

Result

Notes

- Sentiment
- Quantitative characteristics

EMOPIA

- Valence
- Arousal

Shape

- Size
- Speed
- Color
- Shape
- How often a shape is generated

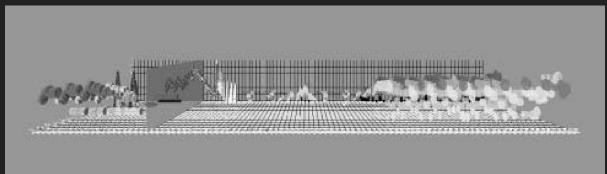
Original concept: Visual overlay on performing musician based on sentiment and quantitative analysis



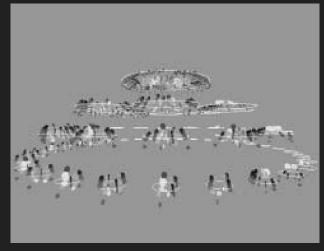
Mappings

Input: music	Output: shapes		
Amplitude	Size of shape		
Note duration/inter-ons et interval	Length of shape (more or less continuous)		
Sentiment: Positive	Warmer, brighter colors		
Sentiment: Negative	Cooler, darker colors		
Time since note onset	Distance from performer		

Current* work in performance visualization



Comprehensive MIDI Player-Interactive virtual space; visualizes multiple channel layers of Valse des Fleurs by Tchaikovsky



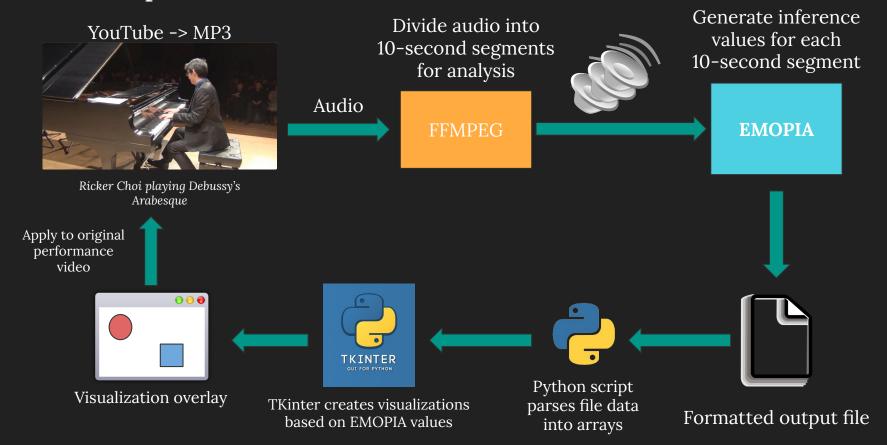
Visualizes musical structure of recursive hierarchy

Rumi Hiraga, Reiko Mizaki, and Issei Fujishiro. 2002. Performance visualization: a new challenge to music through visualization. In Proceedings of the tenth ACM international conference on Multimedia (MULTIMEDIA '02). Association for Computing Machinery, New York, NY, USA, 239–242. https://doi.org/10.1145/641007.641054

Motivation

- Classical music is often hard to approach as a listener
 - Particularly for kids/young children
- Provide a simple visualizer for recordings of classical performances based on the notes currently being played
 - Specifically, use a combination of sentiment and quantitative analysis on the notes being played to influence what the generated visuals look like

Technique



Process: EMOPIA output

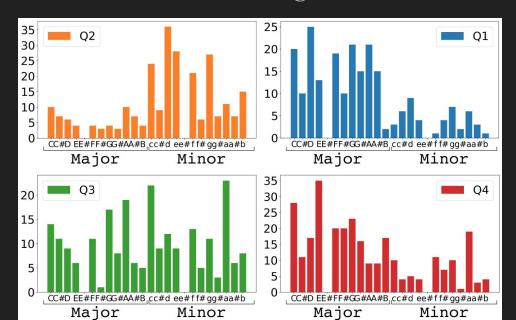
Inferences for every 10 seconds of Ricker Choi playing Debussy's Arabesque No. 1

