

EXTENDED PLAYING TECHNIQUES: THE NEXT MILESTONE IN MUSICAL INSTRUMENT RECOGNITION

First Author

Affiliation1

author1@ismir.edu

Second Author

Retain these fake authors in
submission to preserve the formatting

Third Author

Affiliation3

author3@ismir.edu

ABSTRACT

The expressive variability in which a musical note can be produced conveys some essential information to the modeling of orchestration and style. Yet, although the automatic recognition of a musical instrument from the recording of a single “ordinary” note is now considered a solved problem, the ability of a computer to precisely identify instrumental playing techniques (IPT) remains largely underdeveloped. In this paper, we conduct a benchmark of machine listening systems for query-by-example browsing among 143 instrumental playing techniques, including the most contemporary, for 16 instruments in the symphonic orchestra, thus amounting to 469 triplets of instrument, mute, and technique. We identify and discuss three necessary conditions for significantly outperforming the classical mel-frequency cepstral coefficients (MFCC) baseline: the inclusion of second-order scattering coefficients to account for the presence of amplitude modulations; the inclusion of long-range temporal dependencies; and the resort to large-margin nearest neighbors (LMNN), a supervised metric learning method that reduces intra-class variability in feature space. We report a P@5 of 99.7% for instrument recognition (baseline at 92.5%) and of 61.0% for playing technique recognition (baseline at 50.0%).

1. INTRODUCTION

The gradual diversification of the timbral palette in Western classical music at the turn of the 20th century is reflected in five concurrent trends: the addition of new instruments to the symphonic instrumentarium, either by technological inventions (e.g. theremin) or importation from non-Western musical cultures (e.g. marimba) [50]; the creation of novel instrumental associations, as epitomized by *Klangfarbenmelodie* [51]; the temporary alteration of resonant properties through mutes and other “preparations” [16]; a more systematic usage of extended instrumental techniques, such as artificial harmonics, *col legno batutto*, or flutter tonguing [29]; and the resort to electronics and digital audio effects [61]. The first of these trends has

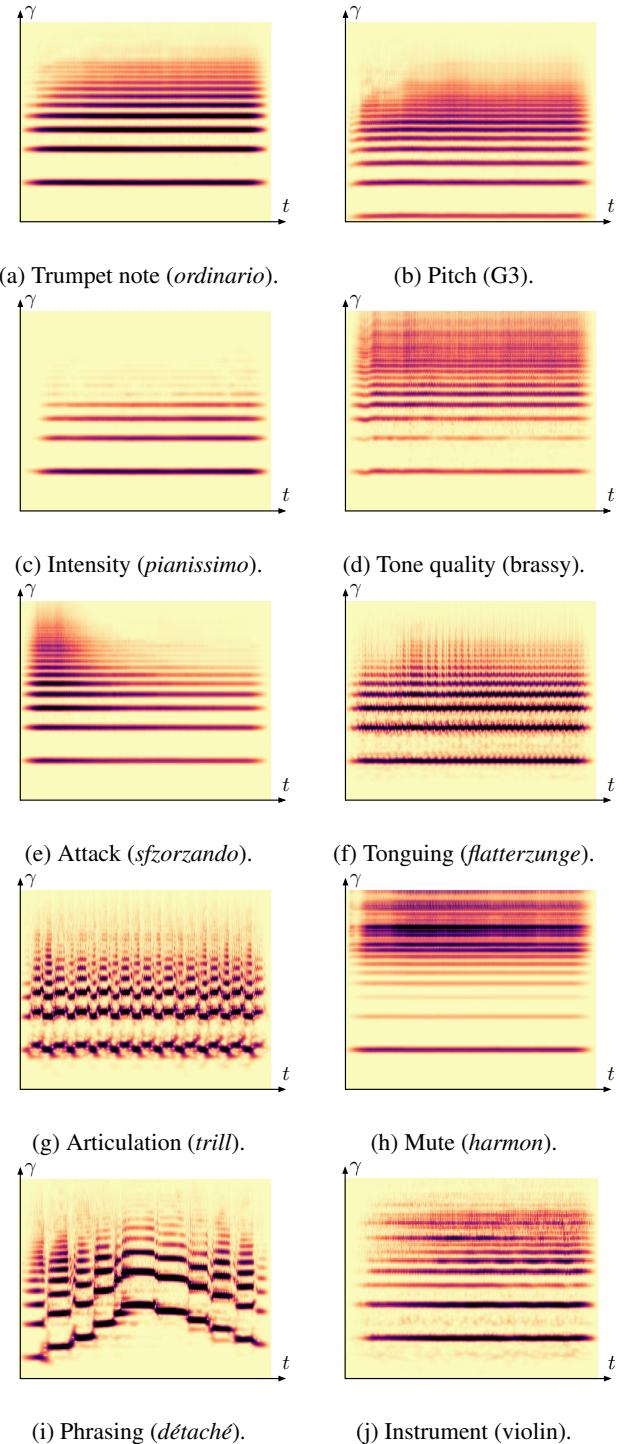


Figure 1: Ten factors of variations of a musical note.

somewhat stalled: to this day, most Western composers rely on an acoustic instrumentarium that is only marginally different from the one that was available in the Late Romantic period. Nevertheless, the latter approaches to timbral diversification were massively adopted into post-war contemporary music. In particular, an increased concern for the concept of musical gesture [22] has liberated many unconventional instrumental techniques from their figurative connotations, thus making the so-called “ordinary” playing style merely one of many compositional – and improvisational – options.

