

# Web Applications for Automatic Audio-to-Score Synchronization with Iterative Refinement: Appendix A

**Adriano Baratè, Goffredo Haus  
Luca A. Ludovico, Giorgio Presti**  
Laboratory of Music Informatics  
Department of Computer Science  
University of Milan  
{first.last}@unimi.it

**Stefano Di Bisceglie  
Alessandro Minoli**  
Department of Computer Science  
University of Milan  
{first.last}@studenti.unimi.it

**Davide A. Mauro**  
Department of Computer  
and Information Technology  
Marshall University  
maurod@marshall.edu

## 1. SCORE-INFORMED BEAT TRACKING

Excerpts of the implementation for the score-informed beat tracking algorithm, including the computation of chroma vectors from the score, and the DTW computation.

### 1.0.1 *chromaFromScore()*

```
1  var chords = xmlDoc.getElementsByTagName("
    chord");
2  transpose = parseInt(document.getElementById(
    "transpose").value);
3  for(var i = 0; i < chords.length; i++) { //for
    every chord
4      var chord = chords[i];
5      var timing_start = timings[
        chord.getAttribute("event_ref")];
6      var dur = getElementDuration(chord);
7      var timing_end = timing_start + (
        quarter_vtu * dur * 4);
8      var pitches = chord.getElementsByTagName(
        "pitch");
9      //for every note in the chord
10     for(var j = 0; j < pitches.length; j++){
11         var pitch_obj = pitches[j];
12         var chroma = 0;
13         var chroma = getChromaFromPitch(
            pitch_obj, transpose);
14         if(chroma == -1)
15             continue;
16         var cv_start = insertIntoChromasVtus(
            chromas_vtus[chroma], timing_start,
            true);
17         var cv_end = insertIntoChromasVtus(
            chromas_vtus[chroma], timing_end, false
        );
18         var dist = deleteBetween(chromas_vtus[
            chroma], cv_start, cv_end);
19         cv_end = cv_end - dist;
20         if(chromas_vtus[chroma][cv_start-1][1] ==
            true){
21             deleteFrom(chromas_vtus[chroma],
                cv_start);
22             cv_end--;
23         }
24         if(chromas_vtus[chroma][cv_end+1][1] ==
            false){
25             deleteFrom(chromas_vtus[chroma], cv_end);
26         }
27     }
28 }
29 for(var chroma = 0; chroma < 12; chroma++){
30     var j = 1;
31     while(j < chromas_vtus[chroma].length-1)
32     {
33         var flag1 = chromas_vtus[chroma][j];
```

```
33     j++;
34     var flag2 = chromas_vtus[chroma][j];
35     insertIntoChvecsScore(flag1[0], flag2[0],
        chroma);
36     j++;
37     }
38 }
```

### 1.0.2 *calcTrackIndexing()*

```
1  //compute new track indexing and add it to
    xmlDoc
2  var track_indexing = xmlDoc.createElement("
    track_indexing");
3  track.appendChild(track_indexing);
4  track_indexing.setAttribute("timing_type", "
    seconds");
5  var spine = xmlDoc.getElementsByTagName("
    spine")[0];
6  var curTiming = 0;
7  for(var i = 0; i < spine.getElementsByTagName(
    "event").length; i++){
8      //for every event
9      var event = spine.getElementsByTagName("
    event")[i];
10     curTiming = curTiming + parseInt(
        event.getAttribute("timing"));
11     var seconds = vtuToSeconds(curTiming);
12     var track_event = xmlDoc.createElement("
        track_event");
13     track_event.setAttribute("event_ref",
        event.getAttribute("id"));
14     track_event.setAttribute("start_time",
        seconds);
15     track_indexing.appendChild(track_event);
16 }
```

