

madmom: a new Python Audio and Music Signal Processing Library

Sebastian Böck^{†*}, Filip Korzeniowski[‡], Jan Schlüter[‡], Florian Krebs[‡], Gerhard Widmer^{†‡}

[†] Department of Computational Perception, Johannes Kepler University Linz, Austria

[‡] Austrian Research Institute for Artificial Intelligence (OFAI), Vienna, Austria

ABSTRACT

In this paper, we present *madmom*, an open-source audio processing and music information retrieval (MIR) library written in Python. *madmom* features a concise, *NumPy*-compatible, object oriented design with simple calling conventions and sensible default values for all parameters, which facilitates fast prototyping of MIR applications. **Prototypes can be seamlessly converted into callable processing pipelines through *madmom*'s concept of *Processors*, callable objects that run transparently on multiple cores.** Processors can also be serialised, saved, and re-run to allow results to be easily reproduced anywhere.

Apart from low-level audio processing, *madmom* puts emphasis on musically meaningful high-level features. Many of these incorporate machine learning techniques and *madmom* provides a module that implements some in MIR commonly used methods such as **hidden Markov models and neural networks**. Additionally, *madmom* comes with several state-of-the-art MIR algorithms for onset detection, beat, downbeat and meter tracking, tempo estimation, and piano transcription. These can easily be incorporated into bigger MIR systems or run as stand-alone programs.

1. INTRODUCTION

Music information retrieval (MIR) has become an emerging research area over the last 15 years. Especially audio-based MIR has become more and more important, since the amount of available audio data in the last years exploded beyond being manageable manually.

Most state-of-the-art audio-based MIR algorithms consist of two components: First, **low-level features are extracted from the audio signal (*feature extraction* stage), and then the features are analysed (*feature analysis* stage) to retrieve the requested information.** Most current MIR systems incorporate machine learning algorithms in the feature analysis stage, with neural networks currently being the most popular and successful ones [17, 11, 2, 3].

Numerous software libraries have been proposed over the years to facilitate research and development of applications in MIR. Some libraries concentrate on low-level feature extraction from audio signals, such as *Marsyas* [18], *YAAFE* [14] and *openSMILE* [9]. Others also include higher level

feature extraction such as onset and beat detection as for example in the *MIRtoolbox* [13], *Essentia* [6] and *librosa* [15]. However, to our knowledge, there exist no library that also includes **machine learning components** (except *Marsyas* [18], which contains two classifiers), although machine learning components are crucial in current MIR applications.

Therefore, we propose *madmom*, a library that incorporates low-level feature extraction and high-level feature analysis based on machine learning methods. This allows the construction of the full processing chain within a single software framework, making it possible to build standalone programs without any dependency on other machine learning frameworks. Moreover, *madmom* comes with several state-of-the-art systems including their trained models, for example for onset detection [17, 8], tempo estimation [3], beat estimation [2, 11], downbeat estimation [4], and piano transcription.

madmom is written in Python, which has become the language of choice for scientific computing for many people due to its free availability and its simplicity to use. The code is released under BSD license and pre-trained models are released under the CC BY-NC-SA 4.0 license.

1.1 Design and Functionality

1.1.1 Object-oriented programming

madmom follows an object-oriented programming (OOP) approach. We encapsulate everything in objects which are often designed as subclasses of *NumPy*'s *ndarray*, offering all array handling routines inherited from *NumPy* [19] with additional functionality. This compactly bundles data and meta-data (e.g., a *Spectrogram* and its *frame rate*) and simplifies meta-data handling for the user.

1.1.2 Rapid prototyping

madmom aims at minimizing the turnaround time from a research idea to a software prototype. To this means, object instantiation is made as simple as possible: E.g., a *log Mel-spectrogram* object can be instantiated with one line of code by only providing the path to an audio file. *madmom* automatically creates all the objects in between using sensible default values.

1.1.3 Simple conversion into runnable programs

Once an audio processing algorithm is prototyped, the complete workflow should be easily transformed into a runnable standalone program with a consistent calling interface. This is implemented using *madmom*'s concept of *Processors*. For details we refer to Section 2.

*To whom correspondence should be addressed; Email: sebastian.boeck@jku.at

1.1.4 Machine learning integration

We aim at a seamless integration of machine learning methods without the need of any third party modules. We limit ourselves to testing capabilities (applying pre-trained models), since it is impossible to keep up with newly emerging training methods in the various machine learning domains. Models that have been trained in an external library should be easily be converted to an internal *madmom* model format.

