

Understanding poetry using natural language processing tools: a survey

Mirella De Sisto^{1,*}, Laura Hernández-Lorenzo ², Javier De la Rosa³, Salvador Ros², Elena González-Blanco⁴

¹Tilburg University, The Netherlands

²National Distance Education University, Spain

³National Library of Norway, Norway

⁴IE University, Spain

*Corresponding author. Tilburg University, Warandelaan 2, 5037 AB Tilburg, The Netherlands. E-mail: m.desisto@tilburguniversity.edu

Abstract

Analyzing poetry with automatic tools has great potential for improving verse-related research. Over the last few decades, this field has expanded notably and a large number of tools aiming at analyzing various aspects of poetry have been developed. However, the concrete connection between these tools and traditional scholars investigating poetry and metrics is often missing. The purpose of this article is to bridge this gap by providing a comprehensive survey of the automatic poetry analysis tools available for European languages. The tools are described and classified according to the language for which they are primarily developed, and to their functionalities and purpose. Particular attention is given to those that have open-source code or provide an online version with the same functionality. Combining more traditional research with these tools has clear advantages: it provides the opportunity to address theoretical questions with the support of large amounts of data; also, it allows for the development of new and diversified approaches.

Keywords: Poetry analysis; Natural Language Processing; Computational Literary Studies.

1. Introduction

The field of automatic poetry analysis is expanding, allowing for the application of automatic tools in verse investigation. This enables new quantitative approaches to address traditional and theoretical questions.

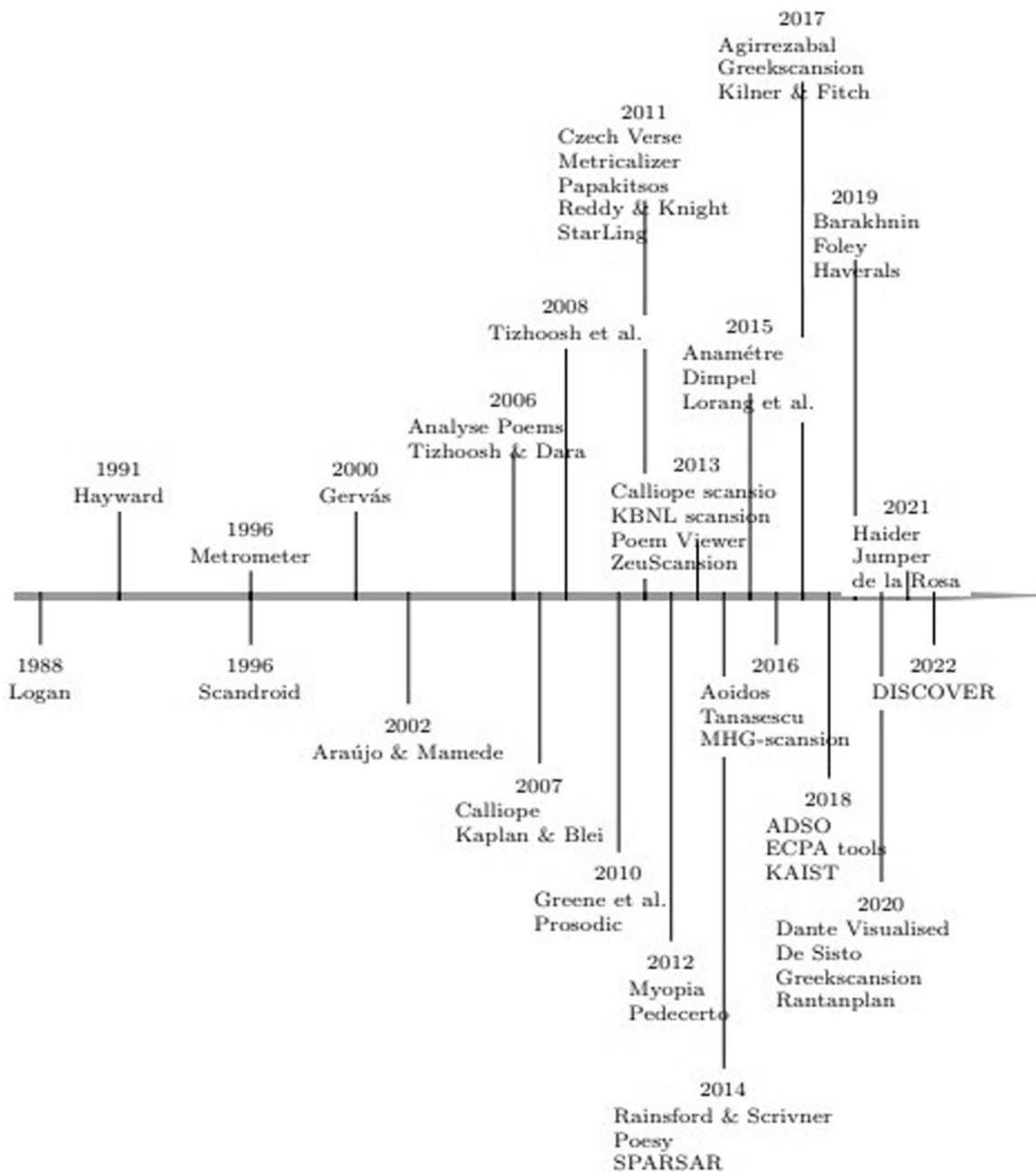
The advancement of Natural Language Processing and text analysis technologies has significantly contributed to the growth of automatic poetry analysis. Multiple and varied tools have been developed, each addressing different aspects of verse with specific goals, such as metrical annotation or the visualization of poetic features. While most tools focus on analyzing English, an increasing number of programs now target other languages and language varieties.

These tools empower scholars to address theoretical questions with the support of large-scale data analysis. However, the development of computational tools often remains disconnected from the theoretical research that could benefit from them. The present article aims to bridge this gap by providing a survey of currently available computational tools for automatic poetry analysis. In this survey, we focus on tools for the

analysis of the structural characteristics of poetry of European languages.

Since most tools are not publicly available, a comprehensive qualitative comparison was not feasible. In addition, such a comparison would require expertise in multiple poetic traditions, which goes beyond the possibilities of a single research group. Therefore, this article aims to offer an exhaustive overview of available tools for European languages, while empirical testing falls beyond its scope.

The article is structured as follows: Section 2 briefly introduces automatic poetry analysis and describes its advantages for poetry-related research. Section 3, along with its subsections, classifies the tools presented in this article based on language (3), purpose and features (4), and algorithmic nature (i.e. rule-based or data-driven) (5). A concluding section summarizes the key points of the review and highlights the importance of bridging the gap between traditional and computational versification studies. [Appendix Table A1](#) provides a detailed overview of the classification and includes the links to the tools' source code and graphical user interface.



tools for poetry visualization, which aim to enhance close-reading by highlighting various details of the line, for example, recurrent sounds in a line.

The following sections provide a survey of the tools based on this classification and offer some details about the specific focus and features of each tool, including a brief explanation of their functioning. Some

tools span two groups and are classified according to their primary function.

4.1 Tools for metrical annotation

The first and the largest macro-group considered in the present survey consists of tools that focus on metrical annotation. These tools are primarily concerned

One of the first metrical annotation tools developed is the tool by [Logan \(1988\)](#), which utilizes sound pattern detection for metrical annotation. In addition to syllabification and stress assignment, this tool identifies the recurring sounds within the lines and their patterns across lines. It also considers the distinctive features of sounds (i.e. + or—labial, + or—coronal, etc.) both on a per-line basis and for the entire poem. By comparing the generated metrical annotation with the metrical pattern of an iambic pentameter, the tool calculates metrical complexity.

Scandroid (Hartman 1996) is a tool for scansion of iambic and anapestic poetry (and has a bias for these meters as a result). If a poem is written in a different meter, the program forces its lines into one of the two possible scansions. When analyzing a line, it first tries to determine whether it is formed by iambic or anapestic feet and, subsequently, calculates line length. Initially, it divides the line into words and then syllabifies words and assigns stress to syllables based on a built-in dictionary and on general principles. The tool can be downloaded from its website, which also provides a tutorial, a manual and the source code.²

Built on Scandroid, Calliope (McAleese 2007) implements metrical scansion by using syntactic information, namely stress assignment derived from syntax, and phonological scansion based on syntax. It uses the Antelope Natural Language Process Parser (Chaumartin 2013) for parsing the text into syntactic elements. Stress is assigned using a lexicon and phonological grouping is performed. The output is parsed by a phonological scansion algorithm, and, subsequently, consistency of the result is assessed. The tool has been tested with iambic and trochaic meters, and with free verse.