### 1.0.3 *getDistance()*

```
1  matrix = [];
2  path_matrix = [];
3  for ( var i = 0; i < ser1.length; i++ )
4  {
5      matrix[i] = [];
6      path_matrix[i] = [];
7      for ( var j = 0; j < ser2.length; j++ )
8      {
9          var c = distFunc(ser1[i], ser2[j]);
10         if(i == 0 && j == 0){
11             matrix[i][j] = c;
12             path_matrix[i][j] = -1;
13         }else if(i > 0 && j == 0){
14             matrix[i][j] = matrix[i-1][j] +
                wh*c;
15             path_matrix[i][j] = 0;
16         }else if(j > 0 && i == 0){
17             matrix[i][j] = matrix[i][j-1] +
                wv*c;
18             path_matrix[i][j] = 1;
19         }else{
20             var h = matrix[i-1][j] + wh*c;
21             var v = matrix[i][j-1] + wv*c;
22             var d = matrix[i-1][j-1] + wd*c;
23             var min = Math.min(v, d, h);
```

Copyright: © 2023 Adriano Baratè, Goffredo Haus et al. This is an open-access article distributed under the terms of the [Creative Commons Attribution 3.0 Unported License](https://creativecommons.org/licenses/by/3.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

```

22         if(min == d){
23             matrix[i][j] = d;
24             path_matrix[i][j] = 2;
25         }else if(min == h){
26             matrix[i][j] = h;
27             path_matrix[i][j] = 0;
28         }else if(min == v){
29             matrix[i][j] = v;
30             path_matrix[i][j] = 1;
31         }
32     }
33 }
34 }
35 return matrix[ser1.length-1][ser2.length-1];

```

## 2. ONSET DETECTION FUNCTIONS

Concerning the reduction of a signal to an ODF, there are mainly two categories: energy-based approaches and phase-based ones. Please note that ODFs detect the physical onsets of a signal. The more the attacks of the notes are slow and/or soft, the more these offsets tend to deviate temporally from the ones perceived during the listening.

*Energy-based approaches* take their cue from the following observation: the occurrence of onsets in the signal is usually accompanied by sudden increments in energy. These changes tend to appear in the spectrum as broadband events. In addition, in the portions of the signal not affected by transients, energy is usually concentrated in the low frequencies, so these changes should be more easily detectable at high frequencies. The ODF called *High-Frequency Content* (HFC) [1] exploits this observation and, for each frame, weighs each frequency bin proportionally to its frequency. As a drawback, it does not consider the temporal evolution of the signal, thus the value of the function at each instant depends exclusively on the current frame.

Other energy-based ODFs consider the spectrum variation between adjacent frames. In particular, they measure the distances between the magnitude of two consecutive frames. For example, by choosing L1-norm as the distance metrics, we obtain the *Spectral Difference* (SD), while L2-norm defines the *Spectral Flux* (SF) [2]. Another approach called *Percussive feature detection* (or, simply, percussive) is based on an algorithm similar to SF but calculated on a logarithmic amplitude scale [3].

In general, energy-based ODFs can be implemented in a simple and efficient way. They give excellent results in the presence of percussive onsets, namely notes with a sharp attack, such as the ones produced by drums or piano. Energy-based ODFs performances get worse in the case of non-percussive signals and when the transients begin to overlap, as in complex signals.

*Phase-based approaches* rely on the idea that, during a transient, the phase of the frequency components changes significantly with respect to the increase expected considering previous time instants. A measure of the temporal instability of the signal phase can be constructed by quantifying the derivative of the phase over time for each frequency component and then adding these quantities to obtain a single time-varying signal [4]. Unlike energy-based ODFs, phase-based ones are also able to detect onsets with loosely defined attacks. However, they are very sensitive to noise, especially for those signal components carrying

very low energy.

In recent years, in order to try to solve the open problem of onset detection, machine learning techniques have been added to the more classic techniques of signal analysis and processing. For example, neural networks can be successfully applied [5].

## 3. FURTHER TESTS

Figures 1 to 4 show the deviations of the original timings and the adjusted ones with respect to the ground truth for 4 of the test tracks.

Then, the solution was tested on excerpts from the *Onset-Leveau* dataset. Compared to the scenario of IEEE 1599 documents, which also contain a *Logic* level with the symbolic description of the events, in this case the score is not available. For this reason, we prepared IEEE 1599 documents which in the *Audio* layer already contained ground-truth timings. In an ideal case, the algorithm should not change them, as they already correspond to the correct time position of the onsets. The best results will be those where *RMSE* and *MAE* are closest to 0.

