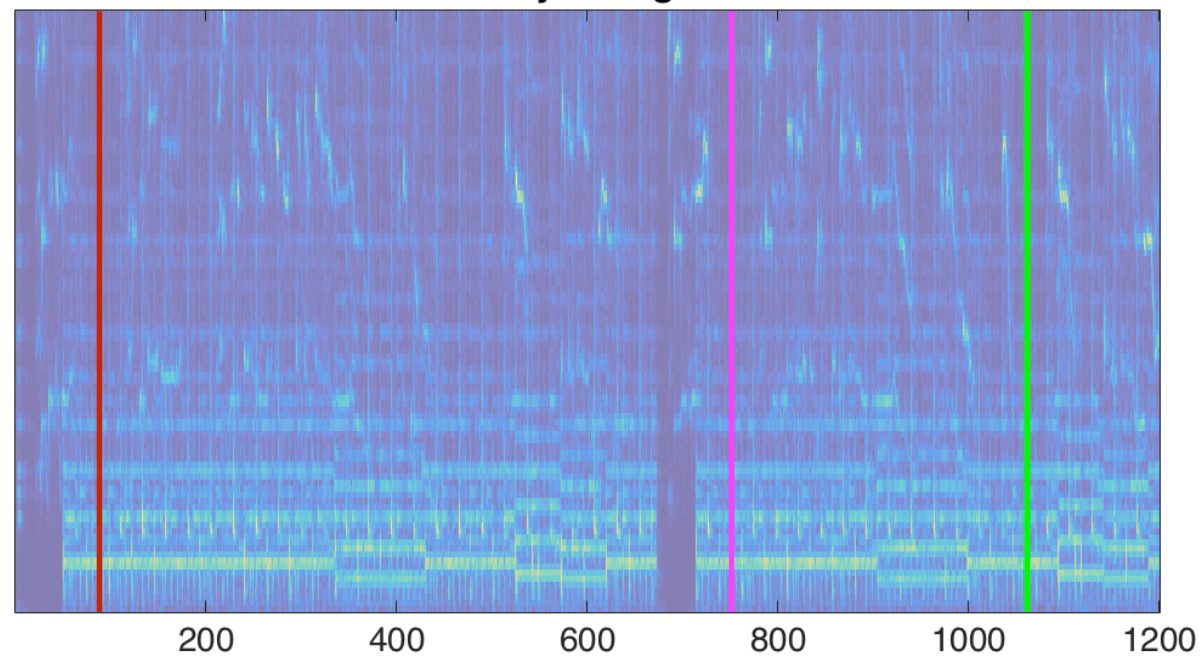
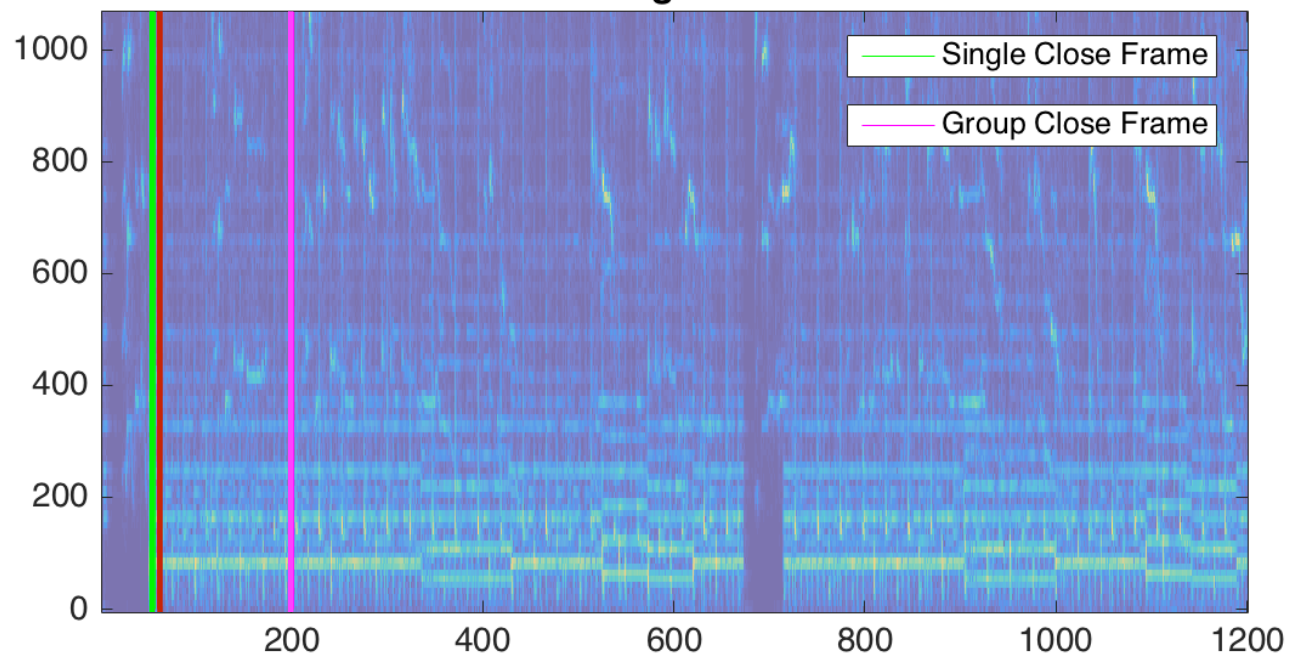