Far from being exclusive to erudite music, extended playing techniques are also commonly found in oral tradition; in some cases, they even stand out as a distinctive component of musical style. Four well-known examples are: the snap pizzicato (“slap”) of the upright bass in rockabilly, the growl of the tenor saxophone in rock’n’roll, the shuffle stroke of the violin (“fiddle”) in Irish folklore, and the glissando of the clarinet in Klezmer music. Consequently, the mere knowledge of organology (the instrumental *what?* of music), as opposed to chironomics (its gestural *how?*), is a rather weak source of information for browsing and recommendation in large music databases.

Yet, past research in music information retrieval (MIR), and especially machine listening, rarely acknowledges the benefits of integrating the influence of performer gestures into a coherent taxonomy of musical instrument sounds. Instead, gestures are either framed as a spurious form of intra-class variability between instruments, without delving into its interdependencies with pitch and intensity; or, symmetrically, as a probe for the acoustical study of a given instrument, without enough emphasis onto the broader picture of orchestral diversity.

One major cause of this gap in research is the difficulty of collecting and annotating data for contemporary instrumental techniques. Fortunately, such obstacle has recently been overcome, owing to the creation of databases of instrumental samples in a perspective of spectralist music orchestration [40]. In this article, we capitalize on the availability of data to formulate a new line of research in MIR, namely the joint retrieval of organological information (“*what* instrument is being played in this recording?”) and chironomical information (“*how* is the musician producing sound?”), while remaining invariant to other factors of variability, which are deliberately regarded as contextual: at what pitches and intensities, but also where, when, why, by whom, and for whom was the music recorded.

Figure 1a shows the constant-*Q* wavelet transform (CQT) of a trumpet musical note, as played with an ordinary technique. Unlike most existing publications on instrument classification, which exclusively focus on pitch (Figure 1b) and intensity (Figure 1c) as the main factors of intra-class variability, this paper aims at accounting for the presence of instrumental playing techniques (IPT), such as changes in tone quality (Figure 1d), attack (Figure 1e), tonguing (Figure 1f), and articulation (Figure 1h), either as intra-class variability (instrument recognition task) or as inter-class variability (IPT recognition task). The analysis

of IPTs whose definition necessarily involves more than a single musical event, such as phrasing (Figure 1i), is beyond the scope of this paper.

Section 2 reviews the existing literature on the topic. Section 3 derives the task of IPT classification from the definition of both a taxonomy of instruments and a taxonomy of gestures. Section 4 describes how two topics in machine listening, namely scattering transforms and supervised metric learning, are relevant to address this task. Section 5 reports the results from an IPT classification benchmark on the Studio On Line (SOL) dataset.

2. RELATED WORK

This section some of the recent MIR literature on the audio analysis of instrumental playing techniques, with a focus on the available datasets afferent to each formulation of the problem.

2.1 Classification of ordinary isolated notes

The earliest works on musical instrument recognition restricted their scope to individual notes played with an ordinary technique – with datasets such as MUMS [47], MIS, RWC [23], and Philharmonia – thus eliminating most factors of intra-class variability due to the performer [6, 11, 18, 25, 27, 41, 57]. These works have culminated with the development of a support vector machine (SVM) classifier trained on spectrotemporal receptive fields (STRF), which are idealized computational models of neurophysiological responses in the central auditory system [14]. Not only did it attain a near-perfect mean accuracy of 98.7% on the RWC dataset, but the confusion matrix of its automated predictions was closely similar to the confusion matrix of human listeners [49]. Therefore, the supervised classification of musical instruments from recordings of ordinary notes could arguably be considered a solved problem; we refer to [8] for a recent review of the state of the art.

2.2 Classification of solo recordings

One straightforward extension of the problem above is the classification of solo phrases, encompassing some variability in melody [30], for which the accuracy of STRF models is around 80% [48]. Since the Western tradition of solo music is essentially limited to a narrow range of instruments (e.g. piano, classical guitar, violin) and genres (sonatas, contemporary, free jazz, folk), datasets of solo phrases, such as solosDb [26], are exposed to strong biases. This issue is partially mitigated by the recent surge of multitrack datasets, such as MedleyDB [9], which has spurred a renewed interest in single-label instrument classification [59]. In addition, the cross-collection evaluation methodology [32] allows to prevent the risk of overfitting caused by the relative homogeneity of these small datasets in terms of artists and recording conditions [10]. To this date, the best classifier of solo recordings is a spiral convolutional network [35] trained on the Medley-solosDB dataset [34], i.e. a cross-collection dataset which aggregates MedleyDB and solosDB following the procedure

of [17]. We refer to [24] for a recent review of the state of the art.

2.3 Multilabel classification in polyphonic mixtures

Because most publicly released musical recordings are polyphonic, the generic formulation of instrument recognition as a multilabel classification task is the most appropriate for large-scale deployment [12, 42]. However, it suffers from two methodological caveats: first, polyphonic instrumentation is not independent from other attributes of information, such as geographical origin, genre, or key; and secondly, the inter-rater agreement decreases with the number of overlapping sources [20, chapter 6]. Such issues are all the more troublesome that there is, to this date, no annotated dataset of polyphonic mixtures that is diverse enough to be devoid of artist bias. The Open-MIC initiative, from the newly created Community for Open and Sustainable Music and Information Research (COSMIR), might contribute to mitigating them in the near future [43].

