

This article was downloaded by: [Mathieu Ramona]

On: 09 February 2012, At: 07:54

Publisher: Taylor & Francis

Informa Ltd Registered in England and Wales Registered Number: 1072954 Registered office: Mortimer House, 37-41 Mortimer Street, London W1T 3JH, UK



## Applied Artificial Intelligence: An International Journal

Publication details, including instructions for authors and subscription information:

<http://www.tandfonline.com/loi/uaai20>

### A PUBLIC AUDIO IDENTIFICATION EVALUATION FRAMEWORK FOR BROADCAST MONITORING

Mathieu Ramona <sup>a</sup>, Sébastien Fenet <sup>b</sup>, Raphaël Blouet <sup>c</sup>, Hervé Bredin <sup>d</sup>, Thomas Fillon <sup>b</sup> & Geoffroy Peeters <sup>a</sup>

<sup>a</sup> IRCAM, Paris, France

<sup>b</sup> Institut Télécom, Télécom ParisTech, CNRS-LTCl, Paris, France

<sup>c</sup> Yacast, Paris, France

<sup>d</sup> LIMSI-CNRS, Orsay Cedex, France

Available online: 06 Feb 2012

To cite this article: Mathieu Ramona, Sébastien Fenet, Raphaël Blouet, Hervé Bredin, Thomas Fillon & Geoffroy Peeters (2012): A PUBLIC AUDIO IDENTIFICATION EVALUATION FRAMEWORK FOR BROADCAST MONITORING, Applied Artificial Intelligence: An International Journal, 26:1-2, 119-136

To link to this article: <http://dx.doi.org/10.1080/08839514.2012.629840>

PLEASE SCROLL DOWN FOR ARTICLE

Full terms and conditions of use: <http://www.tandfonline.com/page/terms-and-conditions>

This article may be used for research, teaching, and private study purposes. Any substantial or systematic reproduction, redistribution, reselling, loan, sub-licensing, systematic supply, or distribution in any form to anyone is expressly forbidden.

The publisher does not give any warranty express or implied or make any representation that the contents will be complete or accurate or up to date. The accuracy of any instructions, formulae, and drug doses should be independently verified with primary sources. The publisher shall not be liable for any loss, actions, claims, proceedings, demand, or costs or damages whatsoever or howsoever caused arising directly or indirectly in connection with or arising out of the use of this material.

## A PUBLIC AUDIO IDENTIFICATION EVALUATION FRAMEWORK FOR BROADCAST MONITORING

**Mathieu Ramona<sup>1</sup>, Sébastien Fenet<sup>2</sup>, Raphaël Blouet<sup>3</sup>, Hervé Bredin<sup>4</sup>, Thomas Fillon<sup>2</sup>, and Geoffroy Peeters<sup>1</sup>**

<sup>1</sup>*IRCAM, Paris, France*

<sup>2</sup>*Institut Télécom, Télécom ParisTech, CNRS-LTCl, Paris, France*

<sup>3</sup>*Yacast, Paris, France*

<sup>4</sup>*LIMSI-CNRS, Orsay Cedex, France*

□ *This paper presents the first public framework for the evaluation of audio fingerprinting techniques. Although the domain of audio identification is very active, both in the industry and the academic world, there is at present no common basis to compare the proposed techniques. This is because corpuses and evaluation protocols differ among the authors. The framework we present here corresponds to a use-case in which audio excerpts have to be detected in a radio broadcast stream. This scenario, indeed, naturally provides a large variety of audio distortions that makes this task a real challenge for fingerprinting systems. Scoring metrics are discussed with regard to this particular scenario. We then describe a whole evaluation framework including an audio corpus, together with the related groundtruth annotation, and a toolkit for the computation of the score metrics. An example of an application of this framework is finally detailed, that took place during the evaluation campaign of the Quaero project. This evaluation framework is publicly available for download and constitutes a simple, yet thorough, platform that can be used by the community in the field of audio identification to encourage reproducible results.*

### INTRODUCTION

Audio identification is a special case of audio event detection that covers the detection and the identification of an audio excerpt (a music track, an advertisement, a jingle, etc.) in an audio recording (either a short excerpt or a broadcast stream). Two techniques are used for that purpose: *audio watermarking*, which relies on embedding meta-information, robust to common alterations, within the audio signal and *audio fingerprinting* (sometimes called *audio hashing*), where audio occurrences are detected through

This work was realized as part of the Quaero Programme, funded by OSEO, French State agency for innovation.

Address correspondence to Mathieu Ramona, IRCAM, 1 place Igor Stravinsky, 75004 Paris, France.  
E-mail: mathieu.ramona@ircam.fr

the recognition of code signatures extracted from short snippets of the signal. These signatures are designed to make a compact representation of the audio content, linked to some perceptually relevant cues, which remains robust to typical distortions observed on audio signals, such as dynamic compression, various encodings, equalization, time scale changes, etc. Because audio watermarking implies an initial processing of the signal source to embed the watermark, it cannot be applied to unknown signals. This paper, hence, focuses on audio fingerprint techniques for audio identification.