Prosodic (Heuser *et al.* 2010) is a tool that focuses on metrical and phonological parsing of poetry. Its scansion process begins with tokenization of the text into words, which are then converted into stressed, syllabified phonetic transcriptions. By default, Prosodic uses the CMU pronunciation dictionary (Weide 1998) to retrieve information about stress patterns and syllable boundaries. For the metrical parsing of words not available in the CMU pronunciation dictionary, a text-to-speech (TTS) engine⁵ is used. Each line is assigned its best available metrical parsing based on a set of constraints, following the principles of Optimality Theory (Prince and Smolensky 1993), which can be

metrical parsing and annotation procedures depend on Prosodic. The tool is available in a code repository.⁹

The combination of metrical annotation and TTS technology has led to the creation of SPARSAR (Delmonte 2014). SPARSAR is a rule-based system aiming to represent and analyze poetic devices in English and Italian. However, it is only evaluated, though, with English corpora. The system utilizes a modified version of VENESSE, a semantically oriented NLP pipeline for deep text understanding. SPARSAR has three modules: the first operates at the sentence level and provides syntactic and semantic analysis; the second is a rule-based system that converts each poem into phonetic characters, divides words into stressed and unstressed syllables, and computes rhyming schemes at line and stanza levels. The third module uses these data to generate a parameterized version of the poem which can be read by a TTS engine. SPARSAR's output includes a visualization and allows for the comparison of poems from different authors using the Pearson correlation coefficient. In addition, it includes a technology for TTS reading of poems. The tool can be downloaded from the project's website,¹⁰ which, besides documentation, also provides examples of analyzed poems and of the generated graphics. However, the platform does not appear to be regularly updated.

Tanasescu, Paget, and Inkpen (2016) developed a tool for poetry classification partially based on Scandroid. The tool uses prosodic and phonetic features to identify meter and rhyme. Lexical stress detection is based on the Scandroid code. For classification, they use a decision tree algorithm (J48) and a bootstrap aggregating machine learning technique (Breiman 1996). The source data of this tool is the poetry available on the online platforms Poetry Foundation¹¹ and sonnets.org.

Similarly, Agirrezabal and his colleagues (Agirrezabal 2017; Agirrezabal, Alegria, and Hulden 2017) use Recurrent Neural Networks (RNN; Schmidhuber 1989) with bidirectional Long Short-Term Memory (bi-LSTM) cells (Hochreiter and Schmidhuber 1997) and Conditional Random Fields (CRF; Lafferty, Abdul-Rahman, and Pereira 2001) to automatically scan poetry in three languages, that is, English, Spanish, and Basque. The tool tokenizes words using the IXA pipeline¹² (Agerri, Bermudez, and Rigau 2014) and then tags POS. Stress is assigned according to an implementation of Groves’ rules (Groves 1998) and using the CMU pronunciation dictionary (Weide 1998) in combination with NETalk (Sejnowski and Rosenberg 1987). In the case that a word is not found in the dictionaries, its stress pattern is guessed by searching for similar entries in the NETalk lexicon and in the CMU dictionary. Finally, syllables are grouped into feet. Like ZeuScansion (Agirrezabal *et al.* 2013), this tool also employed the For

Better For Verse corpus for English for training and the Corpus se Siglo de Oro (Navarro-Colorado 2015) for Spanish.

Kim, Zhao, and Zelalem (2018) from the KAIST research group developed an automatic tool for poetry classification. Their model extracts syllables by referring to the WebCelex dataset and scans lines based on the manually annotated training corpus; aside from special characters (e.g. colons and full stops) and syllables, word boundary is encoded as a character on the syllable. The tool also determines whether a line is metrical or not. The data used for developing this tool is, as for the tool by Tasnasescu *et al.* (2016), the poetry available on Poetry Foundation. The source code is publicly available in a code repository.

Although these tools are all developed to provide metrical annotations of English poetry, they do present some significant differences. The major distinction is between tools developed based on a particular meter (e.g. iambic) and those which are not developed for a specific poetic form. Tools developed by [Logan \(1988\)](#), [Hayward \(1991\)](#), [Hartman \(1996\)](#), and [Tanasescu, Paget, and Inkpen \(2016\)](#) have a bias for iambic meter, making their annotation more accurate with iambic poetry than with other metrical patterns. The same could count for [Greene *et al.* \(Greene, Bodrumlu, and Knight 2010\)](#), since it was trained on Shakespeare’s sonnets. On the other hand, [Prosodic \(Heuser *et al.* 2010\)](#) has a set of constraints which can be adjusted, and the other tools are not based on a specific meter.

Another difference concerns rhyme detection. Only Logan (1988), AnalysePoems (Plamondon 2006), Poesy (Algee-Hewitt *et al.* 2014), and Tanasescu, Paget, and Inkpen (2016) provide information about rhyme patterns.

4.1.2 Spanish

Gervás' tool (Gervás 2000) for Spanish poetry is also one of the earliest automatic annotation tools. It presents a logic programming application for the metrical annotation of Spanish poems. The system uses Definite Clause Grammars to model word syllabification, as well as additional predicates to define synaloepha, syllable count, and rhyme. The logic programming implementation, using SWI Prolog, treats verses as lists of words, and words and syllables as lists of characters. A predicate analysis retrieves from each verse the number of syllables, a list of syllables in the verse, a list of positions of stressed syllables, and the verse's rhyme (which is used to define the rhyme pattern of the stanza). It also includes a system vocabulary, referred to as 'a database of facts' (see Gervás 2000: 1335), containing word information. If data are missing for a given word, a set of rules is used to fill them. After checking if the data are available or need to

be produced, a Word Syllable Parser module is used to parse character types, character groups and syllables, locate the stressed syllables of a word, and extract the rhymes of a word. Additionally, a Metrical Syllable Counter counts the metrical syllables in a verse. The Metrical Syllable Counter consists of a set of Prolog predicates that operate recursively over the list of words in a verse, with another predicate to evaluate if the verse is a hendecasyllable. However, the author acknowledges that this application makes significant errors with the character “y” and has problems with some diphthongs.

On the rule-based side, the ADSO Scansion system (Navarro-Colorado 2018) is a tool for Spanish poetry that uses POS tags for syllabification and stress assignment. First, POS tags (using FreeLing),¹³ Padró and Stanilovsky (2012) are added to the words of every line in a poem. Then, each word is split into syllables, the POS tag is used to classify syllables as stressed or unstressed. As the system is prepared for hendecasyllable detection, a syllable counter is applied. Only when the verse length is shorter or longer than eleven syllables, diaeresis and synaloephas are applied. ADSO’s source code is available in a code repository. The tool was developed on the *Corpus de Siglo de Oro*¹⁴ (Navarro-Colorado 2015), the only available manually annotated corpus of Spanish poetry, which contains sonnets from the 16th and the 17th centuries. The poems are in XML format, meeting TEI standards. Additionally, each line has its metrical pattern indicated in the annotation. The metrical annotation is obtained with the tool developed by Navarro-Colorado and manually verified by experts in some cases.

Rantanplan (De la Rosa *et al.* 2020b) is a rule-based tool for Spanish poetry. It relies on POS tags and its own syllabification algorithm to assign stress to words and lines, in a similar way to the ADSO Scansion system (Navarro-Colorado 2018). Like ADSO, it is based on the *Corpus de Siglo de Oro* (Navarro-Colorado 2015). However, Rantanplan applies all possible

synaloephas and synaeresis at the syllable level before determining the metrical pattern. If the expected metrical length has a greater value than the calculated one, synaloephas and synaeresis are undone in different tests until the discrepancy is resolved. Rantanplan works with different types of verses, not only hendecasyllables, as is the case with ADSO. Rantanplan has its source code available in a code repository. In addition, its scansion can be analyzed and visualized on the web-based interface application PoetryLab (De la Rosa *et al.* 2020a) (see Figure 3). This platform also displays Named Entity Recognition provided by the tool HisMeTag¹⁵ and enjambement detection provided by the Python library JollyJumper.¹⁶ Rantanplan is intended to complement the data repository, Averell,¹⁷ which contributes to reducing the current limitations of quantitative research on poetry by standardizing poetic corpora into an easily downloadable and accessible format.

