

## On the Authenticity of Prose Writings Attributed to Śaṅkara

Ivan ANDRIJANIĆ and Jacek BAŃKOWSKI

**Abstract:** Śaṅkara is traditionally considered the author of an exceptionally large number of works. Indological scholarship has attempted to filter out some of these works within traditional philological and historical frameworks. Many were, however, taken for granted to be authentic, and no serious research into their authenticity has been conducted. This paper attempts a computational stylometric approach to establish the authenticity of prose commentaries attributed to Śaṅkara. The General Imposters (GI) framework appears to be the most suitable existing method developed for the purpose of verifying authorship. The GI calculates the statistical distance between certain texts' features and estimates whether the disputed text is closer to the candidate author than to a set of texts that may not have been composed by him. The paper also presents a machine-based method for separating the words and resolving the sandhi in the Sanskrit text, crucial for the procedure. The success rate in verifying authors of undisputed texts appears to be acceptable enough to proceed to the next step, where 18 prose commentaries traditionally attributed to Śaṅkara are subjected to the GI verification procedure. The result conforms to the most conservative assessments of Śaṅkara's authorship; GI verified the authenticity of the commentaries on the principal Upaniṣads (with the exception of the commentary on the *Śvetāśvataropaniṣad*) and on the *Bhagavadgītā*. Besides these, commentaries on the *Nṛsiṃha-(pūrva)-tāpanīyopaniṣad* and the *Adhyātmapaṭala* were, rather unexpectedly, also successfully verified as genuine works of Śaṅkara.

**Keywords:** authorship, stylometry, Advaita, Vedānta

Ivan ANDRIJANIĆ, University of Zagreb, Croatia; [iandrij@ffzg.unizg.hr](mailto:iandrij@ffzg.unizg.hr);

 0000-0002-1544-585X

Jacek BAŃKOWSKI, Institute of Polish Language, Polish Academy of Sciences, Poland;

[jacek.bakowski@ijp.pan.pl](mailto:jacek.bakowski@ijp.pan.pl);  0000-0003-2480-3396



This article is distributed under a Creative Commons Attribution 4.0 licence (<https://creativecommons.org/licenses/by/4.0/>).

## Introduction

Within the broader field of digital humanities, contemporary computational stylometry represents a particularly interesting and exciting area. Broadly outlined, stylometry implies the measurement of textual stylistic affinities in order to address questions like authorship and chronology. Advancements in computing power have made it increasingly feasible to carry out complex operations that involve extensive statistical calculations, which were considered unachievable until recently. One of the most studied stylometric disciplines is authorship attribution, where features of a text of unknown authorship are compared to the determined profiles of known authors in order to find a matching candidate.<sup>1</sup> However, in the history of Indian philosophy, a different setup might be of greater interest, where features of a text of disputed authorship are compared to undisputed texts of the candidate author. Such a setup is referred to as authorship verification. In Indian philosophy, this might be important because many spurious works were traditionally attributed to certain famous authors often without credible verification. Such is a case with Śaṅkara (8th cent. CE)<sup>2</sup>, to whom a vast number of texts is ascribed in manuscript colophons and by monastic tradition.<sup>3</sup>

In two articles, ANDRIJANIĆ (2020a, 2020b) experimented with an authorship verification method named the General Imposters (GI) framework in order to assess the accuracy of the method on Sanskrit philosophical texts.<sup>4</sup> As the method gave satisfactory results in verifying authorship of undisputed texts,<sup>5</sup> ANDRIJANIĆ (2020a; 2020b) verified traditional attribution of the *Kāthopaniṣadbhāṣya* (KaUBh), *Īsopaniṣadbhāṣya* (ĪUBh) and *Chāndogyopaniṣadbhāṣya* (ChUBh) to Śaṅkara. However, two serious shortcomings are visible in this experiment. The first problem is that rather small text samples were used in both experiments.<sup>6</sup>

<sup>1</sup> The term “stylometry” was coined by Wincenty LUTOSLAWSKI in 1898. More on stylometry, its history and methods one can find in HOLMES 1994; JUOLA 2006; KOPPEL et al. 2009 and STAMATOS 2009. For a more general introduction to the authorship problem, we recommend LOVE 2002.

<sup>2</sup> For an overview and evaluation of previous attempts to date Śaṅkara see HARIMOTO 2006, who narrows the date of *Brahmasūtrabhāṣya* between 756 and 772.

<sup>3</sup> According to BELVALKAR 1930: 241 about 435 works are ascribed to Śaṅkara in manuscript colophons. Belvalkar made his estimation according to Aufrecht’s *Catalogus Catalogorum* and Reports and Descriptive Catalogues of the Government Library in Madras.

<sup>4</sup> Imposters method is originally proposed by KOPPEL and WINTER 2014: 5–6 and further developed by SEIDMAN 2014 and POTHA and STAMATOS 2017. KESTEMONT et al. 2016 employed the method on the disputed writings of Julius Caesar, while variations of the method won first prize at the PAN-2013 and PAN-2014 evaluation lab on uncovering plagiarism, authorship and social software misuse.

<sup>5</sup> The success rate of the GI procedure applied to Sanskrit philosophical texts reached in certain setups up to 80% of successful attributions.

<sup>6</sup> In fact, KOPPEL and WINTER 2014: 8 have shown that GI method accuracy increases as the

Another problem is that the number of texts used was quite limited. Therefore, it is possible to doubt statistical reliability of results that came from such small samples. Due to the utilisation of manually segmented Sanskrit text corpus in both of Andrijačić's studies, it was not feasible to compile a substantial quantity of texts.

Now, let us clarify the importance of the segmentation issue. The GI method relies on a feature vector, usually consisting of word or  $n$ -gram<sup>7</sup> frequencies. At first glance, character  $n$ -grams might seem uninformative, meaningless and counter-intuitive. However, according to JUOLA 2006, they have turned out to be the best performing feature type in the sophisticated authorship attribution, although they carry little information or meaning. One of the reasons for the effectiveness of this measure is that these units tend to capture "a bit of everything", being sensitive to both the content and form of a text (HOVARDAS and STAMATATOS 2006; KOPPEL et al. 2009; STAMATATOS 2009).<sup>8</sup> Admittedly, some have expressed caveats regarding their use, since many of them are "closely associated to particular content words and roots" (KOPPEL et al. 2009: 13). However, the use of  $n$ -grams increases the amount of measurement data to be observed, as in a text there is more  $n$ -grams than entire words, which is worth noticing from the strictly quantitative point of view (STAMATATOS 2009; DAELEMANS 2013).

---

length of the input documents increases. However, they took into consideration rather small texts. This means that the method is successful even when such short texts of 1,500 words are used. However, the problem is here in the selection, because Andrijačić used smaller sections taken from voluminous works. Using a randomly selected smaller set of features from a larger set is expected to yield more reliable results compared to utilising only a small fragment of text. For the problem of text size and sampling in stylometry, see LUYCKX and DAELEMANS 2011 and EDER 2015.

<sup>7</sup> Character  $n$ -grams are adjoining and partially overlapping sequences of  $n$ -letters. E.g. as a character  $n$ -gram sequence (with  $n = 3$ ), the Sanskrit phrase *tattvamasi* "thou art that" will be analysed as "tat" "att" "ttv" "tva" "vam" "ama" "mas" "asi" (cf. ANDRIJANIĆ 2020b: 107). However, if the phrase is segmented into words, the phrase reads *tad tvam asi* and the character 3-gram sequence would explicitly catch spaces between words. The sequence would be analysed as: "tad" "ad∅" "d∅t" "∅tv" "tva" "vam" "am∅" "m∅a" "∅as" "asi". In authorship studies, character  $n$ -grams are recognised as a powerful alternative to words (word unigrams). Cf. KESTEMONT et al. 2016: 87.

<sup>8</sup> To some extent they are therefore similar to function words. We understand function words as a small closed-class category set of words which contribute to sentence meaning only indirectly, such as articles, prepositions, particles and determiners (MORROW 1986: 423). The prevailing opinion is that function words, being heavily grammaticalised, do not carry meaning in isolation but are instead used much more frequently than content words (ZIF 1949). Unlike content words, function words might not be so influenced by the topic of the text. Their high frequency of use makes them interesting to study quantitatively, and they are universally employed by authors in a given language. Most importantly, it is often considered that their usage is not under an author's conscious control during the writing process. Thus, they are a reliable basis for textual comparisons (KESTEMONT 2014).

Furthermore, since it has been proved that authorship attribution based on word frequencies provides poorer results with highly inflective languages, *n*-grams' ability to function independently of a language constitutes a crucial argument for their use (RYBICKI and EDER 2011: 319–320). Indeed, in weakly inflected languages much of their functional linguistic information is expressed through minimal units of meaning or grammatical morphemes, usually in the form of individual words such as prepositions or articles (MORROW 1986). On the other hand, it has been proved that highly inflected (and agglutinative) languages show a greater susceptibility to analysis by *n*-grams – which has been attested with languages such as Latin, Polish or Hungarian (RYBICKI and EDER 2011: 319–320). Sanskrit follows it to no less extent, extensively using the case endings, as well as other forms of inflection – and thus is closer to such languages like Latin and Polish, highly inflected in comparison with English.