The audio identification technology underlies several key applications, including broadcast monitoring, internet content identification, copyrights control, and interactive behavioral targeted advertising. This explains a great effort of contributions in the community during the last decade, despite the relative novelty of the domain, mostly from industrial actors, such as Philips (Haitsma and Kalker 2002), Shazam (Wang 2003), Google (Weinstein and Moreno 2007), (Mohri, Moreno, and Weinstein 2008), Fraunhofer (Allamanche et al. 2001), (Herre, Allamanche, and Hellmuth 2001), or Microsoft (Burges, Platt, and Jana 2003), (Burges, Platt, and Jana 2002). Many propositions also emerge from the academic research area: Ircam owns a technology based on a double-nested Fourier transform (Rodet, Worms, and Peeters 2003), NTT Basic Research Lab has proposed the *active search* algorithm (Smith, Murase, and Kashino 1998), and the Pompeu Fabra University owns a technology based on the so-called *audio DNA* model (Neuschmied, Mayer, and Batlle 2001), (Cano et al. 2002).

However, it remains impossible at present to compare the different systems described in the literature, because no common framework or corpus has been proposed, apart from the TRECVID evaluation on video copy detection (Smeaton, Over, and Kraaij 2006). Indeed, most of the evaluations are applied to private corpuses, of which volume and nature vary greatly between the articles. Also, the evaluation protocols, as well as the scoring metrics, are generally based on different use-cases and reflect different underlying priorities for the authors, which induce very different insights and conclusions on the algorithms.

Moreover, because the key point of audio fingerprinting is the detection of audio events under common distortions, evaluations are often based on the application of controlled synthetic audio distortions applied to audio items, which do not necessarily reflect the constraints of a real-world use-case.

We thus propose in this article the first public evaluation framework focused on audio identification, based on a scenario involving the detection of audio excerpts in broadcast radio streams. The framework consists of a public corpus and an evaluation toolkit named PxAFE. This corpus is not based on artificial audio degradations but on the real-world degradations induced by the radio broadcast production, which implies a wide

variety of combined distortions. This corpus hence makes a challenging and realistic task for audio identification methods. Relevant metrics, related to the use-case, are also provided, together with a discussion about their respective characteristics.

These contributions define a consistent evaluation framework that is made publicly available to the community in order to encourage benchmarking in the field of audio identification instead of private evaluations. We will then thoroughly describe the process and the results of the first evaluation campaign of the Quaero project<sup>1</sup> on audio identification, held in September 2010, which is based on this evaluation framework.

This paper is structured as follows. First, the various evaluation schemes in the literature will be briefly examined, in order to give an outline of the common protocols, as well as the synthetic audio degradations generally used to assess the robustness of the fingerprint codes. A new evaluation framework will then be proposed, which includes a corpus and an evaluation toolkit. The latter includes the implementation of the scoring metrics discussed earlier. Then, the next section will detail an example of the application of this framework on the evaluation campaign of the Quaero project, together with the results of the participants. Some comments and short-term perspectives on the framework will finally be given in the concluding section.

## AUDIO IDENTIFICATION EVALUATION

### Audio Identification in the Context of Event Detection

Audio event detection covers a wide range of scenarios of audio analysis. Depending on the type of events that must be recognized, their duration, and their acoustical variability, varying techniques are employed.

Frame-based classification methods are generally used when each event is defined by an acoustic source, such as gun shots (Clavel, Ehrette, and Richard 2005), applause or cheers (Cai et al. 2003). More general acoustic classes, such as speech or music, can also be considered as audio events, and involve a very large literature (Lin et al. 2005; Ramona and Richard 2009, etc.). This kind of problem implies the use of statistical methods to cover the whole range of variability of the acoustic sources, and focuses on the possible confusion among the classes.

On the other hand, events may not denote sources, but audio signals themselves, as in the detection of jingles (Pinquier and André-Obrecht 2004), advertisements, or musical tracks. This scenario dramatically reduces the scope of variations of the event occurrences, which is restricted to typical audio degradations that affect the signal while keeping it perceptually recognizable.

This case covers what is denoted by *audio identification* in this paper. It is characterized by the fact that a single example is available for each event in the training process and involves specific techniques: typically audio fingerprinting and audio watermarking. Another typical aspect of this scenario is the very large number of events generally involved. For instance, a music track identification task can scale up to several million different classes. The confusion among classes is thus critical, because the slightest overlap might lead to a very large number of false alarms. The audio identification scenario thus focuses on the compromise between discrimination and robustness to audio degradations.

Figure 1 illustrates the typical workflow of an audio fingerprinting system. As all machine-learning systems, audio fingerprinting requires two modes. The first is the *database creation*, where the set of target audio items is processed by the system for the extraction of their fingerprints. All fingerprints are stored in a database and allow to link a given content to tags or to metadata. The other mode allows *audio identification*, based on the fingerprints computed from the audio stream.

Although the framework presented here is originally dedicated to the evaluation of works on audio fingerprinting, it can indeed be used for any audio identification task, with no restriction on the technique employed.

The following sections detail the typical audio degradations found in the literature, as well as the evaluation protocols for audio identification systems.

### Considering Usual Audio Degradations

A very diverse panel of audio degradations can be found in the literature, designed to reproduce most of the audio effects that can be applied to an audio signal, affecting its quality, without changing its semantic content (i.e., what is perceptually received by the listener). Most of them

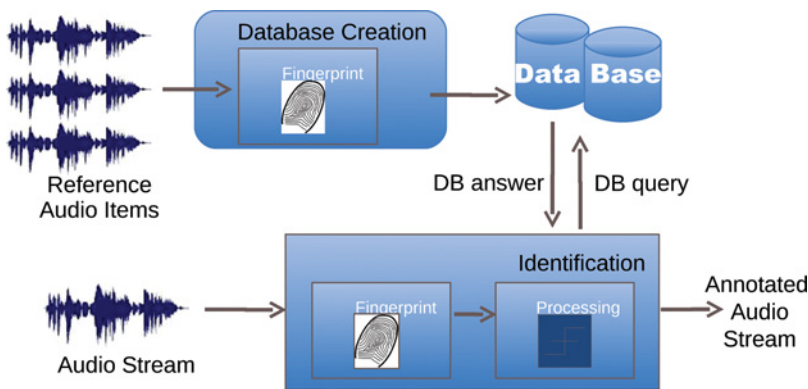


FIGURE 1 Illustration of the audio fingerprinting workflow. (Figure is provided in color online.)

are inspired by the studies on robustness of Haitsma and Kalker (2002) and Allamanche et al. (2001).