2.4 Single-instrument playing technique classification

Lastly, there is a growing interest for studying the role of the performer in musical acoustics, from both perspectives of sound production and sound perception. Besides its interest in audio signal processing, this topic is connected to other disciplines, such as biomechanics and gestural interfaces [45]. The majority of the available literature focuses on the range of IPTs afforded by a single instrument: recent examples include clarinet [37], percussion [53], piano [7], guitar [13, 19, 52], violin [60], cello [15, chapter 6], and erhu [58]. Some publications frame timbral similarity in a polyphonic setting, yet do so according to a purely perceptual definition of timbre – with continuous attributes such as brightness, warmth, dullness, roughness, and so forth – yet without connecting these attributes to the discrete latent space of IPTs [3].

In this paper, we formulate the retrieval of expressive parameters of musical timbre at the scale of the symphonic orchestra at large, while expliciting these parameters in terms of sound production (i.e. through a discrete set of instructions, readily interpretable by the performer) rather than by means of perceptual epithets only. We refer to [31] for a recent review of the state of the art.

3. TASKS

In this section, we distinguish taxonomies of musical instruments from taxonomies of musical gestures.

3.1 Taxonomies

The Hornbostel-Sachs taxonomy (H-S) strives to organize the diversity of musical instruments according to their manufacturing characteristics only, and is purposefully unaffected by sociohistorical background [46]. Because it offers an unequivocal way of describing any acoustic instrument without any prior knowledge on its afferent IPTs, it serves as a *lingua franca* in ethnomusicology and museology, especially for ancient or rare instruments which may

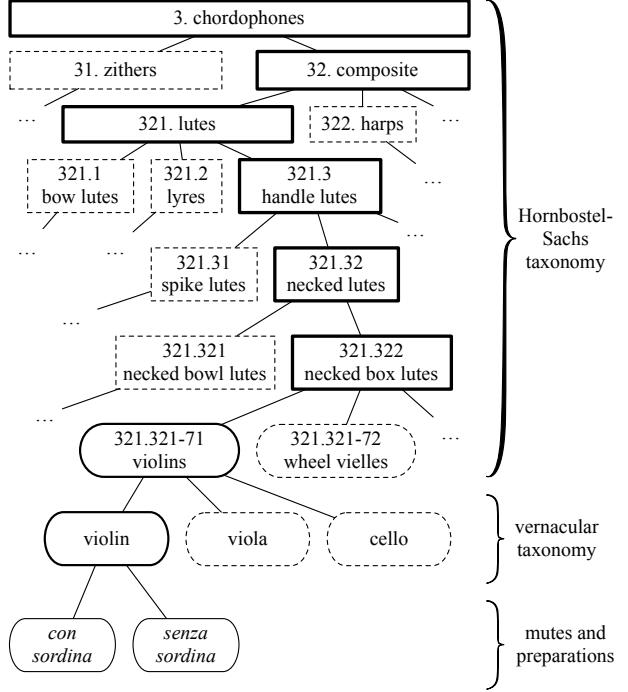


Figure 2: Taxonomy of musical instruments.

lack available informants. The location of the violin in H-S (321.321-71), as depicted in Figure 2, also encompasses the viola and the cello in addition to the violin. This is because these three instruments, viewed as inert objects, share a common morphology, despite differences in posture for the performer: both violin and viola are usually played under the jaw whereas the cello is held between the knees. Accounting for these differences begs to refine H-S by means a vernacular taxonomy. Most instrument taxonomies in music signal processing, including MedleyDB and AudioSet [21], reach the vernacular level rather than conflating all instruments belonging to the same H-S node. In some cases, an even finer level of granularity is attained by the listing of potential alterations to the instrument – be them permanent or temporary, at the time scale of more than a single note – that affect its resonant properties after the end of the conventional manufacturing process, e.g. mutes and other preparations [16]. The only example of node in the MedleyDB taxonomy reaching this level is *tack piano* [9].

Unlike musical instruments, which are approximately amenable to a hierarchical taxonomy of resonating objects, IPTs result from a complex synchronization between multiple gestures, which may involve both hands and arms, as well as diaphragm, vocal tract, and sometimes the whole body. As a result, there is no immediate way to interface them with H-S, or indeed any tree-like structure [28]. Instead, every playing technique is described by a finite collection of categories, each belonging to a different “namespace”; Figure 3 illustrates such namespaces in the case of the violin. It therefore appears that, rather than aiming for a mere increase in granularity with respect to H-S, a coherent research program around extended playing techniques should formulate them as belonging to a meronomy,

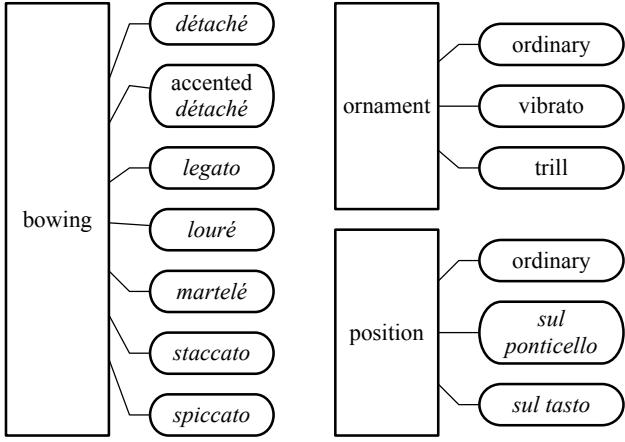


Figure 3: Namespaces of violin playing techniques.

i.e. a modular entanglement of part-whole relationships, in the fashion of the Visipedia initiative in computer vision [5]. In recent years, some publications have attempted to lay the foundations of such a modular approach, with the aim of making H-S relevant to contemporary music creation [38, 56]; yet, such considerations are still in large part speculative, and offer no definitive procedure for evaluating, let alone training, information retrieval systems.

