

Figure 2. Excerpt of the first movement of Mozart’s Piano Sonata n.2 (K. 279, m.7-10) as seen on the web application used to annotate and review the dataset. Each annotated label can be put in doubt by the reviewer: once a *feedback* is opened over one given measure, comments can be added to exchange with the annotator. Feedbacks have three states: *raised* (in white – just submitted by a reviewer), *in conflict* (red – waiting for a consensus) and *resolved* (green, as in the figure).

Corpus ID: corpus:MIR:mozartpianosonatas:texture:2022:version1.0
Raw Corpus Definition: existing corpus of real symbolic items. Digital Scores of Mozart’s Piano Sonatas following the <i>Neue Mozart-Ausgabe (NME)</i> from Mozart Annotated Sonatas [14]. 9 annotated movements in 3 sonatas (K. 279, K. 280, K. 283). Sampling: Western classical style piano music, available data, all movements in Sonata Form. Type of media diffusion: original files in .mscx and .tsv (Tab Separated Values), annotations in .tsv, .txt and .dez (Dezrann format).
Annotations Origin: manual Concepts definition: concepts and syntax for texture annotations from [7]. Annotation rules: annotation guidelines provided with the dataset. Annotators: one annotator (expert knowledge in music and piano), 2 reviewers with the same background. Validation/reliability: review, one for each movement. Annotation tools: Dezrann web-interface [24].
Documents and Storing Identifier and storage: Köchel number, scores from NMA (Neue Mozart Ausgabe) reference edition, Git repository of the dataset by [14] in musescore and TSV format, annotation files on a Git repository ¹ and online on Dezrann ² .

Table 1. Description of the annotated dataset.

2.2 Annotated corpus

Following the syntax presented in Section 2.1, we release manual annotations of texture at the bar level in the 3 movements of the 3 sonatas K. 279, K. 280 and K. 283 by W. A. Mozart, totaling a set of 1164 bar annotations in the 9 movements. Those early sonatas were all composed in the end of the year 1774, which presumes a style consistency and limits the shift of compositional practice that could be induced by the evolution of piano manufacture. Moreover, these movements present an interesting diversity in tonality, rhythmic signature and tempo, while covering a large variety of textures. All those movements are in Sonata Form, which opens perspectives to pursue research on the links between texture and form [22].

Table 1 synthesizes the properties of the dataset following the conventions proposed in [23].

2.3 Annotation procedure

2.3.1 Granularity of the annotations

Textural segments can have highly variable duration, which tends to complicate their annotation. As a compromise to facilitate statistics and computational processing of the annotations, we propose in this work to annotate a single texture label for each bar, as indicated in the proposition of syntax [7]. In some cases where the texture strongly shifts in the middle of a bar, a comma (,) is used to segment the label in two parts. In particular, this situation may occur on boundaries between musical phrases. By contrast, a texture can also remain unchanged over several consecutive bars, which will lead to a repetition of a same texture label over these bars.

2.3.2 Annotation and reviewing methodology

This section describes our annotation protocol which was inspired by [25] and illustrated in Figure 3. One expert *E* annotated all the measures of the 9 movements with texture labels on Dezrann², a web platform dedicated to the annotation of musical analyses on scores [24]. The annotations were committed to a Git repository to trace the different stages of annotation and reviews. A Python script³ was used to systematically check that the labels were consistently formed before to proceed to the review phase, relieving the reviewers from syntax checking to concentrate on the textural decomposition. Two reviewers *R*₁ and *R*₂ (distinct from the expert), with strong musical background and piano music knowledge, as well as knowledge of the texture syntax used, reviewed those annotations on Dezrann. They added a *feedback* label each time they disagreed or were unsure of the annotator label choice, detailing the reasons of the doubt and possibly providing a new proposition for the label. The expert *E* then studied all the doubts, resolved the obvious ones and discussed with the two reviewers to reach a consensus on the remaining labels.

Over the 1164 initial annotations of the annotator, 31% raised a reviewer feedback and 22% (256) were updated in

¹ <http://algonus.fr/data>

² <http://dezrann.net>

³ <http://algonus.fr/code>

the end. This substantial number emphasizes the importance of the reviewing phase as well as the complexity of texture annotation due to its variety and the high expressivity of the syntax. Besides correcting possible omissions or typos in the labels, the majority of conflicts ensured a more homogeneous and precise use of the syntax in the whole dataset – in each movements and between them. Hence, most of the discussions between the annotator and the reviewers involved several measures: consecutive ones sharing similar texture, cyclic resurgence of thematic materials (for example in both exposition and recapitulation of sonata form) or comparable textures across pieces. Note that in some cases (like in music analysis in general), it is possible to have different interpretations of the textural layers in a musical passage. The goal of the review here was not to blur those distinct possible analyses but to provide consistent labels that made sense for all.