To sum up, the *n*-grams approach combines all advantages of both word functions and *n*-grams: “high frequency, good dispersion, content-independence [and] unconscious use” and is often able to capture more refined grammatical patterns (KESTEMONT 2014).<sup>9</sup> Furthermore, and this will be a fundamental concept for the rest of our reasoning here, there is a subtle usage of the presence of whitespaces by *n*-grams, namely, it allows for more observation-per-word, but what is more, due to its explicit encoding, it makes a representation sensitive to inflectional information – which is simply ignored in a word-level approach (KESTEMONT 2014) – and which is predominant in Sanskrit. It also allows one to highlight the important status of words' first letters, which are particularly important in how words are cognitively accessed in the mental lexicon (RUBIN 1995: 74).<sup>10</sup>

The problem is that Sanskrit words available in electronic and printed texts are connected to one another due to sandhi and Devanāgarī writing conventions. ANDRIJANIĆ 2020b showed that unsegmented and unsandhied Sanskrit texts analysed as *n*-grams do not yield satisfactory results with GI even when large text-samples are used, which tends to confirm the above observations. Text segmentation can also, to a certain extent, isolate and bring to an equal form some functors and retrieve some functional and stylistic information from them. Thus, by breaking up our Sanskrit sandhied words into smaller units we

<sup>9</sup> A very special attention should be given to grammatical morphemes, also named “functors” by KESTEMONT, which broaden and extend the concepts of function words to include all grammatical morphemes realised either as individual words or phrases (KESTEMONT 2014).

<sup>10</sup> We can operate here an interesting parallel with art history research. In the 19th century Giovanni MORELLI (1816–1891) suggested that the attribution of Italian master's paintings should be based on frequent, functional, inconspicuous (and maybe even unconscious) details rather than content-related elements (KESTEMONT et al. 2012: 61–62).

were able to harvest more and better information from texts. Furthermore, this approach allows us to isolate the previously mentioned first word letters. All this brings us to the conclusion that Sanskrit texts where words are separated work much better.

Also, to obtain more reliable results, a larger body of text samples is needed. Recently, a solution to this problem came to hand when a reliable automatic text segmentation method was introduced by HELLWIG and NEHRDICH 2018.

Therefore, in the first part of this paper we shall describe the GI method and the machine-learning text segmenter developed by HELLWIG and NEHRDICH 2018. Then, we shall evaluate whether the GI method accurately attributes machine segmented texts of undisputed authorship to their authors. If the results turn out to be satisfactory, we will move on to the final phase in which we will evaluate whether a body of prose writings, traditionally attributed to Śaṅkara, can actually be recognised as his works.

### Imposters method

The GI algorithm depends on measuring the distance between a feature vector representing the disputed text and text(s) that belong to a candidate author on one hand, and the distance between the same disputed text and the set of “imposters”, that is texts composed by authors that cannot be authors of the disputed text, on the other. In our experiment, feature vectors (that represent a certain text) consist of relative frequencies of words (word unigrams) or character trigrams. Let “*D*” stand for a vector of features representing the disputed text; “*C*” for one or more texts by the target author (candidate texts). “*I*” stands for the set of imposter texts that could not have been composed by the candidate author. The method measures in a number of iterations whether “*D*” (disputed text) is closer to “*C*” (candidate) than to the “*I*” (imposters set).

All calculations in this paper are made by the function `imposters()`, a part of the `stylo` package (EDER et al. 2016), an open source stylometric script written in the statistical programming environment R (cf. EDER 2018). Function `imposters()` is by default set to 100 iterations; in each of these iterations a random subset of 10% of features from “*D*” and “*C*” is selected, and compared to one half of the imposter set. The result (from 0 to 1) indicates a proportion of iterations where “*D*” is closer to the set of candidates “*C*” than to the imposters set “*I*”.

At this point, the question arises as to what result could indicate a successful verification. If the result would be e.g. 0.5 (in which half of the iterations were closer to the candidate and half closer to the imposters), would this mean

that the result is positive or not? For this purpose, function `imposters.optimize()` is designed to find optimal parameters.<sup>11</sup> The optimizer calculates values that set the threshold for successful and unsuccessful verifications. In our machine segmented corpus, the threshold (calculated with the Cosine Delta distance measure) for the word unigrams is 0.66, which means that any higher score indicates higher probability of successful attribution. A score below 0.34 indicates that the candidate author is unlikely the author of the disputed text. Everything between 0.34 and 0.66 represents a “grey area”, a zone of uncertainty where the classifier refrained from reaching a decision. For character trigrams, threshold is similarly at 0.66 and above for successful and at 0.32 and below for unsuccessful verification.

## Distance metrics employed in the GI

Distance metrics, as indicated in the GI description, play a crucial role in the algorithm. Both distance and its measurement seem to be absolutely intuitive concepts. Quite naturally, in everyday life the distance between two points is based on the Euclidean measure, e.g. the straight line between them.

The same will occur with the much less intuitive notion of the distance between two completely different texts of different length and made up of different words. The optimal measurement method will again depend on the most suitable criteria to apply in our case.

We will then approach the problem of measuring the distance between a given pair of documents  $A$  and  $B$ . Those documents will be represented by two document vectors  $a$  and  $b$  consisting of  $n$  features in some fixed order;  $a_i$  and  $b_i$  will represent the value of the  $i$ -th feature in both of these documents, respectively, which means that each different word corresponds to a different dimension – see the Vector Space Models representation (KESTEMONT et al. 2016: 4–5).

In our experiment, we use two distance measures that have yielded consistently good results in stylometric studies. The first measure is MinMax, which has been shown to be more successful than Manhattan and Cosine (KESTEMONT et al. 2016). The MinMax measure is defined as follows:<sup>12</sup>

$$\text{minmax}(a, b) = 1 - \left( \frac{\sum_{i=1}^n \min(a_i, b_i)}{\sum_{i=1}^n \max(a_i, b_i)} \right) \quad (\text{KESTEMONT et al. 2016: 5}).$$

<sup>11</sup> Based on the “score shifter” from KESTEMONT et al. 2016. The `c@1` measure of classifier’s performance (PEÑAS and RODRIGO 2011) is applied to identify a “grey zone” where the classifier is not able to make a decision.

<sup>12</sup> The MinMax measure was developed by M. RUŽIČKA 1958 for use in the field of phytogeography.

The second one is Cosine Delta,<sup>13</sup> which consist of a Cosine Distance function, but applied on z-score normalised features:

$$\text{cosine}(D, D') = 1 - \frac{\vec{f}(D) \cdot \vec{f}(D')}{\|\vec{f}(D)\|_2 \|\vec{f}(D')\|_2} \quad (\text{JANNIDIS et al. 2015: 9})$$

$$\text{with z-score: } \frac{f_i(D) - \mu_i}{\sigma_i} \quad (\text{JANNIDIS et al. 2015: 9}).$$

The cosine operates on vectors projected in a multi-dimensional space, and therefore is really useful as it can easily establish how the two documents are similar regardless of their size and words stock. Indeed, the angle between the two vectors is independent of their length in the same way that the angle between two segments is also independent of their length. It is also easier to interpret as it is a value of the interval  $[0, 1]$ ; the smaller the angle, the higher the similarity of the two texts (MOISL 2015: 95, 96, 200).

Word/ $n$ -gram frequencies follow Zipf's law of distribution. In other words, the frequency of any word is inversely proportional to its rank in the frequency table (ZIPF 1935). Therefore, the distance between two texts would be affected by a few top-scoring words. The z-score, introduced by BURROWS 2002, standardises word frequencies to overcome this problem inherent to the nature of language. For each word  $i$  in a given document  $D$ , it normalises the word's frequency over the whole corpus, so that the mean for each word is 0 and the standard deviation is 1 by subtracting the population mean  $\mu_i$  from the individual word's score and then dividing the difference by the standard deviation  $\sigma_i$  (EVERT et al. 2017: 6). The profile of the most frequent words' frequencies as a whole is more meaningful than some specific words (EVERT et al. 2017: 14), which means that the focus is more on many weak discriminators than on a small number of strong ones (BURROWS 2002: 268). We can consider this as a global approach on the whole words set.

On the other hand, the MinMax measure, reliant on counting common words/ $n$ -grams between documents, is size-dependent. First, the number of features will tend to increase with the length of the texts (MOISL 2015: 76), even if their topics are different. And it will perform worse in the case of big disproportion in the size of the compared documents.

<sup>13</sup> Developed in JANNIDIS et al. 2015 and EVERT et al. 2017, who have also demonstrated that this measure produces very good results compared to other distance metrics.

## Text segmentation

Due to its various linguistic peculiarities, even preliminary tasks such as word segmentation are non-trivial in Sanskrit. Not only because of the lack of white spaces between words, but also because of loose syntax, which gives weak indications of the presence of sentence boundaries (HELLWIG 2016). But Sanskrit text segmentation is made even more complex on account of a set of phonetic changes (sandhi) that occur at adjacent word boundaries. The contact phonemes of neighbouring words are changed and sometimes even merged. In that way, Sanskrit sentences appear as unseparated strings, incorporating multiple lexemes in forms that differ from their standard dictionary forms, making them difficult to recognise. Therefore, a simple maximum matching algorithm (PALMER 2010: 20) based just on lexical analysis is ineffective. Furthermore, sandhi resolution is non-deterministic, which means that different combinations of unsandhi words can result in the same merged sequence.<sup>14</sup> As a result, the same text can be segmented into several different sets of words. Thus, sandhi resolution in many cases depends on the semantic context of the full sentence. Until recently, this constituted a major obstacle to the automatic analysis of large corpora of Sanskrit texts.