The results achieved on the dataset tracks are shown in Table 1. The tracks in the upper area perceptively present sharp note attacks, and those in the lower area smooth attacks.

By weighing each track by its number of onsets, the means and standard deviations were calculated, grouping results by track family and by the characteristics of their onsets. These measures are respectively denoted by symbols  $\mu_{RMSE}$ ,  $\mu_{MAE}$ ,  $\sigma_{RMSE}$  and  $\sigma_{MAE}$ , and reported in Table 2. Figures 5 to 8 show their trend as the size of the neighborhood changes; the bars represent the standard deviation of the timings from their mean. Finally, Figures 9 and 10 show the deviations of adjusted timings with respect to the ground truth for some tracks in Table 1.

## 4. IEEE 1599 STRUCTURE

The standard IEEE layered structure is depicted in Fig. 11

## 5. REFERENCES

- [1] P. Masri, “Computer modelling of sound for transformation and synthesis of musical signals,” Ph.D. dissertation, University of Bristol, 1996.
- [2] C. Duxbury, M. Sandler, and M. Davies, “A hybrid approach to musical note onset detection,” in *Proc. Digital Audio Effects Conf.(DAFX’02)*. sn, 2002, pp. 33–38.
- [3] D. Barry, D. Fitzgerald, E. Coyle, and B. Lawlor, “Drum source separation using percussive feature detection and spectral modulation,” *IEE Irish Signals and Systems Conference (ISSC)*, pp. 13–17, 2005.
- [4] P. Brossier, J. P. Bello, and M. D. Plumbley, “Real-time temporal segmentation of note objects in music signals,” in *Proceedings of ICMC 2004, the 30th Annual International Computer Music Conference*. University of Surrey, 2004.

WAV file	duration	# onsets	type	narrow neighborhood		large neighborhood	
				<i>RMSE</i>	<i>MAE</i>	<i>RMSE</i>	<i>MAE</i>
guitar3	15 s	56	P	8 ms	7 ms	8 ms	7 ms
jazz2	14 s	47	CM	20 ms	15 ms	21 ms	16 ms
piano1	15 s	19	P	12 ms	10 ms	12 ms	10 ms
pop1	15 s	27	CM	13 ms	11 ms	13 ms	11 ms
rock1	15 s	59	CM	12 ms	10 ms	24 ms	15 ms
synthbass	7 s	23	M	6 ms	5 ms	21 ms	9 ms
techno2	6 s	56	CM	17 ms	15 ms	17 ms	15 ms
cello1	14 s	61	P	23 ms	17 ms	43 ms	32 ms
clarinet1	30 s	35	M	22 ms	17 ms	86 ms	80 ms
classic2	20 s	38	CM	23 ms	19 ms	40 ms	29 ms
classic3	14 s	4	CM	28 ms	28 ms	28 ms	28 ms
distguit1	6 s	19	P	10 ms	7 ms	10 ms	7 ms
guitar2	15 s	36	P	18 ms	12 ms	34 ms	22 ms
jazz3	11 s	53	CM	21 ms	16 ms	35 ms	24 ms
sax1	12 s	9	M	28 ms	23 ms	64 ms	55 ms
trumpet1	14 s	56	M	21 ms	19 ms	21 ms	19 ms
violin2	15 s	73	P	20 ms	14 ms	41 ms	29 ms

Table 1: Results achieved on *Onset\_Leveau*.

category	narrow neighborhood				large neighborhood			
	$\mu_{RMSE}$	$\sigma_{RMSE}$	$\mu_{MAE}$	$\sigma_{MAE}$	$\mu_{RMSE}$	$\sigma_{RMSE}$	$\mu_{MAE}$	$\sigma_{MAE}$
M	19 ms	3.91 ms	16 ms	3.52 ms	43 ms	9.54 ms	37 ms	9.44 ms
P	17 ms	2.15 ms	12 ms	1.52 ms	29 ms	4.71 ms	21 ms	3.36 ms
CM	18 ms	1.27 ms	14 ms	0.99 ms	25 ms	2.22 ms	19 ms	1.46 ms
sharp onsets	13 ms	1.08 ms	11 ms	0.90 ms	17 ms	1.43 ms	12 ms	1.02 ms
smooth onsets	21 ms	1.30 ms	17 ms	0.99 ms	41 ms	2.84 ms	32 ms	2.32 ms

Table 2: Aggregated results achieved on *Onset\_Leveau*.