2.3.3 Syntax knowledge and annotation reproducibility

The annotator followed the syntax recalled in Section 2.1 with the help of a catalog of common textural configurations provided in [7]. As shown in Section 2.3.2, the texture of a bar can sometimes be interpreted in various ways, leading to different labels. The consideration of the neighboring bars can also impact the texture estimation. To limit inconsistency in the labels, we provide annotation guidelines in addition to the dataset. These guidelines typically indicate the order of consideration of the different textural elements when estimating a label. They aim at encouraging a sense of normalization in the annotations, although such a normalization has not been strictly formalized. This aims at helping futur annotators to understand and reproduce the annotation labels. Since the web platform used allows to share several analyses of the same musical piece, divergent labels could be further used ultimately to propose alternative analyses and diversify the dataset for machine learning applications.

3. CORPUS ANALYSIS

3.1 Common textures and layers

Among the 1164 annotated bars, 16.6% include two labels instead of one due to a substantial shift of texture in the middle of the bar (see Section 2.3.1), resulting in 1357 textural configurations. From this set, we can extract 2317 non-empty textural layers. Among the most common layers, we find simple single-voice melodies (M1) which totaled 24.9% of all written layers, followed by melodies with scale motives (M1s, 7.0%, see Figure 1.b) and with parallel motions (M2p, 3.7%), generally doubled at the third (see Figure 1.a) or at the sixth.

Another important family of textural layers encompasses repetitive accompaniment figures like arpeggios, in which notes of the current harmony are played one by one (HS1, here). A famous example is the *Alberti bass*, an idiomatic pattern which alternates notes in low-high-middle-high order (see Figure 1.c and the first half of Figure 2) and is typical of Mozart’s piano music. The combination

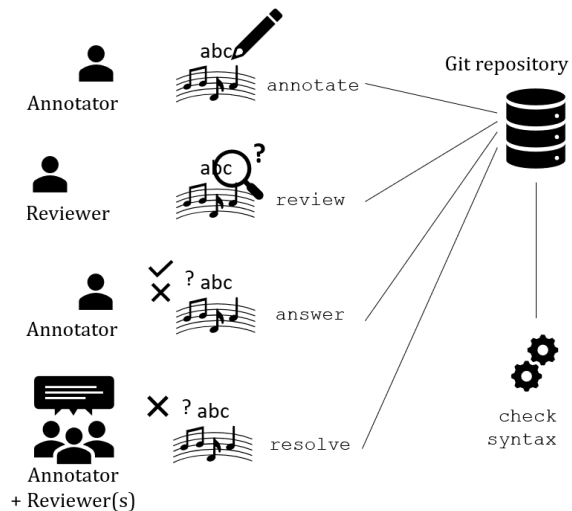


Figure 3. Annotation protocol. Once each measure of a given piece has been annotated (1), these labels are read by a reviewer (2) who can emit a ‘doubt’ on questionable annotations, along with a justification or a proposal of correction. Then, the annotator treats each doubt, by either resolving it or arguing towards another solution (3). Finally, conflicts that remained unresolved are discussed between the annotator and reviewer(s) until a consensus is found (4). At each step of the process, the correctness of the syntax is automatically checked and the changes are committed to a Git repository.

of basic melody M1 and single-voice accompaniments HS1 is the most represented in the dataset, with 7.3% of all annotated labels.

3.2 Textural elements in labels

Symbolic texture labels are highly expressive and rich in information. We call *textural elements* the set of unary attributes, each one notated with a dedicated character, that can appear in our labels. We consider that a bar includes a textural element if it appears at least once in the whole label. For example, a bar includes the textural element *h* (homorhythmy), if the whole texture, or at least one of its layers, is homorhythmic (and therefore annotated with *h*). In the case of two successive textural configurations described in one label (*A*, *B*), the presence of the element on one side is sufficient to consider its presence in the whole label. The textural elements are all listed in Table 2.