A different approach is employed by Jumper (Marco *et al.* 2021), which addresses the metrical annotation of Spanish poetry without performing syllabification but by counting syllables just by counting the vowels of each word. Stress is assigned using the same method: the position of graphically accented vowels is defined by counting vowels, and, in case of no graphical stress, Spanish orthographic rules are applied. The verse level module is composed by three submodules: the first module counts syllables by counting vowels and assigns stress on the line level; the second module applies line compensation (e.g. synaloepha) to lines which are longer than eleven syllables; finally the third module addresses metric ambiguities. The poem level module analyzes the scansion performed by the line level module for each line evaluates whether only one line length is attested in the poem (in the case of a fixed-meter poem), and whether some line lengths are more frequent than others. Jumper’s code is available in a code repository.



Figure 4. Metrical annotation with Metricalizer.

Finally, [De la Rosa et al. \(2021\)](#) conducted metrical pattern annotation of Spanish, English, and German poetry using transformer-based models. They utilize monolingual and multilingual BERT-base and RoBERTa models ([Wolf et al. 2019](#)) to predict line stress patterns. The findings demonstrate that these models effectively capture sufficient structural information for poetry scansion. Specifically, they perform well for Spanish in a monolingual setting, and the cross-lingual transfer, when training the model on the three languages, yields to particularly promising results for English and German. The training, test, and evaluation datasets consist of the following: for Spanish, a subset of the *Corpus de Siglo de Oro* ([Navarro-Colorado 2015](#)); for English, a subset of the For Better For Verse corpus; for German, the manually annotated corpus from [Haider et al. \(2020\)](#). The German corpus comprises 158 poems from 1575 to 1936 which are annotated in terms of syllable stress, foot boundaries, caesuras, and line main accent. The original lines are available on the online platform *Antikoerperchen Lyrik Datenbank*.¹⁸

4.1.3 Ancient Greek

A tool for Ancient Greek hexameter, developed by [Papakitsos \(2011\)](#) is one of the scansion machines that emerged in the early 2010s. The primary goal of the tool is to divide the line in six parts (feet or meters) and determine their length. The tool has four modules: the first module detects syllables; the second module assigns length based on the number of syllables; the third module identifies foot structure; and the final module returns a Boolean (true or false) value indicating the success of the scansion. When the Boolean is true, the ‘identity of the syllables and feet is known’ ([Papakitsos 2011](#)).

Another tool developed for Ancient Greek is `Greek_scansion`, created by [Conser \(2017\)](#). This tool analyzes prosody in Ancient Greek poetry and prose. It offers two versions: a generic one not tied to a specific meter, and another specifically designed for iambic trimeter. The tool provides line scansion as a list of syllable markers, namely “long,” “short,” and “unknown.” The syllabification module is based on James Tauber’s Greek accentuation Python library¹⁹ and can potentially work with a different syllabifier. Unfortunately, the source code, available in a code repository, does not appear to be updated regularly.

Greek_scansion (Schumann 2020; Schumann *et al.* 2022), homonym of the tool by Conser (Conser 2017), was implemented based on the syllabification method proposed by Papakitsos (Papakitsos 2011). Like the tool by Papakitsos this rule-based tool also focuses on annotating Ancient Greek hexameter. It uses a class ruleset containing linguistic rules for syllabification

and syllable count, and a verse class that contains the verse model, including information about the poetic text such as syllabification, syllable count resulting from hyphenation, and the hexameter scheme. The source code for this tool is available in a code repository.

4.1.4 Czech

In the expanding field of poetry scansion, the *Versologie* research group at the Institute of Czech Literature has been particularly prolific, developing a set of tools for poetry analysis. These data are available in the Database of Czech Metres (DCM; Ibrahim and Plecháč 2011), which provide metrical and stanzaic descriptions of poems contained in the Corpus of Czech Verse. Poetry is analyzed in terms of prosody and metrical structure through a set of algorithms. With DCM, poetry can be statistically evaluated in terms of verse meter, verse length, line ending, metrical pattern, stanza and rhyme scheme, and types of fixed metrical forms (e.g. sonnet). In addition to DCM, the *Versologie* research group also developed a tool called *Gunstick* for automatic analysis of rhyme. Its annotations are available in the Database of Czech rhymes.²⁰ The set of tools to analyze the material available in the Corpus of Czech Verse and the DCM are available on the Institute of Czech Literature's website.

4.1.5 Dutch, Middle Dutch, and Renaissance Dutch

The KBNL scansion generator (Koppelaar and Van Oostendorp 2013; Van Oostendorp 2014) is an annotation tool for Dutch poetry. It uses information about stressed syllables from the CELEX Lexical database for Dutch²¹ to determine which syllables are stressed in the words that form a line. Subsequently, the algorithm attempts to match the line with an ideal template (in their example, a fixed number of iambs). An online platform is for the KBNL scansion generator (Koppelaar and Van Oostendorp 2013; Van Oostendorp 2014), offering comprehensive information about the tool, a demo for users to scan poetry, and a link to the code repository.

A Middle Dutch syllabifier (Haverals, Karsdorp, and Kestemont 2019) has been developed using a bi-LSTM RNN. The tool provides syllabification of Middle Dutch texts. Each word is inputted to the model at the character level, and the model provides the syllable grouping. In Haverals (2020) word stress is implemented. The model was trained on the Corpus Van Reenen-Mulder (Reenen and Mulder 1993). This tool can be applied to both verse and prose, and in Haverals (2020), it was used to define the metrical pattern of Dutch medieval poetry. The source code is available in a code repository. In addition, the input data for the machine learning algorithm in various

formats and a pdf-file with further information about the dataset are also publicly accessible on Zenodo.

De Sisto (2020) developed a tool similar to that by Haverals, Karsdorp, and Kestemont (2019), for the metrical annotation of a Dutch Renaissance poetry corpus. Like the Middle Dutch syllabifier and other data-driven tools, the tool for Renaissance Dutch uses a bi-LSTM RNN. The tool is based on two models: one is trained on a list of the most frequent words syllabified and annotated in terms of stress (with a focus on polysyllabic words); the other on random lines from the corpus with stress pattern annotations. These two models are then combined to extract the metrical pattern of a line based on the lexical stress assignment of the words it contains. Line annotation only considers syllabification and disregards word boundary. The Dutch Renaissance poetry corpus annotated by De Sisto (2020) and the source plain texts are available in a code repository.²² The original source texts used to build the corpus are available at *De Digitale Bibliotheek voor de Nederlandse Letteren* (DBNL).²³

4.1.6 French

Metrometer (Beaudouin and Yvon 1996) is a tool for French poetry which was developed in the mid-nineties. As mentioned earlier, in the case of the French poetic tradition, metrical annotation is based on the number of syllables per line. Metrometer uses phonemic transcriptions, phonological and syntactic information to calculate the number of syllables and of metrical positions per line while considering the occurrence of metrical phenomena such as synaloepha and liaison. The input is converted into phonemic transcriptions enriched with information such as syntactic components. The tool combines a syntactic analyzer with a system of grapheme-to-phoneme transcription. The rules of the tool were developed based on alexandrines by P. Corneille and J. Racine, hence it is more suitable for annotating this specific form of versification.

Reddy and Knight (2011) propose an unsupervised and language-independent model for rhyme schemes detection. The model infers rhyme schemes based on word repetition within stanzas. Their study tests this model on English and French data. English and French poetry data were used in the rhyme recognition process. The English data are from the open access part of the corpus used by Sonderegger (2011) consisting of 12,000 stanzas from the 15th to 20th centuries. The French data are from the ARTFL project²⁴ containing 3,000 stanzas.