- [5] J. Schlüter and S. Böck, “Improved musical onset detection with convolutional neural networks,” in *2014 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2014, pp. 6979–6983.

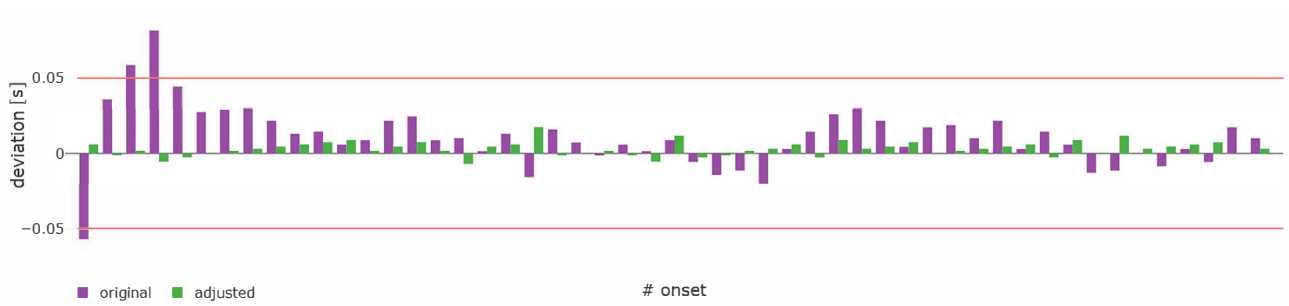


Figure 1: Results for piano1.mp3.



Figure 2: Results for violin1.mp3.

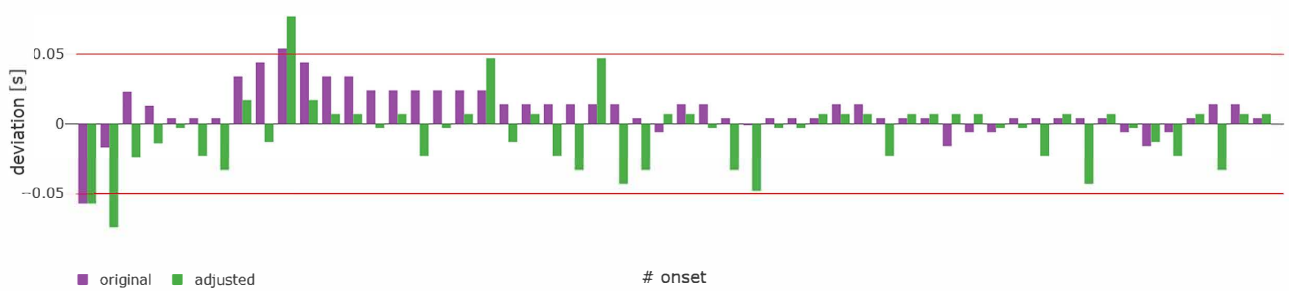


Figure 3: Results for string\_4et.mp3.



Figure 4: Results for jazz\_4et.mp3.

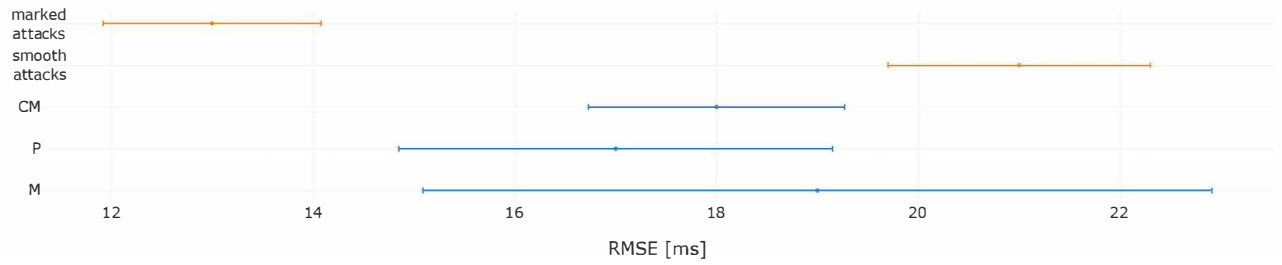


Figure 5: *RMSE* of the adjustments with a narrow neighborhood.