In 2018, a new model designed to solve the sandhi problem was released by Oliver Hellwig and Sebastian Nehrlich, based on the character-level approach, as well as Neural Network and Deep Learning (HELLWIG and NEHRDICH 2018). They introduced innovative character-based models for Sanskrit word splitting (SWS) that outperform previous models by large margins, which was achieved by using as a base a new dataset for SWS made of sentences with manually validated splits. The model has been written in Python programming language and is based on TensorFlow, a symbolic math library dedicated to machine learning, developed by the Google Brain Team in 2015 and based on data flow and differentiable programming.<sup>15</sup> As with all machine learning systems, the purpose is to learn – based on a sample data – a desired behaviour in order to imitate it. In other words, machine learning systems can learn, on the basis of a sufficient number of examples, which we call a training set, a desired behaviour and then reproduce it.

Hellwig and Nehrlich released a new dataset based on the Digital Corpus of Sanskrit (DCS). Each sentence of the DCS has been re-analysed with the help of the `SanskritTagger` software. Lastly, the dataset is made up of the surface forms of sentences in the DCS to which we add the split points and sandhi rules proposed by the `Tagger`. According to KITAGAWA and KOMACHI 2017, the input can be enriched with multinomial split probabilities extracted from the training data.

<sup>14</sup> E.g. *tattvamasi* can, besides *tat tvam asi* “thou art that”, be tentatively separated as *tattvam asi* “thou art a *tattva* (principle)”.

<sup>15</sup> More information about TensorFlow can be found here: <https://www.tensorflow.org/>.

For almost all deep learning methods, the size of the training dataset is crucial. The one used by Hellwig and Nehrdich contains 561,596 sentences made up of 4,171,682 tokens. Of course, no less important is the quality and variety of the input (that is the training set). As the system will learn from the examples contained in the training set, its quality will directly impact the performance of the system. Thus, the data stream must be versatile and varied enough in order to obtain results that meet our needs, that is a system correctly reproducing the desired behaviour. In order to provide a sufficient variety of vocabulary, most sentences came from epic and scientific domains. Indeed, while most epic texts are composed in an easy, plain Sanskrit, the scholarly works tend to be much more elaborated. Furthermore, selecting both of the domains ensures a large enough coverage of the vocabulary necessary to finally obtain a system which will provide statistically reliable results – that is, in our case, the one that will correctly perform the operation on text with resolved sandhi.

For this authorship analysis, we gathered 82 texts made up of 1,307,610 word-strings before segmentation. We had to deal with two flaws – firstly, the system operates only on properly coded IAST (International Alphabet of Sanskrit Transliteration) words. As the system is operated on the character-level, any character incorrectly coded will not only be misinterpreted, but will also influence the results for the following characters and, finally, can impact the whole final result for a given sentence. Secondly, the maximum length of sentences to be segmented at one time is 128 characters. To overcome these limits, we wrote a basic Python script to ensure the pre-processing of the text by dividing it into smaller sequences of 128 characters and detecting any character not compatible with the IAST standard. Some word-sentences were even longer than 128 characters and, therefore, were not segmented correctly because the exceeding part of the word was skipped. As this type of problem is exceptionally uncommon, it should not have any impact on the final result of the authorship verification process. Finally, the computation was performed with Python v3.5.2, TensorFlow v1.8.0 and produced 2,287,451 words after segmentation. The estimated error ratio is about 15% on the level of text lines, which means that about 85% of all lines processed with the model do not contain wrong sandhi resolutions.

## **Texts preparation**

As indicated by KOPPEL and WINTER (2014: 5–6), imposters have to be chosen carefully. Imposters have to be in the same language, conceptually and temporarily close to the candidate author and to the disputed text. If imposters belong to radically different genres, false positive results might appear. Two web pages contain a sizeable number of Sanskrit texts that can be used as imposters. The first one is Göttingen Register of Electronic Texts in Indian Languages

(GRETIL), the second is *Advaitasāradā* (AŚ), which contains a number of texts in the Devanāgarī script attributed to Śaṅkara and to later Advaita Vedānta authors. For the purpose of this experiment, however, a number of other important texts that do not exist in electronic form were also prepared. Vimuktātman’s *Iṣṭasiddhi*, Sureśvara’s *Naiṣkarmyasiddhi* and a part of his *Bṛhadāraṇyakopaniṣadbhāṣya-vārtika* was prepared by performing OCR on scanned Devanāgarī texts that were further transliterated into the IAST standard. Also, some of the texts that will be used in the second part of the experiment, where prose texts attributed to Śaṅkara will be examined, do not exist in electronic form. Therefore, we prepared in the same way the *Adhyātmapaṭalavivarāṇa*, *Hastāmalakastotrābhāṣya* (HastBh), *Nṛsiṃha-(pūrva)-tāpanīyopaniṣadbhāṣya* (NṛsTBh), *Sanatsujātīyabhāṣya* (SanatBh), *Śvetāśvataropaniṣadbhāṣya* (ŚvUBh), *Viṣṇusahasranāmabhāṣya* and a rather small segment of the *Pātañjalayogaśāstravivarāṇa* (PātŚVi). The result was in some ways “noisy” because of mistakes that appear during OCR, especially where the Sanskrit text is scanned in lower resolution or the image is blurred. Some words might also be wrongly separated due to hyphenation when OCR fails to recognise it. However, “unclean” texts behaved well in the first experiment and were attributed correctly to their authors by the GI classifier.

### Manually segmented vs. automatically segmented corpus

In this part of the paper, we will first evaluate the method; the first results will suggest removing text-comment pairs and working with the most successful settings.

Regarding Śaṅkara, in this part of the experiment we use four works for which we have best indications that they were composed by Śaṅkara himself. The first is the *Brahmasūtrabhāṣya* (BSBh), which can be taken as a standard for determining Śaṅkara’s authorship. At the beginning of his *Pañcapādīkā*, Padmapāda mentions Śaṅkara by name as the author of the BSBh.<sup>16</sup> Sureśvara, who mentions Śaṅkara by name in *Naiṣkarmyasiddhi* 4.74 and 4.76, composed a commentary on the *Bṛhadāraṇyakopaniṣadbhāṣya* (BĀUBh) in which he mentions Śaṅkara as his teacher (commentary on BĀUBh 6,5.25).<sup>17</sup> Sureśvara also composed a commentary on the *Taittirīyopaniṣadbhāṣya* (TaittUBh) that may fall into the same category.<sup>18</sup> The fourth must be the *Upadeśasāhasrī*

<sup>16</sup> For the BSBh, the GRETIL edition will be used. It is not clear on what printed edition the GRETIL e-text was based. Also, the GRETIL e-text does not contain the introduction that we prepared for this experiment according to *Works of Śaṅkarācārya in original Sanskrit*. Vol. 1. Delhi: Motilal Banarsidass, 1964, reprint 2007.

<sup>17</sup> For the BĀUBh, the GRETIL edition will be used. It is not clear what printed edition served as a basis for the GRETIL e-text.

<sup>18</sup> For the TaittUBh, we used the GRETIL edition based on *Works of Śaṅkarācārya in original Sanskrit*. Vol. 1: *Ten Principal Upaniṣads with Śaṅkara-bhāṣya*. Delhi: Motilal Banarsidass, 1964, reprint 2007.

(Upad),<sup>19</sup> which is cited 20 times in Sureśvara's *Naiṣkarmyasiddhi* (MAYEDA 2006, vol. I: 45). Given the fact that Sureśvara explicitly mentions Śaṅkara as his teacher, it is quite safe to claim that Śaṅkara authored Upad.<sup>20</sup> For these four works – besides external evidence for Śaṅkara's authorship – internal evidence of similarity in teachings and terminology have already been presented in Indological scholarship.<sup>21</sup>

Also, in this part we shall assess which setup with regard to the distance metrics and choice of feature vectors yields best results. The text corpus we used was more than ten times larger than the corpus used in ANDRIJANIĆ 2020b.<sup>22</sup> Texts range from very short treatises, such as Nāgārjuna's *Yuktiṣaṣṭikakārikā* with 899 words, to voluminous works, such as Vācaspati Miśra's *Nyāyavārttikatātparyāṭikā*, the largest treatise on our list, with 167,357 words. However, we should keep in mind that short texts might be a problem, since they behave very unstably in multivariate calculations and tend to group with other small texts. The table in the appendix presents our complete corpus with word count. Most works are complete, except for ones marked with asterisk.

The manually segmented corpus used by ANDRIJANIĆ (2020a: 276 and 2020b: 110) yielded in its best setup a quite acceptable 83% of successful verification.<sup>23</sup> In the first step, we segmented automatically more or less the same corpus as used in ANDRIJANIĆ 2020a and 2020b. In the automatically segmented corpus of the same size and features, the rate of successful verifications dropped from 83% to 60%. However, the level of mistaken attribution (10%) remained the same. This is because the classifier did not make a decision in 20% cases with automatically segmented text-corpus. The reason for the lower success may need to be sought in the fact that the process of separating the sandhis is done with a 15% error rate, although it is questionable whether sandhi errors should have such an influence on the higher level task.