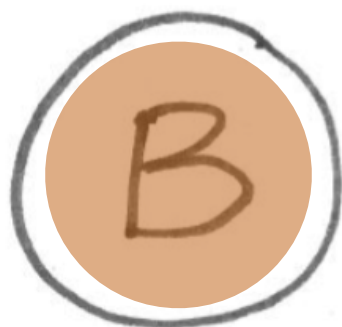
DISTANCE BETWEEN GROUP OF FRAMES

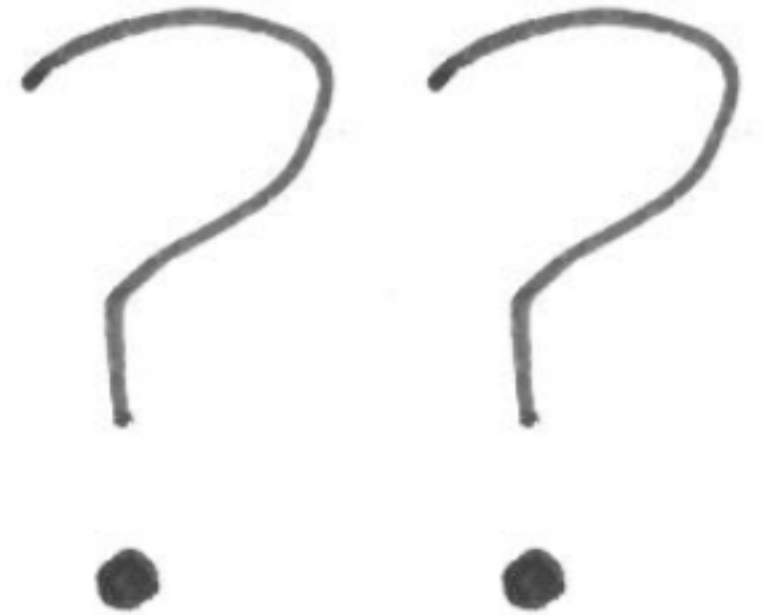


Frame 60 - Background and vocals

Frame 87 - Only background music







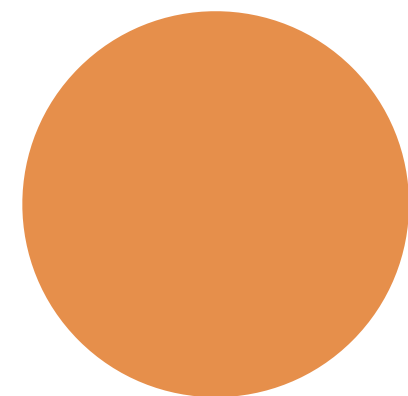
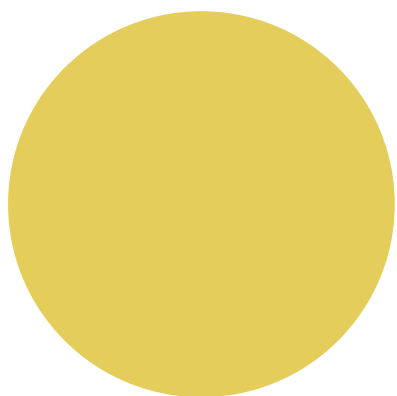




SINGLE FRAME
SIMILARITY



GROUPS OF FRAMES
SIMILARITY



EXPERIMENTS

EXPERIMENTS

Demixing Secrets Dataset 100 (DSD100)

2016 Signal Separation Evaluation Campaign (SiSEC)

EXPERIMENTS

Demixing Secrets Dataset 100 (DSD100)

2016 Signal Separation Evaluation Campaign (SiSEC)