The main issue of this discussion is the distinction among three kinds of degradations:

- *Numerical degradations*, by far the most convenient to apply, because they can be simulated numerically.
- *Acoustic degradations*, which involve somehow a conversion to acoustic waves. Their simulation requires more equipment (microphones, loudspeakers, etc.), but remains possible.
- “*Real-world*” *degradations* are a much more complex blend that combines numerous degradations and requires a whole sound production chain, e.g., broadcast radio production and transmission.

We detail here most of the degradations found in the literature, which fall in the first two categories:

**Audio encodings (numerical):**

- *MP3 or WMA encoding/decoding*, from very low (8 kbps) to usual bitrates (128 kbps),
- *Real Media encoding/decoding*,
- *GSM encoding/decoding*, at full rate, with a controlled carrier to interference (C/I) ratio,
- *Resampling*, down to half the sample rate and up again.

**Filtering (numerical):**

- *All-pass filtering*, using an IIR filter,
- *Equalization*, with a 10 to 30-band equalizer,
- *Band-pass filtering*,
- *Telephone band-pass*, between 135 and 3700 Hz,
- *Echo filter*, simulating old time radio.

**Noise addition, of controlled SNR (numerical):**

- *White noise addition*, using a uniform or gaussian white noise,
- *Real-world noise addition*, adding a capture of noisy environment (e.g., a crowded pub),
- *Speech addition*.

**Dynamics changes (numerical):**

- *Amplitude dynamic compression*,
- *Multiband companding*, specifically defined in the TRECVID evaluation,

- *Volume change*, affecting the global volume with a constant or slightly varying factor.

#### **Temporal changes (numerical):**

- *Time scale modification*, up to +4% and -4% without affecting the pitch,
- *Linear speed change*, up to +4% and -4%, with both tempo and pitch affected,
- *Time shift*, the signal is slightly shifted in order to affect the alignment of the temporal frames. (This will be discussed thoroughly in the next section.)

#### **Acoustic conversions (acoustic):**

- *D/A A/D conversion*, using a commercial analog tape recorder,
- *Re-recording*, through a loudspeaker/microphone chain. Possibly in a noisy environment.

These reproducible degradations are shared by almost all the experiments in the literature (Belletini and Mazzini, 2010; Liu, Yun, and Kim 2009; Jang et al. 2009), but very few examples of real-world degradations are found (Betser, Collen, and Rault 2007; Cano et al. 2002; Piquier and André-Obrecht 2004), and all of them are based on radio broadcast recordings. In fact, such recordings include most of the artificial degradations detailed here: the signal is generally numerically encoded, and all the filtering processes, such as all-pass and band-pass filtering and equalization, are very common effects in radio broadcast production. This is also true for time scale, pitch shift, and linear speed modifications, which are used especially on musical tracks. Finally, real-world noise addition is also observed, because most radio-show hosts speak during the instrumental introduction of the songs.

Synthetic distortions are strictly controlled and studied independently, whereas real-world radio broadcast signals provide a varied set of complex combinations among all these distortions. Finally, the audio streaming constraint induces the loss of alignment between the original audio excerpts and the observed audio frames, which will be discussed later on.

These remarks motivate our proposal of an evaluation framework based on a radio broadcast corpus.

### **Existing Evaluation Protocols**

As stated in the previous section, a large majority of the past contributions rely on synthetic audio degradations; this constrains the evaluation protocol and the corpus. Indeed, in most cases, the corpus consists of a collection of unrelated audio (mostly music) tracks. The queries are subsets

of the original tracks, on which various degradations are applied. The dominant evaluation protocol thus consists in detecting in the queries the original tracks learnt from the corpus.

However, it remains focused on the false rejects (also called false negatives or deletion errors). A slightly different protocol involves a collection of matching and nonmatching pairs of audio excerpts. The matching pairs include original clean tracks and their degraded versions, and the nonmatching pairs are arbitrary. Through the number of matching and nonmatching pairs, the balance between false rejects and false alarms (also called false positives or insertion errors) can be controlled.

Another interesting approach deals with the distances between fingerprint codes themselves. It cannot evaluate directly the performance of an algorithm, but the comparison of the distributions of matching and nonmatching fingerprint pairs gives very useful insights on the discriminativity of a fingerprint code. This was initiated by Haitsma and Kalker (2002), who measure the bit error rate (BER, i.e., the Hamming distance) on the binary Philips fingerprint.

As stated in the previous section, the last protocol encountered in the literature is based on real broadcast recordings that include occurrences of the corpus's audio items. The drawback is a reduced control over the number of occurrences, but the "real-world" combinations of degradation and the presence of long sections with no occurrence enforce the realism of the evaluation.

Another major argument in favor of broadcast-oriented evaluation is the arbitrariness of the occurrences positions in the audio streams, which imply random desynchronization between the original item signal and the occurrence signal. Indeed, as stated in a previous publication (Ramona and Peeters 2011), a slight time-shift between the original audio and the degraded audio induces distortions on the fingerprint that are more important than most degradations. Many evaluations skip this issue because they apply degradations directly on the original audio sample and thus preserve the temporal alignment. A real-world corpus implicitly induces random time-shifts in the occurrences.