1.1.5 State-of-the-art features

Many existing libraries provide a huge variety of low-level features but few musically meaningful high-level features. *madmom* tries to close this gap by offering high-quality state-of-the-art feature extractors for downbeats, beats, onsets, tempo, piano transcription, etc..

1.1.6 Reproducible research

In order to foster reproducible research, we want to be able to save and load the specific settings used to obtain the results for a certain experiment. In *madmom* this is implemented using Python’s own *pickle* functionality which allows to save an entire processing chain (including all settings) to a file.

1.1.7 Few dependencies

madmom is built on top of three excellent and wide-spread libraries: *NumPy* [19] provides all the array handling sub-routines for *madmom*’s *data classes*. *SciPy* [10] provides optimised routines for the fast Fourier transform (FFT), linear algebra operations and sparse matrix representations. Finally, *Cython* [1] is used to speed up time critical parts of the library by automatically generating C code from a Python-like syntax and then compiling and linking it into extensions which can be transparently used from within Python. These libraries are the only installation and runtime dependencies of *madmom* besides the *Python* standard library itself, supported in version 2.7 as well as 3.3 and newer.

1.1.8 Multi-core capability

We designed *madmom* to be able to exploit the multi-core capabilities of modern computer architectures, by providing function to run several programs or *Processors* in parallel.

1.1.9 Extensive documentation

All source code files contain thorough documentation following the *NumPy* format. The complete API reference, instruction on how to build and install the library, as well as interactive *IPython* [16] notebooks can be found online at <http://madmom.readthedocs.io>. The documentation is build automatically with *Sphinx*¹.

1.1.10 Open development process

We follow an open development process and the source code and documentation of our project is publicly available on GitHub: <http://github.com/CPJKU/madmom>. To maintain high code quality, we use continuous integration testing via TravisCI², code quality tests via QuantifiedCode³, and test coverage via Coveralls⁴.

¹<http://www.sphinx-doc.org>

²<http://www.travis-ci.org>

³<http://www.quantifiedcode.com>

⁴<http://www.coveralls.io>

2. LIBRARY DESCRIPTION

In this section, we will describe the overall architecture of *madmom*, its packages as well as the provided standalone programs.

madmom’s main API is composed of classes, but much of the functionality is implemented as functions (in turn used internally by the classes). This way, *madmom* offers the ‘best of both worlds’: concise interfaces exposed through classes, and detailed access to functionality through functions. In general, the classes can be split in two different types: the so called *data classes* and *processor classes*.

Data classes represent data entities such as audio signals or spectrograms. They are implemented as subclasses of *NumPy*’s *ndarray*, and thus offer all array handling routines inherited directly from *NumPy* (e.g. transposing or saving the data to file in either binary or human readable format). These classes are enriched by additional attributes and expose additional functionality via methods.

Processor classes exclusively store information on how to process data, i.e., how to transform one data class into another (e.g., from an (audio-) *Signal* into a *Spectrogram*). In order to build chains of transformations, each data class has its corresponding processor class, which implements this transformation. This enables a simple and fast conversion of algorithm prototypes to callable processing pipelines.

2.1 Packages

The library is split into several packages, grouped by functionality. For a detailed description including examples of usage please refer to the library’s documentation.

2.1.1 *madmom.processors*

Processors are one of the fundamental building blocks of *madmom*. Each *Processor* accepts a number of processing parameters and must provide a *process* method, which takes the data to be processed as its only argument and defines the processing functionality of the *Processor*. An *OutputProcessor* extends this scheme by accepting a second argument which defines the output and can thus be used to write the output of an algorithm to a file. All *Processors* are callable, making it easy to use them interchangeably with normal functions. Further, the *Processor* class provides methods for saving and loading any *Processor* to a file – including all parameters – using Python’s own *pickle* library. This facilitates the reproducibility of an experiment.

Multiple *Processors* can be combined into a processing chain using a ***SequentialProcessor* or *ParallelProcessor***, which either execute the chain sequentially or in parallel, using multiple CPU cores if available.

2.1.2 *madmom.audio*

The *madmom.audio* package includes basic audio signal processing and “low-level” functionality.

The *Signal* and *FramedSignal* classes are used to load an audio signal and chop it into (overlapping) frames. Following *madmom*’s automatic instantiation approach, both classes can be instantiated from any object up the instantiation hierarchy – including a simple file name. *madmom* supports almost any existing audio and video formats – provided *ffmpeg*⁵ is installed – and transparently converts sample rates and number of channels if needed.