100 SONGS

POLYPHONIC

$F_s = 44.1 \text{ KHz}$

30s LONG

EXPERIMENTS

Demixing Secrets Dataset 100 (DSD100)

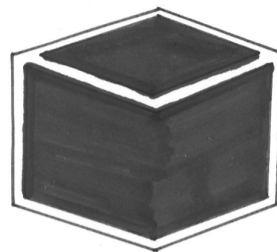
2016 Signal Separation Evaluation Campaign (SiSEC)

100 SONGS

POLYPHONIC

$F_s = 44.1 \text{ KHz}$

30S LONG



BSS Eval toolbox 3.0

Signal to Distortion Ratio (SDR)

Source Image to Spatial Distortion Ratio (ISR)

Source to Interference Ratio (SIR)

Source to Artifacts Ratio (SAR)

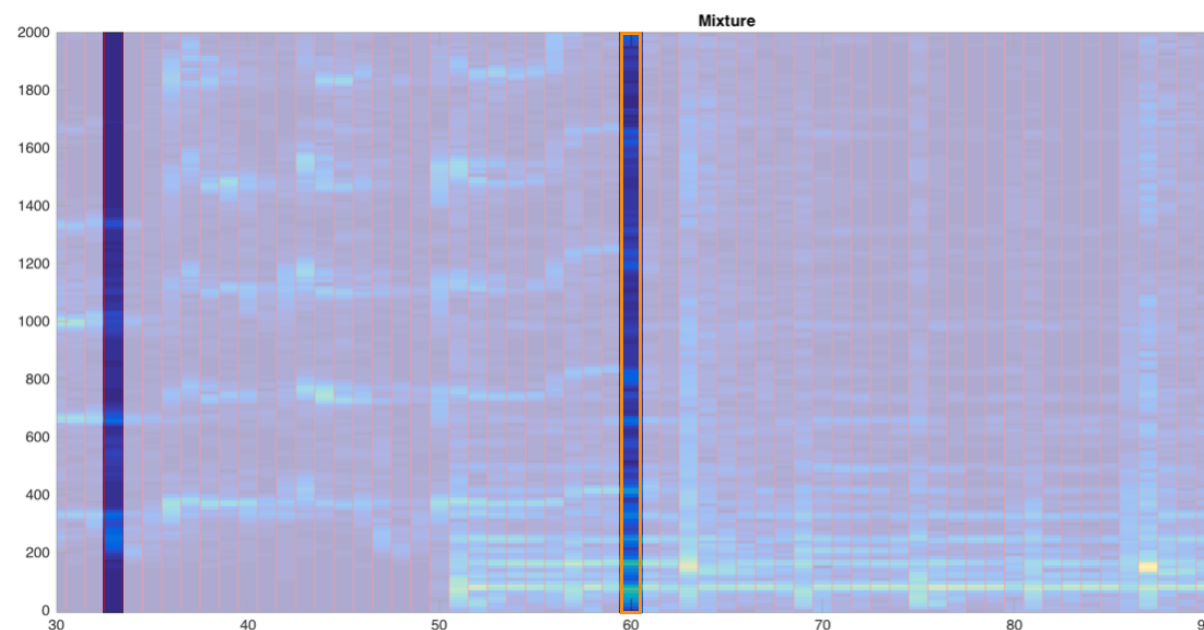
FFT size = 4096

Hopsize = 2048

BASELINE :