Unsurprisingly, a vast majority of the dataset (94.6%⁴) contains a layer with at least one melodic function (*M*). Following on function combinations, 20.2% of labels only contain melodic layers (presence of the element *M* and absence of *H* and *S*). The coincident presence of the three functions concerns 23.5% of the labels and the most common combination is made of melodic and/or harmonic layers without any static (*S*) ones (34.1%) as it is the case in a typical melody plus accompaniment section. Finally, the proportion of measures with harmonic layers (textural

⁴ The proportions of all annotated textures and textural elements are provided in the dataset repository: <http://algonus.fr/data>.

element H) varies between annotated movements, notably according to the tempo: this percentage is 26% higher in slower movements (the second of each full sonata) than in the average of the 6 others (84.4% versus 58.2%).

3.3 Density and diversity

We compute the diversity and the global density of each annotated textural configuration. As detailed in Section 2.1, the diversity corresponds to the number of stacked layers while the global density corresponds to the approximate number of monophonic voices that can be heard simultaneously. Figure 4 puts in relation these two values for the set of annotated labels. The diagonal is assimilated to polyphony where each voice sounds like a new distinct musical idea, an individual layer. On the contrary, the bottom row, with diversity of value 1, corresponds to monophonic texture; thus, any note is contributing to the same unique musical entity. This case is common at the end of structural parts and cadences, typically in homorhythmy (see Figure 1.d). Homophonic textures, for example combinations of a main melody and accompaniment, are found between these two areas. Among the annotated non-empty textural configurations, 29.2% lay in the 2|2 combination of diversity|density, including the common case of melody and single-voice accompaniment, presented earlier. Textures in 2|3 and 2|4 (three or four voices merged into two layers) illustrate denser homophonic variants, in which a melody may be doubled at the third or at the sixth, or an accompaniment made of simultaneous notes (in homorhythmy). The combinations above the diagonal correspond to situations where the number of layers is higher than the number of voices (2|1 or 3|2). This can happen when two distinct lines alternate in antiphony (call & response): several textural layers are perceived whereas their notes do not overlap.

The diversity, which corresponds to the number of layers, rarely exceeds 2. However, Figure 4 shows that a systematic separation into two layers, which follows an intuitive organization of the score content between the two hands of the pianist, would not be sufficiently representative of the corpus.

In other repertoires, the number of voices – the vertical density – could be globally higher. Bach three-voice inventions would be mainly categorized as a continuous 3|3 polyphony, and some works of the Romantic Era including very large chords would be extremely dense vertically. Hence, it is easy to imagine that this textural space, similarly to the one described in [4], is prone to convey strong stylistic content. Moreover, the evolution of texture throughout a single piece of music could be modeled as a trajectory in this space, which offers promising future analytical perspectives.

4. APPLICATION: PREDICTION OF TEXTURAL ELEMENTS

We present in this section preliminary experiments of texture prediction.

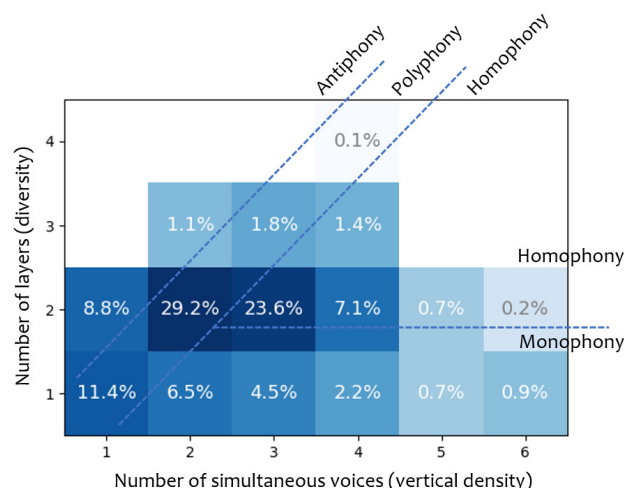


Figure 4. Repartition of textural configurations of the dataset according to their density and diversity. We divided this textural space into areas corresponding to the main texture types evoked in Section 3.3: monophony, homophony, polyphony and antiphony [4,26]. Empty textures (silence) are not taken into account in this figure.

4.1 Symbolic textural descriptors

To facilitate the prediction of textural elements in musical scores, we propose a set of 62 high level features computed on the raw musical score. Some of these were inspired from previous works including [27] and were complemented with original ones especially elaborated for the modeling of high level textural concepts in polyphonic piano scores. The Python implementation of these descriptors is publicly released⁵. Provided code enables to compute the descriptors in *Stream* objects from the Music21 Python library [28] as well as note lists directly stored in TSV (Tab Separated Values) files, as in [14]⁶.