Lastly, another tool for French poetry is *Anamètre* (Delente and Renault 2015). The tool is used for the analysis of French poetry and drama from the early seventeenth to the early twentieth century. The meter

calculation program employs a deterministic approach based on Cornulier's metrical method (Cornulier 1982). Developed in Python, this tool identifies vowels which are part of the metrical pattern and from those calculates the meter based on the prosodic and morpho-syntactic constraints of the French language, which vary diachronically and depend on the context. Additionally, *Anamètre* uses a lexical resource of grammatical syllables and employs algorithms to address issues of ambiguity by considering left and right context.

4.1.7 German and Middle High German

Metricalizer (Bobenhausen 2011; Bobenhausen and Hammerich 2015) is a rule-based tool for metrical analysis of German poetry. It detects words, line, rhyme, and stanza structures in a text and provides word syllabification. The metrical analysis is based on a rule-based algorithm that incorporates prosodic and morphological information. In addition, the tool can identify metrical complexity by detecting discrepancies between prosodic and metrical structures. The results produced by Metricalizer can be saved in TEI-XML format. An online platform is available for accessing Metricalizer. In addition to the user interface, the website offers an analyzed poetry corpus from the *Freiburg Anthologie*.²⁵ The corpus analysis provides a rich overview of the data that can be generated with the tool. The poems, organized by author, are initially displayed as plain text and can be analyzed and presented with their corresponding annotations (according to metrical pattern, prosodic level, word boundary, and complexity). The database allows for other options, such as visualizing the transcription and comparing poems. An annotated sample from the corpus from the *Freiburg Anthologie* available on the Metricalizer website is shown in Fig. 4.

Haider (2021) conducted experiments to test how different models learn meter on two large corpora of English and German poetry. After preprocessing the corpora, different models were compared. The results demonstrated that bi-LSTM-CRF models with syllable embeddings outperform the baseline CRF model and BERT models. Notably, the experiments yielded interesting findings regarding caesura: the use of POS information supported the correct identification of caesura, which highlights its connection to poetic rhythm and its dependency on syntax. The bi-LSTM-CRF models also perform better in terms of emotion prediction. Specifically, emotion prediction seems to play a role in learning POS and verse measure, leading Haider to propose a relationship between meter and emotion. The two large corpora used for this study contain approximately 3 million lines of English poetry and around 2 million lines of German poetry. The English

corpus contains all poetry available in the English Project Gutenberg.²⁶ The German corpus contains data from digitized corpora such as the German Text Archive,²⁷ the Digital Library of Textgrid²⁸ and the German version of Project Gutenberg. The data are available in both JSON and TEI P5 XML format.

The tool *Automatische Mittelhochdeutsche Metrik* 2.0 (Dimpel 2015) provides metrical annotation for Middle High German. Originally designed for authorship detection to identify an author's metrical profile, the tool annotates and statistically evaluates line features, such as stress placement, by using the ErMaStat package, which follows the *Kadenz-Taxonomie* "rhythm taxonomy" by Heusler (1927). Not all features were statistically significant in improving the performance of the tool. This was the case, for instance, of elision, which did not significantly improve line annotation. Consequently, these features have been excluded from the tool. However, they can be activated manually by the user in the user interface in case they are of interest to the researcher.

Estes and Hench (2016) use machine learning technology for annotating Middle High German poetry. Their supervised learning model, namely a bi-LSTM with CRF neural network, predicts the metrical value of syllables. The meter analyzed is both stress-based and quantity sensitive, therefore both stress- and weight-related features are predicted by the tool. Other features considered are syllable position within the line, word boundary, syllable length in characters, syllable constitution in characters, and elision. To obtain the most likely scansion of a line, additional metrical rules are considered and marginal probabilities calculated. The tool was trained on material from the *Mittelhochdeutsche Begriffsdatenbank*.²⁹ The database can be consulted online. However, most of its material still has copyright, hence, the whole texts are not publicly available. The Middle High German scansion tool (MHG-scansion model) by Estes and Hench (2016) has a graphical user interface and its source code available in a code repository.

4.1.8 Latin

Also dating from the early 2010s, Pedecerto (Colombi et al. 2012) has been developed within the project *Traditio patrum* FIRB of the *Università di Udine*. It is used to annotate approximately 244,000 dactylic lines from the *Musisque Deoque* digital archive. Unfortunately, neither the source code nor the algorithm used for developing the tool seems to be available. The website presenting Pedecerto offers a user-friendly online demo of the tool which can be used to scan verse of various meters. In addition, the platform provides a corpus of annotated poetry³⁰ and the possibility of doing multiple types of searches, such as

searching for word or verse type, based on line metrical pattern or based on specific metrical sequences. The plain text can be instead browsed at the *Musisque Deoque*³¹ archive.

Calliope scansio (Jacquet 2013) dates back to the same year. This tool for Latin poetry detects the number of syllables per line and groups them into feet. Syllables are looked up in a dictionary, and syllabification occurs on a character-by-character basis, following specific rules. Once the syllables are defined, their quantity is assigned based on their position in the words. The resulting line annotation is matched to a metrical template, for example hexameter, by considering all possible combinations and groupings of the material, until the best combination is found. The code repository of Calliope scansio contains some annotated and non-annotated material, however, the folder contains a limited amount of material and, just like the whole tool, has not been updated since 2014.³² The code repository of Calliope scansio mentions a version of the tool with graphical user interface, but it does not appear to be currently available. In addition, its repository does not seem to be regularly updated.

4.1.9 Old Occitan

Slightly different from the other tools introduced in this section, the metrical annotation machine by Rainsford and Scrivner (2014) serves the purpose of adding metrical annotation to treebanks (i.e. syntactically annotated corpora). The process of annotation of Old Occitan follows three steps: first, a syllable-level tokenizer algorithm syllabifies words in a line; second, the results are manually corrected, and lexical stress is added; third, a scansion algorithm, equipped with accepted rules of syllabification, assigns metrical positions to each syllable. The material used by Rainsford and Scrivner is available—although without annotation—on the Old Occitan corpus *Le Roman de Flamenca Corpus*.³³ The online platform appears to be under construction; therefore, while the plain text of the poems can be browsed, the metrical annotations do not seem to be available.

4.1.10 Portuguese

The first tool for metrical annotation of Portuguese verse also dates back to the early 2000s. It was developed by Araújo and Mamede (2002) and analyzes poems in terms of line length, rhyme, and the number of lines in each stanza. Its architecture is inspired by Galaxy-II (Seneff et al. 1998) with a central process that coordinates the various modules: the external interface, which generates phonetic transcriptions and syllabic division; the lexical interface, which connects the information generated by the external interface to the lexicon; the system interface, which includes

first case study considers three types of features, namely rhyme, rhythm (the number of words and the number of syllables in a line and throughout the poem), and frequency of adjectives in comparison to prose. Their results suggest that rhythmical features, followed by rhyme, were the most accurate in distinguishing. The frequency of adjectives, on the other hand, was found to be not significant for poetry recognition. In the second case study, the classifier was trained on multiple classes representing poetic forms and one representing prose; nine features were considered, such as the number of lines and the number of words in a line, and in the whole texts, the number of stanzas, the number of alliteration and rhymes. This method performs considerably well, albeit using a small dataset for training. In the third case study, the Festival Text-to-Speech Synthesis System (Black, Taylor, and Caley 2002) was used to extract various text features, particularly those relevant to semantics and prosody, that were then employed to train a decision tree classifier. As in the previous case studies, the employment of the extracted features gave high-accuracy results.

Subsequently, Tizhoosh, Sahba, and Dara (2008) explored various techniques of poetry recognition. They focused on extracting genre-specific poetic features, for example, rhyme, meter, in order to define new poem classification techniques. Two classifiers were used, namely naïve Bayesian and Multilayer Perceptron. They evaluated five different approaches: (1) a traditional text classification method employing Information gain and Mutual Information (Yang and Pederesen 1997); (2) a method that considered shape, that is, structural characteristics of poetry such as line length; (3) a method which combined shape and rhyme features; (4) a method which considered shape, rhyme, and metrical structure; and (5) a method which combined poetic features of method 3 and word frequency. Despite all approaches performing well, the highest accuracy was achieved with method 5.

