

Retrieving Music Semantics from Optical Music Recognition by Machine Translation

Martha E. Thomae
McGill University
martha.thomaeelias@mail.mcgill.ca

Antonio Ríos-Vila
University of Alicante
arios@dlsi.ua.es

Jorge Calvo-Zaragoza
University of Alicante
jcalvo@dlsi.ua.es

David Rizo
University of Alicante
drizo@dlsi.ua.es

José M. Iñesta
University of Alicante
inesta@dlsi.ua.es

Abstract

In this paper, we apply machine translation techniques to solve one of the central problems in the field of optical music recognition: extracting the semantics of a sequence of music characters. So far, this problem has been approached through heuristics and grammars, which are not generalizable solutions. We borrowed the seq2seq model and the attention mechanism from machine translation to address this issue. Given its example-based learning, the model proposed is meant to apply to different notations provided there is enough training data. The model was tested on the PriMuS dataset of common Western music notation incipits. Its performance was satisfactory for the vast majority of examples, flawlessly extracting the musical meaning of 85% of the incipits in the test set—mapping correctly series of accidentals into key signatures, pairs of digits into time signatures, combinations of digits and rests into multi-measure rests, detecting implicit accidentals, etc.

Introduction

We present a machine learning-based approach to retrieve the semantics of a sequence of (graphic) music symbols, which constitutes a central problem in the field of Optical Music Recognition (OMR). OMR is the process of converting the digital image of a score into a symbolic file encoding the music content of that score. The traditional OMR workflow consists of four stages: preprocessing, symbol recognition, music reconstruction, and music encoding. The third stage, music reconstruction, must retrieve the actual musical meaning of the graphical symbols recognized in the previous stage. So far, the models proposed to solve this problem are based on rules [1] or grammars [2, 3], which prevents their use in other contexts (e.g., in notation systems other than the one for which they were implemented). For high scalability, we propose a machine learning-based approach which learns the semantics of a particular notation system by providing the model with enough training examples.

We use an encoding introduced by [4] to represent the graphical and semantic information obtained by the second and third stages of the OMR workflow, respectively. This encoding provides an intermediate representation that, during the music encoding stage of the OMR process, becomes a well-established music format, such as MusicXML, MEI, or `**kern`.

Background

Agnostic and Semantic Encodings of Sequences

In MEC'2017, [5] presented the concept of agnostic and semantic sequential representations of a music score. The agnostic encoding represents the output of the music symbol recognition stage of the OMR, where we only have the graphical information about the symbols (their shapes and positions) and no musical meaning. The agnostic representation is a sequential encoding of the graphical symbols in the score (Figure 1b). Each token in the sequence encodes two types of information: the label of the symbol (e.g., C clef, quarter note, half note, sharp) and its vertical position within the staff (e.g., third line, fourth space). On the other hand, the semantic

representation is a sequential encoding of symbols in a score, which includes their musical meaning (Figure 1c). Translating an agnostic sequence into a semantic one involves several tasks, including re-interpreting a series of accidentals into a key signature and parsing the position of the notes in the staff into pitch values.