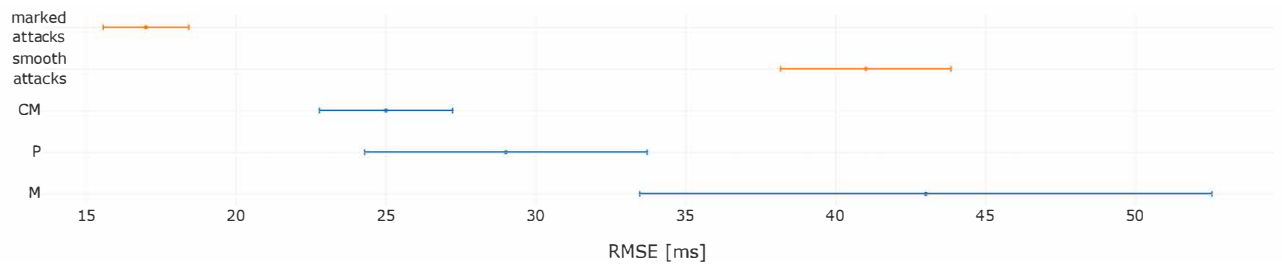


Figure 6: *RMSE* of the adjustments with a large neighborhood.

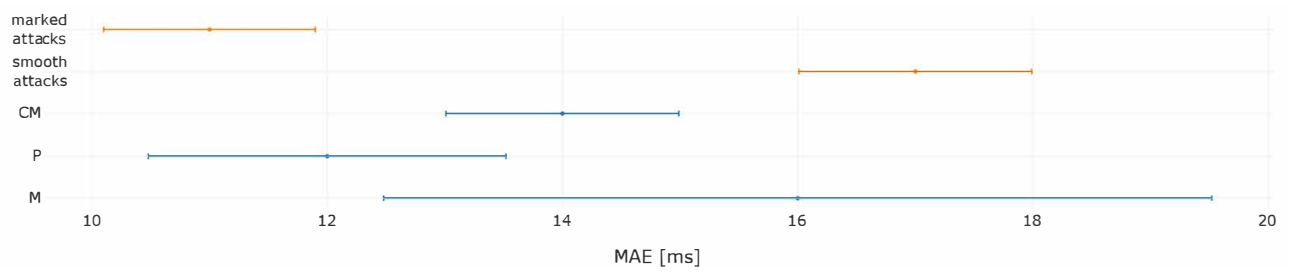


Figure 7: *MAE* of the adjustments with a narrow neighborhood.

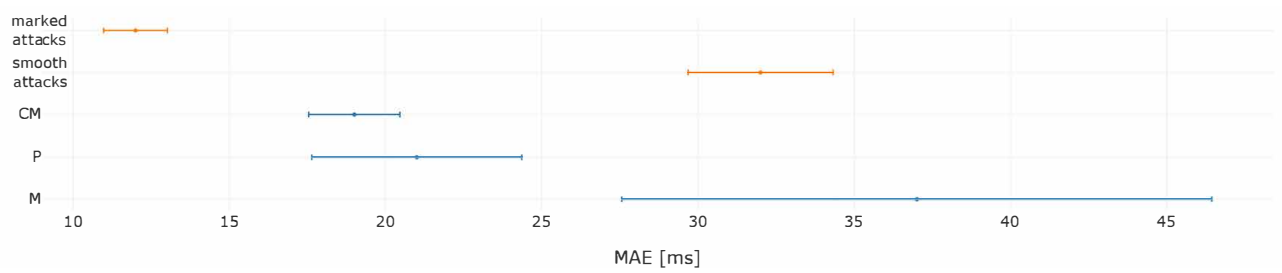


Figure 8: *MAE* of the adjustments with a large neighborhood.