In the next step, we proceed with larger corpus in hope that a larger dataset might statistically compensate for flawed segmentation. Therefore, a corpus of 64 works (including Śaṅkara's works that are used as candidate texts) belonging to 36 authors was measured by MinMax and Cosine Delta distance measures. According to KESTEMONT et al. (2016: 90–91), the MinMax metric works better

<sup>19</sup> For the Upad, we used the GRETEL edition based on Mayeda's critical edition (MAYEDA 2006).

<sup>20</sup> Cf. MAYEDA 2006: 44–49 for further detailed argumentation.

<sup>21</sup> For the BĀUBh, see MARSCHNER 1933; for the Upad, see MAYEDA 2006: 23–44.

<sup>22</sup> In ANDRIJANIĆ 2020b corpus consisted of 25 works of 11 authors with altogether 157,592 words.

<sup>23</sup> In two Andrijanić's studies slightly different text corpus of known authors was used. The best performing setup included measurement of a feature vector consisting of word unigrams, and the best distance measure was MinMax.

than Manhattan and Cosine (not to be confused with Cosine Delta). EVERT et al. 2017 showed that Cosine Delta produces very good results compared to other distance metrics, although they did not compare it to the MinMax metrics.<sup>24</sup> In ANDRIJANIĆ 2020a and 2020b, MinMax performed slightly better than Cosine Delta, while both significantly outperformed Burrows' Delta (BURROWS 2002). In our experiment with trigrams (large automatically segmented corpus) measured with Cosine Delta, we obtained only 61% of successful verifications, 10% mistakes and for the rest (29%) classifier did not reach a verdict.

By inspecting these results more carefully, a few strange issues have arisen. For example, Śaṅkara's BSBh ended up in a grey zone; the classifier failed to attribute it to Śaṅkara. But in the manually segmented corpus it is correctly verified. Thus, let us analyse what might have been the problem – a bad segmentation or something else? When we scanned the whole corpus with the GI classifier, it turned out that the GI recognises Śaṅkara's BSBh and Vācaspati Mīśra's *Bhāmatī* as works of the same author. As *Bhāmatī* is a commentary on the BSBh, the *Bhāmatī* reiterates or glosses over a significant amount of words, and this must have interfered in the classification process. Thus, when the *Bhāmatī* was excluded from the imposters list, the BSBh was correctly attributed to Śaṅkara. On the other hand, when the BSBh was excluded from the imposters list, the *Bhāmatī* was correctly attributed to Vācaspati. The same problem appeared with all the other pairs of commentaries: Śaṅkara's TaittUBh and Sureśvara's commentary TaittUBhV; Nāgārjuna's *Mūlamadhyamakakārikās* (MMK) and Candrakīrti's commentary *Prasannapadā*; Udayana's *Nyāyavārttikatātparyapariśuddhi* and Vācaspati's *Nyāyavārttikatātparyāṭikā*. This indicates that the classifier is very sensitive when it comes to recognition of related works. Indeed, when the *Bhāmatī* and TaittUBhV were excluded from the list of imposters, the classifier attributed both the BSBh and TaittUBh correctly to Śaṅkara and vice versa; when the BSBh and TaittUBh were removed from the list of imposters, the *Bhāmatī* and TaittUBhV were correctly attributed to Vācaspati and Sureśvara.<sup>25</sup> The same happened for the MMK, which was verified as Nāgārjuna's work when the *Prasannapadā* was taken out of the imposters list; when the *Prasannapadā* was in the list, the classifier did not reach a decision.<sup>26</sup> A notable example comes from Maṇḍana Mīśra, whose works at first resisted correct attribution. However,

<sup>24</sup> See also EDER 2018.

<sup>25</sup> The same happened with Vācaspati and Udayana; when the authorship of Udayana's *Nyāyavārttikatātparyapariśuddhi* is examined, Vācaspati's *Nyāyavārttikatātparyāṭikā* should be removed from the imposters list and vice versa.

<sup>26</sup> It did not work so well the other way around with trigrams measured with Cosine Delta. While the MMK was on the imposters list, the *Prasannapadā* was classified as not authored by Candrakīrti. However, when Nāgārjuna's MMK was removed from the imposters list, the *Prasannapadā* reached a score of 0.39, meaning the classifier could not make a decision. Nevertheless, this is better than reaching a wrong decision. However, MinMax in both setups (trigrams and unigrams) and Cosine Delta with unigrams confirmed Candrakīrti's authorship.

when we take a closer look at the *Vibhramaviveka* and *Brahmasiddhi*, we notice that the *Vibhramaviveka* is a short metrical work and that the *Brahmasiddhi* is a voluminous mixture of prose and metrical material. Thus, we experimentally divided the *Brahmasiddhi* into the metrical and prose parts and the classifier managed to attribute all three samples correctly to the same author.

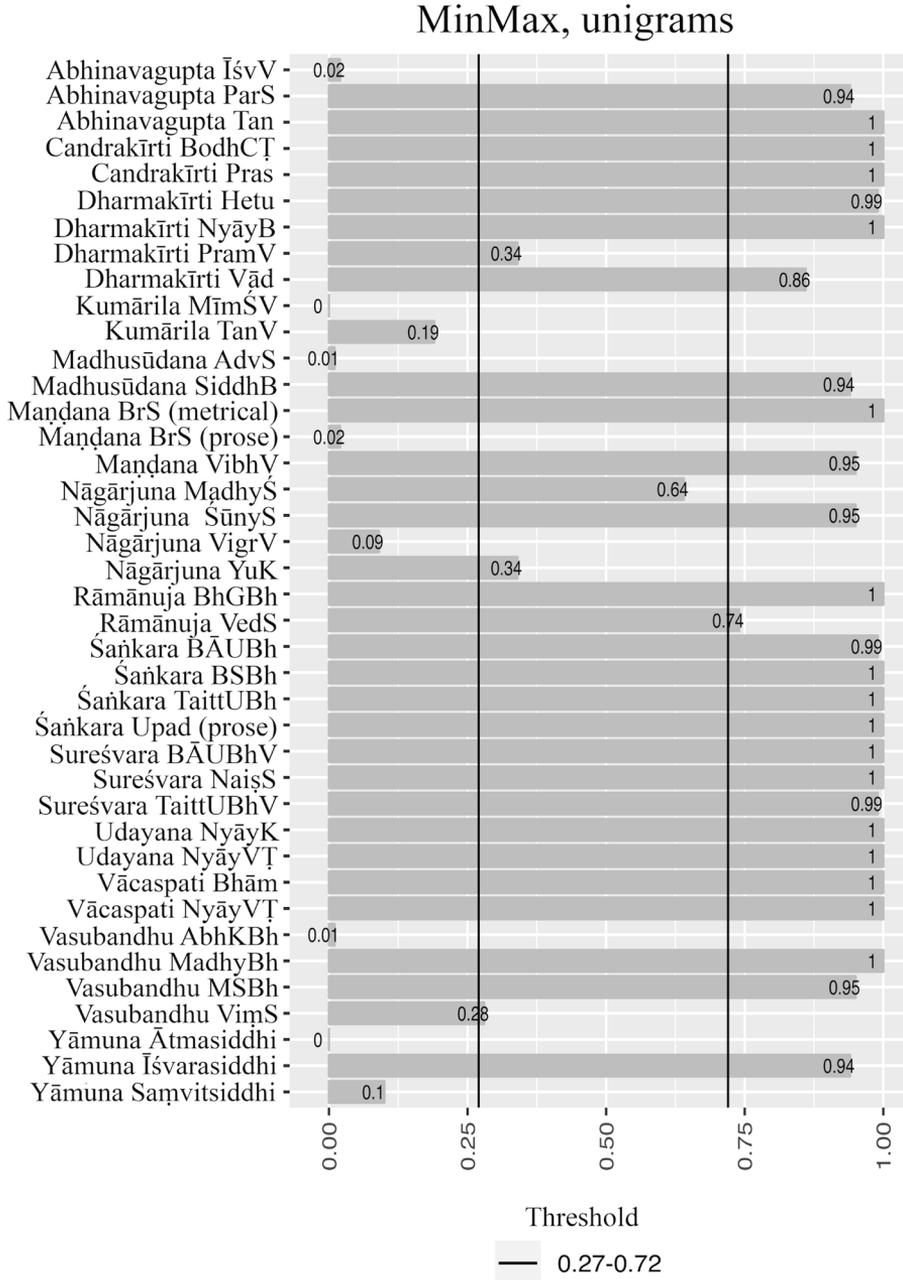
The results of the experiment conducted on texts of undisputed authorship show a significant improvement of results with these adjustments (exclusion of commentaries from the imposters list and the distinction between the metrical and prose texts). We tried four different setups; two different distance measures: MinMax and Cosine Delta; and two types of text segmentations: word unigrams (words) and trigrams.<sup>27</sup>

The Cosine Delta obviously outperformed MinMax in our experiment, most probably because the corpus contains texts of very different sizes. On the other hand, in all four setups, all four of Śaṅkara's works were correctly attributed, thus confirming Śaṅkara's strong authorship signal.