Instance of K.A.M

for vocal separation

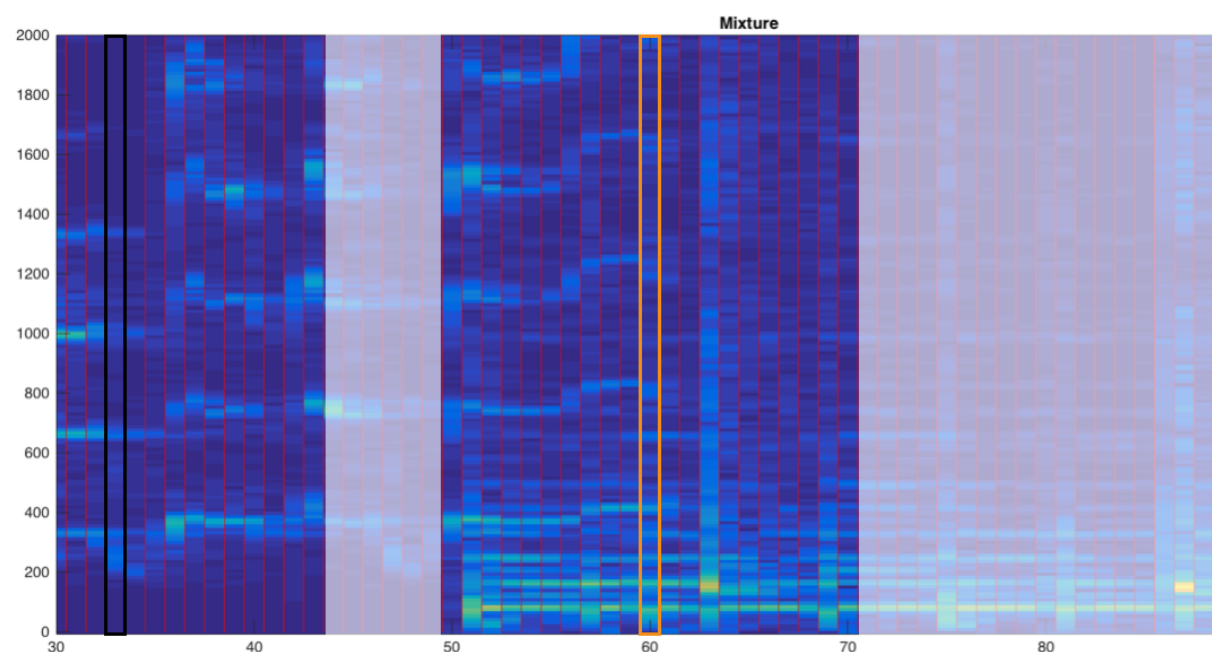


DISTANCE BETWEEN FRAMES

PROPOSED EXTENSION :