A further advancement for poetry recognition was represented by the tool developed by [Lorang and *et al.* \(2015\)](#). This tool used an image-based approach to identify poetry in historical newspapers. The purpose of the tool was to detect and make available for further investigation poetry which appears in non-poetic historical texts. In order to do so, supervised machine learning was used to create a classifier able to distinguish snippets that contain poetry from those that do not. A preprocessing phase is followed by feature extraction. The preprocessing algorithm performs blurring on the snippet from a newspaper, then, the blurred image undergoes bi-Gaussian binarization and pixel consolidation. The feature extraction algorithm

saves a computed set of attributes to the data structure representing the processed image.

Similarly, the project by [Kilner and Fitch \(2017\)](#) addressed poetry detection in digitized historical Australian newspapers from the AusLit database ([Kilner 2002](#)). The authors trained a naïve Bayesian classifier adapting the code by Shiffman (Bayesian Filtering,⁴¹ [2008](#)) and that by Paul Graham (A Plan for Spam,⁴² [Graham 2002](#)), and by using the overProof ([Evershed and Fitch 2014](#)) correction algorithm to improve the quality of OCR. The features considered for recognizing poetry were based on Tizhoosh *et al.* ([Tizhoosh, Sahba, and Dara 2008](#)).

Finally, the most recent tool for poetry identification was developed by [Foley \(2019, 2020\)](#). The tool is an automatic, learned model for poetry extraction from digitally scanned books. Its purpose is to identify poetry in scanned documents and in particular in scanned pages which do not necessarily only contain poetry. The model (using bi-LTSM RNN) has been trained on a few thousand labeled pages and detects anomalous formatting in the page, such as differences in punctuation, margins, and capitalization. This tool has its code, model, and datasets available and downloadable.

4.3 Tools for poetry visualization⁴³

4.3.1 English

In this section, we describe a range of tools capable of visualizing different aspects of a poem. Most of these tools have emerged in recent years, indicating an increasing interest in poetry from the perspective of distant reading (Moretti 2013), macroanalysis (Jockers 2013), and quantitative approaches.

The first research-oriented method was presented by [Kaplan and Blei \(2007\)](#). It visualizes stylistic similarities and relationships between poems by embedding them in a vector space and using Principal Component Analysis. Style detection is used to distinguish poetic texts and cluster poems by poet. A set of features is used to differentiate different styles and classify poems. These features are categorized as orthographic (e.g. word count, average line length), syntactic (e.g. frequencies of POS), and phonemic (e.g. sound patterns, assonance, and consonance). Each feature was assigned a weight. The CMU Pronouncing Dictionary for North American English ([Weide 1998](#)) is used for word-to-phoneme translation. The method is applied to a miscellaneous corpus of American poets from various sources and different centuries. The authors conclude that their algorithm is able to distinguish between poems from poets with different styles, verifying whether a single long poem follows the same style, and providing evidence of known mentoring relationship between some poets.

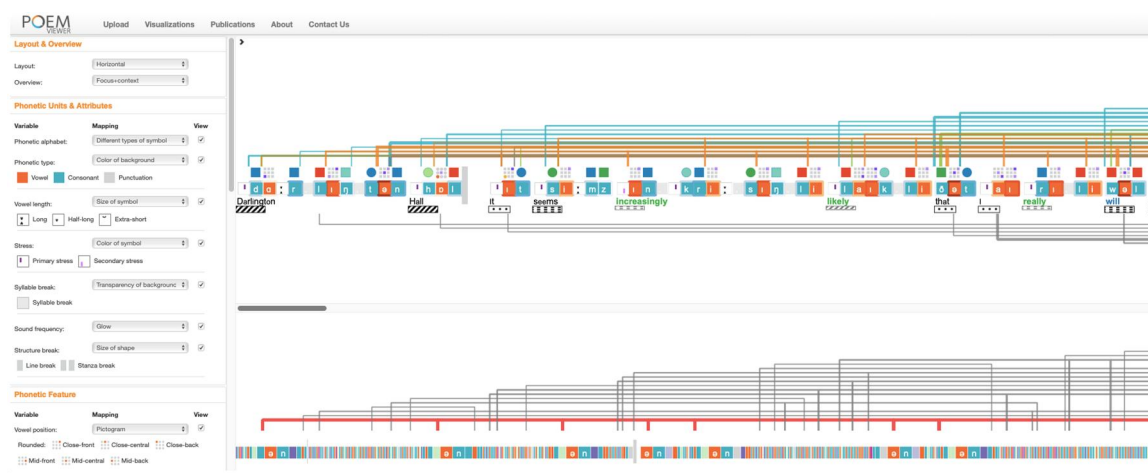


Figure 5. Sample of poetry visualization with PoemViewer.

Myopia (Chaturvedi *et al.* 2012) is a software tool implemented in Python for poetry visualization in the context of close reading. It was developed based on the poems available in the Poetess Archive,⁴⁴ which contains 4,000 entries of poetry written by British and American poets between 1750 and 1900. Myopia uses a schema extended from TEI, where a poem’s characteristics are structured as tags and attributes according to XML standards. It is designed for comparing literary attributes of different poems and to improve the display of a multidimensional representation of poetry encoded in TEI. The tool has a visual interface and currently exists as a desktop tool.

Poem Viewer ([Abdul-Rahman et al. 2013](#)) is part of the Visual Informatics Lab at Oxford and is a web-based tool for visualizing poems which aims of enhancing close reading.⁴⁵ It was developed to create a user-centered tool that allows the visualization of different features. The features that can be displayed pertain to the phonic aspect of the poem: the tool shows consonants, vowels, their phonetic features, and phonetic relations (rhyme, assonance, and alliteration). It also visualizes the vocalic movement within the vocalic triangle, depicting the transition from one vowel to another while reading the line. In addition, word repetition, word sentiments, and semantic relations can be visualized. Poem Viewer has an online graphical user interface platform where poetry can be analyzed. The platform gives many customizable options regarding the visualization, such as displaying vowel length or word repetition. An example is given in [Fig. 5](#).

The Eighteenth-Century Poetry Archive (ECPA; Huber 2018) provides the community with a range of new web-based poetry visualization tools for this corpus. For each digital edition of a particular poem, users can apply analysis tools to extract information on the

phonological, morphological, syntactic, semantic, and pragmatic (mainly, named entities) layers. Additionally, it integrates some original and adapted visualization tools, including the aforementioned Poem Viewer, as well as Phonemia—which visualizes the phonemic makeup of a poem—and DoubleTreesJS—which aims to provide a more compact view of words in their context of use than traditional concordances (Key Word in Context). These tools are available on the ECPA’s graphical user interface platform and are connected to the digital edition site of each poem. The analysis tools are displayed by default, while Phonemia, Poem Viewer, and DoubleTreeJS must be launched by the user.

4.3.2 Spanish

One tool for visualization of Spanish poetry is the project DISCOVER (Bermúdez-Sabel, Ruiz Fabo, and Martínez Cantón 2023). This tool adapts existing resources for the analysis of rhymes in the DISCO corpus, featuring Spanish sonnets (Ruiz Fabo *et al.* 2017a). This corpus was annotated using ADSO Scansion system (Navarro-Colorado 2018), ANJA by Pablo Ruiz Fabo (Ruiz Fabo *et al.* 2017b), and RhymeTagger⁴⁶ by Petr Plecháč (2018) for rhyme detection. The interface, implemented in PHP and Javascript, allows users to filter results by corpus and prosodic elements and view rhyme trends. DISCOVER has an online graphical user platform and allows the user to limit the corpus by origin, century, and even gender of the author. It also provides filters for stress, rhyme, the presence of enjambment, and even visualize rhyme trends. However, the code is not available.

4.3.3 Italian

Dante Visualized (Ferraro 2020) was developed to enhance distant reading of Dante’s *Divina Commedia*,

and has the potential to be applied to other texts in future research. The tool analyzes both stylistic and semantic features, offering a schematic representation of the structure and rhythm, a visual representation of the sentiment analysis, and the distribution of keywords in the three parts of the work. The tool was developed by using machine learning techniques and the Naive Bayes Classifier (Perkins 2010). The tool was implemented in Python. A JavaScript library (Bostock D3.js⁴⁷) was used for producing the visualizations. The code repository of the Dante Visualized (Ferraro 2020) includes the tool itself as well as the results of the case study on the *Divina Commedia* by Dante.