The selected descriptors are computed at different levels of the musical content: pitches, onsets or temporal slices. Temporal slicing is equivalent to the action of `Stream.chordify()` method in the Music21 library, also called salami-slicing [29]. For elements followed by the symbol ‘*’, we compute average, deviation, median, and extremal values.

Pitches. We compute the total number of distinct pitches, the number of pitch classes, the notes duration*, the notes MIDI pitch* and an indicator of pitch reuse.

Onsets. We compute the total number of onsets, the number of simultaneous pitches by onset*, the regularity*, the harmonic intervals and the number of gaps* between pitches – where a *gap* is detected when the interval between two simultaneous notes is larger than a fourth [3].

Slices. We compute the number of slices, the number of simultaneous pitches*, the number of gaps*, the distance pitch ambitus*, the proportion of consonant intervals (minor and major third, perfect fourth and perfect fifth), the proportion of rests and the longest rest.

⁵ <http://algomus.fr/code>

⁶ See https://github.com/DCMLab/mozart_piano_sonatas/tree/main/notes

Textural element	Log. Reg.	Random	All True
M (melodic)	0.912	0.665	0.973
H (harmonic)	0.744	0.545	0.799
S (static)	0.616	0.457	0.599
h (homorhythmy only)	0.673	0.396	0.453
p (parallel motions)	0.572	0.315	0.401
o (octave motions)	0.538	0.211	0.244
h+ (h or p or o)	0.810	0.538	0.708
p+ (p or o)	0.602	0.393	0.479
s (scale motives)	0.363	0.282	0.332
t (sustained notes)	0.669	0.161	0.193
b (oscillations)	0.098	0.116	0.103
r (repeated notes)	0.501	0.183	0.200
_ (sparsity)	0.587	0.258	0.306
, (sequential)	0.520	0.198	0.291

Table 2. F1-scores of the logistic regression models for the prediction of each textural element, compared to – respectively – a uniform random model and a model that always predict the presence of the textural element. F1-scores are averaged on the 9 folds of the cross-validation.

4.2 Description of the models

Different machine learning models were compared to predict the presence of textural elements presented in Section 3.2 from the set of 62 descriptors detailed in Section 4.1. Descriptor values are given to the model as vectors of 62 floats extracted from each measure. The prediction of the presence of each textural element was formulated as a binary classification task leading to a dedicated model for each of the 14 textural elements. Whereas modern data-driven machine learning approaches tend to favor neural networks for their power of abstraction, we stick to simple Logistic Regression models applied on pre-processed high level features as they seemed more adapted to the limited size of our training set and are more easily interpretable. We used Scikit-learn [30] implementation with L-BFGS solver, a maximum of 100 iterations and L2-regularization. Decision Tree classifiers and Support Vector Machine with linear kernel were also tested without significant improvements. In all our models, output classes weights are balanced with respect to their proportion in the dataset. The evaluation criterion is the F1-score, using cross-validation with *leave-one-piece-out* strategy to avoid overfitting due to similarities and repetitions inside movements.

4.3 Results and interpretations

The results are presented in Table 2. We observe that the presence of melody (M) is very well detected (F1-score of 0.912), despite being highly unbalanced in the dataset (94.6% of labels contains at least one layer with a melodic function). Harmonic or static functions are quite more difficult to predict. They were also more difficult to annotate: the determination of these functions also involves a part of subjectivity, despite the efforts made in the reviewing process to ensure the consistency of annotations.

The predictions of notes simultaneities and semblant motions (h, p, o) obtain fair results. We could have ex-

pected better from them, as well as for t (sustained notes) and r (repeated notes), since they only involved rare divergences during reviews, their determination being more straightforward. Furthermore, they are closer to the use of defined descriptors – notably those which are related, for instance, to the number of notes played at the same time or the presence of certain harmonic intervals in onsets or slices. The success of predicting h+ compared to specific relationships h, p or o is meaningful: considering the fact that parallel motions are specific cases of homorhythmy, it seems more practical and sound to focus on this more general case.