3.2 Application setting and evaluation

In what follows, we adopt a middle ground position between the two aforementioned approaches: neither a supervised classifier (as in a hierarchical taxonomy), nor a caption generator (as in a meronomy), our system is a query-by-example search engine in a large database of isolated notes. This system is meant to provide a small number k of nearest neighbors in the dataset of musical instrument samples to any user-defined audio query $x(t)$. In the context of contemporary music creation, this $x(t)$ may be an instrumental or vocal sketch; a sound event recorded from the environment; a computer-generated waveform; or any mixture of the above [40]. Upon inspecting the k nearest neighbors returned by the search engine, the composer may decide to retain one of the retrieved notes, in which case its attributes (pitch and intensity, but also the exact playing technique) are readily available and can be included into the musical score to approximate the query.

Faithfully evaluating such a system is a difficult procedure, and ultimately would rest on its practical usability, as judged by the composers themselves. Nevertheless, a useful quantitative metric for this task is the precision at k ($P@k$) of the test set with respect to the training set, both under a instrument taxonomy and an IPT taxonomy.

3.3 Studio On Line dataset (SOL)

The Studio On Line dataset (SOL) was recorded at Ircam in 2002 and is freely downloadable as part of the Orchids software for computer-assisted orchestration¹. It comprises 16 musical instruments playing 25444 isolated notes

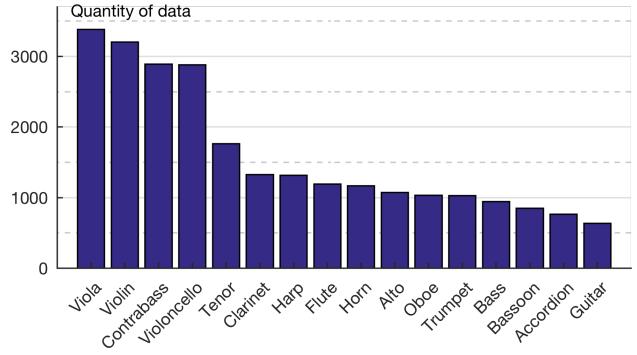


Figure 4: Instruments in the SOL dataset.

in total. The distribution of these notes, shown in Figure 4, spans the full combinatorial diversity of applicable intensities, pitches, preparations (i.e. mutes), as well as all applicable playing techniques – whose most common are shown in Figure 3 – is heavy-tailed (average 178, standard deviation 429): this is because some playing techniques are shared between many instruments (e.g. *tremolo*) whereas other are instrument-specific (e.g. *xylophonic* which is specific to the harp). The SOL dataset has 143 IPTs in total, and 469 applicable instrument-mute-technique triplets.

4. METHODS

In this section, we describe the scattering transform and supervised metric learning used to implement all query-by-example systems in our benchmark.

4.1 Scattering transform

The scattering transform is a cascade of two wavelet modulus operators: the first layer extracts the average spectral envelope of $x(t)$ at frequencies λ_1 , whereas the second layer $S_2x(t, \lambda_1, \lambda_2)$ extracts amplitude modulations of this spectral envelope at rates λ_2 . The set of frequencies discretizes the auditory range according to the mel scale, with $Q_1 = 12$ bins per octave at topmost frequencies; whereas rates λ_2 follow a geometric sequence between λ_1 and some minimal rate T^{-1} , with $Q_2 = 1$ bin per octave. We refer to [2] for a general introduction to scattering transforms in audio classification, and to [33] for a discussion on its application to musical instrument classification in solo recordings, as well as its close connections with STRF. The scattering transform is theoretically suited to model extended playing techniques, since various values of the rate λ_2 provably characterize some of the most common nonstationarities in sound production, including tremolo, vibrato, and dissonance [1]. In the following, we denote by $Sx(t, \lambda)$ the concatenation of all scattering coefficients, whether the generic scattering path λ corresponds to a singleton (λ_1) or a pair (λ_1, λ_2).