Evaluation protocols in the literature also cover a wide range of score metrics. Generally used with queries as subsets of the corpus, the accuracy rate (the number of queries correctly identified) is by far the most common criterion. However, it does not cover false alarms (false positives). A more thorough approach is to use the false negative/false positive (FN/FP) pair as a measure of performance. The TRECVID evaluation plan for copy detection (Smeaton, Over, and Kraaij 2006) also defines a refined metric that specifically balances false negatives and positives.

When the method involves a threshold or a tunable parameter, a receiver operating characteristic (ROC) curve is used to illustrate the evolution of the FP/FN measures with regard to the parameter. The precision and recall metrics (derived from the FP/FN measures) are also



**TABLE 1** Comparative List of the Evaluation Protocols in the Literature, Specifying Corpus and Queries Sizes, Protocols and Score Metrics

Articles	Corpus	Queries	Protocol	Metric
(Allamanche et al. 2001) (Fraunhofer)	15 k	NA	Subsets	Acc + Top 10
(Cano et al. 2002)	50 k	12 h (104)	Subsets + Broad	FP/FN
(Haitsma and Kalker 2002) (Philips)	4	4	Subsets + BER	Nb hits
(Wang 2003) (Shazam)	10 k	250	Subsets	Acc + FP
(Pinquier and Andre-Obrecht 2004)	200	10 h (132)	Broad	Acc
(Seo et al. 2006)	8 k	NA	Subsets	ROC (Pre/Rec)
(Covell and Baluja 2007) (Google)	10 k	1000	Subsets	Acc + ROC (FP/FN)
(Betser, Collen, and Rault 2007)	30	18 h (243)	Broad	Recall
(Kim and Yoo 2007)	NA	7 M	Pairs	ROC (FP/FN)
(Mohri, Moreno, and Weinstein 2008)	15 k	3,600	Subsets	Acc
(Liu, Yun, and Kim 2009)	13 k	2400	Pairs	Acc + ROC (FP/FN)
(Jie, Gang, and Jun 2009)	500	NA	Subsets	Top 1,5,10,20,50
(Jang et al. 2009)	100	44 k	Pairs	Acc + ROC (FP/FN)
(Belletini and Mazzini 2010)	100 k	NA	Subsets	Acc
(Li, Liu, and Xue 2010)	1,822	100	Subsets + BER	Top 1,5,10 + FP/FN
(Smeaton, Over, and Kraaij 2006) (TRECVID)	400 h	12000	Subsets	Balanced FP/FN

used in a similar way. Finally, since audio identification is often based on a nearest neighbor scheme, instead of computing the accuracy on the first result, a tolerance can be set to accept detections that are within the  $N$  best ranked. This measure is denoted by “Top- $N$ .”

Table 1 summarizes the different evaluation protocols presented here, together with the score metrics, the size of the corpuses, and the number of queries. The last line describes the TRECVID copy detection evaluation task, mentioned earlier. It is the only other evaluation campaign related to our subject, but it is mostly dedicated to video indexing and does not propose a specific task for audio identification.

It can be noted that, even though audio identification is typically considered a large-scale problem, most of the evaluations are limited to a few thousand, or a even a few hundred, tracks in the corpuses or used as queries. We hope the framework provided here, and described in the next section, offers a larger scale than most evaluations mentioned in the table.

## PROPOSED EVALUATION FRAMEWORK

### Broadcast-Oriented Corpus

The evaluation corpus described in this paper comes from a basic media monitoring use-case. Given a set of target musical tracks, it

determines if and when an audio item has been broadcasted, i.e., the time of broadcasting and the identity of the target track occurrences.

The evaluation corpus has been drawn in the framework of the subtask *Audio Identification/Fingerprint* of the Quaero project. It provides a corpus of target audio items, for the database building, and several continuous radio broadcast streams, for the identification. The annotation process was done semi-automatically and entirely checked by human operators. For the evaluation within the Quaero project, the corpus is characterized by around 8,000 audio items of one minute in length. These audio excerpts correspond to audio segments previously broadcasted and manually annotated and extracted. There is one excerpt per audio item, available in RAW format, little-endian, with 16 bits per sample, at a sampling rate of 11025 Hz.

The test audio streams consists in full days of radio stream, from different stations, captured and saved in succeeding five-minutes chunks, encoded in AAC, with a bit-rate of 64 kbps and a sampling rate of 11025 Hz. The item signatures may not be completely broadcasted in the streams. The beginning of a signature is not known.

For legal reasons, we are not allowed to distribute the whole stream associated to a media. We hence have built an *artificial* stream made of a concatenation of short audio excerpts (from 3 s to 45 s) coming from different media. The stream is made up of around 8,000 audio items provided in 4 files of 4 hours each. As the proposed corpus includes real radio broadcasted items from a set of 15 media, it is likely to cover all the challenges of audio identification for radio monitoring. For each target audio item, 30 seconds of audio data are provided to compute the fingerprint signature.

The provided corpus is partly synthetic, because it relies on a concatenation of excerpts. Nevertheless, it is important to note that the excerpts themselves come from real-world signals, and thus provide realistic degradations, as stated earlier, considered as the main issue for a serious study on robustness.