The poor results of models to predict oscillation b can be justified by the limited proportion of positive examples in the dataset (around 5.1%). Finally, the difficulty of annotating scale motives (s) in practice is reflected in their prediction result (F1-score of 0.36). From slow descending sequences of neighboring pitches to cells of short-duration notes that share the same repeated contour: many variants of this element can be found in the corpus, some of which were making consensus difficult to reach between annotator and reviewers. This issue can be addressed in the definition of the syntax as the need to refine the criteria and guidelines for the annotation of targeted concepts. The scale patterns sometimes span over wider regions than the duration of a bar, making the descriptors inadequate in this case. Providing a larger context to our models therefore appears as a promising perspective to improve the detection of this textural element.

5. CONCLUSION AND FURTHER WORKS

We provide an open dataset of texture annotations that specifically handles piano classical music. This dataset contains annotations of 9 movements of Mozart’s piano sonatas, all in Sonata Form, and annotation guidelines are provided in order to allow for an easy extension of the dataset by the community. It could be completed with other movements of Mozart’s piano sonatas, but also with more diverse piano music from the Classical Era. An extension of the texture syntax to consider specific textures from other periods would also certainly bring new insights for the study of composition styles.

The dataset can be used for musicological purposes, for instance to study the links between texture and a variety of annotations – harmony, modulations, cadences, phrases and section boundaries – that are common subjects of interest in the MIR community. It also opens perspectives in MIR for computer-aided analysis and composition. As a first application, we implemented a set of textural descriptors and used them for the prediction of textural elements, obtaining encouraging results. An in-depth analysis of the relations between the textural descriptors and textural elements could help improving the prediction, especially on harmonic, static layer detection, or scale motives. This work will also allow to study the evolution of texture in the musical pieces and to better integrate this dimension for automatic generation of music following textural scenarios.

6. ACKNOWLEDGEMENT

We would like to thank Nathalie Hérold for fruitful discussions during the initial phase of the project. We also thank the Algomus team and the anonymous reviewers for their helpful feedback. This research is partly funded by Région Hauts-de-France.

7. REFERENCES

- [1] D. Moreira, “Textural design: a compositional theory for the organization of musical texture,” Ph.D. dissertation, Centro de Letras e Artes, Universidade Federal do Estado do Rio de Janeiro, Rio de Janeiro), 2019.
- [2] N. Hérold, “Timbre et analyse musicale : les possibilités d’intégration du timbre dans l’analyse formelle des œuvres pour piano du dix-neuvième siècle,” in *L’interprétation musicale*. Delatour France, 2012, pp. 79–103.
- [3] Q. R. Nordgren, “A measure of textural patterns and strengths,” *Journal of Music Theory*, vol. 4, no. 1, pp. 19–31, 1960.
- [4] D. Huron, “Characterizing musical textures,” in *International Computer Music Conference (ICMC 1989)*, 1989, pp. 131–134.
- [5] B. Duane, “Texture in eighteenth- and early nineteenth-century string-quartet expositions,” Ph.D. dissertation, Northwestern University, 2012.
- [6] C. P. d. Silva, “Por uma definição unificada de textura musical,” Master’s thesis, Escola de Música, Universidade Federal da Bahia, 2018.
- [7] L. Couturier, L. Bigo, and F. Levé, “Annotating Symbolic Texture in Piano Music: a Formal Syntax,” in *Sound and Music Computing Conference (SMC 2022)*, 2022.
- [8] O. Cífka, U. Şimşekli, and G. Richard, “Groove2Groove: one-shot music style transfer with supervision from synthetic data,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 28, pp. 2638–2650, 2020.
- [9] C. Weiß, M. Mauch, S. Dixon, and M. Müller, “Investigating style evolution of western classical music: A computational approach,” *Musicae Scientiae*, vol. 23, no. 4, pp. 486–507, 2019.
- [10] L. Soum-Fontez, M. Giraud, N. Guimard-Kagan, and F. Levé, “Symbolic textural features and melody/accompaniment detection in string quartets,” in *International Symposium on Computer Music Multidisciplinary Research (CMMR 2021)*, 2021.
- [11] D. Régnier, N. Martin, and L. Bigo, “Identification of rhythm guitar sections in symbolic tablatures,” in *International Society for Music Information Retrieval Conference (ISMIR 2021)*, 2021.
- [12] Z. Wang, D. Wang, Y. Zhang, and G. Xia, “Learning interpretable representation for controllable polyphonic music generation,” in *International Society for Music Information Retrieval Conference (ISMIR)*, 2020.
- [13] J. Zhao and G. Xia, “Accomontage: Accompaniment arrangement via phrase selection and style transfer,” *arXiv preprint arXiv:2108.11213*, 2021.
- [14] J. Hentschel, M. Neuwirth, and M. Rohrmeier, “The annotated mozart sonatas: Score, harmony, and cadence,” *Transactions of the International Society for Music Information Retrieval*, vol. 4, no. 1, 2021.
- [15] J. Devaney, C. Arthur, N. Condit-Schultz, and K. Nisula, “Theme and variation encodings with roman numerals (tavern): A new data set for symbolic music analysis,” in *International Society for Music Information Retrieval Conference (ISMIR 2015)*, 2015.
- [16] M. Giraud, R. Groult, E. Leguy, and F. Levé, “Computational fugue analysis,” *Computer Music Journal*, vol. 39, no. 2, 2015.
- [17] D. R. W. Sears, M. T. Pearce, W. E. Caplin, and S. McAdams, “Simulating melodic and harmonic expectations for tonal cadences using probabilistic models,” *Journal of New Music Research*, vol. 47, no. 1, pp. 29–52, 2017. [Online]. Available: <https://doi.org/10.1080/09298215.2017.1367010>
- [18] P. Allegraud, L. Bigo, L. Feisthauer, M. Giraud, R. Groult, E. Leguy, and F. Levé, “Learning Sonata Form Structure on Mozart’s String Quartets,” *Transactions of the International Society for Music Information Retrieval*, vol. 2, no. 1, pp. 82–96, 2019.
- [19] M. Neuwirth, D. Harasim, F. C. Moss, and M. Rohrmeier, “The Annotated Beethoven Corpus (ABC): A dataset of harmonic analyses of all Beethoven string quartets,” *Frontiers in Digital Humanities*, vol. 5, p. 16, 2018. [Online]. Available: <https://www.frontiersin.org/article/10.3389/fdigh.2018.00016>
- [20] D. Tymoczko, M. Gotham, M. S. Cuthbert, and C. Ariza, “The romantext format: A flexible and standard method for representing roman numeral analyses,” in *International Society for Music Information Retrieval Conference (ISMIR)*, 2019.
- [21] M. Giraud, F. Levé, F. Mercier, M. Rigaudière, and D. Thorez, “Towards modeling texture in symbolic data,” in *International Society for Music Information Retrieval Conference (ISMIR 2014)*, 2014, pp. 59–64.
- [22] A. Tenkanen and F. Gualda, “Detecting changes in musical texture,” in *International Workshop on Machine Learning and Music*, 2008.
- [23] G. Peeters and K. Fort, “Towards a (better) definition of annotated mir corpora,” in *International Society for Music Information Retrieval Conference (ISMIR)*, 2012.