In order to match a decibel-like perception of loudness, we apply the path-adaptive, quasi-logarithmic compression

$$\tilde{S}x_i(\lambda) = \log \left(1 + \frac{Sx_i(\lambda)}{\varepsilon \times \mu(\lambda)} \right) \quad (1)$$

¹ Link to SOL dataset: <http://forumnet.ircam.fr/product/orchids-en/>

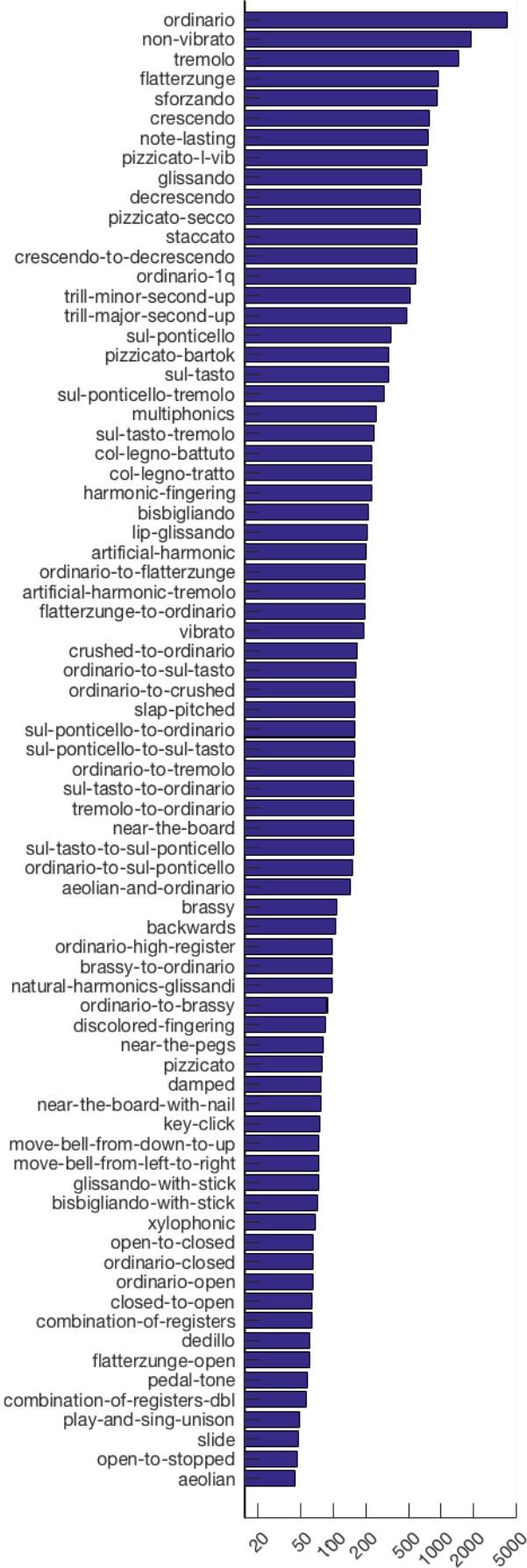


Figure 5: Playing techniques in the SOL dataset.

where $\varepsilon = 10^{-3}$ and $\mu(\lambda)$ is the median value of the scattering coefficient $\mathbf{S}\mathbf{x}_i(\lambda)$ for path λ across samples i .

4.2 Metric learning

Linear metric learning algorithms generate a matrix \mathbf{L} such that the Mahalanobis distance

$$D_{\mathbf{L}}(\mathbf{x}_i, \mathbf{x}_j) = (\tilde{\mathbf{S}}\mathbf{x}_i - \tilde{\mathbf{S}}\mathbf{x}_j)^T \mathbf{L}^T \mathbf{L} (\tilde{\mathbf{S}}\mathbf{x}_i - \tilde{\mathbf{S}}\mathbf{x}_j) \quad (2)$$

between all pairs of samples $(\mathbf{x}_i, \mathbf{x}_j)$ optimizes some objective function. Observe that the Euclidean distance is a particular case of the equation above, in which case \mathbf{L} boils down to identity. We refer to [4] for a review of the state of the art. In particular, the large-margin nearest neighbors (LMNN) algorithm aims at bringing all k nearest neighbors \mathbf{x}_j of every \mathbf{x}_i closer than the canonical Euclidean distance $D(\mathbf{x}_i, \mathbf{x}_j) = \|\tilde{\mathbf{S}}\mathbf{x}_i - \tilde{\mathbf{S}}\mathbf{x}_j\|$ if \mathbf{x}_i and \mathbf{x}_j belong to the same class, and further apart otherwise [54, 55].

As compared to a class-wise generative model (such as Gaussian mixtures), a global linear model ensures some robustness to minor alterations of the taxonomy, which is important in the context of IPT: e.g. what some instrumentists may call *slide*, others will call *glissando*. Furthermore, although one of its major drawback relies on its strong dependency on the Euclidean neighbors determine intra-class variability [44], this drawback is alleviated in the case of a future space based on scattering transform coefficients, whose Euclidean metric provably approximates the extent of elastic deformation needed to shear $\mathbf{x}_i(t)$ into $\mathbf{x}_j(t)$ in the time-frequency domain [39, Theorem 2.16].

5. EXPERIMENTAL RESULTS

In this section, we apply the aforementioned methods to instrument and IPT query-by-example retrieval in the Studio On Line (SOL) dataset.

To evaluate the performance of the alternative implementations of the query by example system, the precision at rank k ($p@k$) is considered. It is computed as the number of relevant items among the k items retrieved by the system divided by k . Recall here is evaluated as the seed and the retrieved item having the same label in the reference partition. Each item of the database is considered as a seed, and the precision of the system under evaluation is computed for this seed. The precision for this system is the precision averaged over all the seeds.

In all subsequent experiments, we set the number of retrieved items to $k = 5$.