The groundtruth is determined by a set of XML files, one for each media file. Of course the XML annotation specifies only the items that are present in the corpus delivered. The annotation, as stated earlier, comes from an audio identification engine manually checked, with a precision of about 1 s. Each file lists occurrences of target audio items, with the following XML structure:

```
<MusicTrack>
  <id>123456</id>
  <idMedia>548</idMedia>
  <title>Some kind of wonderful</title>
  <artist>Grand Funk Rail road</artist>
  <album>Caught in the act</album>
  <genre>Pop/Rock International</genre>
```

```
<startDate>2010-07-05 00:03:43</startDate>
<endDate>2010-07-05 00:03:55</endDate>
</MusicTrack>
```

where:

- `<id>` identifies the audio item on air between `<startDate>` and `<endDate>`,
- `<idMedia>` identifies the source media, in order to extract identification scores for specific media,
- `<title>`, `<artist>`, `<album>` and `<genre>` are various metadata on the target item. This makes it possible, for instance, to extract identification scores by genre,
- `<startDate>` and `<endDate>` respectively indicate the start and end time of the elements denoted by the ID. However, since the item signatures provided for the database build are subsets of the original songs (as stated earlier), these time limits are larger than the actual occurrence of the item signature alone. This issue will be thoroughly discussed in the next section.

Audio streams and signature files are freely available for academic research use.<sup>2</sup>

## Scoring

The present use-case implies the following constraints:

- Occurrences of known audio items are to be detected in an audio stream.
- The audio items are only known through short snippets called *audio signatures*.
- The position of these signatures within the original tracks is unknown.

Please note that if audio tracks in the stream are not occurrences of any item in the corpus, they are not considered as “occurrences.” They are part of the “no-item” areas described later, like any unknown signal, because they cannot be recognized.

### *Scope of Evaluation*

Because only a part of the audio items (the signature) is taken into account in the training process, one could consider the sole signature area being the track itself, as in Figure 2(a), where an occurrence of a target item is shown in medium gray, between two areas with no item in light gray; the signature snippet is shown in dark gray. In this context, detections timestamps can be evaluated precisely.

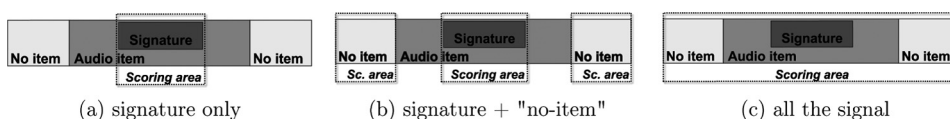


FIGURE 2 Possible evaluation scopes.

The exclusion of the no-item areas greatly reduces the scope of evaluation. The evaluation should rather imply the whole signal. Because music tracks generally imply repetitive structures (chorus, verses, etc.), the areas of the song outside the signatures are likely to present high correlations with the signatures themselves.<sup>3</sup> It would therefore be preferable to exclude the latter from the evaluation scope, as shown in Figure 2(b), because the notions of missed detection and false alarms are ambiguous outside the signatures.

However, as stated earlier, the position of the signatures in the tracks is unknown; it is thus impossible to define the scope of evaluation according to it. Therefore, the case described in Figure 2(c), implying the whole signal, is the only one applicable here. Although this configuration is theoretically less reliable than the previous schemes, this issue is answered by not considering the temporal positions of the detections in the items, as discussed in the next paragraph on score metrics.

### Score Metrics

The metric for the evaluation of audio identification is based on a punctual event detection scheme, which means that only instants of decision are taken into account, instead of segments. The possible segmentation of the signal into frames is exclusively related to the audio identification process, and is not relevant to the following score metrics.

A previous section exposed the main score metrics found in the literature. Most evaluations are based on the accuracy (i.e., number of correctly detected occurrences over the number of occurrences), which is here equivalent to the Precision measure, and inversely proportional to the false reject rate (or deletion error rate). The Recall measure is similarly related to the false alarm rate (or insertion error rate). Note that in the context of audio identification, no distinction is made by the community between substitutions (i.e., mistaking an item for another one) and insertions. ROC curves are also commonly used, but not adapted to an evaluation framework designed to compare different algorithms. Indeed, in a proper benchmark, each system is evaluated as a stand-alone application that does not require any parameter tuning.

This framework is thus restricted to false reject and false alarm measures. The TRECvid evaluation plan (Smeaton, Over, and Kraaij 2006) defines several balances between the two measures, but the present use-case

has no preference for a specific error. The counting of the false alarms is discussed in this section.

Let us denote a collection of  $N$  audio occurrences of items. Each occurrence  $n$  is between time boundaries  $t_n^{sta}$  and  $t_n^{end}$  in the signal and is related to an item of index  $i_n$  in the database. The evaluation considers a set of  $D$  punctual detections. Each detection  $d$  is related to an item index  $j_d$  and instant  $\tau_d$ . As stated before, the signature is not precisely located, so that a correct detection involves only the observation of at least one detection of the correct item within the scope of the occurrence (the medium gray "Audio item" scope in Figure 2). The number of correct detections (Accuracy) is thus defined as:

$$S_{OK} = \text{Card}\{n \in [1 \dots N], \exists d, j_d = i_n \text{ and } t_n^{sta} \leq \tau_d \leq t_n^{end}\}. \quad (1)$$

The number of false rejects is straightforward and is implicit in the Accuracy definition ( $S_{FR} = N - S_{OK}$ ). The global score is defined as the following rate:

$$R = 1/N \cdot (S_{OK} - [S_{FA} + S_{FA}^{out}]), \quad (2)$$

where  $S_{FA}$  and  $S_{FA}^{out}$ , respectively, denote the false alarm rate within and outside the occurrences.