- [24] M. Giraud, R. Groult, and E. Leguy, “Dezrann, a web framework to share music analysis,” in *International Conference on Technologies for Music Notation and Representation (TENOR 2018)*, 2018, pp. 104–110.
- [25] J. Hentschel, F. C. Moss, M. Neuwirth, and M. Rohrmeier, “A semi-automated workflow paradigm for the distributed creation and curation of expert annotations,” in *International Society for Music Information Retrieval Conference (ISMIR)*, 2021.
- [26] B. Benward and M. Saker, *Music in Theory and Practice (Volume 1)*. McGraw-Hill Professional, 2008 (8th ed.).
- [27] C. McKay, J. Cumming, and I. Fujinaga, “JSymbolic 2.2: Extracting features from symbolic music for use in musicological and MIR research,” in *International Society for Music Information Retrieval Conference (ISMIR 2018)*, 2018, pp. 348–354.
- [28] M. S. Cuthbert and C. Ariza, “music21: A toolkit for computer-aided musicology and symbolic music data,” in *International Society for Music Information Retrieval Conference (ISMIR 2010)*, 2010, pp. 637–642.
- [29] C. W. White and I. Quinn, “The yale-classical archives corpus,” *Empirical Musicology Review*, vol. 11, no. 1, 2016.
- [30] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, “Scikit-learn: Machine learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