	Low Valence	High Valence
High Arousal	Q2	Q1
Low Arousal	Q3	Q4

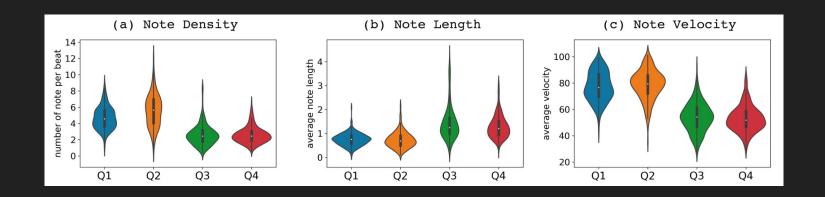
```
. .
                                              outputs
          is emotion Q4, Inference values:
                                             [0.05236754 0.03710493 0.05777654 0.9899472 ]]
           is emotion Q4, Inference values:
                                              [0.02216598 0.0195886 0.02014354 0.99775684]]
abs 20.wav
           is emotion Q4, Inference values:
                                              [0.01405534 0.03163977 0.02138974 0.99696064]]
abs 30.wav
           is emotion Q1, Inference values:
                                              [0.6617407 0.01601456 0.01004033 0.50370824]]
abs 40.wav
           is emotion Q3, Inference values:
                                              [0.25627008 0.39471266 0.65608525 0.3135395 ]]
           is emotion Q4, Inference values:
abs 50.wav
                                              [0.18802476 0.2107108 0.0188778 0.96072435]]
abs 60.wav
           is emotion Q4, Inference values:
                                              [0.03207375 0.07148175 0.0592729 0.9922824 ]]
           is emotion Q4, Inference values:
                                              [0.18821049 0.04989502 0.28378034 0.9648091 ]]
abs 80.wav is emotion Q4, Inference values:
                                              [0.46439102 0.32005224 0.02169288 0.61169475]]
           is emotion Q4, Inference values:
                                              [0.0207687 0.0260953 0.51218325 0.60367954]]
abs 100.wav is emotion Q4, Inference values:
                                               [0.02896808 0.0095521 0.02661344 0.9969424 ]]
abs_110.wav is emotion Q4, Inference values:
                                               [0.43077004 0.09250324 0.05918495 0.85090196]]
abs_120.wav is emotion Q1, Inference values:
                                               [0.52503
                                                           0.01079857 0.01055894 0.50426686]]
abs_130.wav is emotion Q4, Inference values:
                                               [0.4978638
                                                          0.03958806 0.16938588 0.5492348 ]]
abs 140.wav is emotion Q4, Inference values:
                                                          0.01146556 0.49717066 0.7328177 ]]
abs_150.wav is emotion Q4, Inference values:
                                               [0.08625619 0.08709763 0.19441879 0.9777598 ]]
abs 160.wav
            is emotion Q2, Inference values:
                                               [0.6067872  0.6470002  0.01204914  0.04304875]]
abs 170.wav is emotion Q3, Inference values:
                                               [0.02242045 0.02241473 0.9891635 0.11356956]]
abs 180.wav is emotion 03. Inference values:
                                               [0.01853796 0.02256764 0.63111085 0.5197379 ]]
abs_190.wav is emotion Q4, Inference values:
                                               [0.02922424 0.01407809 0.01523031 0.99783885]]
abs_200.wav is emotion Q4, Inference values:
                                               [0.01203746 0.05413383 0.04454004 0.9952735 ]]
abs_210.wav is emotion Q1, Inference values:
                                               [0.99676216 0.02690898 0.00879209 0.04750074]]
abs_220.wav is emotion Q1, Inference values:
                                               [0.7161766 0.5723025 0.0310793
abs 230.wav is emotion Q4, Inference values:
                                               [0.15028392 0.01569371 0.02630714 0.99106795]]
                                               [0.01715324 0.0232556 0.02378375 0.997772
abs 240.wav
            is emotion Q4, Inference values:
abs_250.wav is emotion Q4, Inference values:
                                               [0.15640613 0.03024157 0.16956954 0.949735
```

Valence: is the audio positive or negative?

- High Valence -> Major tonalities -> Positive sentiment -> Warm colors
- Low Valence -> Minor tonalities -> Negative sentiment -> Cool colors



Arousal: "...approximately maps to frequency of note occurrences and their strength...we measure them by <u>note density</u>, <u>length</u>, and <u>velocity</u>"



Arousal: "...approximately maps to frequency of note occurrences and their strength...we measure them by <u>note density</u>, <u>length</u>, and <u>velocity</u>"

- High Arousal: high note density, short note lengths, high note velocities
 - High note density -> Shapes come out more often
 - Short note lengths -> Shapes move away very fast from the performer
 - High note velocities -> Shapes are larger
 - o Generally higher arousal -> Shape is pointier
- Low Arousal: low note density, longer note lengths, smaller note velocities
 - Low note density -> Shapes come out less often
 - Longer note lengths -> Shapes are slower to move away from the performer
 - Smaller note velocities -> Shapes are smaller
 - Generally lower arousal -> Shape is rounder