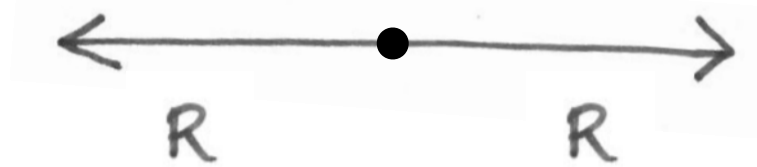
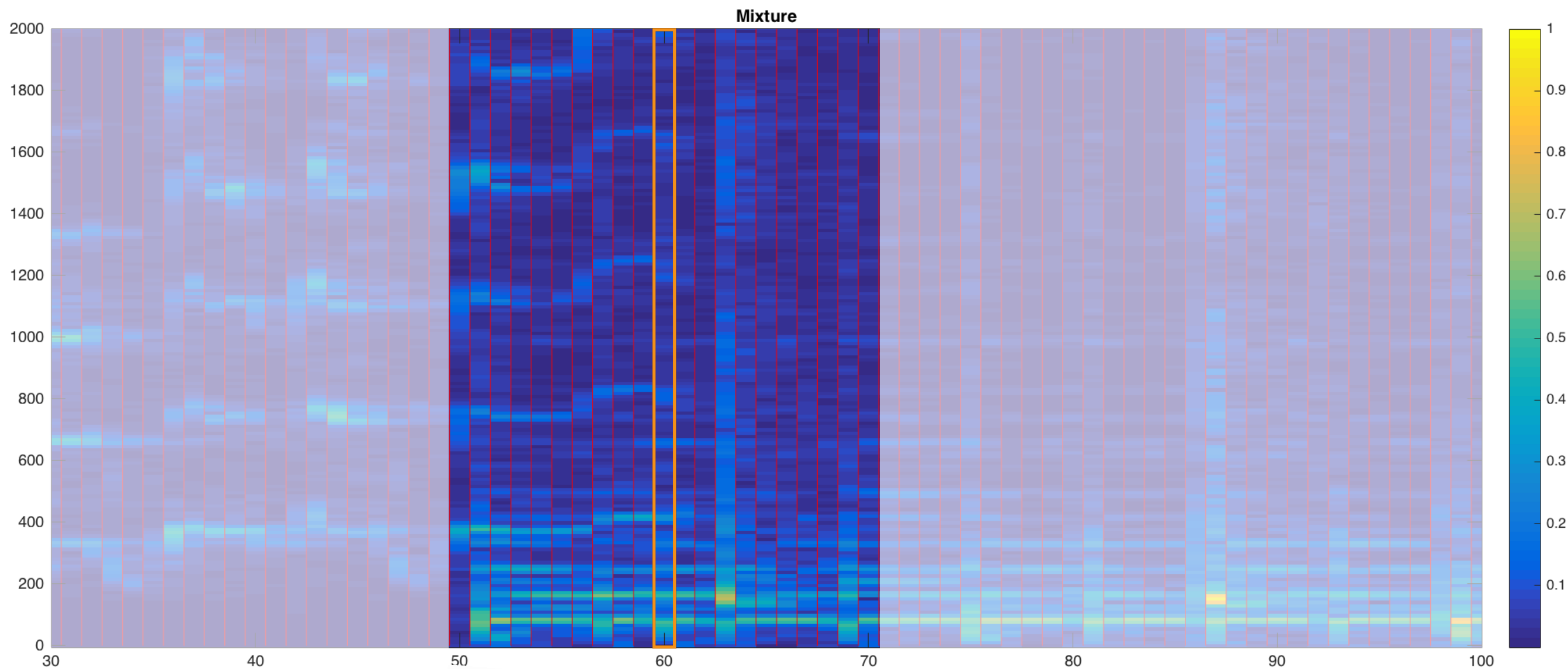
Baseline +
temporal context