VIOLIN ETUDES: A COMPREHENSIVE DATASET FOR F_0 ESTIMATION AND PERFORMANCE ANALYSIS

Nazif Can Tamer Pedro Ramoneda Xavier Serra

Music Technology Group, Universitat Pompeu Fabra, Barcelona

nazifcan.tamer@upf.edu

ABSTRACT

Violin performance analysis requires accurate and robust f_0 estimates to give feedback on the playing accuracy. Despite the recent advancements in data-driven f_0 estimators, their application to performance analysis remains a challenge due to style-specific and dataset-induced biases. In this paper, we address this problem by introducing *Violin Etudes*, a 27.8-hours violin performance dataset constructed with domain knowledge in instrument pedagogy and a novel automatic f_0 -labeling paradigm. Experimental results on unseen datasets show that the CREPE f_0 estimator trained on *Violin Etudes* outperforms the widely-used pre-trained version trained on multiple manually-labeled datasets. Further preliminary findings suggest that (i) existing data-driven f_0 estimators may overfit to equal temperament, and (ii) iterative re-labeling regularized by our novel *Constrained Harmonic Resynthesis* method can simultaneously enhance datasets and f_0 estimators. Our dataset curation methodology is easily scalable to other instruments owing to the quantity of pedagogical data online. It also supports a range of MIR research directions thanks to the performance difficulty labels from educational institutions.

1. INTRODUCTION

Accurate f_0 tracking is fundamental for violin performance analysis due to the prime role of intonation in violin mastery. Musicians regard intonation and pitch accuracy as the most important criteria for assessing a string performance [1], and most of the previous work on violin performance analysis also focus on vibrato and intonation [2–5]. A study on intonation patterns of artist-level violinists [4] found that highly-regarded musicians deviate significantly from equal-temperament while remaining coherent in their intonation preferences in the close vicinity of just-noticeable difference (95% confidence intervals within just 6 cents). Another study found that violinists’ intonation can be better approximated by other tuning systems, e.g., Pythagorean, rather than the standard equal-temperament [2]. From an engineering perspective, into-

nation analysis is reliable only if the f_0 estimates are more precise than the intonation consistency of the player. Thus, these studies imply that an f_0 estimator suitable for violin performance analysis needs to conform to a frequency precision higher than 6 cents and remain consistent irrespective of the player’s deviation from equal temperament.

Alongside frequency precision, an f_0 estimator has to fulfill two more requirements for reliable performance analysis: (i) robustness to octave errors that result from the complex frequency response of the violin body and (ii) high temporal precision that can handle fast string crossings common in violin performance. However, in the current paradigm, there is a trade-off in complying with these two necessities: Most of the f_0 estimators leverage temporal post-processing stages in order to eliminate octave errors at the expense of temporal precision (e.g., Viterbi for the f_0 estimator of PRAAT [6], pYIN [7], and CREPE [8]; a custom filtering for Melodia [9]). Moreover, the Viterbi implementations of some f_0 estimators are even more restrictive: pYIN does not allow a jump bigger than 2.5 semitones between consecutive frames, and that number is 2.4 semitones for CREPE. These restrictions are detrimental to violin performance analysis, as it is common to see very fast and abrupt string crossings in violin repertoire (e.g., 21 semitone jumps in Figure 2).

Despite the above-mentioned problems, monophonic f_0 estimation is considered a mature task in MIR literature, mainly owing to the high accuracies reported in benchmark datasets. Data-driven f_0 estimators claim 96-99% accuracies in datasets such as MIR-1k [10] and MDB-stem-synth [11], which are all highly restricted in terms of performance virtuosity. Although these numbers seem promising in theory, this is not what we found when we use these f_0 estimators in the real-life analysis of advanced-level violin performances. In this paper, to better address the real-world needs of performance analysis using data-driven methods, we introduce the *Violin Etudes*¹, a 27-hour large-scale monophonic dataset comprised of pedagogical violin performances by professional violin players. We also provide our methodologies for dataset curation and automatic f_0 labeling and show the strength of our approach by outperforming the pre-trained version of CREPE f_0 estimator on its own train data.



© N. C. Tamer, P. Ramoneda, and X. Serra. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0).
Attribution: N. C. Tamer, P. Ramoneda, and X. Serra, “Violin Etudes: A Comprehensive Dataset for f_0 Estimation and Performance Analysis”, in *Proc. of the 23rd Int. Society for Music Information Retrieval Conf.*, Bengaluru, India, 2022.