6 6 6 6 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7					
Quadrant	`	Q2: Low Valence, High Arousal	· '	Q4: High Valence, Low Arousal	

Larger

Faster

Cooler

More often

rectangle)

Pointier (line or

Larger

Faster

Warmer

More often

rectangle)

Pointier (line or

Smaller

Slower

Cooler

Less often

Rounder (circle)

Smaller

Slower

Warmer

Less often

Rounder (circle)

Shape

Size

Speed

Color

How often

shapes are generated

What kind of

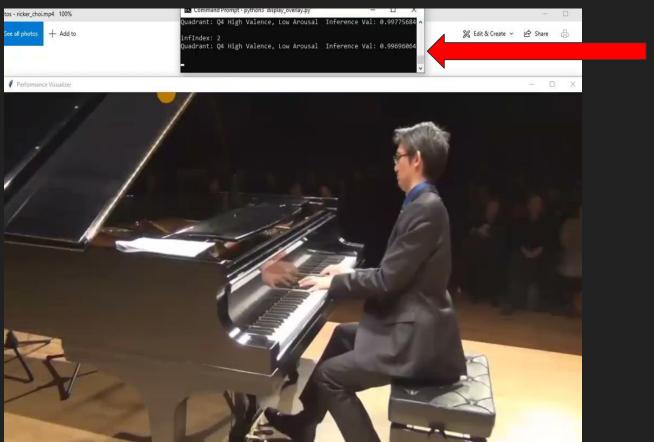
shape?

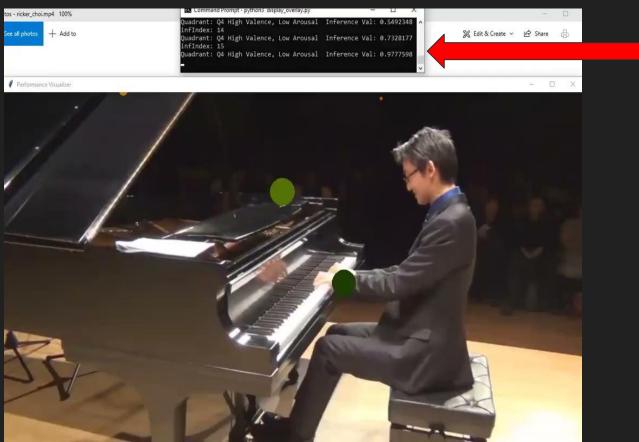
Characteristics

```
# 01: High Valence, High Arousal
# Shapes should be BIGGER, FASTER, WARM-COLORED and should come out MORE OFTEN; shapes should be POINTIER
if Q array[data.infIndex] == 'Q1':
                                                                            Shape characteristics are NOT
   base radius = base radius + 20*quad
                                                                            binary! The inference value acts as a
   base speed = base speed + 20*quad
   base color = color gen("warm", guad)
                                                                            scaling factor to determine how fast,
   data.circleTime = round((5 + round(10*quad)) / 2)
   base shape = random.choice([MyLine, MyRectangle])
                                                                            how big, how warm/cool, etc. of the
# Q2: Low Valence, High Arousal
# Shapes should be BIGGER, FASTER, COOL-COLORED and should come out MORE OFTEN;
                                                                           shape uld be POINTIER
elif Q array[data.infIndex] == 'Q2':
   base_radius = base_radius + 20*quad
   base_speed = base_speed + 20*quad
   base color = color gen("cool", guad)
   data.circleTime = round((5 + round(10*quad)) / 2)
   base_shape = random.choice([MyLine, MyRectangle])
# Q3: Low Valence, Low Arousal
# Shapes should be SMALLER, SLOWER, COOL-COLORED and should come out LESS OFTEN; shapes should be ROUNDER
elif Q_array[data.infIndex] == 'Q3':
   base radius = base radius - 20*quad
   base speed = base speed - 20*quad
   base_color = color_gen("cool", quad)
   data.circleTime = round((15 + round(10*quad)) / 2)
   base shape = Circle
# Q4: High Valence, Low Arousal
# Shapes should be SMALLER, SLOWER, WARM-COLORED and should come out LESS OFTEN; shapes should be ROUNDER
elif Q array[data.infIndex] == 'Q4':
   base_radius = base_radius - 20*quad
   base speed = base speed - 20*quad
   base color = color gen("warm", guad)
   data.circleTime = round((15 + round(10*quad)) / 2)
   base shape = Circle
```