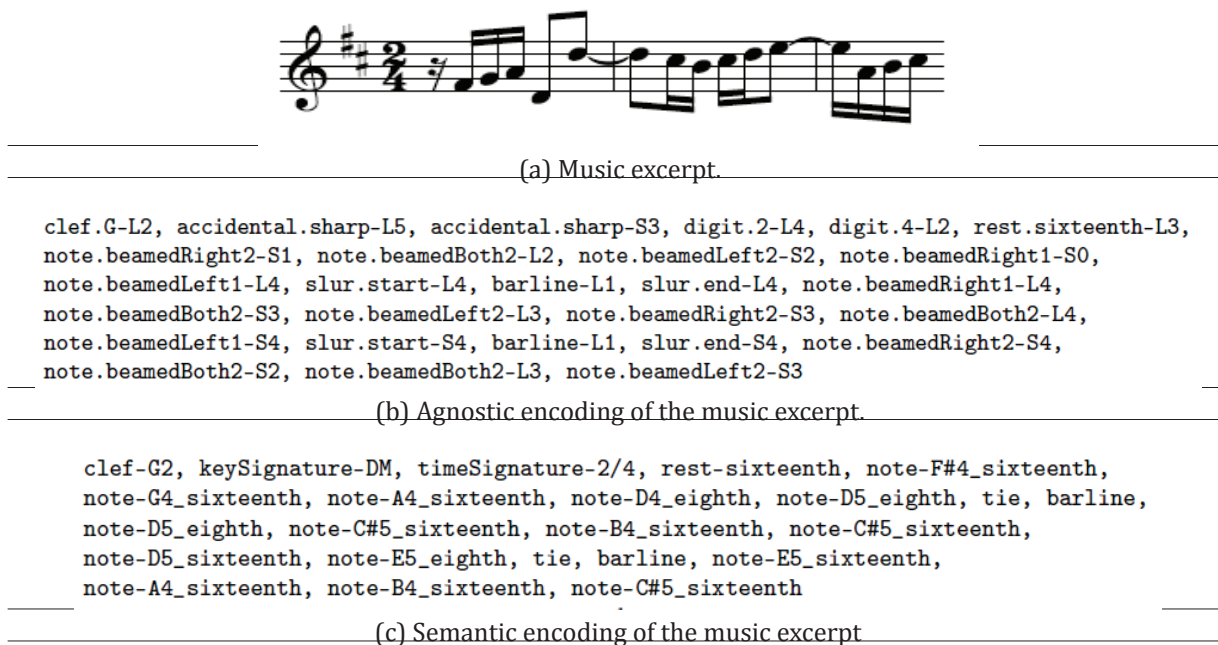


Figure 1: Example of the agnostic and semantic encoding of a musical excerpt [4].

In [4], it was shown that sequential encodings are suitable for converting a digital image into either an agnostic or a semantic representation without human-encoded rules, with more robust results in the agnostic [6]. In this paper, we implement a machine translator that takes the agnostic representation of a sequence of notes in the staff and generates its corresponding semantic representation, in order to take advantage of the performance of the agnostic case for OMR.

Translation Model Description

The main task of the model is to translate an agnostic sequence into its corresponding semantic sequence. We used a “seq2seq” model, first introduced by [7] for machine translation. A seq2seq model consists of two parts, an encoder that maps the input sequence (in this case, the agnostic sequence) onto a fixed-dimension vector, and a decoder that builds the target sequence (here, the semantic sequence) from that vector. We added an attention mechanism, which has been used to improve the translation results by selectively focusing on parts of the input sentence during translation [8]. The attention mechanism allows us to visualize which tokens (graphical symbols) of the agnostic sequence affect the translated tokens of the semantic sequence. In other words, it can show us what the model is paying attention to when translating (Figure 2).