⁵<http://www.ffmpeg.org>

Signal is a subclass of *ndarray* with additional attributes like *sample rate* or *number of channels*. *FramedSignal* supports float hop sizes, making it possible to build systems with an arbitrary frame rate – independently of the signal’s sample rate – and ensures that all frames are temporally aligned, even if computed with different frame sizes.

The *ShortTimeFourierTransform* and *Spectrogram* classes represent the complex valued STFT and magnitudes respectively. They are the key classes for spectral audio analysis and provide windowing, automatic circular shifting (for correct phase) and zero-padding. Both are independent of the data type (integer or float) of the underlying *Signal*, resulting in spectrograms of the same value range. A *Spectrogram* can be filtered with a *Filterbank* (e.g. Mel, Bark, logarithmic), which in turn can be parametrised to reduce the dimensionality or transform the spectrogram into a logarithmically spaced pitch representation closely following the auditory model of the human ear. *madmom* also provides standard *MFCC* and *Chroma* features.

2.1.3 madmom.features

The *madmom.features* package includes “high-level” functionality which are related to certain MIR tasks, such as beat tracking or onset detection. *madmom*’s focus is to provide musically meaningful and descriptive features rather than a vast number of low to mid-level features. At the time of writing this paper, *madmom* contains state-of-the-art features for onset detection, beat and downbeat tracking, rhythm pattern analysis, tempo estimation and piano transcription.

All features are implemented as *Processors* without a corresponding *data class*. Users can thus use the provided functionality and build algorithms on top of these features. For most of the features, *madmom* also provides stand-alone programs with a consistent calling interface to process audio files (see Section 2.2).

2.1.4 madmom.evaluation

All features come with code for evaluation. The implemented metrics are those commonly found in the literature of the respective field.

2.1.5 madmom.ml

Most of today’s top performing music analysis algorithms incorporate machine learning, with neural networks being the most universal and successful ones at the moment. *madmom* includes Python implementations of commonly used machine learning techniques, namely *Gaussian Mixture Models* (GMM), *Hidden Markov Models* (HMM), and different types of *neural networks* (NN), including feed forward, recurrent, and convolutional layers, various activation functions and special purpose units such as *long short-term memory* (LSTM).

madmom provides functionality to use these techniques without any dependencies on third-party modules, but does not contain training algorithms. This decision was made on purpose since the library’s main focus is on applying machine learning techniques to MIR, rather than providing an extensive set of learning techniques. However, trained models can be easily converted to be compatible with *madmom*, since neural network layers usually are simply defined as a set of weights, biases and an activation function they apply to the input data.

2.1.6 madmom.models

madmom comes with a set of pre-trained models which are distributed under a Creative Commons attribution non-commercial share-alike license, i.e. they can be freely used for research purposes as long as derivative works are distributed under the same license. *madmom* uses the exact same mechanism to save and load the models it uses for *Processors* to be pickled.

2.2 Standalone Programs

madmom comes with a set of standalone programs, covering many areas of MIR. These programs are simple wrappers around the functionality provided by the *madmom.features* package, and provide a simple and easy to use command line interface. They are implemented as *Processors* and can operate either in *single* or *batch* mode, processing single or multiple input files, respectively. Additionally all programs can be *pickled*, serialising all parameters in a way that the program can be executed later on with the exact same settings.

Table 1 lists selected programs included in the library with the performance achieved at the annual *Music Information Retrieval Evaluation eXchange* (MIREX)⁶, where MIR algorithms are compared on hidden test datasets. We aggregated the results of all years (2006-2015), i.e., a rank 1 means that the algorithm is the best performing one of all submissions from 2006 until present. The outstanding results in Table 1 highlight the state-of-the-art features *madmom* provides.

Table 1: **Ranks of the programs included in *madmom* for the MIREX evaluations, results aggregated over all years (2006-2015).** Asterisks indicate pending submissions.

Program	Task	Year	Rank
CNNOnsetDetector [17]	onset	2013	1
OnsetDetector [8]	onset	2013	2
BeatTracker [5]	beat MCK	2015	1
DBNBeatTracker [2]	beat SMC	2015	1
CRFBeatDetector [11]	beat MAZ	2015	1
GMMPatternTracker [12]	downbeat	2015	2
DBNDownBeatTracker [4]	downbeat	2016	*
TempoDetector [3]	tempo	2015	1
CRFChordTranscriptor	chord	2016	*
PianoTranscriptor	transcription	2016	*

3. CONCLUSION

This paper gave a short introduction into *madmom*, its design principles and library structure. Up-to-date information on functionality can be found in the project’s online documentation at <https://madmom.readthedocs.io> and source code repository at <https://github.com/CPJKU/madmom>.