in proximity kernel



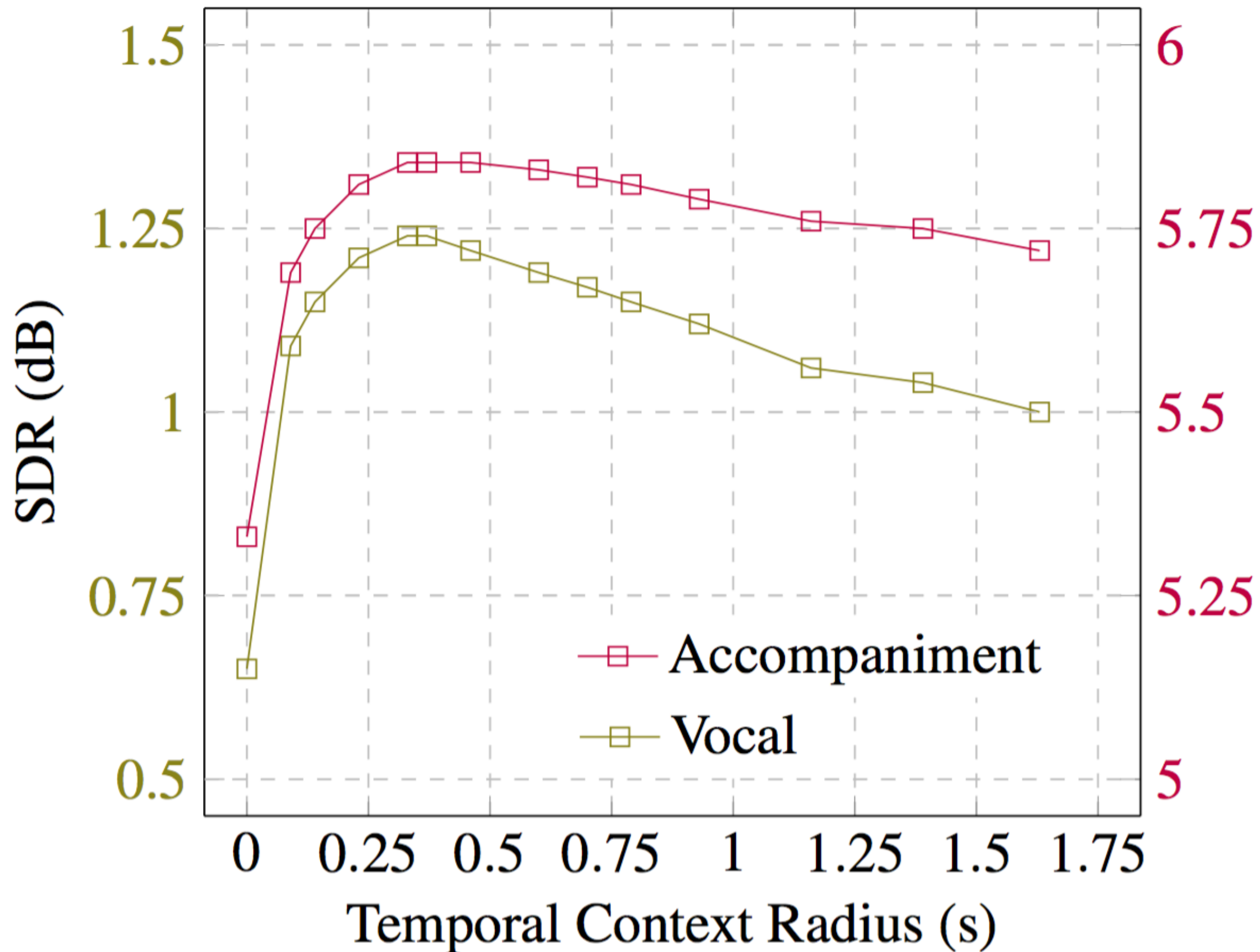
DISTANCE BETWEEN GROUP OF FRAMES

PROPOSED EXTENSION

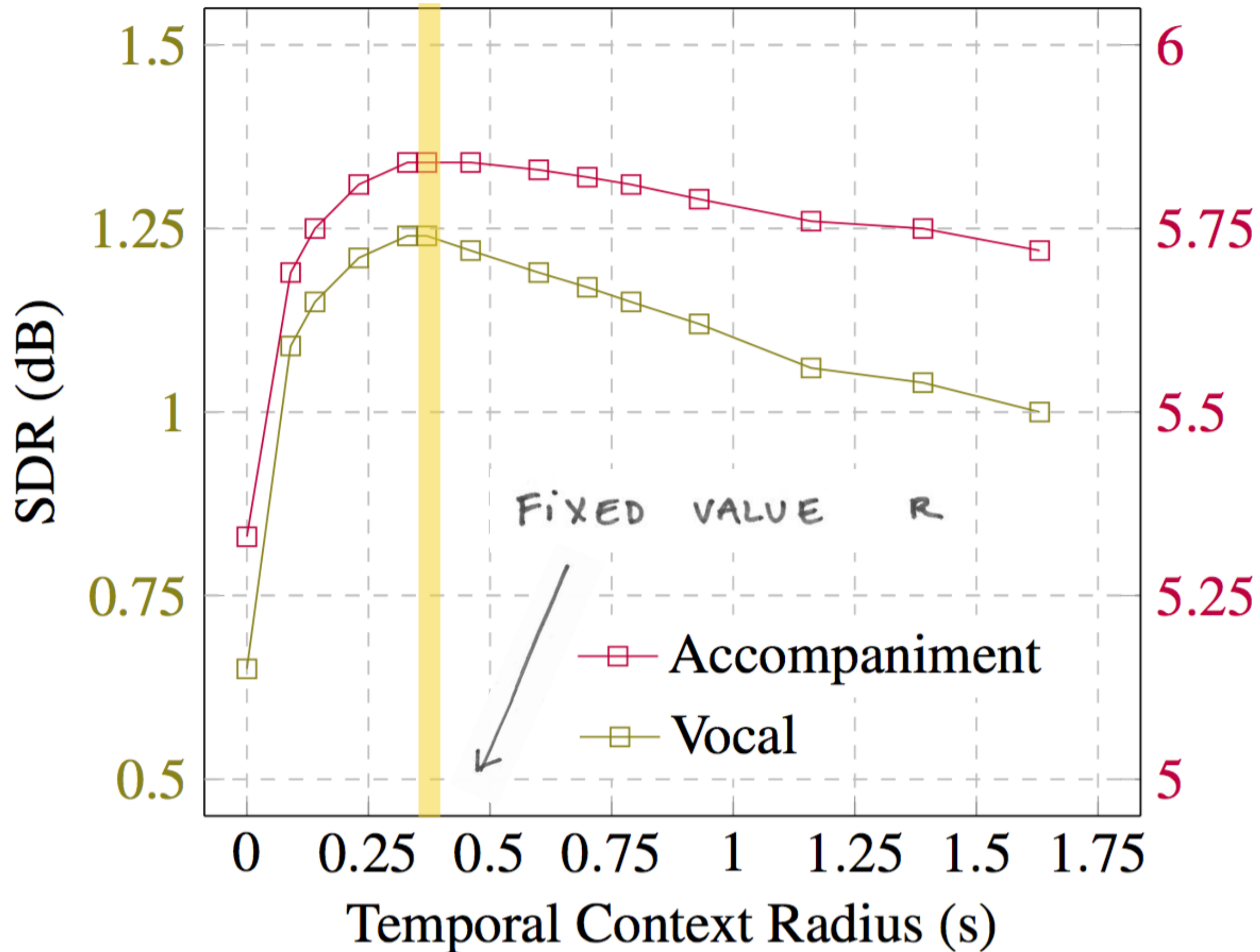


TEMPORAL CONTEXT

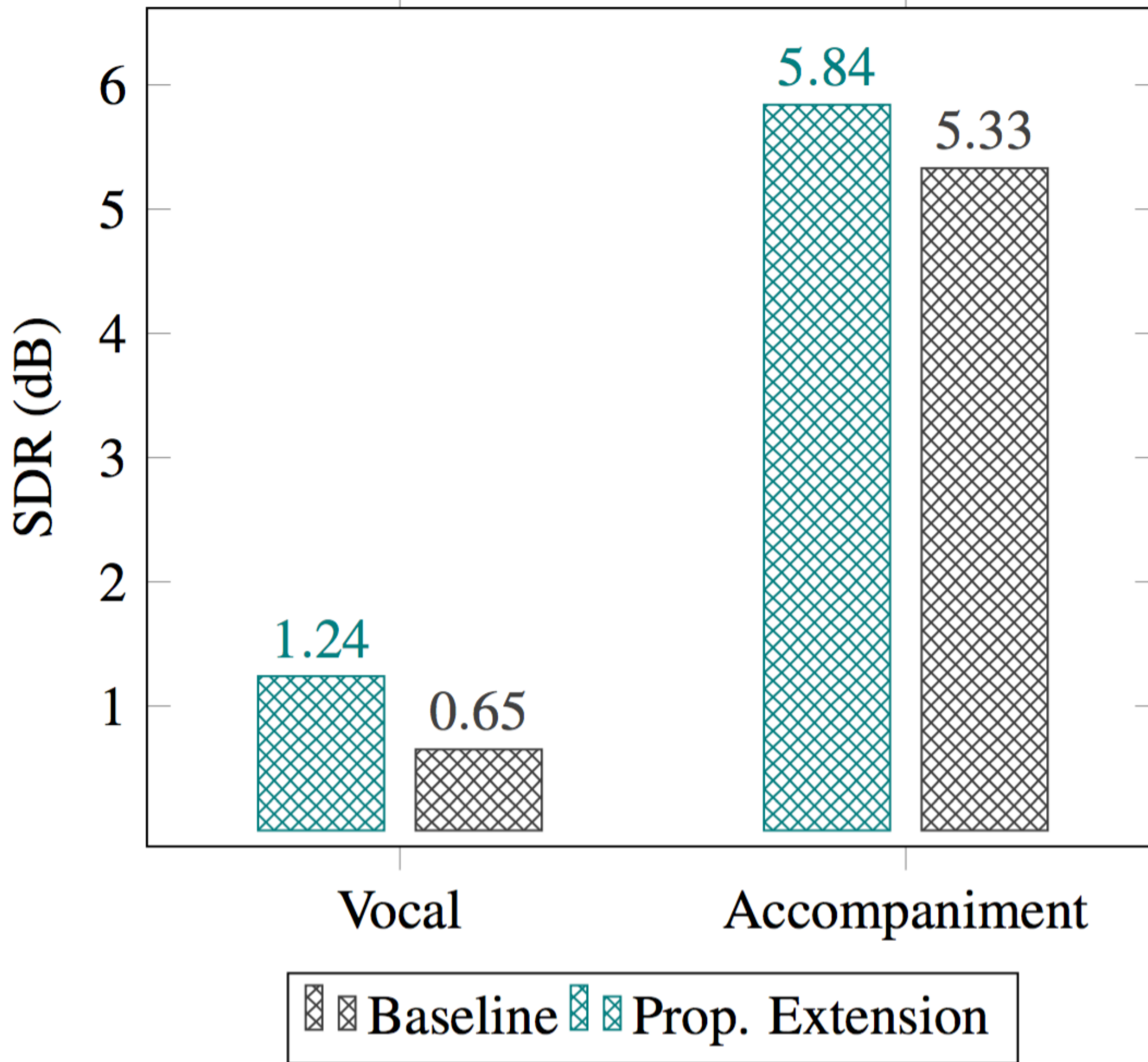
Temporal Context impact on the SDR



Temporal Context impact on the SDR



SDR for DSD100 Test dataset



WE PROPOSE A SIMPLE EXTENSION:

WE PROPOSE A SIMPLE EXTENSION:

INTRODUCE A TEMPORAL CONTEXT
IN PROXIMITY KERNELS

WE PROPOSE A SIMPLE EXTENSION:

INTRODUCE A TEMPORAL CONTEXT
IN PROXIMITY KERNELS

⊕ STABILISES THE SOURCE ESTIMATES

WE PROPOSE A SIMPLE EXTENSION:

INTRODUCE A TEMPORAL CONTEXT
IN PROXIMITY KERNELS

⊕ STABILISES THE SOURCE ESTIMATES

⊕ IMPROVES (A BIT) THE SEPARATION
PERFORMANCE

ANY
QUESTIONS ?

Delia Fano Yela

Sebastian Ewert

Derry Fitzgerald

Mark Sandler

d.fanoyela@qmul.ac.uk