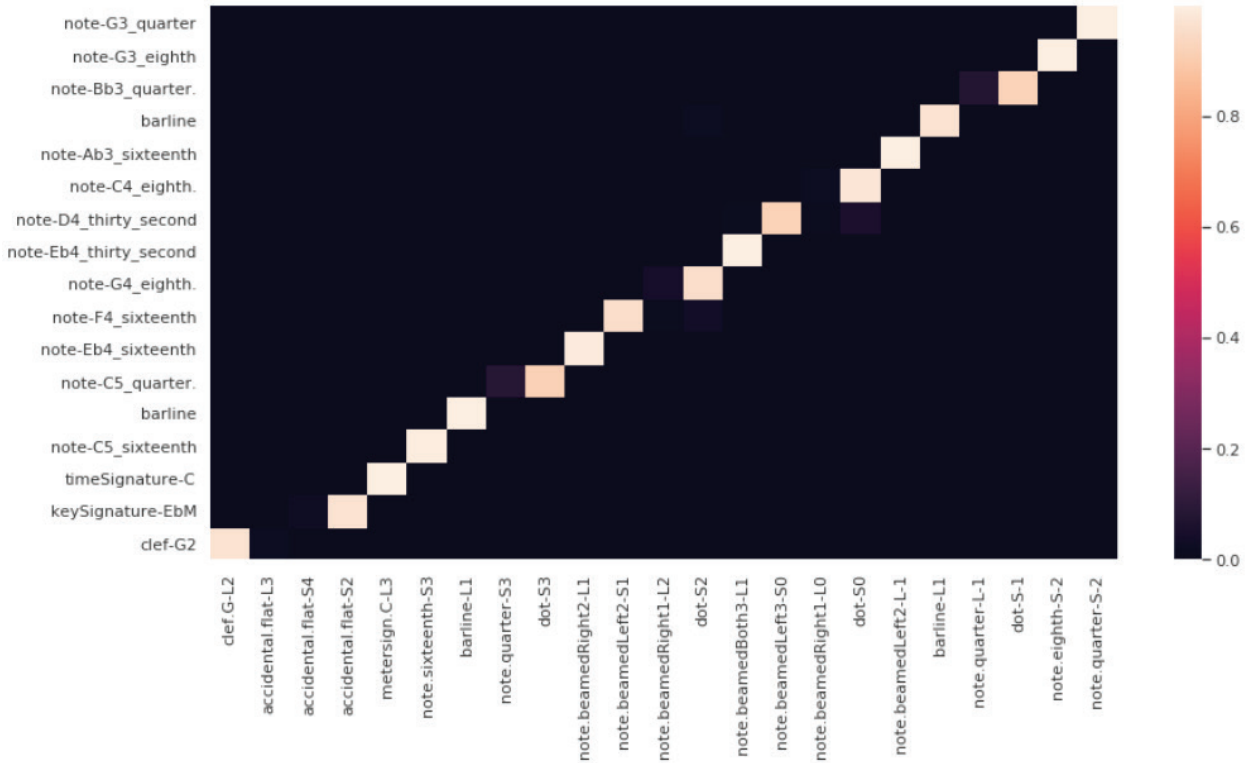


Figure 2: Attention matrix of the model when translating an agnostic sequence (horizontal axis) into a semantic sequence (vertical axis).

Experiment and discussion

We tested this model’s performance on the Printed Images of Music Staves (PrIMuS) dataset [4]. The PrIMuS dataset consists of 87,678 music incipits from RISM encoded in a variety of formats, including the agnostic and semantic representations mentioned above. We used an 80:10:10 split for training, validation, and testing.

We evaluated the model using the edit distance, which measures the number of operations (in terms of insertion, deletion, and substitution of tokens) needed for two strings to match. Given an agnostic sentence, the edit distance was computed between the corresponding semantic sequence in the dataset and the translated sequence obtained. The model flawlessly extracted the musical meaning of 85% of the agnostic sentences in the test set, correctly identifying key signatures, time signatures, multi-measure rests, dotted notes, and notes affected by notated or implied accidentals (see Figures 3 and 4).



(a) Beginning of one of the incipits in the test set.

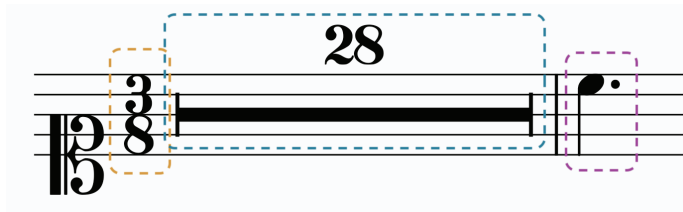
clef.G-L2 accidental.flat-L3 accidental.flat-S4 accidental.flat-S2
 metersign.C/-L3 note.half-L1

(b) Agnostic encoding of the music excerpt.

clef-G2 keySignature-EbM timeSignature-C/ note-Eb4_half

(c) Semantic encoding generated by the model.

Figure 3: Example of the translation of key signatures (green) and implicit accidentals (purple) by the model.



(a) Beginning of one of the incipits in the test set.

clef.C-L1 digit.3-L4 digit.8-L2 digit.2-S5 digit.8-S5 multirest-L3
barline-L1 note.quarter-S4 dot-S4

(b) Agnostic encoding of the music excerpt.

clef-C1 timeSignature-3/8 multirest-28 barline note-C5_quarter.

(c) Semantic encoding generated by the model.

Figure 4: Example of the translation of time signatures (yellow), multi-measure rests (blue), and dotted notes (purple) by the model.

According to the edit distance values obtained, for 7% of the test sentences, only one error was made in the translation. One example of this is the sequence shown in Figure 2, where the last dotted note (coming from the agnostic tokens “note.quarter-L-1 dot-S-1”) is wrongly translated into a Bb instead of Ab. As seen in the attention matrix of Figure 2, when translating dotted notes, the translator pays more attention to the dot token than to the preceding note token. Similar to dotted notes, the model also pays considerably more attention to the last accidental in a series of accidentals at the moment of parsing the key signature.

As can be seen from Figure 6, most error-free sentences lie on the average-sentence-length region (the 18–33 interval with the highest data concentration in Figure 5). Analyzing some of the examples with the highest edit distance values, some of the patterns found are the presence of a clef change, after which the translator’s performance consistently drops for all following tokens; and long sentences (of more than 35 tokens, lying in the right end of the distribution shown in Figure 5).