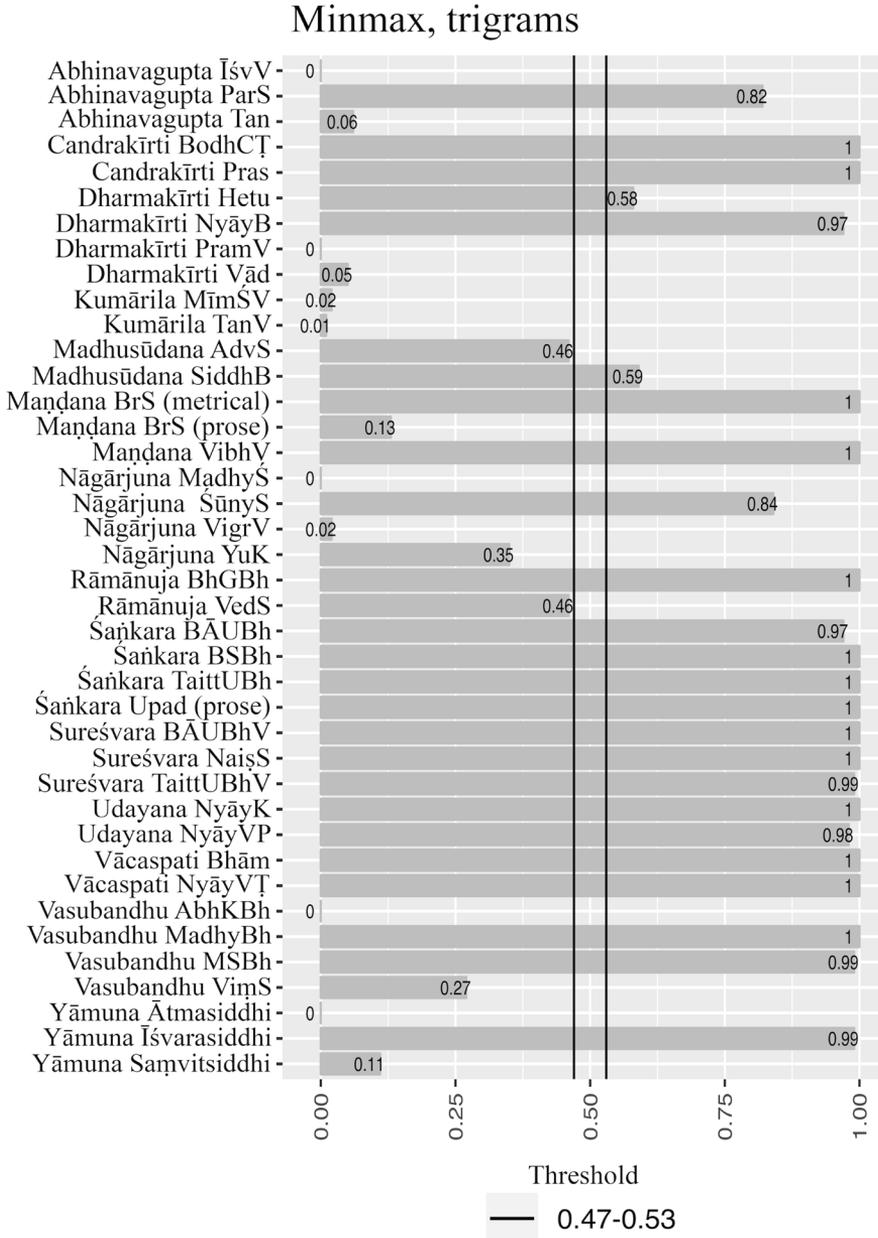
With a success rate of 77.5–80% obtained on a large text corpus, we can be quite satisfied. However, the mistake ratio should also be taken into account. Trigrams measured using the Cosine Delta have an error ratio of 7.5%, while with word unigrams it is 10%. Therefore, both trigram and unigram frequency vectors measured with the Cosine Delta appear to be the most successful setups.

---

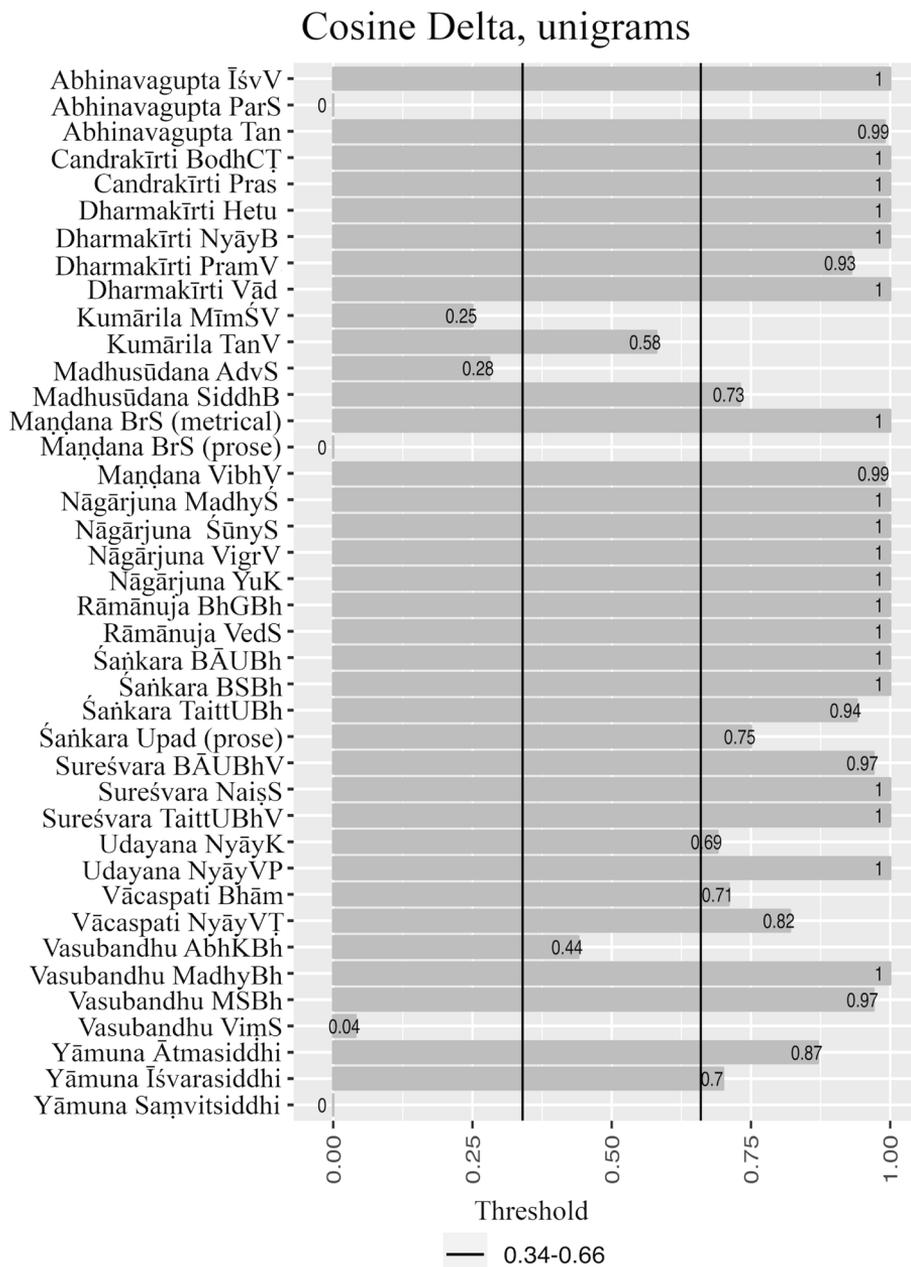
<sup>27</sup> These trigrams are different than in ANDRIJANIĆ 2020b, where trigrams were made out of raw unsegmented text corpus. In this paper, trigrams are executed on segmented texts with resolved sandhis, thus catching spaces between words.