The expression of  $R$  depends on the definition of the false alarm rates, which depends on the tolerance accepted. The most straightforward definition counts each false alarm as one error, as shown in Figure 3(a):

$$\begin{cases} S_{FA,1} = \text{Card}\{d \in [1 \dots D], \exists n, t_n^{sta} \leq \tau_d \leq t_n^{end} \text{ and } j_d \neq i_n\}, \\ S_{FA,1}^{out} = \text{Card}\{d \in [1 \dots D], \nexists n, t_n^{sta} \leq \tau_d \leq t_n^{end}\}. \end{cases} \quad (3)$$

However, this measure is strongly biased because false alarms are upper bounded by  $D$  (i.e., can be arbitrarily numerous), whereas correct detections are bounded by the number of occurrences  $N$ . The nonhomogeneity

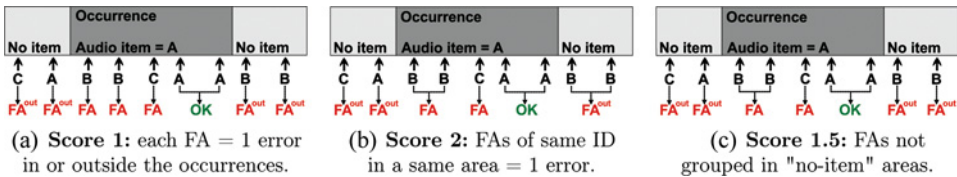


FIGURE 3 Comparison of the false alarms counting methodologies. (Figure is provided in color online.)

between the two measures can lead to negatives scores, especially with a high number of detections.

The scheme presented in Figure 3(b) corrects this balance by grouping as a single error the false alarms of a same wrong item in a given area. However, false alarms of distinct items are still counted separately, otherwise the balance would be strongly biased in favor of correct detections. In order to unify the evaluation scheme, the areas between the occurrences (No-item areas in Figure 3) are considered as occurrences where only false alarms are counted (no item can be detected). Hence, the evaluation does not consider individual detections, but only the items detected within areas. False alarms are expressed as follows:

$$\begin{cases} S_{FA,2} = \text{Card}\{n \in [1 \dots N], \exists d, t_n^{sta} \leq \tau_d \leq t_n^{end} \text{ and } j_d \neq i_n\}, \\ S_{FA,2}^{out} = \text{Card}\{n \in [1 \dots N], \nexists d, t_n^{sta} \leq \tau_d \leq t_n^{end}\}. \end{cases} \quad (4)$$

The last proposition, illustrated in Figure 3(1c), is a compromise between the previous two. On a musical radio station, the “no-item” areas are rather short, and only contain transitions between musical tracks. Other stations, though, can produce mostly talk shows, lasting several hours. The grouping of false alarms of the same item separated by long durations is therefore not relevant. The metric 1.5 counts the false alarms by items within the occurrences, and by detections outside them:

$$\begin{cases} S_{FA,1.5} & = & S_{FA,2} \\ S_{FA,1.5}^{out} & = & S_{FA,1}^{out} \end{cases} \quad (5)$$

The three metrics we have detailed are jointly provided and used in the present evaluation framework. None is favored *a priori* over the others. The global scores stand as follows:

$$R_1 = 1/N \cdot \left( S_{OK} - [S_{FA,1} + S_{FA,1}^{out}] \right) \quad (6)$$

$$R_2 = 1/N \cdot \left( S_{OK} - [S_{FA,2} + S_{FA,2}^{out}] \right) \quad (7)$$

$$R_{1.5} = 1/N \cdot \left( S_{OK} - [S_{FA,2} + S_{FA,1}^{out}] \right) \quad (8)$$

### The PyAFE Evaluation Toolkit

As simple as it may seem at first, the actual implementation of these scoring metrics can be complex. Researchers should not have to implement their own version. This would increase the risk of getting several (yet differing)

implementations of the same scoring metric, therefore compromising the fair comparison paradigm for which we aim. Moreover, as state-of-the-art audio fingerprinting systems are getting really close to perfection, it becomes crucial that the performance of two systems can be compared accurately. That is why we introduce the **PyAFE**<sup>4</sup> toolkit:

### Python Audio Fingerprinting Evaluation

PyAFE was developed in the framework of the *Evaluation* work-package of the Quaero project and is made freely available as open-source, downloadable software.<sup>5</sup>

It was designed as a modular piece of software, in order to be easily extended in the future:

- Two modules provide the necessary functions to parse the groundtruth and detection XML files (whose formats are described on the **PyAFE** website).
- The core module includes the implementation of score metrics  $R_1$ ,  $R_2$ , and  $R_{1.5}$  described previously. The number of correct detections  $S_{OK}$ , false rejects  $S_{FR}$ , and false alarms  $S_{FA}$  and  $S_{FA}^{out}$  are provided.

Included in the **PyAFE** toolkit, is an all-in-one command line evaluation tool. It provides an easy-to-use, straightforward way of getting evaluation results.

This tool actually browses all subdirectories of the groundtruth base directory looking for annotation files. For each of them, the corresponding detection file in the detection base directory is evaluated. Depending on the requested level of verbosity, it can output one single score for the whole set of files, one score per file, or even the detailed list of errors made for every single file. Another useful option allows performing the evaluation using only a subset of audio targets, by providing the list of their identifiers.

The full documentation can be found on the **PyAFE** website, together with sample groundtruth and detection files.