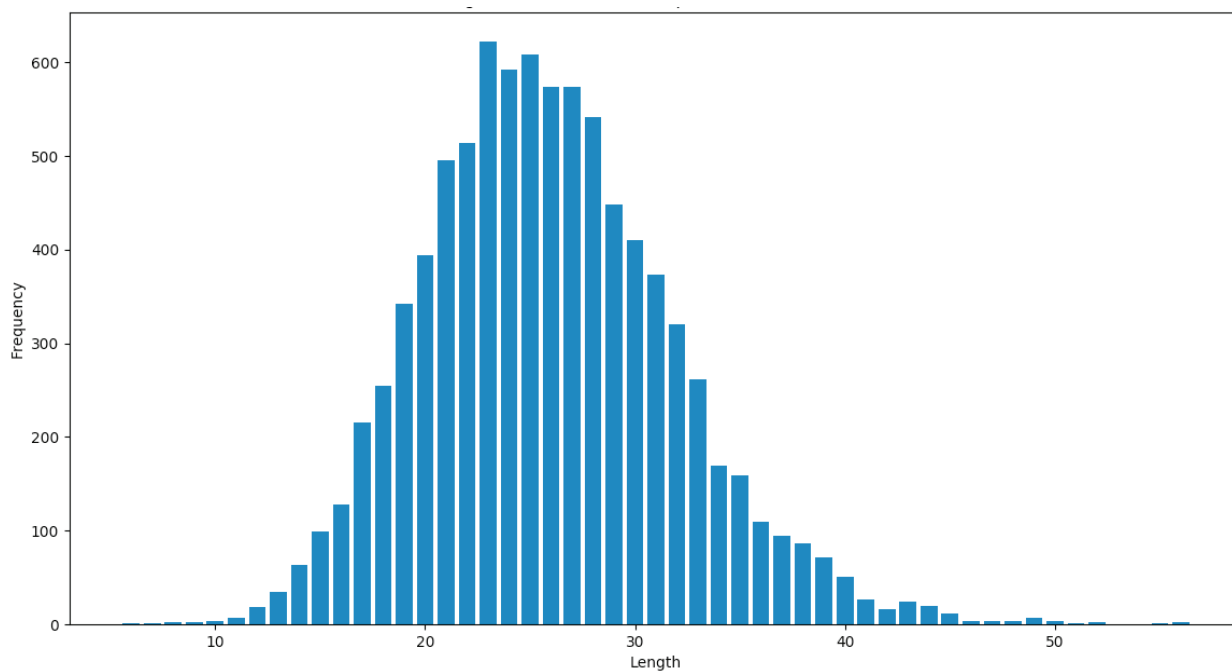


Figure 5: Length of the semantic sequences in the test set.