¹ *Violin Etudes* dataset is available for research purposes on <https://doi.org/10.5281/zenodo.6564408>

2. RELATED WORK

2.1 Monophonic f_0 Estimation

F_0 estimation is an essential step for many tasks in music and speech processing. Early f_0 estimators were handcrafted methods [6, 7, 12–19], but today data-driven methods [8, 20–24] are preferred following the general trend of deep learning. Both in handcrafted and data-driven techniques, f_0 estimation literature can be classified into two main approaches: time-domain and frequency-domain. Time-domain approaches used to utilize techniques such as auto-correlation during the times of handcrafted signal processing [6, 7, 12–15], while nowadays end-to-end deep learning methods directly use the time domain signal as their input [8, 20, 23, 24]. Frequency-domain approaches, on the other hand, were used to apply spectral analysis and select spectral peaks with signal processing [16–19]; but recently data-driven frequency-domain approaches require a spectrogram derivative as their input signal [21, 22]. A final important concept to mention here is self-supervised f_0 estimation, adopted by SPICE [22] and DDSP-inv [24]. Although the f_0 estimation strategy used in this paper is based on supervised learning, our combined procedure of iterative f_0 -labeling is analogous to a self-supervision applied post-training.

2.2 Automatic f_0 -labeling

To guarantee the f_0 label correctness for semi-automatically labeled datasets, Salamon et al. created an analysis-synthesis method [11] which forces the audio to represent any f_0 error in the annotation. The method was used for creating the resynthesized MDB-mf0-synth and Bach10-mf0-synth datasets. However, applying this method to unlabeled datasets is yet to be explored.

2.3 Large Scale Performance Datasets

Collecting large-scale datasets from community platforms has long been an important data source for research on action recognition and multimodal learning. The data collection procedure most often involves searching for keywords on YouTube, generally without collecting further metadata other than the query term itself. To our knowledge, the three main instrument performance datasets collected in this fashion are aimed at self-supervised source separation and spatial localization: MUSIC dataset [25] consists of solo and duet performances spanning over 11 instrument categories, including 53 solo violin performances. They extended the dataset for the task of video-to-audio synthesis with additional solo performances and constructed the MUSIC-Extra-Solo dataset [26] which include 213 solo violin performances, some of which include backing track. In a similar but more constrained scenario, the Solos dataset [27] is formed by searching for 13 classical instruments and the word 'audition' on YouTube and includes 66 solo violin performances in approximately 400 minutes. However, none of these datasets provide any label or metadata on the musical content: whether they are monophonic, include different renditions of the same score, the player or

recording conditions, or the supposed difficulty of the performance. The only label is the instrument name.

A more constrained and informed large-scale data collection endeavor is the GiantMIDI-piano dataset [28] where the authors queried YouTube with the piano repertoire they collected from International Music Score Library Project (IMSLP) and later transcribed the audio into MIDI. By including more meta-data such as the composers, work titles, and style, the dataset is more suitable for musical analysis. However, it is also susceptible to MIDI transcription errors. The *Violin Etudes* dataset introduced in this work has similarities to this controlled approach, but in an even more restricted scenario: By collecting query words from the pedagogical repertoire, we have control over the performance difficulty. By keeping track of the performer and providing multiple renditions for the same etudes, we control the expressivity and recording conditions. Last but not least, by manually curating and removing the works, including double stops and chords, we ensure monophony and control the timbre.

3. VIOLIN ETUDES DATASET

From the 16th century onwards, music pedagogy has been creating teaching curricula that guide the students from the very early stages of their journey to the professional level. Inspired by how humans learn an instrument, we present the *Violin Etudes* dataset, the first large-scale MIR dataset rooted in instrument pedagogy. Etudes and caprices form the backbone of the traditional violin method and are defined as study pieces presenting "a technical problem or challenge in the context of a musical setting" [29]. Alongside being vital for education, these pedagogical materials have an unparalleled potential as datasets for intelligent systems, especially MIR applications. They most often come with inherent difficulty labels and are organized in human curricula, which can be used for autonomous learning. Composers often provide textual descriptions on the purpose of the study and techniques involved (e.g., Figure 3), and they are still actively researched by instrument pedagogues on their technical content e.g. [30–32]. They are most often for solo instruments, which is favorable for signal processing. And most importantly, there are hundreds of instrument teachers actively recording their reference performances of their teaching material, which supports the much-needed data for deep learning applications.

3.1 Data Collection

While previous large-scale YouTube data collection endeavors mainly focus on data for self-supervision, we opted for a more controlled approach in curating the material and stored metadata that would be used as ground truth in many topics such as expressive performance analysis and performance difficulty analysis. The individual violin methods are selected from the standard violin curriculum by a trained violinist, but interested readers are encouraged to search for '*violin (or flute/trumpet/piano...) etudes*' to see how easily they can create similar lists, e.g., [33–36].