5. The role of rules and data

Poetry analysis tools can be classified based on the type of technology they are built upon and, more specifically, the role played by rules and data in their structure.

A large number of tools are rule-based, especially, those for metrical annotation. These tools are constructed based on a set of predefined rules that govern the annotation process. For metrical annotation, explicit annotation rules are outlined, and the tool follows these rules to annotate new material.

With the advancements of machine and deep learning techniques, the use of these technologies has spread in the field of automatic poetry analysis, leading to different types of systems. In these systems, rules are not explicitly provided to the tool; instead, the model learns implicit generalizations from large amounts of data. Machine learning algorithms are employed to extract information and create a model, which is then used to analyze new data (Patternson and Gibson 2017). This process can be either supervised or unsupervised. In supervised learning, the model is trained on manually annotated data to make generalizations for annotating new data, while unsupervised learning involves using raw data and fine-tuning the model with manually annotated material. This type of tool, where the model makes generalizations from data, is called data-driven or inference-based. Although the terms may seem contradictory, they refer to the same aspect: data-driven refers to the fact that the model receives data and learns from them; inference-based refers to the process of making generalizations and making predictions from the data.

In the present article, both rule-based and inference-based tools are considered. As mentioned above, within the group of metrical annotation tools, most tools are rule-based. Defining metrical and prosodic rules for verse annotation has proven to be an effective technique. These tools generally rely on rules related to linguistic features, such as syllabification, stress assignment, phonetic representation, as well as metrical features of the

line, such as line length, number of stressed syllables per line, rhyme, and the number of lines in a stanza.

Nevertheless, in the last few years, the use of machine learning technologies in the field of automatic metrical annotation has significantly increased. The most commonly employed approach involves RNNs (Schmidhuber 1989) with bi-LSTM cells (Hochreiter and Schmidhuber 1997) and CRF (Lafferty, Abdul-Rahman, and Pereira 2001). These techniques have demonstrated high accuracy in metrical annotation. Some of the metrical annotation tools developed using machine learning techniques include the MHG-scansion model (Estes and Hench 2016), the tool by Tanasescu, Paget, and Inkpen (2016), the tool by Agirrezabal (Agirrezabal 2017; Agirrezabal, Alegria, and Hulden 2017), the tool by KAIST (Kim, Zhao, and Zelalem 2018), the Middle Dutch Syllabifier (Haverals, Karsdorp, and Kestemont 2019) and Stresser (Haverals 2020), the tool by De Sisto (2020) and that by Haider (2021). The manually annotated material used as training data of these models is generally annotated according to syllable structure, stress, rhyme, word boundary, and, in the case of quantitative verse (i.e. Middle High German for MHG-scansion model), in terms of syllable quantity. These tools typically employ a similar architecture based on RNN with bi-LSTM. A different method is proposed in the study by De la Rosa *et al.* (2021), where transformer-based models (Wolf *et al.* 2019) are employed for detecting line stress patterns.

Similar to data-driven metrical annotation tools and, actually, a predecessor of them, the tool by [Hayward \(1991, 1996\)](#) also uses Artificial Neural Networks for evaluating line metricality. In addition, [Greene, Bodrumlu, and Knight \(2010\)](#) use unsupervised learning for extracting word stress from a corpus as part of their tool which annotates, translates, and generates poetry.

Machine learning techniques are also utilized in the field of poetry visualization, as demonstrated by the ongoing project Dante Visualized (Ferraro 2020). This tool uses a Naive Bayes Classifier (Perkins 2010) trained on a manually created dataset of subsets of texts from other authors that are linguistically and temporally similar (Ferraro 2020).

Finally, inference-based methods are highly common among tools of poetry recognition; in fact, all tools considered in the present survey employ them: [Tizhoosh and Dara \(2006\)](#), [Tizhoosh, Sahba, and Dara \(2008\)](#), [Lorang *et al.* \(2015\)](#), and [Foley \(2019, 2020\)](#).

6. Conclusion

In this article, we have presented a survey of automatic tools for poetry analysis, with a focus on those

23. Available at: <https://github.com/ChristopheJacquet/Calliope/tree/master/test> [accessed 18 Jan. 2022].
24. See at: https://cl.indiana.edu/~obscrivn/F1_495.html [accessed 21 Jun. 2023].
25. <https://aoidos.ufsc.br> [accessed 21 Jun. 2023].
26. StarLing is a software package for various types of linguistic text and database analyses. For more details: <https://starlingdb.org/program.php?lan=en> [accessed 21 Jun. 2023].
27. We would like to thank an anonymous reviewer for pointing this out to us.
28. Available at: <https://ruscorpورا.ru/new/search-poetic.html> [accessed 6 Jan. 2022].
29. We would like to thank an anonymous reviewer for providing us information about this corpus and about research based on it. Given that the website and other resources are only available in Russian, it would not have been possible for us to include it in our review without the reviewer's contribution.
30. See [Plecháč, Zelenkov, and Sela \(2020\)](#), p. 136, footnote 7.
31. For relevant literature, see: [Manevitz and Yousef 2007](#); [McCallum and Nigam 1998](#); [Ruiz and Srinivasan, 1998](#); [Yang 1999](#).
32. Unfortunately, the link to this blog post does not exist anymore: http://shiffman.net/teaching/a2z_2008/bayesian/ [last accessed 27 Oct. 2015]
33. <http://www.paulgraham.com/spam.html> [accessed 22 Feb. 2022].
34. Apart from the tools mentioned in this section, there are some studies of rhyme through network analysis. In this sense, we would like to highlight the panel "Understanding Rhyme Through Network Analysis" presented at Digital Humanities 2020 Conference ([Houston et al. 2020](#)).
35. See its website at: <http://www.poetessarchive.org> [accessed 18 Oct. 2021].
36. According to its authors, their goal was to produce "a suitable tool to help poets to explore the information spaces of poems. Following the nested design process" ([Abdul-Rahman et al. 2013](#)).
37. <https://github.com/versotym/rhymetagger>
38. <https://d3js.org/>
39. See the poem 'I dwell in Possibility' (<https://www.poetryfoundation.org/poems/52197/i-dwell-in-possibility-466>).

References

- Abdul-Rahman, A. *et al.* (2013) 'Rule-based Visual Mappings—with a Case Study on Poetry Visualization', *Computer Graphics Forum*, 32: 381–90. <https://doi.org/10.1111/cgf.12125>
- Agerri, R., Bermudez, J., and Rigau, G. (2014) 'Ixa Pipeline: Efficient and Ready to Use Multilingual nlp Tools', *LREC*, 2014: 3823–8.
- Agirrezabal, M. (2017) *Automatic Scansion of Poetry*. Donostia, Spain: University of Basque Country.
- Agirrezabal, M., Alegria, I., and Huldén, M. (2017) A *Comparison of Feature-Based and Neural Scansion of Poetry*. in: *Proceedings of the International Conference Recent Advances in Natural Language Processing*, RANLP 2017, pp. 18–23. Varna, Bulgaria. INCOMA Ltd.
- Agirrezabal, M. *et al.* (2013) 'ZeuScansion: A Tool for Scansion of English Poetry', *Journal of Language Modelling*, 4: 3. <https://doi.org/10.15398/jlm.v4i1.102>
- Algee-Hewitt, M.R. *et al.* (2014) *The Stanford Literary Lab Transhistorical Poetry Project Phase ii: Metrical Form*. in: *9th Annual International Conference of the Alliance of Digital Humanities Organizations*, {DH} 2014, Lausanne, Switzerland, 8-12 July 2014, *Conference Abstracts*. Alliance of Digital Humanities Organizations (ADHO). <http://dharchiv.org/paper/DH2014/Paper-788.xml>.