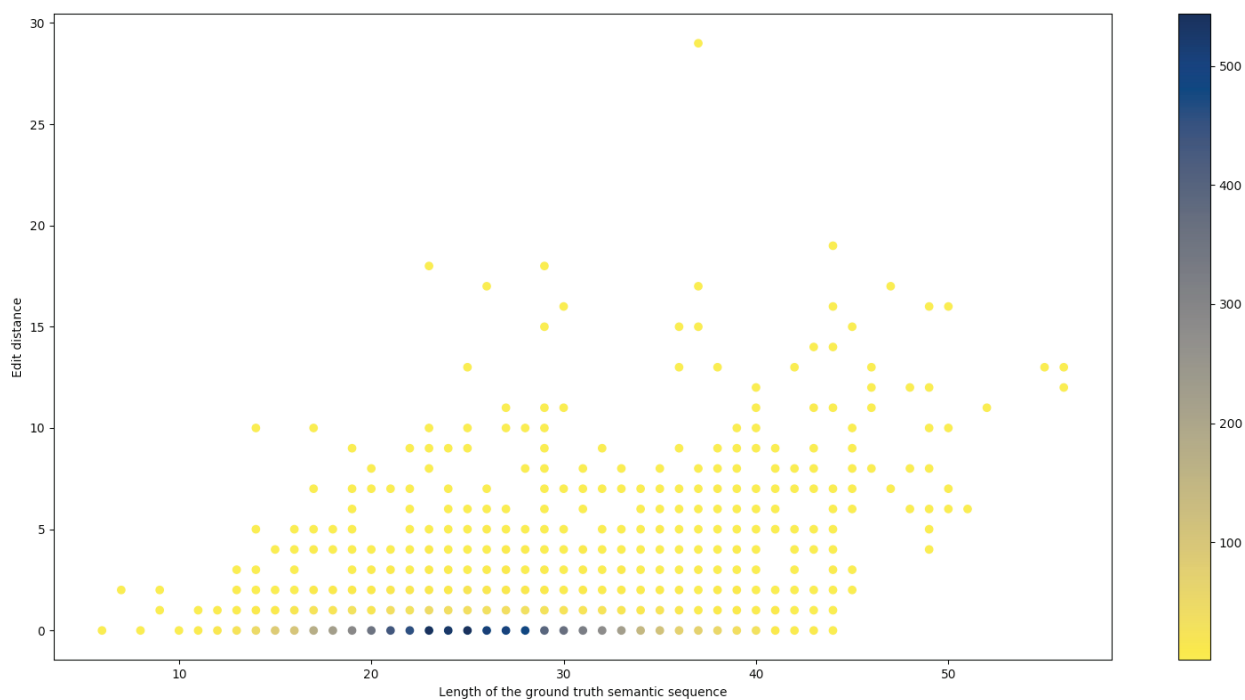


Figure 6: Color density plot of the edit distance of all sentences in the test set. The color bar indicates the frequency of a particular (sentence length, edit distance) pair.

Conclusion

Given its example-based learning, the model we propose is meant to apply to different notation systems provided there is enough training data. The performance in the PrIMuS dataset was satisfactory for the vast majority of examples. However, we plan to improve the attention mechanism to enhance its performance before tackling notation systems with more complex semantics (e.g., mensural notation). Other future work includes the substitution of the semantic representation by ****kern**, a well-established music encoding format that also encodes the music symbols sequentially for each staff. The advantages of ****kern** over the semantic encoding are that the former allows for rendering the encoded sequence in Verovio, and that there is technology already available to obtain more complex formats (e.g., MEI or MusicXML) from ****kern** [9].

Acknowledgements

This work is supported by the Spanish Ministry HISPAMUS project TIN2017-86576-R, partially funded by the EU, and by CIRMMT's Inter-Centre Research Exchange Funding and McGill's Graduate Mobility Award.

Works cited

- [1] Rossant, Florence, and Isabelle Bloch. "Robust and Adaptive OMR System Including Fuzzy Modeling, Fusion of Musical Rules, and Possible Error Detection" *EURASIP Journal on Advances in Signal Processing*, no. 1 (2007), <https://doi.org/10.1155/2007/81541>.
- [2] Couasnon, Bertrand. "DMOS: A Generic Document Recognition Method, Application to an Automatic Generator of Musical Scores, Mathematical Formulae and Table Structures Recognition Systems" in *Proceedings of the 6th International Conference on Document Analysis and Recognition (ICDAR 2001)*, 215–20, <https://doi.org/10.1109/ICDAR.2001.953786>.
- [3] Szwoch, Mariusz. "Guido: A Musical Score Recognition System" in *Proceedings of the 9th International Conference on Document Analysis and Recognition (ICDAR 2007)*, vol. 2, 809–13, <https://doi.org/10.1109/ICDAR.2007.4377027>.
- [4] Calvo-Zaragoza, Jorge, and David Rizo. "End-to-End Neural Optical Music Recognition of Monophonic Scores" *Applied Sciences* 8, no. 4 (2018), 606–29, <https://doi.org/10.3390/app8040606>.
- [5] Rizo, David, Jorge Calvo-Zaragoza, José M. Iñesta, and Ichiro Fujinaga. "About Agnostic Representation of Musical Documents for Optical Music Recognition" presented at the Music Encoding Conference, Tours, France, May 16-19, 2017.
- [6] Calvo-Zaragoza, Jorge, and David Rizo. "Camera-PrIMuS: Neural End-to-End Optical Music Recognition on Realistic Monophonic Scores" in *Proceedings of the 19th International Society for Music Information Retrieval Conference (ISMIR 2018)*, 248–55.
- [7] Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. "Sequence to Sequence Learning with Neural Networks" in *Proceedings of the 27th International Conference on Neural Information Processing Systems (NIPS'14)*, 3104–12, <https://www.arxiv-vanity.com/papers/1409.3215/>
- [8] Luong, Minh-Thang, Hieu Pham, and Christopher D. Manning. "Effective Approaches to Attention-Based Neural Machine Translation" *ArXiv:1508.04025* (2015), <http://arxiv.org/abs/1508.04025>
- [9] Sapp, Craig Stuart. "Verovio Humdrum Viewer" presented at the Music Encoding Conference, Tours, France, May 16-19, 2017.