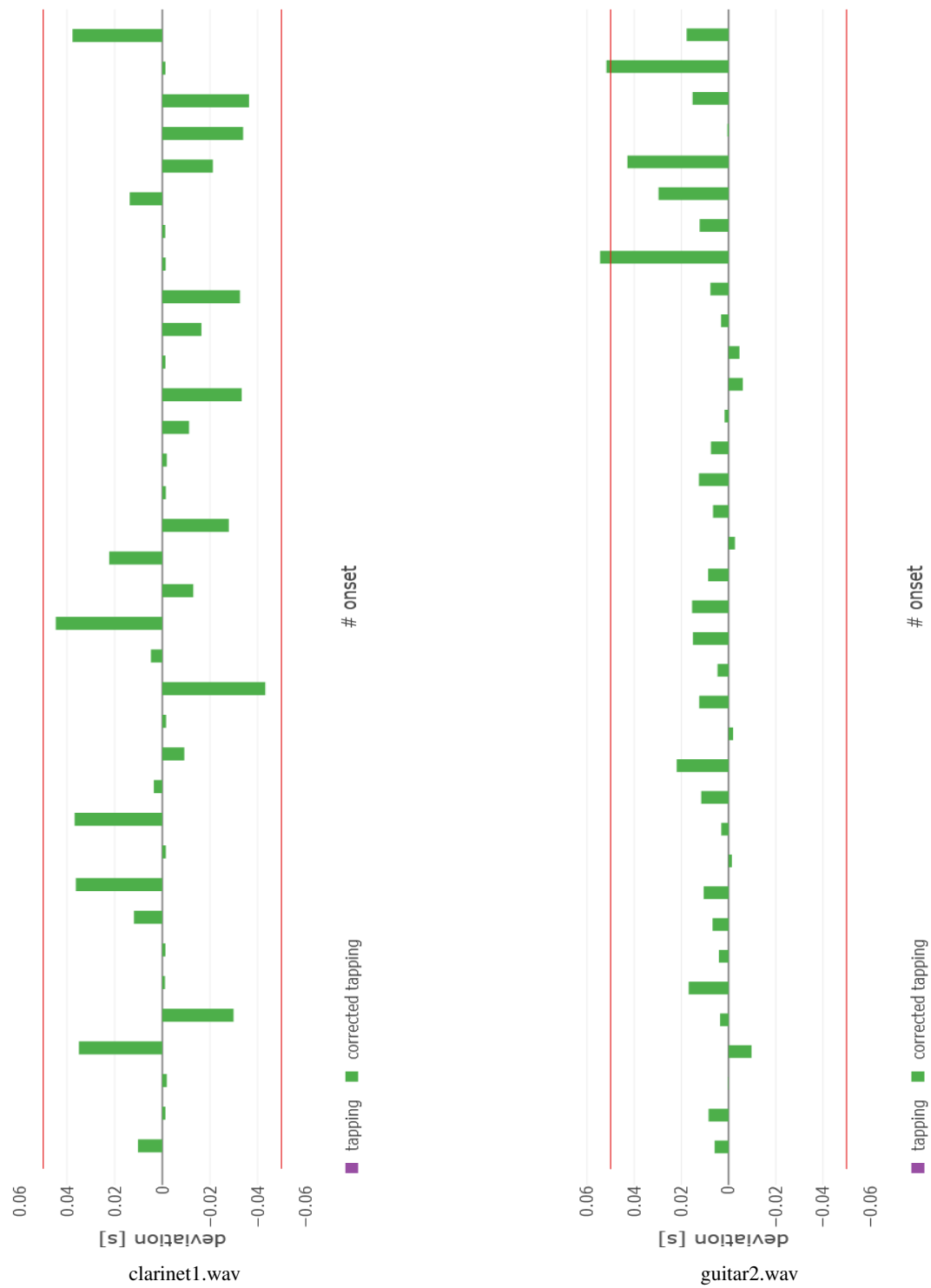


Figure 9



Figure 10