5.1 Evaluation of instrument recognition

In the task of instrument recognition, each of the k elements \mathbf{x}_j returned by the system is considered relevant to the query \mathbf{x}_i if and only if \mathbf{x}_i and \mathbf{x}_j correspond to the same instrument, regardless of pitch, intensity, mute, and playing technique.

We compare scattering features to a baseline of mel-frequency cepstral coefficients (MFCC), corresponding to the 13 lowest frequencies after applying a discrete cosine

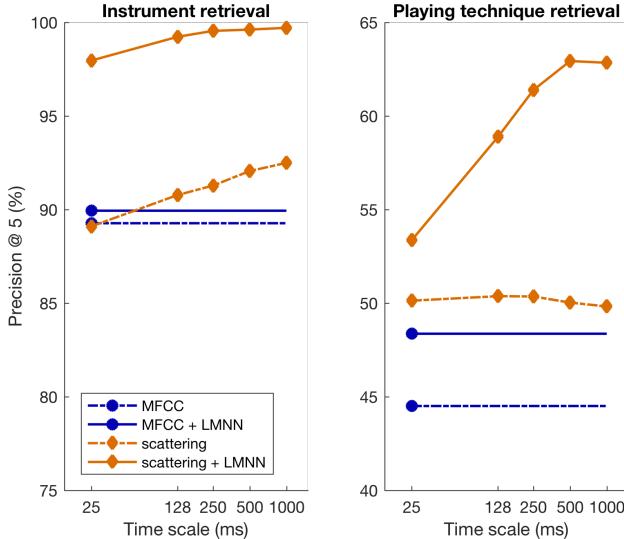


Figure 6: Summary of results on the SOL dataset.

transform (DCT) on the logarithm of the 40-band mel-frequency spectrum. We evaluate the influence of LMNN on both MFCC and scattering coefficients/ In addition, we vary the maximum time scale T of amplitude modulation between 25 ms and 1 s. In the case of MFCC, $T = 25$ ms corresponds to the inverse of the lowest audible frequency ($T = 40$ Hz). Therefore, increasing the frame duration T has no effect on the value of MFCC, because the mel-spectrogram is equivalent to a local averaging of the wavelet scalogram at the time scale T , leaving unchanged the result of the global averaging of $\mathbf{S}(t, \lambda)$ at the time scale of the whole musical note.

Figure 6 summarizes our results. We find that MFCC reach a relatively P@5

It which achieves a p@5 of 89%. Keeping all 40 MFCCs coefficients rather than the lowest 13 degrades accuracy down to 84%. This result is in line with the fact that the low frequency modulation of the spectrum are most useful to model the resonant cavity, be it the vocal tract for speech and the instrument body for musical instruments.

Scattering: 89%. Disabling median renormalization – i.e. setting 84%. Disabling logarithmic compression altogether: 76%. These results are consistent with [36].

Enlarging the frame size do not influence the performance of the MFCC baseline. On contrary the p@5 benefit from an increase of the frame size of the scattering, see Figure 6(left).

Considering the LMNN algorithm to project the features only marginally improves the performance for the MFCC baseline (p@5 of 90%), whereas the scattering features highly benefit from it. The gain decreases while the framesize is increased, probably due to a glass ceiling effect.

In order to control if the LMNN artificially benefits from the higher degree of freedom provided by the higher dimensionality of the scattering features (494 for the 25 ms. frame size), another set of baseline features based on monomials of the MFCCs with the same dimensionality (494) is considered. Considering the LMNN over this new

set of features gives a slight increase of performance (p@5 of 91%), we thus assume that the better performance is due to the higher level of expressivity of the scattering features.

5.2 Evaluation of playing technique recognition

As expected, considering the playing technique partition as a reference gives a harder task to perform as the MFCC baseline achieves a p@5 of 45% and the scattering a p@5 of 50%. When considering a frame size of 25 ms., the LMNN projection gives an equivalent

Again, enlarging the frame size do not influence the performance of the MFCC baseline, see Figure 6(right). More surprisingly, the performance of the scattering features stalls with respect of the frame size. An increase would have been expected, since many playing techniques involves long term modulation of energy through frequency bands which are expected to be nicely captured by the scattering features.

It is only when considering the LMNN projection that the p@5 strongly increases with respect to the frame size. We believe that the fact that the system benefit from a richer input only if a supervised projection is considered is due to the need to select *a priori* among the large range of modulations modeled by the scattering features the ones that are relevant for the task at hand.

6. CONCLUSION

playing technique is important to fully describe the interaction between a performer and a instrument.

playing technique modeling

need to capture long term modulations

scattering useful for this end

benchmark results demonstrate the need to have a supervised step to specialize the representation

The results described in this paper validate the processing pipeline on the reference partitions of instruments and playing techniques. Future work will focus on the definition of a reference similarity among couples of instruments and playing techniques that is defined perceptually and the evaluation of the proposed processing pipeline against this perceptually defined reference.

7. REFERENCES

- [1] Joakim Andén and Stéphane Mallat. Scattering representation of modulated sounds. In *Proc. DAFX*, 2012.
- [2] Joakim Andén and Stéphane Mallat. Deep scattering spectrum. *IEEE Trans. Sig. Proc.*, 62(16):4114–4128, 2014.
- [3] Aurélien Antoine and Eduardo R. Miranda. Musical acoustics, timbre, and computer-aided orchestration challenges. In *Proc. ISMA*, 2018.
- [4] Aurélien Bellet, Amaury Habrard, and Marc Sebban. A survey on metric learning for feature vectors and structured data. *arXiv preprint arXiv:1306.6709*, 2013.