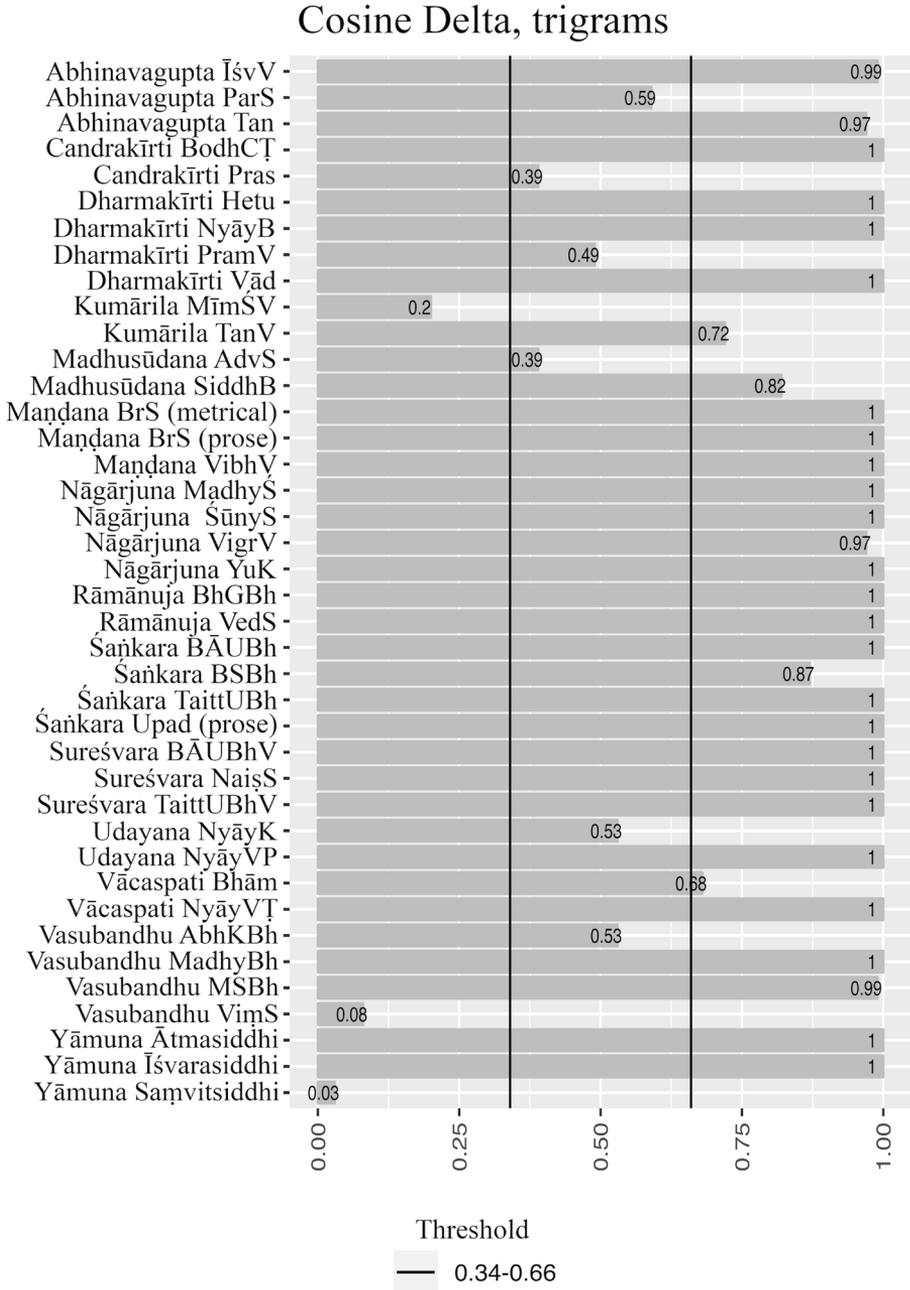
**Fig. 1.** Shows the results of the experiment in which texts were segmented into the word unigrams and measured with MinMax; that setup yielded 67.5% successful attributions.



**Fig. 2.** Shows the results of the experiment in which texts were segmented in the word trigrams and measured with MinMax; that setup yielded only 60% successful attributions.



**Fig. 3.** Shows the results of the experiment in which texts were segmented in the word unigrams and measured with the Cosine Delta; that setup yielded 80% successful attributions.



**Fig. 4.** Shows the results of the experiment in which texts were segmented in the word trigrams and measured with the Cosine Delta; that setup yielded 77.5% successful attributions.

## Authenticity of prose texts attributed to Śaṅkara

In this section, we evaluate prose works traditionally attributed to Śaṅkara. The candidate set consists of the same four works (BSBh, BĀUBh, TaittUBh and the prose part of the Upad) used in the previous experiment.

However, the question arises as to which prose works attributed to Śaṅkara for which we have no external evidence should be evaluated for Śaṅkara's authorship. The large body of works that are attributed to Śaṅkara has already been filtered by editors of Śaṅkara's complete works, and then by a number of scholars. Our choice of works to be tested is a kind of concurrence of these previous attempts. BELVALKAR (1929: 218) pointed out that in addition to the previously mentioned four works, these works probably come from Śaṅkara himself: *Aitareyopaniṣadbhāṣya* (AiUBh),<sup>28</sup> ChUBh, *Bhagavadgītābhāṣya* (BhGBh), ĪUBh, KaUBh, *Kenopaniṣad-(pada)-bhāṣya* (KeUBh), *Muṇḍakopaniṣadbhāṣya* (MuUBh), *Praśnopaniṣadbhāṣya* (PraśUBh). All these works are included in collected works of Śaṅkara<sup>29</sup> and they all pass HACKER's criteria (HACKER 1978) of being attributed to Śaṅkara-*bhagavat(pūjya)-pāda* in colophons. Moreover, ĪUBh, KaUBh, BhG and ChUBh already passed two stylometric tests on limited corpus (ANDRIJANIĆ 2020a, 2020b). ĪUBh, KeUBh, KaUBh, BhGBh and *Gauḍapādīyabhāṣya* (GauBh) also pass HACKER's (1950) terminological criteria (MAYEDA 1965a, 1965b, 1967, 1967–1968; ANDRIJANIĆ 2020a). On the other hand, we selected a number of questionable works: ŚvUBh, SanatBh, NṛsTBh, HastBh, *Viṣṇusahasranāmabhāṣya*, *Adhyātmapaṭalavivarāṇa* and *Pātañjalayogaśāstravivarāṇa*. In the *Śaṅkaradigvijaya* (ŚDV), Śaṅkara's biography composed between 1650 and 1789 (BADER 2000: 55), Śaṅkara's writings are enumerated in vs. 6,61–63. These include the Upad, BSBh,

<sup>28</sup> All printed editions of Śaṅkara's commentary on the AiU include commentaries on three *adhyaṅgas* of the second *āranyaka* of the *Aitareyāranyaka* (2,4–6) that can be understood as *Aitareyopaniṣad* proper. However, in a number of manuscripts, a larger commentary is preserved, that comprises a running commentary on full second and third *āranyaka*. BELVALKAR 1930: 242 considers this larger commentary authentic. For a comprehensive overview of the problem see DAVID 2017, who also argues in favour of the authenticity of the “longer” Bhāṣya (DAVID 2017: 733–745). At this moment we shall evaluate only the shorter text, at least until the critical edition of the “longer” version, being prepared by Hugo David, will be available.

<sup>29</sup> First collection of Śaṅkara's works appears to be *Sri Sankaracharya's Miscellaneous Works* in 4 vols., ed. by A. Mahadeva SASTRI and K. RANGACHARYA (Mysore: Government Branch Press, 1898–1899). The *Works of Sri Sankaracharya* (Memorial edition) (Srirangam: Sri Vani Vilas Press, 1910) was printed in 20 vols. It was retyped and printed in 11 volumes in Śrīraṅgam as *Śrīśaṅkaragrānthāvalīh*. The 1910 edition was rearranged in 10 vols. in the *Complete Works of Sri Sankaracharya in the Original Sanskrit*, Madras: Samanta Books, 1981–1983. Widely used Motilal Banarsidass edition *Works of Śaṅkarācārya in Original Sanskrit* in 3 vols. (1964–1985) is based on the four-volume edition edited by Hari Raghunath SASTRI (Poona: Ashtekar & Co.). See REIGLE and REIGLE 2005.

commentaries on the Upaniṣads,<sup>30</sup> BhGBh, SanatBh and NṛsTBh. Cidvilāsa's *Śaṅkaravijayavilāsa* 10,2–3<sup>31</sup> mentions BSBh, BhGBh, commentaries on ten Upaniṣads, the *Viṣṇu-* and *Rudrasahasranāma*. For ŚvUBh ANDRIJANIĆ 2019 presented arguments that the work is several centuries later than Śaṅkara. Nevertheless, we conducted the GI test to see whether it will confirm Andrijanić's arguments.<sup>32</sup> ŚvUBh and HastBh meet Hacker's colophon criteria, while SanatBh and NṛsTBh partly meet Hacker's colophon criteria as they are sometimes attributed to Śaṅkarācārya and sometimes to Śaṅkarabhagavat. To these works we also added the *Lalitātrisatistotrabhāṣya* because it is included in the VVP 18 edition of Śaṅkara's collected works. PātŚVi is not included in any collection of Śaṅkara's works, but it is included in the experiment because a number of scholars have argued in favour of its authenticity. We used only the critically edited text from PātŚVi 1.1 (HARIMOTO 2014: 171–183) and 1,23–28 (HARIMOTO 2014: 47–84).

The two tables below list works attributed to Śaṅkara that we have examined. In the first column is the title of the work together with the edition on the basis of which the test was made. The second column contains brief remarks about previous scholarship on authorship. The third column contains GI results obtained in the most successful setup (trigrams measured with the Cosine Delta metric). If the result is above 0.66, the GI classifies the work as authentic (i.e. the classifier considers that the author is the same as the author of BSBh, BĀUBh, TaittUBh and the prose portion of the Upad). If the result is below 0.34, the GI renders it inauthentic. Numbers between 0.34 and 0.66 indicate a “grey zone”, where the classifier did not reach a verdict. As words measured with the Cosine Delta reached a similar result as the trigrams, we indicate the result obtained with word unigrams in brackets.

<sup>30</sup> Dhanapati Sūri in *Ḍiṇḍima* 6,61, a commentary on the ŚDV from 1798, enumerates the Upaniṣads that were commented by Śaṅkara: the ĪUBh, KeUBh, KaUBh, PraśUBh, MuUBh, AiUBh, ChUBh, BĀUBh and TaittUBh. Acyuta, another commentator on the ŚDV, in his *Advaitarājyalakṣmī* from 1805 (information on Acyuta's date is from HACKER 1951: 28), adds the *Viṣṇusahasranāmabhāṣya* and the *vākya* and *pada* versions of the KeUBh. It is worth noting that both do not mention the ŚvUBh.

<sup>31</sup> Between the 14th and 18th cent. (BADER 2000: 24).