## APPLICATION OF THE FRAMEWORK FOR THE QUAERO PROJECT

The Quaero project includes a subproject focused on *Audio Identification and Fingerprinting*. Starting in 2010, an audio fingerprinting evaluation campaign is organized every year during summer, throughout the existence of the Quaero project. The pilot evaluation took place in September 2010, for which a dedicated corpus was collected.

## Annotated Radio Broadcast Corpus

This corpus consists in the recording of 5 weeks of the French radio station RTL, captured and saved on disk in 5-minute chunks. Therefore, the total duration reaches 840 hours. Similarly to what was described in Section 3.1, the whole dataset was annotated by manually checking the output of an audio identification engine (with precision around 1 second). The whole corpus is divided into two parts: four weeks make the training set and the remaining week constitutes the test set. A set of 7, 309 1-minute-long target audio signatures was gathered to build the database.

## Pilot Evaluation

A few months before the submission deadline, participants were provided with the training set, the corresponding annotations, and the target signatures. They were also provided with the **PyAFE** toolkit, knowing that this very tool would be used by the evaluation coordinator for the actual evaluation. The test set—obviously free of any annotation—was then distributed to participants, and they submitted back the output (XML files) of their audio fingerprinting systems. Three participants submitted at least one run for the 2010 Quaero evaluation campaign:

- **Participant 1** provided a system based on double-nested FFT, combined to a k-NN search among the database codes, and a post-processing that correlates succeeding detection timestamps.
- **Participant 2** has developed a fingerprinting system based on the Shazam algorithm (Wang 2003).
- **Participant 3** has developed an audio fingerprinting system that implements some slight modification over the system described by Haitsma and Kalker (2002). It is based on quantizing differences of energy measures from overlapped short-term power spectra.

**TABLE 2** Pilot Quaero Evaluation Results

System	$S_{OK}/N$	$S_{FA,1} (S^{out})$	$R_1$	$S_{FA,1.5} (S^{out})$	$R_{1.5}$	$S_{FA,2} (S^{out})$	$R_2$
Participant 1	445/459	0 (2)	96.5%	0 (2)	96.5%	0 (2)	96.5%
Participant 2	381/459	0 (0)	83.0%	0 (0)	83.0%	0 (0)	83.0%
Participant 3	442/459	0 (2)	95.9%	0 (2)	95.9%	0 (2)	95.9%

The false alarm scores are indicated inside and outside the occurrences (respectively, before and between the parentheses).



After a necessary adjudication phase during which the test corpus annotations were corrected when necessary, the final results were published within the Quaero consortium. Table 2 reports the results of the various submitted runs, as provided by the **PyAFE** toolkit.

## CONCLUSION & FUTURE WORK

Audio fingerprint is one of the main industrial challenges of the last years, related to the diffusion of music. Although many systems have been proposed to perform this task, no comparison among existing technologies has been performed because of the lack of unified evaluation frameworks. In this paper we described a proposal for the evaluation of audio fingerprint algorithms in the case of broadcasted music. This framework contains the definition of score metrics, their public implementation and a public test-set corresponding to the use-case of broadcast monitoring of music. Because this test-set contains radio streams, it naturally allows representing several degradation types artificially created in previous evaluations. The whole framework is accessible online. As an example, we presented the use of this framework and the results obtained during the first Quaero audio fingerprint evaluation.

The current framework focuses on the punctual detection of music tracks (“when has this music track been broadcasted?”) in a corpus, given a short signature of each track. Further scenarios will include the detection of the exact boundaries of music track diffusion (“when did this broadcast radio or TV start playing the song and end it?”) or boundaries within the music tracks themselves (“which part of the track has been played?”). Further works will concentrate on extending the framework to the detection of advertisement and jingles in audio streams, as well as blind recurrent patterns detection in audio streams (therefore, without previous knowledge of signatures).

In a further step, it could be worth defining distortion measures between the reference signature and the broadcasted audio. This could lead to an objective distortion measure for each corpus.

## NOTES

1. Please consult <http://quaero.org> for more information on the project.
2. <http://pyafe.niderb.fr>.
3. This correlation cannot be quantified, since it highly depends on the fingerprint code design, but it suffices to say that it is used as a basic assumption for automatic musical structure retrieval (see Peeters, Burthe, and Rodet 2002), for an example based on fingerprint techniques).
4. PyAFE is pronounced like the last name of Edith Piaf the famous French singer.
5. <http://pyafe.niderb.fr>