- [5] Serge Belongie and Pietro Perona. Visipedia circa 2015. *Pattern Recognition Letters*, 72:15 – 24, 2016.
- [6] Emmanouil Benetos, Margarita Kotti, and Constantine Kotropoulos. Musical instrument classification using non-negative matrix factorization algorithms and subset feature selection. In *Proc. IEEE ICASSP*, 2006.
- [7] Michel Bernays and Caroline Traube. Expressive production of piano timbre: touch and playing techniques for timbre control in piano performance. In *Proc. SMC*.
- [8] DG Bhalke, CB Rama Rao, and Dattatraya S Bormane. Automatic musical instrument classification using fractional Fourier transform based-MFCC features and counter propagation neural network. *Journal Int. Inform. Syst.*, 46(3):425–446, 2016.
- [9] Rachel M Bittner, Justin Salamon, Mike Tierney, Matthias Mauch, Chris Cannam, and Juan Pablo Bello. Medleydb: A multitrack dataset for annotation-intensive mir research. In *Proc. ISMIR*, 2014.
- [10] Dmitry Bogdanov, Alastair Porter, Perfecto Herrera Boyer, and Xavier Serra. Cross-collection evaluation for music classification tasks. In *Proc. ISMIR*, 2016.
- [11] Judith C Brown. Computer identification of musical instruments using pattern recognition with cepstral coefficients as features. *J. Acoust. Soc. Am.*, 105(3):1933–1941, 1999.
- [12] Juan José Buried, Axel Robel, and Thomas Sikora. Polyphonic musical instrument recognition based on a dynamic model of the spectral envelope. In *Proc. IEEE ICASSP*, pages 173–176. IEEE, 2009.
- [13] Yuan-Ping Chen, Li Su, Yi-Hsuan Yang, et al. Electric guitar playing technique detection in real-world recording based on f0 sequence pattern recognition. In *Proc. ISMIR*, 2015.
- [14] Taishih Chi, Powen Ru, and Shihab A Shamma. Multiresolution spectrotemporal analysis of complex sounds. *J. Acoust. Soc. Am.*, 118(2):887–906, 2005.
- [15] Magdalena Chudy. *Discriminating music performers by timbre: On the relation between instrumental gesture, tone quality and perception in classical cello performance*. PhD thesis, Queen Mary University of London, 2016.
- [16] Tzenka Dianova. *John Cage's Prepared Piano: The Nuts and Bolts*. PhD thesis, U. Auckland, 2007.
- [17] Patrick J Donnelly and John W Sheppard. Cross-dataset validation of feature sets in musical instrument classification. In *Proc. IEEE ICDMW*, pages 94–101. IEEE, 2015.
- [18] Antti Eronen and Anssi Klapuri. Musical instrument recognition using cepstral coefficients and temporal features. In *Proc. IEEE ICASSP*, volume 2, pages II753–II756. IEEE, 2000.
- [19] Raphael Foulon, Pierre Roy, and François Pachet. Automatic classification of guitar playing modes. In *Proc. CMMR*. Springer, 2013.
- [20] Ferdinand Fuhrmann. *Automatic musical instrument recognition from polyphonic music audio signals*. PhD thesis, Universitat Pompeu Fabra, 2012.
- [21] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *Proc. IEEE ICASSP*, 2017.
- [22] R.I. Godøy and M. Leman. *Musical Gestures: Sound, Movement, and Meaning*. Taylor & Francis, 2009.
- [23] Masataka Goto, Hiroki Hashiguchi, Takuichi Nishimura, and Ryuichi Oka. RWC music database: music genre database and musical instrument sound database. 2003.
- [24] Yoonchang Han, Jaehun Kim, Kyogu Lee, Yoonchang Han, Jaehun Kim, and Kyogu Lee. Deep convolutional neural networks for predominant instrument recognition in polyphonic music. *Proc. Trans. Audio Speech Lang. Process.*, 25(1):208–221, 2017.
- [25] Perfecto Herrera Boyer, Geoffroy Peeters, and Shlomo Dubnov. Automatic classification of musical instrument sounds. *J. New. Mus. Res.*, 32(1):3–21, 2003.
- [26] Cyril Joder, Slim Essid, and Gaël Richard. Temporal integration for audio classification with application to musical instrument classification. *IEEE Trans. Audio Speech Lang. Process.*, 17(1):174–186, 2009.
- [27] Ian Kaminskyj and Tadeusz Czaszejko. Automatic recognition of isolated monophonic musical instrument sounds using kNNC. *J. Intell. Inf. Syst.*, 24(2-3):199–221, 2005.
- [28] Sefki Kolozali, Mathieu Barthet, György Fazekas, and Mark B Sandler. Knowledge representation issues in musical instrument ontology design. In *Proc. ISMIR*, pages 465–470, 2011.
- [29] Stefan Kostka. *Materials and Techniques of Post Tonal Music*. Taylor & Francis, 2016.
- [30] A.G. Krishna and Thippur V. Sreenivas. Music instrument recognition: from isolated notes to solo phrases. In *Proc. IEEE ICASSP*. IEEE, 2004.
- [31] Marc Leman, Luc Nijs, and Nicola Di Stefano. *On the Role of the Hand in the Expression of Music*, pages 175–192. Springer International Publishing, Cham, 2017.
- [32] Arie Livshin and Xavier Rodet. The importance of cross database evaluation in sound classification. In *ISMIR 2003*, 2003.