<sup>32</sup> The ŚvUBh does not meet Hacker's terminological criteria. Terms and concepts such as *saccidānanda* that appear in later Advaita Vedānta are used, together with long purāṇic quotations. The second important problem is a quotation from the *Bṛhatsaṃhitā* dated to the 12th cent. For further details and a review of previous views on the authenticity of the ŚvUBh see ANDRIJANIĆ 2019.

**Table 1.** Works verified as authentic in comparison to BSBh, BĀUBh, TaittUBh and the prose section of the Upad. Feature vector consists of relative frequencies of trigrams (unigrams in brackets), distance measure is the Cosine Delta.

Title	Number of words	Remarks	GI result (0.34–0.66)
<i>Adhyātmapaṭala-vivarāṇa</i> (TSS 41)	3,460	“More or less debatable” (BELVALKAR 1929: 219). HACKER 1968–1969: 147 considers it authentic. NAKAMURA 1983: 306 considers it possible that Śaṅkara is the author. PANDE 1994: 109–110, 113 and LEGGET 1978: 218–228 argue for its authenticity.	0.88 (0.84)
<i>Aitareyopaniṣad-bhāṣya</i> (GRETIL)	6,904	The longer, unpublished, commentary is, according to BELVALKAR 1930: 242, authentic. Meets Hacker’s colophon criterion (HACKER 1978: 46).	1 (0.97)
<i>Bhagavadgītā-bhāṣya</i> (GRETIL)	28,624	Meets Hacker’s terminological criteria (MAYEDA 1965a). Meets also Hacker’s colophon criterion (HACKER 1978: 46).	1 (1)
<i>Chāndogypopaniṣad-bhāṣya</i> (GRETIL)	49,930	Verified by the GI method as genuine against a limited corpus (ANDRIJANIĆ 2020b). Meets also Hacker’s colophon criterion (HACKER 1978: 46).	1 (1)
<i>Gauḍapādīya-bhāṣya</i> (GRETIL)	18,507	“More or less debatable” (BELVALKAR 1929: 218). Meets Hacker’s colophon criterion (HACKER 1978: 46). VETTER 1968/69 argues for its authenticity. Hacker considers it authentic (HACKER 1968–1969, 1972), noting few cautious remarks (1968–1969: 115–117, fn. 2). Also meets Hacker’s terminological criteria (MAYEDA 1967–1968).	1 (1)
<i>Īsopaniṣadbhāṣya</i> (GRETIL)	2,232	According to ANDRIJANIĆ 2020a, meets most of Hacker’s terminological criteria, while GI also verifies it as genuine against the limited imposter corpus. Meets Hacker’s colophon criterion (HACKER 1978: 45).	0.93 (0)
<i>Kaṭhopaniṣad-bhāṣya</i> (GRETIL)	11,237	According to ANDRIJANIĆ 2020a, meets Hacker’s terminological criteria, while GI also verifies it as genuine against the limited imposter corpus. Meets Hacker’s colophon criterion (HACKER 1978: 46).	0.98 (1)
<i>Kenopaniṣad-(pada)-bhāṣya</i> (GRETIL)	6,048	Meets Hacker’s terminological criteria (MAYEDA 1967). Also meets Hacker’s colophon criterion (HACKER 1978: 46).	1 (1)

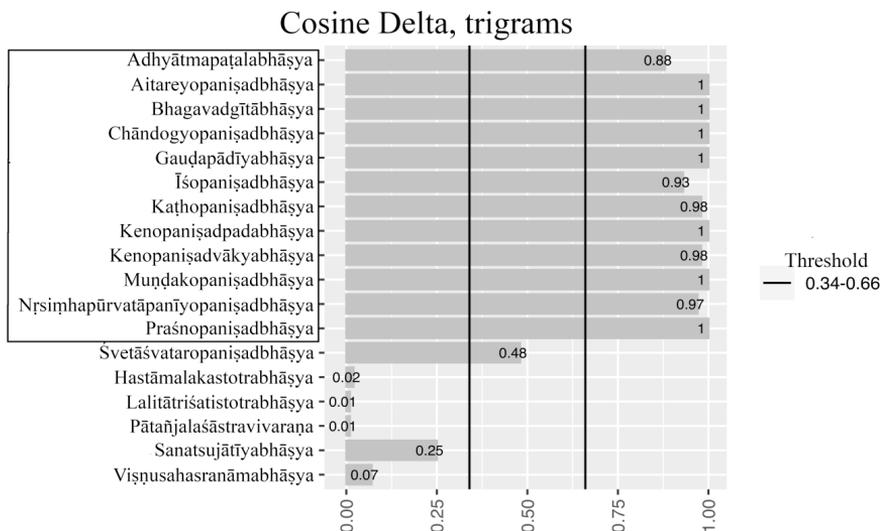
Title	Number of words	Remarks	GI result (0.34–0.66)
<i>Kenopaniṣad-(vākya)-bhāṣya</i> (AŚ)	4,990	“More or less debatable” (BELVALKAR 1929: 218). Meets Hacker’s terminological criteria (MAYEDA 1967). Also meets Hacker’s colophon criterion (HACKER 1978: 46).	0.98 (1)
<i>Muṇḍakopaniṣad-bhāṣya</i> (GRETIL)	5,857	“Most probably” authentic (BELVALKAR 1929: 218). Meets Hacker’s colophon criterion (HACKER 1978: 46).	1 (1)
<i>Nṛsiṃha-(pūrva)-tāpanīyopaniṣad-bhāṣya</i> (VVP 10)	21,777	JACOB 1886: 70 emphatically denies Śaṅkara’s authorship. According to BELVALKAR 1929: 218, “More or less debatable”. Attributed in colophons both to Śaṅkara- <i>ācārya</i> and - <i>bhagavat</i> (HACKER 1978: 48).	0.79 (0.93)
<i>Praśnopaniṣad-bhāṣya</i> (GRETIL)	8,117	“Most probably” authentic (BELVALKAR 1929: 218). Meets Hacker’s colophon criterion (HACKER 1978: 46).	1 (1)

**Table 2.** Works not verified as authentic in comparison to BSBh, BĀUBh, TaittUBh and the prose section of the Upad. Feature vector consists of relative frequencies of trigrams, distance measure is the Cosine Delta.

Title	Number of words	Remarks	GI result (0.34–0.66)
<i>Hastāmālakastotra-bhāṣya</i> (AŚ)	3,491	“More or less debatable”, according to BELVALKAR 1929: 218. According to PANDE 1994: 110, it should be “confidently excluded” from the list of Śaṅkara’s writings.	0.02 (0.03)
<i>Lalitātriśatistotra-bhāṣya</i> (VVP 18)	21,345	“Certainly spurious” according to BELVALKAR 1929: 219. According to SANDERSON 2017: 7 fn. 7, the attribution to Śaṅkara- <i>bhagavat</i> from the colophon is surely false.	0.01 (0)
<i>Śvetāśvataropaniṣad-bhāṣya</i> (ĀāSS 17)	17,287	“More or less debatable” (BELVALKAR 1929: 218). According to ANDRIJANIĆ 2019 cannot be ascribed to Śaṅkara.	0.48 (0)
<i>Sanatsujātīya-bhāṣya</i> (VVP 13)	18,707	“More or less debatable” (BELVALKAR 1929: 219). HACKER 1978: 50–51 raised a number of arguments against Śaṅkara’s authorship. In colophons it is attributed both to Śaṅkara- <i>ācārya</i> and - <i>bhagavat</i> (HACKER 1978: 48). PANDE 1994: 109, 113 argues against Śaṅkara’s authorship.	0.25 (0)

Title	Number of words	Remarks	GI result (0.34–0.66)
<i>Viṣṇusahasranāmabhāṣya</i> (VVP 13)	22,306	“More or less debatable” (BELVALKAR 1929: 219). SASTRY 1980: xxi–xxii argued for its authenticity. PANDE 1994: 109, 113 argues against Śaṅkara’s authorship.	0.07 (0)
<i>Pātañjalayogaśāstravivarāṇa</i> (HARIMOTO 2014)	8,228	“Certainly spurious” according to BELVALKAR 1929: 218. Meets Hacker’s terminological criteria (HARIMOTO 2014: 244–247) and the colophon criterion “but not without some caveats” (HARIMOTO 2014: 243). PātSVi is not included in any complete works of Śaṅkara.	0.01 (0.01)

**Fig. 5.** Results from Tables 1 and 2. The table shows the results for works attributed to Śaṅkara. For the works outlined with a dash, Śaṅkara’s authorship has been confirmed, while it has not been confirmed for the others.