- Araújo, J. P., and Mamede, N.J. (2002) *Classificador de Poemas*. Lisbon: CCTE.
- Barakhnin, V.B., Kozhemyakina, O. Y., and Borzilova, Y. S. (2019) ‘The Development of the Information System of the Representation of the Complex Analysis Results for the Poetic Texts’, *Vestnik NSU: Information. Technologies*, 17: 5–17. <https://doi.org/10.25205/1818-7900-2019-17-1-5-17>
- Beaudouin, V., and Yvon, F. (1996) ‘The Metrometer: A Tool for Analysing French Verse’, *Literary and Linguistics Computing*, 11: 23–31. <https://doi.org/10.1093/lilc/11.1.23>
- Bermúdez Sabel, H., Ruiz Fabo, P., and Martínez Cantón, C. (2023) ‘DISCOOvering Spanish Sonnets: A Circular Reading Experience’, in: A. S. Bories, P. Plecháč, and P. Ruiz Fabo (eds.) *Computational Stylistics in Poetry, Prose and Drama*, pp. 67–86. Berlin, Germany: De Gruyter. <https://doi.org/10.1515/9783110781502-004>
- Black, A.W., Taylor, P., and Caley, R. (2002) *The Festival Speech Synthesis System: System Documentation*. Edinburgh, Scotland: University of Edinburgh, The Centre for Speech Technology Research.
- Bobenhausen, K. (2011) ‘The Metricalizer2—automated Metrical Markup of German Poetry’, *Current Trends in Metrical Analysis*. Bern, Switzerland: Peter Lang.
- Bobenhausen, K., and Hammerich, B. (2015) ‘Métrique Littéraire, Métrique Linguistique et Métrique Algorithmique de l’allemand Mises en Jeu Dans le Programme Metricalizer2’, *Langages*, 3(199): 67–88.
- Breiman, L. (1996) ‘Bagging Predictors’, *Machine Learning*, 24: 123–40.
- Chaturvedi, M. et al. (2012) *Myopia: A Visualization Tool in Support of Close Reading*. in: *Digital Humanities Conference*. <https://www-archiv.fdm.uni-hamburg.de/dh2012/conference/programme/abstracts/myopia-a-visualization-tool-in-support-of-close-reading.1.html>.
- Chaumartin, F. R. (2013) *Antelope, Une Plate-forme de TAL Permettant d’Extraire les Sens du Texte: Théorie et Applications de l’interface Syntaxe-sémantique* 228. Université Paris-Diderot - Paris VII. PhD Dissertation. <https://theses.hal.science/tel-00803531/file/TheseFRC.pdf>.
- Colombi, E. et al. (2012) ‘Pedecerto [WWW Document]’, <http://www.pedecerto.eu/public/index>. Date accessed 21 June 2023.
- Conser, A. (2017) ‘Greek_scansion [WWW Document]’, https://github.com/aconser/greek_scansion. Date accessed 24 September 2020.
- Cornulier, B. (1982) *Théorie du Vers: Rimbaud, Verlaine, Mallarmé*. Paris: Seuil.
- De la Rosa, J. et al. (2021) ‘Transformers Analyzing Poetry: Multilingual Metrical Pattern Prediction with Transformer-based Language Models’, *Neural Computing and Applications*, 35: 1–6. <https://doi.org/10.1007/s00521-021-06692-2>
- De la Rosa, J. et al. (2020a) ‘PoetryLab as infrastructure for the analysis of Spanish poetry’ in: Navarretta, C. and Eskevich, M. (Eds.), *Selected Papers from the CLARIN Annual Conference 2020. Presented at the CLARIN Annual Conference 2020*. Linköping Electronic Conference Proceedings 180 pp. 82–7. <https://doi.org/10.3384/ecp1809>.
- De la Rosa, J. et al. (2020b) ‘Fast and Accurate Syllabification and Scansion of Spanish Poetry’, *Review Process Length Natural*, 65: 83–90.

- Jockers, M.L. (2013) *Theme, Im: Macroanalysis: Digital Methods and Literary History*, pp. 118–53. Urbana: University of Illinois Press.
- Kaplan, D.M., and Blei, D.M. (2007) ‘A computational approach to style in American poetry’, *Seventh IEEE International Conference on Data Mining (ICDM 2007). Presented at the Seventh IEEE International Conference on Data Mining (ICDM 2007)*, pp. 553–8. Omaha, NE: IEEE. <https://doi.org/10.1109/ICDM.2007.76>
- Kilner, K. (2002) ‘AustLit: The Australian Literature Resource [WWW Document]’, www.austlit.edu.au. Date accessed 9 January 2020.
- Kilner, K., and Fitch, K. (2017) ‘Searching for My Lady’s Bonnet: Discovering Poetry in the National Library of Australia’s Newspapers Database’, *Digital Scholarship in the Humanities*, 32: i69–83. <https://doi.org/10.1093/llc/fqw062>
- Kim, H., Zhao, Y., and Zelalem, B. (2018) *Automatic Classification of Poetry by Neural Scansion*. School of Electrical Engineering and Computer Science, Faculty of Engineering, University of Ottawa. PhD Dissertation. https://ruor.uottawa.ca/bitstream/10393/37309/1/Kesarwani_Vaibhav_2018_thesis.pdf.
- Koehn, P. et al. (2007) ‘Moses: Open source toolkit for statistical machine translation’, *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions. Association for Computational Linguistics*, Prague, Czech Republic. Association for Computational Linguistics., pp. 177–80.
- Koppelaar, H., and Van Oostendorp, M. (2013) ‘The Scansion Generator [WWW Document]’, KB Lab, <http://lab.kb.nl/tool/scansion-generator>. Date accessed 31 October 2023.
- Kozmin, A.V. (2006) *Automated Analysis of Verse with Starling Software Package. Automatic Verse Analysis in Starling*. in: *Computational Linguistics and Intellectual Technologies. International Conference "Dialogue 2006" Proceedings*, pp. 265–268.
- Lafferty, J., Abdul-Rahman, A., and Pereira, F.C. (2001) *Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data*. In *Proceedings of the Eighteenth International Conference on Machine Learning (ICML 2001)*, Williams College, Williamstown, MA, USA, June 28 - July 1, 2001. 282–289.
- Logan, H.M. (1988) ‘Computer Analysis of Sound and Meter in Poetry’, *College Literature*, 15: 19–24.
- Lorang, M.L. et al. (2015) ‘Developing an Image-Based Classifier for Detecting Poetic Content in Historic Newspaper Collections’, *D-Lib Magazine*, 21, 1–24.
- Mamede, N. et al. (2004) ‘An electronic assistant for poetry writing’, *Proceedings of the 9th Ibero-American Conf. on Artificial Intelligence. Presented at the LNCS*. Berlin, Heidelberg: Springer. https://doi.org/10.1007/978-3-540-30498-2_29.
- Manevitz, L., and Yousef, M. (2007) ‘One-class Document Classification via Neural Networks’, *Neurocomputing*, 70: 1466–81.
- Marco, G. et al. (2021) ‘Automated Metric Analysis of Spanish Poetry: Two Complementary Approaches’, *IEEE Access*, 9: 51734–46. <https://doi.org/10.1109/ACCESS.2021.3069635>
- Maron, M. E. (1961) ‘Automatic Indexing: an Experimental Inquiry’, *Journal of the ACM*, 8: 404–17.
- McAless, W. (2007) *Improving Scansion with Syntactic Analysis: An Investigation into the Effectiveness of a Syntactic Analysis of Poetry by Computer Using Phonological Scansion Theory*. Technical Report, Department of Computing, The Open University. <https://doi.org/10.21954/ou.ro.00016055>.
- McCallum, A., and Nigam, K. (1998) *A comparison of Event Models for Naive Bayes Text Classification*. in: *AAAI Conference on Artificial Intelligence*. https://courses.washington.edu/ling572/papers/mccallum1998_AAAI.pdf.
- McClelland, J.L., and Rumelhart, D.E. (1988) *Explorations in Parallel Distributed Processing: A Handbook of Models, Programs, and Exercises*. Cambridge, MA: The MIT Press.
- Minsky, M. (1961) ‘Steps towards Artificial Intelligence’, *Proceedings of the IRE*, 49: 8–30.
- Mittmann, A., von Wangenheim, A., and dos Santos, A.L. (2016) ‘Aoidos: A system for the automatic scansion of poetry written in Portuguese’, *Computational Linguistics and Intelligent Text Processing. Presented at the 17th International Conference*. Konya, Turkey: CICLing.
- Moretti, F. (2013) *Distant Reading*. London: Verso.
- Navarro-Colorado, B. (2015) ‘A computational linguistic approach to Spanish Golden Age Sonnets: metrical and semantic aspects’, *Proceedings of the Fourth Workshop on Computational Linguistics for Literature*, pp. 105–13. Denver, CO: Association for Computational Linguistics. <https://doi.org/10.3115/v1/W15-0712>
- Navarro-Colorado, B. (2018) ‘A Metrical Scansion System for Fixed-metre Spanish Poetry’, *Digital Scholarship in the Humanities*, 33, 112–127. <https://doi.org/10.1093/llc/fqx009>
- Padr  L. I., and Stanilovsky, E. (2012) *FreeLing 3.0: Towards Wider Multilinguality Proceedings of the Language Resources and Evaluation Conference (LREC 2012) ELRA*. Istanbul, Turkey. European Language Resources Association (ELRA). http://www.lrec-conf.org/proceedings/lrec2012/pdf/430_Paper.pdf.
- Papakitsos, E.C. (2011) ‘Computerized Scansion of Ancient Greek Hexameter’, *Literary and Linguistics Computing*, 26: 57–69. <https://doi.org/10.1093/llc/fqq015>
- Patternson, J., and Gibson, A. (2017) *Deep Learning*. Sebastopol, CA: O’Reilly Media.
- Perkins, J. (2010) ‘Text Classifier for Sentiment Analysis—Naive Bayes Classifier’, *StreamHacker*, <https://streamhacker.com/2010/05/10/text-classification-sentiment-analysis-naive-bayes-classifier/>. Date accessed 10 September 2020.
- Pilshchikov, I., and Starostin, A. (2011) ‘Automated Analysis of Poetic Texts and the Problem of Verse Meter’, in *Current Trends in Metrical Analysis*. Lausanne, Switzerland: Peter Lang.
- Plamondon, M.R. (2006) ‘Virtual Verse Analysis: Analysing Patterns in Poetry’, *Literary and Linguistics Computing*, 21: 127–41. <https://doi.org/10.1093/llc/fql011>
- Plech c, P. (2018) ‘A Collocation-Driven Method of Discovering Rhymes (in Czech, English, and French Poetry)’, *Taming the Corpus: From Inflection and Lexis to Interpretation*, pp. 79–95. Cham, Switzerland: Springer.
- Plech c, P. et al. eds. (2019) *Quantitative Approaches to Versification*. Prague: Institute of Czech Literature of the Czech Academy of Sciences.
- Plech c, P., Zelenkov, J., and  ela, A. (2020) ‘The Batenkov Phenomenon and the Problem of Authorship Verification: A