Color

```
# Helper function: given a mode and a scaling number, generates a hex color for the shape
def color gen(mode, scale):
    if mode == "warm":
        r lim = round(255*scale)
       b \lim = 255 - round(255*scale)
                                                                  Scale positively affects red value and
       g \lim = round((r \lim + b \lim) / 2)
       red = lambda: random.randint(0, r lim)
                                                                  negatively affects blue value
       green = lambda: random.randint(0,g lim)
       blue = lambda: random.randint(0,b lim)
        return '#%02X%02X%02X' % (red(),green(),blue())
    elif mode == "cool":
        r \lim = 255 - round(255*scale)
       b lim = round(255*scale)
                                                                   Scale negatively affects red value and
       q \lim = round((r \lim + b \lim) / 2)
                                                                   positively affects blue value
       red = lambda: random.randint(0,r lim)
       green = lambda: random.randint(0,g lim)
       blue = lambda: random.randint(0,b lim)
        return '#%02X%02X%02X' % (red(), green(), blue())
```



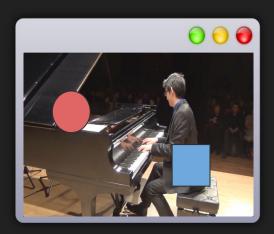


- Method to make the window transparent only exists for Windows, not macOS
 - Had to use CMU's virtual machines
- Very old and bad documentation
- Very limited functionality for common-sense things
 - You can play a video in Tkinter with external libraries
 - But you can't use this video as a background or a frame
 - And you can't draw shapes directly on top of this video



Normal video playing

- Method to make the window transparent only exists for Windows, not macOS
 - Had to use CMU's virtual machines
- Very old and bad documentation
- Very limited functionality for common-sense things
 - You can play a video in Tkinter with external libraries
 - But you can't use this video as a background or a frame
 - And you can't draw shapes directly on top of this video



Place transparent TKinter window on top of original video (this only works on Windows OS)

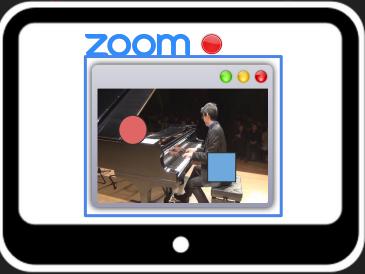
- Method to make the window transparent only exists for Windows, not macOS
 - Had to use CMU's virtual machines
- Very old and bad documentation
- Very limited functionality for common-sense things
 - You can play a video in Tkinter with external libraries
 - But you can't use this video as a background or a frame
 - And you can't draw shapes directly on top of this video



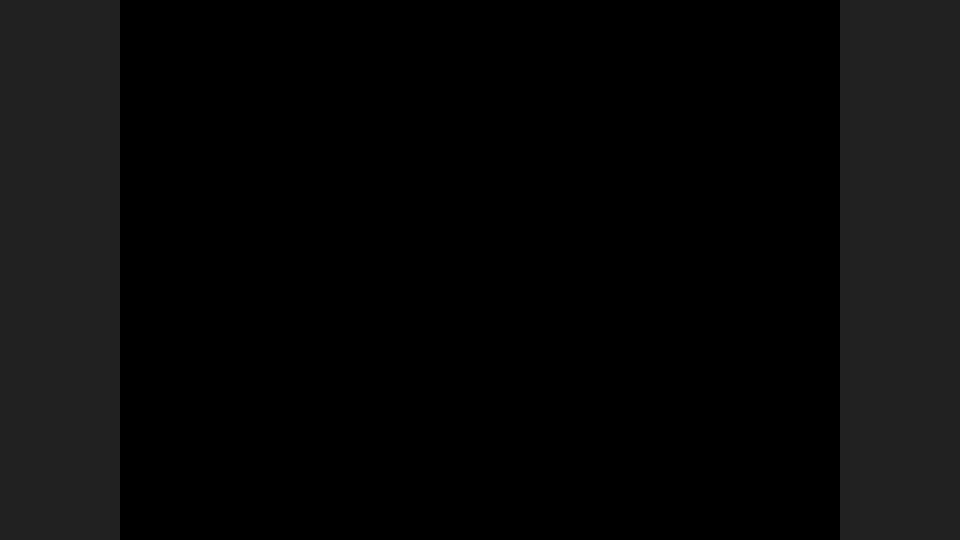
Record this with Zoom

- Method to make the window transparent only exists for Windows, not macOS
 - Had to use CMU's virtual machines
- Very old and bad documentation
- Very limited functionality for common-sense things
 - You can play a video in Tkinter with external libraries
 - But you can't use this video as a background or a frame
 - And you can't draw shapes directly on top of this video