Future work aims at including a streaming mode, i.e. providing online real-time processing of audio signals in a memory efficient way instead of processing whole audio files at a

⁶<http://www.music-ir.org/mirex/wiki/>

time. In addition, we will gradually extend the set of features and algorithms, as well as add tools to automatically convert models that have been trained with popular machine learning libraries such as Lasagne [7].

4. ACKNOWLEDGMENTS

This work is supported by the European Union Seventh Framework Programme FP7 / 2007-2013 through the GiantSteps project (grant agreement no. 610591) and the Austrian Science Fund (FWF) project Z159.

5. REFERENCES

- [1] S. Behnel, R. Bradshaw, C. Citro, L. Dalcin, D. Seljebotn, and K. Smith. Cython: The Best of Both Worlds. *Computing in Science Engineering*, 13(2), 2011.
- [2] S. Böck, F. Krebs, and G. Widmer. A Multi-model Approach to Beat Tracking considering Heterogeneous Music Styles. In *Proc. of the 15th Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2014.
- [3] S. Böck, F. Krebs, and G. Widmer. Accurate Tempo Estimation based on Recurrent Neural Networks and Resonating Comb Filters. In *Proc. of the 16th Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2015.
- [4] S. Böck, F. Krebs, and G. Widmer. Joint Beat and Downbeat Tracking with Recurrent Neural Networks. In *Proc. of the 17th Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2016.
- [5] S. Böck and M. Schedl. Enhanced Beat Tracking with Context-Aware Neural Networks. In *Proc. of the 14th Int. Conf. on Digital Audio Effects (DAFx)*, 2011.
- [6] D. Bogdanov, N. Wack, E. Gómez, S. Gulati, P. Herrera, O. Mayor, G. Roma, J. Salamon, J. Zapata, and X. Serra. Essentia: an open source library for sound and music analysis. In *In Proc. of ACM Multimedia*, 2013.
- [7] S. Dieleman, J. Schlüter, C. Raffel, E. Olson, S. K. Sønderby, D. Nouri, E. Battenberg, A. van den Oord, et al. Lasagne: First release., 2015.
- [8] F. Eyben, S. Böck, B. Schuller, and A. Graves. Universal Onset Detection with Bidirectional Long Short-Term Memory Neural Networks. In *Proc. of the 11th Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2010.
- [9] F. Eyben, F. Weninger, F. Gross, and B. Schuller. Recent Developments in openSMILE, the Munich Open-Source Multimedia Feature Extractor. In *In Proc. of ACM Multimedia*, Barcelona, Spain, 2013.
- [10] E. Jones, T. Oliphant, P. Peterson, et al. SciPy: Open source scientific tools for Python, 2001–. [Online; accessed 2016-05-20].
- [11] F. Korzeniowski, S. Böck, and G. Widmer. Probabilistic Extraction of Beat Positions from a Beat Activation Function. In *Proc. of the 15th Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2014.
- [12] F. Krebs, S. Böck, and G. Widmer. Rhythmic Pattern Modeling for Beat and Downbeat Tracking in Musical Audio. In *Proc. of the 14th Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2013.
- [13] O. Lartillot and P. Toivainen. A Matlab toolbox for musical feature extraction from audio. In *Proc. of the 10th Int. Conf. on Digital Audio Effects (DAFx)*, 2007.
- [14] B. Mathieu, S. Essid, T. Fillon, and J. Prado. YAAFE, an easy to use and efficient audio feature extraction software. In *Proc. of the 11th Int. Society for Music Information Retrieval Conf. (ISMIR)*, 2010.
- [15] B. McFee, C. Raffel, D. Liang, D. Ellis, M. McVicar, E. Battenberg, and O. Nieto. librosa: Audio and Music Signal Analysis in Python. In *Proc. of the 14th Python in Science Conf. (SCIPY)*, 2015.
- [16] F. Pérez and B. E. Granger. IPython: A System for Interactive Scientific Computing. *Computing in Science Engineering*, 9(3), 2007.
- [17] J. Schlüter and S. Böck. Musical Onset Detection with Convolutional Neural Networks. In *Proceedings of the 6th International Workshop on Machine Learning and Music*, Prague, Czech Republic, 2013.
- [18] G. Tzanetakis and P. Cook. MARSYAS: a framework for audio analysis. *Organised Sound*, 4, 2000.
- [19] S. van der Walt, S. C. Colbert, and G. Varoquaux. The NumPy Array: A Structure for Efficient Numerical Computation. *Computing in Science Engineering*, 13(2), 2011.