## REFERENCES

- Allamanche, E., J. Herre, O. Hellmuth, B. Fröba, T. Kastner, and M. Cremer. 2001. Content-based identification of audio material using MPEG-7 low level description. In *Proceeding of the international symposium on music information retrieval (ISMIR '01)*. Bloomington, Indiana, USA.
- Belletini, C., and G. Mazzini. 2010. A framework for robust audio fingerprinting. *Journal of Communications* 5:409–424.
- Betsler, M., P. Collen, and J.-B. Rault. 2007. Audio identification using sinusoidal modeling and application to jingle detection. In *Proceeding of the international symposium on music information retrieval (ISMIR '07)*. Vienna, Austria.
- Burges, C. J., J. C. Platt, and S. Jana. 2002. Extracting noise-robust features from audio data. In *Proceeding of the international conference on acoustics, speech and signal processing (ICASSP '02)*, vol. 1, 1021–1024. Orlando, Florida, USA.
- Burges, C. J., J. C. Platt, and S. Jana. 2003. Distortion discriminant analysis for audio fingerprinting. *IEEE Transactions on Speech and Audio Processing* 11:165–174.
- Cai, R., L. Lu, H.-J. Zhang, and L.-H. Cai. 2003. Highlight sound effects detection in audio stream. In *Proceeding of the IEEE international conference on multimedia and expo (ICME '03)*, vol. 3, 37–40. Baltimore, Maryland, USA.
- Cano, P., E. Battle, H. Mayer, and H. Neuschmied. 2002. Robust sound modeling for song detection in broadcast audio. In *Proc. 112th AES Convention*, 1–7. Munich, Germany.
- Clavel, C., T. Ehrette, and G. Richard. 2005. Events detection for an audio-based surveillance system. In *Proceeding of the IEEE International Conference on Multimedia and Expo (ICME '05)*, 1306–1309. Amsterdam, The Netherlands.
- Covell, M., and S. Baluja. 2007. Known-audio detection using waveprint: Spectrogram fingerprinting by wavelet hashing. In *Proceeding of the international conference on acoustics, speech and signal processing (ICASSP '07)*, 237–240. Honolulu, Hawaii, USA.
- Haitisma, J., and T. Kalker. 2002. A highly robust audio fingerprinting system. In *Proceeding of the international symposium on music information retrieval (ISMIR '02)*. Paris, France.
- Herre, J., E. Allamanche, and O. Hellmuth. 2001. Robust matching of audio signals using spectral flatness features. In *Proceeding of the IEEE workshop on applications of signal processing to audio and acoustics (WASPAA '01)*, 127–130. New Paltz, New York, USA.
- Jang, D., C. D. Yoo, S. Lee, S. Kim, and T. Kalker. 2009. Pairwise boosted audio fingerprint. *IEEE Transactions on Information Forensics and Security* 4:995–1004.
- Jie, T., L. Gang, and G. Jun. 2009. Improved algorithms of music information retrieval based on audio fingerprint. In *Proceedings of the 3rd international symposium on intelligent information technology application workshops (IITAW '09)*.
- Kim, S., and C. D. Yoo. 2007. Boosted binary audio fingerprint based on spectral subband moments. In *Proceeding of the international conference on acoustics, speech and signal processing (ICASSP '07)*, vol. 1, 241–244. Honolulu, Hawaii, USA.
- Li, W., Y. Liu, and X. Xue. 2010. Robust audio identification for MP3 popular music. In *Proceeding of the 33rd international ACM SIGIR conference on research and development in information retrieval*, 627–634. Geneva, Switzerland.
- Lin, C. C., S. H. Chen, T. K. Truong, and Y. Chang. 2005. Audio classification and categorization based on wavelets and support vector machine. *IEEE Transactions on Speech and Audio Processing* 13:644–651.
- Liu, Y., H. S. Yun, and N. S. Kim. 2009. Audio fingerprinting based on multiple hashing in DCT domain. *IEEE Signal Processing Letters* 16:525–528.
- Mohri, M., P. Moreno, and E. Weinstein. 2008. Efficient and robust music identification with weighted finite-state transducers. *IEEE Transactions on Audio, Speech and Language Processing* 18: 197–207.
- Neuschmied, H., H. Mayer, and E. Battle. 2001. Identification of audio titles on the internet. In *Proceedings of the international conference on web delivering of music (Wedelmusic '01)*. Florence, Italy.
- Peeters, G., A. L. Burthe, and X. Rodet. 2002. Toward automatic music audio summary generation from signal analysis. In *Proceedings of the international symposium on music information retrieval (ISMIR '02)*. Paris, France.

- Pinquier, J., and R. André-Obrecht. 2004. Jingle detection and identification in audio documents. In *Proceedings of the international conference on acoustics, speech and signal processing (ICASSP '04)*, vol. 4, 329–332.
- Ramona, M., and G. Peeters. 2011. Audio identification based on spectral modeling of bark-bands energy and synchronisation through onset detection. In *Proceeding of the international conference on acoustics, speech and signal processing (ICASSP '11)*, 477–480. Prague, Czech Republic.
- Ramona, M., and G. Richard. 2009. Comparison of different strategies for a SVM-based audio segmentation. In *Proceedings of the 17th European signal processing conference (EUSIPCO '09)*. Glasgow, Scotland.
- Rodet, X., L. Worms, and G. Peeters. 2003. Brevet FT R&D/03376: Procédé de caractérisation d'un signal sonore—Patent 20050163325 Method for characterizing a sound signal.
- Seo, J. S., M. Jin, S. Lee, D. Jang, S. Lee, and C. D. Yoo. 2006. Audio fingerprinting based on normalized spectral subband moments. *IEEE Signal Processing Letters* 13:209–212.
- Smeaton, A. F., P. Over, and W. Kraaij. 2006. Evaluation campaigns and trecvid. In *Proceedings of the 8th ACM international workshop on multimedia information retrieval (MIR '06)*, 321–330. Santa Barbara, California, USA.
- Smith, G., H. Murase, and K. Kashino. 1998. Quick audio retrieval using active search. In *Proceeding of the international conference on acoustics, speech and signal processing (ICASSP '98)*, vol. 6, 3777–3780. Seattle, Washington, USA.
- Wang, A. L. C. 2003. An industrial-strength audio search algorithm. In *Proceedings of the International Symposium on Music Information Retrieval (ISMIR '03)*. Washington, D.C., USA.
- Weinstein, E., and P. Moreno. 2007. Music identification with weighted finite-state transducers. In *Proceedings of the international conference on acoustics, speech and signal processing (ICASSP '07)*, vol. 2, 689–692. Honolulu, Hawaii, USA.