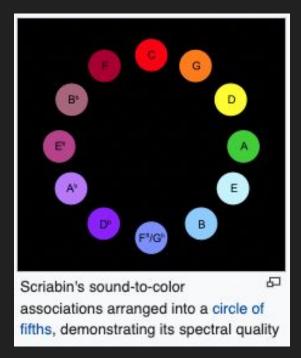


Reflection

- Non-cross-platform software is really bad!
- Old software is really bad!
- How objective are these visualizations? For example, there could be people who associate low valence/minor tonalities with *warm* colors rather than cooler colors
 - Should these visualizations be customizable, then?
 - o Do these visualizations enforce a specific type of viewing specific types of music?
 - Like mild propaganda
- Does it make sense to use randomized values to influence the shapes generated as well? Or should all values only be determined by analyzed data?
- How does this compare to chromesthesia?*

Reflection

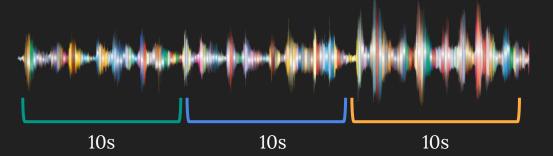
How does this compare to chromesthesia?



- All-in-one pipeline to do this with any uploaded video
- More direct analysis functions rather than all-encompassing valence/arousal values
- Should inference values from non-dominant quadrants also influence the visuals?
- Sliding windows*
- Use better/more modern interfaces than Tkinter to directly show visualizations on top of input performance video
 - OpenCV or even PyGame are alternatives
- Improved visualizations
 - For example: growing/shrinking shapes rather than static ones
 - Customizable visualizations? Users could input what mappings they prefer
- Real time? Bit tricky, but technically possible!*
 - Requires lots of computer science techniques that I am bad at
 - What implications does this have?
 - Use in performance AR/VR
- Do this with other instruments besides just piano pieces (EMOPIA is tailored for piano pop music specifically)
- Use Detectron to actually detect where the performer is rather than just making shapes come out from near the center
- Brightness of colors?
- More advanced graphics manipulation than just object generation (for example, modifying the color scheme of the entire performance frame)

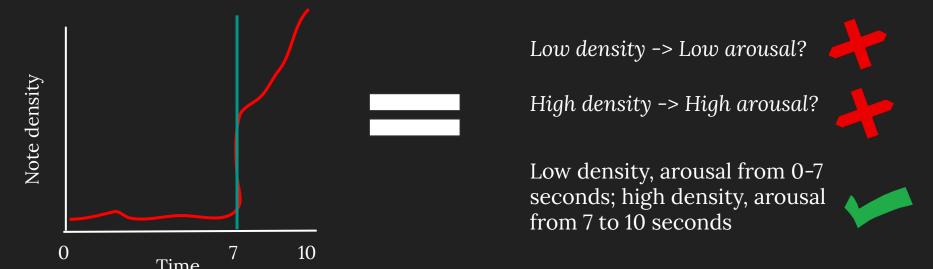
Sliding windows

 Currently, within a 10-second segment of the performance audio, drastic changes within a small part in the segment can wildly skew inference results



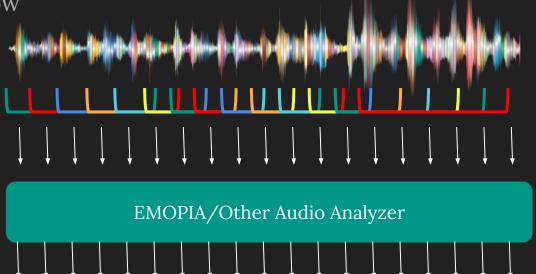
Sliding windows

 Currently, within a 10-second segment of the performance audio, drastic changes within a small part in the segment can wildly skew inference results



Sliding windows

 A better solution would be to have overlapping windows that sample the original audio to be analyzed, where inference values are generated for each window



Sliding windows

- Already a very common technique in any audio/signal processing context
- Obviously more expensive computationally and storage-wise
- But would get you much more accurate data & visuals

Real time? Bit tricky, but technically possible!

- Requires lots of computer science techniques (dynamic programming) that I am bad at
- What implications does this have?
 - Use in performance AR/VR
 - For both the performer and listeners
 - 3D visual objects instead of 2D



Kresge VR Project @ CMU

Real time? Bit tricky, but technically possible!

- Requires lots of computer science techniques (dynamic programming) that I am bad at
- What implications does this have?
 - Use in performance AR/VR
 - For both the performer and listeners
 - 3D visual objects instead of 2D