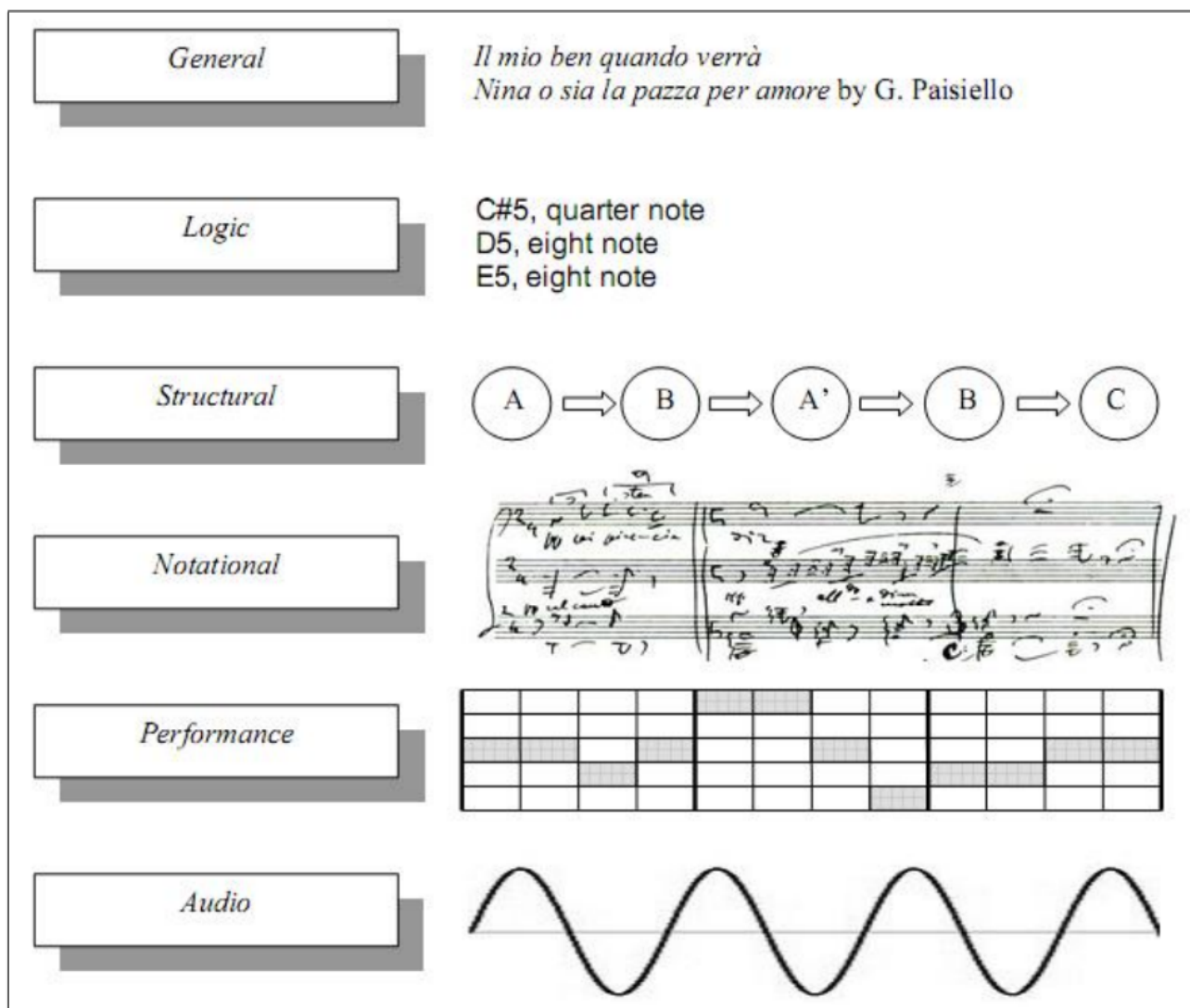


Figure 11: IEEE1599 Layers