- Multivariate Statistical Approach to an Unsolved Question', *Acta Slavica Estonica*, 12: 131–65.
- Prince, A., and Smolensky, P. (1993) *Optimality Theory: Constraint Interaction in Generative Grammar*. Manuscript, Rutgers University and University of Colorado, Boulder.
- Quinlan, J.R. (1986) 'Induction of Decision Trees', *Machine Learning*, 1: 81–106. <https://doi.org/10.1007/BF00116251>
- Rainsford, T.M., and Scrivner, O. (2014) 'Metrical annotation for a verse treebank', *Proceedings of the Thirteenth International Workshop on Treebanks and Linguistic Theories (TLT13)*, Tübingen, Germany. University of Tübingen. <https://doi.org/10.5281/zenodo.10054982>
- Reddy, S., and Knight, K. (2011) 'Unsupervised discovery of rhyme schemes', *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, Portland, Oregon, USA. Association for Computational Linguistics, 77–82.
- van Reenen, P., and Mulder, M. (1993) 'Een Gegevensbank van 14de-eeuwse Middelnederlandse Dialecten op Computer', *Lexikos*, 3: 259–81. <https://doi.org/10.5788/3-1-1110>
- Ruiz Fabo, P. (2020) 'The Diachronic Spanish Sonnet Corpus: TEI and Linked Open Data Encoding, Data Distribution, and Metrical Findings', *Digital Scholarship in the Humanities*, 36: i68–80. <https://doi.org/10.1093/lc/fqaa035>.
- Ruiz Fabo, P. *et al.* (2017a) *Diachronic Spanish Sonnet Corpus (DISCO)*. Madrid: UNED.
- Ruiz Fabo, P. *et al.* (2017b) 'Enjambment detection in a large diachronic corpus of Spanish sonnets', *Proceedings of the Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage*. in: *Social Sciences, Humanities and Literature*, Vancouver, Canada. Association for Computational Linguistics.. 27–32.
- Ruiz, M.E., and Srinivasan, P. (1998) 'Automatic text categorization using neural networks', *Proceedings of the 8th Workshop on Advances in Classification Research*, University of Washington, 59–72.
- Schmidhuber, J. (1989) 'The Neural Bucket Brigade: A Local Learning Algorithm for Dynamic Feedforward and Recurrent Networks', *Connection Science*, 403: 412.
- Schumann, A. K. (2020) *Automatische Versmaßannotation in Werken der Homerischen Hexameter-Dichtung*. Berlin, Germany: Freie Universität Berlin.
- Schumann, A. K. *et al.* (2022) 'Using Finite-state Machines to Automatically Scan Ancient Greek Hexameter', *Digital Scholarship in the Humanities*, 37: 242–53. <https://doi.org/10.1093/lc/fqab071>
- Sejnowski, T.J., and Rosenberg, C.R. (1987) 'Parallel Networks that Learn to Pronounce English Text', *Complex Sys.*, 1: 145–68.
- Seneff, S. *et al.* (1998) 'Galaxy-II: A reference architecture for conversational system development', in: *The 5th International Conference on Spoken Language Processing, Incorporating The 7th Australian International Speech Science and Technology Conference*, Sydney Convention Centre, Sydney, Australia, 30th November - 4th December 1998. ISCA. 10.21437/ICSLP.1998-478.
- Sonderegger, M. (2011) 'Application to Graph Theory to an English Rhyming Corpus', *Computer Speech and Language*, 25: 655–78.
- Tanasescu, C., Paget, B., and Inkpen, D. (2016) *Automatic Classification of Poetry by Meter and Rhyme 6. Proceedings of the twenty-Ninth International Florida Artificial Intelligence Research Society Conference (FLAIRS 2026)*. Association for the Advancement of Artificial Intelligence. <https://cdn.aaai.org/ocs/12923/12923-60593-1-PB.pdf>
- Tizhoosh, H.R., and Dara, R.A. (2006) 'On Poem Recognition', *Pattern Analysis and Applications*, 9: 325–38. <https://doi.org/10.1007/s10044-006-0044-8>
- Tizhoosh, H.R., Sahba, F., and Dara, R. (2008) 'Poetic Features for Poem Recognition: A Comparative Study', *Journal of Pattern Recognition Research*, 3: 24–39. <https://doi.org/10.13176/11.62>
- Van Oostendorp, M. (2014) 'Introducing a Scansion Machine for Dutch Poetry and Prose', *Loquens*, 1: e002. <https://doi.org/10.3989/loquens.2014.002>
- Weide, R. (1998) *The CMU Pronunciation Dictionary*. <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>, (accessed 17 Jan. 2024).
- Wolf, T. *et al.* (2019) 'HuggingFace's Transformers: State-of-the-art Natural Language Processing'. <https://doi.org/10.48550/ARXIV.1910.03771>, (accessed 17 Jan. 2024)
- Yang, Y. (1999) 'An Evaluation of Statistical Approaches to Text Categorization', *Information Retrieval*, 1: 67–88.
- Yang, Y., and Pedersen, J.O. (1997) 'A comparative study on feature selection in text categorization', *Proceedings of the Fourteenth International Conference on Machine Learning. Presented at the ICML'97, Nashville*, 8-12 July 1997, 412–20.

Table A1 Summary of the reviewed tools including availability of source code, GUI, websites and code repositories.

(continued)