- [33] Vincent Lostanlen. *Convolutional operators in the time-frequency domain*. PhD thesis, 'Ecole normale supérieure, 2017.
- [34] Vincent Lostanlen, Rachel Bittner, and Slim Essid. Medley-solos-DB: a cross-collection dataset of solo musical phrases, 2018.
- [35] Vincent Lostanlen and Carmine Emanuele Cella. Deep convolutional networks on the pitch spiral for musical instrument recognition. In *Proc. ISMIR*, 2016.
- [36] Vincent Lostanlen, Grégoire Lafay, Joakim Andén, and Mathieu Lagrange. Relevance-based quantization of scattering features for unsupervised mining of environmental audio. *Submitted to EURASIP J. Audio Speech Music Process.*, 2018.
- [37] Mauricio A Loureiro, Hugo Bastos de Paula, and Hani C Yehia. Timbre classification of a single musical instrument. In *Proc. ISMIR*, 2004.
- [38] Thor Magnusson. Musical organics: A heterarchical approach to digital organology. *J. New Music Research*, 46(3):286–303, 2017.
- [39] Stéphane Mallat. Group invariant scattering. *Comm. Pure Applied Math.*, 65(10):1331–1398, 2012.
- [40] Yan Maresz. On computer-assisted orchestration. *Contemp. Mus. Rev.*, 32(1):99–109, 2013.
- [41] Keith D. Martin and Youngmoo E. Kim. Musical instrument identification: A pattern recognition approach. In *Proc. ASA*, 1998.
- [42] Luis Gustavo Martins, Juan José Burred, George Tzanetakis, and Mathieu Lagrange. Polyphonic instrument recognition using spectral clustering. In *Proc. ISMIR*, 2007.
- [43] Brian McFee, Eric J. Humphrey, and Julián Urbano. A plan for sustainable mir evaluation. In *Proc. ISMIR*, 2016.
- [44] Brian McFee and Gert R. Lanckriet. Metric learning to rank. In *Proc. ICML*, 2010.
- [45] Cheryl D Metcalf, Thomas A Irvine, Jennifer L Sims, Yu L Wang, Alvin WY Su, and David O Norris. Complex hand dexterity: a review of biomechanical methods for measuring musical performance. *Front. Psychol.*, 5:414, 2014.
- [46] Jeremy Montagu. It's time to look at Hornbostel-Sachs again. *Muzyka (Music)*, 1(54):7–28, 2009.
- [47] Frank J Opolko and Joel Wapnick. McGill University Master Samples (MUMS), 1989.
- [48] Kailash Patil and Mounya Elhilali. Biomimetic spectro-temporal features for music instrument recognition in isolated notes and solo phrases. *EURASIP J. Audio Speech Music Process.*, 2015(1):27, 2015.
- [49] Kailash Patil, Daniel Pressnitzer, Shihab Shamma, and Mounya Elhilali. Music in our ears: the biological bases of musical timbre perception. *PLOS Comput. Biol.*, 8(11):e1002759, 2012.
- [50] Curt Sachs. *The History of Musical Instruments*. Dover Books on Music. Dover Publications, 2012.
- [51] Arnold Schoenberg. *Theory of Harmony*. University of California, 100th anniversary edition edition, 2010.
- [52] Li Su, Li-Fan Yu, and Yi-Hsuan Yang. Sparse cepstral, phase codes for guitar playing technique classification. In *Proc. ISMIR*, 2014.
- [53] Adam R Tindale, Ajay Kapur, George Tzanetakis, and Ichiro Fujinaga. Retrieval of percussion gestures using timbre classification techniques. In *Proc. ISMIR*, 2004.
- [54] Kilian Q Weinberger, John Blitzer, and Lawrence K Saul. Distance metric learning for large margin nearest neighbor classification. In *Proc. NIPS*, pages 1473–1480, 2006.
- [55] Kilian Q Weinberger and Lawrence K Saul. Distance metric learning for large margin nearest neighbor classification. *J. Mach. Learn. Res.*, 10(Feb):207–244, 2009.
- [56] Stéphanie Weisser and Maarten Quanten. Rethinking musical instrument classification: towards a modular approach to the Hornbostel-Sachs system. *Yearb. Tradit. Music*, 43:122–146, 2011.
- [57] Alicja A Wieczorkowska and Jan M Żytkow. Analysis of feature dependencies in sound description. *J. Intell. Inf. Syst.*, 20(3):285–302, 2003.
- [58] Luwei Yang, Elaine Chew, and Sayid-Khalid Rajab. Cross-cultural comparisons of expressivity in recorded erhu and violin music: Performer vibrato styles. 2014.
- [59] Hanna Yip and Rachel M Bittner. An accurate open-source solo musical instrument classifier. In *Proc. ISMIR, Late-Breaking / Demo (LBD) session*.
- [60] Diana Young. Classification of common violin bowing techniques using gesture data from a playable measurement system. In *Proc. NIME*. Citeseer, 2008.
- [61] Udo Zölzer. *DAFX: Digital Audio Effects*. Wiley, 2011.