## Concluding observations

- a) The GI result confirmed Belvalkar’s intuition (1929: 218) and verified all 11 titles from his list of works that most likely come from Śaṅkara himself. Almost the same result was obtained when word frequency vectors were measured, with the only exception of *Īsopaniṣadbhāṣya*. However, the ĪUBh was confirmed by different setups in ANDRIJANIĆ 2020a and 2020b, and in our study trigrams measured using the Cosine Delta and MinMax, together with word unigrams measured with MinMax confirmed Śaṅkara’s authorship. To summarise, the *Adhyātmapaṭalavivarāṇa*, *Aitareyopaniṣadbhāṣya*,

*Bhagavadgītābhāṣya*, *Chāndogyopaniṣadbhāṣya*, *Gauḍapādīyabhāṣya*, *Īśopaniṣadbhāṣya*, *Kāṭhopaniṣadbhāṣya*, *Kenopaniṣadbhāṣya* (*pada* and *vākya*), *Muṇḍakopaniṣadbhāṣya*, *Nṛsimha-(pūrva)-tāpanīyopaniṣadbhāṣya* and *Praśnopaniṣadbhāṣya* are verified by most GI setups as written by the same author who composed the *Brahmasūtrabhāṣya*, *Bṛhadāraṇyakopaniṣadbhāṣya*, *Taittirīyopaniṣadbhāṣya*, and the prose part of the *Upadeśasāhasrī*. It is indeed notable that the list is almost the same as Hacker's list of authentic works (HACKER 1968–1969: 147), which also includes the *Adhyātmapaṭalavivarāṇa*. The only exception from Hacker's list is the NṛTBh, which is verified as Śaṅkara's by GI in all setups. All works that Mayeda and Andrijačić subjected to Hacker's terminological analysis were also confirmed. In this way, the GI analysis largely confirmed traditional philological analysis, with an exception of the PātŚVi. For the PātŚVi there is no evidence against Śaṅkara's authorship, and some arguments even speak in favour of its authenticity. It should be noted that only a small part of the PātŚVi was examined in our analysis and that it is not impossible that, if a larger text sample was used, the result might be different.

- b) The experiment with the GI authorship verification framework conducted on Sanskrit philosophical texts showed that the classifier is quite reliable in identifying authors of undisputed texts and confirms the superiority of analysis based on *n*-grams over the content-words based one. Moreover, it seems that text segmentation is a prerequisite for this kind of stylometric Sanskrit analysis as the sandhi rules tend to decrease the stylometric signal. The classifier appears to be highly sensitive when it attributed commentaries on the same works, in which many words glossed over from the original text are repeated, to the same authors. This shows sensitivity, but also calls for caution when choosing imposters and candidate authors. We conclude that commentaries by different authors on the same works and works that comment on each other should be excluded from the test.
- c) The third important issue is that the GI classifier is sometimes confused in verifying prose and metrical works that belong to the same author. The reason for this is that authors possibly had to choose words differently in order to fit the metrical scheme. On the other hand, Sureśvara's works, which are all in the *śloka* meter, were verified by the classifier as authored by the same hand. It is also important to note that GI did not confuse them with other authors who composed their texts in the *śloka* meter. Therefore, if we try to establish the authorship of a prose text, it might be better in some cases to take only those candidate texts which are also in prose, and vice versa. This is important for the future evaluation of the numerous metrical works attributed to Śaṅkara. It would be less reliable to take Śaṅkara's prose commentaries as candidate texts. In this case, the questionable metrical

works attributed to Śaṅkara should probably be judged only in relation to the metrical part of the *Upadeśasāhasrī*.

- d) As computers' power is growing and every day more complex operations become easier to perform, we are witnessing a big change in the field of author studies. Automatic segmentation and sandhi are no longer a problem to computer-assisted Sanskrit texts analysis, and we are now able to analyse huge texts' corpora. Thus, in the future we will see many breakthroughs in the field of computational stylometry to assess authorship verification and attribution, potentially throughout entire literature.

## Supplementary Material

All additional material needed to recreate the experiment can be found at: <https://github.com/JacekBakowski/stylometry/tree/main/papers/2024-otao> (accessed 18 January 2024).

## Appendix

Table with the texts used in the first experiment. Most of the texts are complete, except the texts marked with an asterisk.

Author	Work	Number of words
Abhinavagupta	<i>Īśvarapratyabhijñāvimarśinī</i>	43,031
	<i>Paramārthasāra</i>	1,739
	<i>Tantrāloka</i>	88,351
Annambhatta	<i>Tarkasaṅgraha</i>	1,974
Asaṅga	<i>Abhidharmasamuccaya</i>	24,736
Bhāskara	<i>Bhagavadgītābhāṣya</i>	27,195
Candrakīrti	<i>Bodhisattvayogācāracaṭuḥśatakaṭikā</i>	23,224
	<i>Prasannapadā</i>	78,235
Dharmakīrti	<i>Hetubindu</i>	6,289
	<i>Nyāyabindu</i>	2,359
	<i>Pramāṇavārttika</i>	16,255
	<i>Vādanyāya</i>	10,104
Gaṅgeśa	<i>Tattvacintāmaṇi</i>	34,249
Jayarāśi	<i>Tattvopaplavasimha</i>	14,453
Kavirājayati	<i>Sāṃkhyatattvaprādīpa</i>	4,924

Author	Work	Number of words
Kumārila	<i>*Mīmāṃsāslokaṽrttika</i>	7,289
	<i>*Tantravārttika</i>	5,614
Madhusūdana	<i>Advaitasiddhi</i>	133,946
	<i>Siddhāntabindu</i>	8,560
Madhva	<i>Anuvyākhyāna</i>	29,255
Maṇḍana Mīśra	<i>Brahmasiddhi</i>	40,018
	<i>Vibhramaviveka</i>	2,002
Māṭhara	<i>Māṭharavṛtti</i>	17,918
Nāgārjuna	<i>Madhyamakaśāstra</i>	6,566
	<i>Śūnyatāsaptati</i>	1,170
	<i>Vigrahavyāvartanī</i>	6,052
	<i>Yuktiṣaṣṭikakārikā</i>	899
Padmapāda	<i>Pancapādikā</i>	28,574
Praśastapāda	<i>Pādārthadharmasaṁgraha</i>	11,073
Rāmānuja	<i>Bhagavadgītābhāṣya</i>	40,026
	<i>Vedārthasaṁgraha</i>	18,830
Śabara	<i>Mīmāṃsāsūtrabhāṣya</i>	123,358
—	<i>Sāṁkhyaparibhāṣā</i>	3,714
Śaṅkara	<i>Brahmasūtrabhāṣya</i>	109,993
	<i>Bṛhadāraṇyakopaniṣadbhāṣya</i>	101,952
	<i>Taittirīyopaniṣadbhāṣya</i>	17,195
	<i>Upadeśasāhasrī (Gadya)</i>	5,415
—	<i>Sarvamataṣaṁgraha</i>	7,716
Sthiramati	<i>Triṁśikāvijñāptibhāṣya</i>	8,727
Sureśvara	<i>*Bṛhadāraṇyakopaniṣadbhāṣyavārtika</i>	31,146
	<i>Naiṣkarmyasiddhi</i>	13,391
	<i>Taittirīyopaniṣadbhāṣyavārtika</i>	15,499
Toṭaka	<i>Śrutisārasamuddharaṇa</i>	3,781
Udayana	<i>Nyāyakusumāñjali</i>	34,547
	<i>Nyāyavārtikatātparyapariśuddhi</i>	86,988
Vācaspati Mīśra	<i>Bhāmatī</i>	152,511
	<i>Nyāyavārtikatātparyatīkā</i>	167,357

Author	Work	Number of words
Vasubandhu	<i>Abhidharmakośabhāṣya</i>	7,711
	<i>Madhyāntavibhāgabhāṣya</i>	6,825
	<i>Mahāyānasūtrālamkārahāṣya</i>	23,432
	<i>Viṃśatikasiddhi</i>	2,252
Vātsyāyana	<i>Nyāyasūtrabhāṣya</i>	42,189
Dharma-rājadhvarīndra	<i>Vedāntaparibhāṣā</i>	12,119
Sadānanda	<i>Vedāntasāra</i>	3,809
Veṅkaṭanātha	<i>Nyāyapariśuddhi</i>	27,503
Vijñānabhikṣu	<i>Sāṃkhyasāra</i>	7,994
Vimuktātman	<i>*Iṣṭasiddhi</i>	35,123
Yāmuna	<i>Samvitsiddhi</i>	3,366
	<i>Īśvarasiddhi</i>	2,324
	<i>Ātmasiddhi</i>	10,799
—	<i>Yuktidīpikā</i>	54,988

## Funding

This research was funded in part – covering entire J. B. contribution – by the National Science Centre, Poland, grant number 2021/43/O/HS2/02392.

## Abbreviations and primary sources

ĀāSS 17 *Kṛṣṇayajurvedīyaśvetāśvataropaniṣacchāmkarabhāṣyopetā, tathā Śaṃkarānandakṛtā Śvetāśvataropaniṣaddīpikā, Nārāyaṇakṛtā Śvetāśvataropaniṣaddīpikā, Vijñānabhagavatkrtaṃ Śvetāśvataropaniṣadvivaraṇam*. Ed. by V. G. Āpaṭe. Ānandāśrama-saṃskṛtagranthāvalīḥ 17, 1890.

AbhKBh *Abhidharmakośabhāṣya*

AdvS *Advaitasiddhi*

AiU *Aitareyopaniṣad*

AiUBh *Aitareyopaniṣadbhāṣya*

AŚ *Advaitaśāradā*. <https://advaitasharada.sringeri.net> (accessed 18 January 2024).

BĀUBh *Bṛhadāraṇyakopaniṣadbhāṣya*

BĀUBhV	<i>Bṛhadāranyakopaniṣadbhāṣyavārtika</i>
Bhām	<i>Bhāmatī</i>
BhG	<i>Bhagavadgītā</i>
BhGBh	<i>Bhagavadgītābhāṣya</i>
BodhCT	<i>Bodhisattvayogācāracaṭuṣṭakaṭīkā</i>
BrS	<i>Brahmasiddhi</i>
BS	<i>Brahmasūtra</i>
BSBh	<i>Brahmasūtrabhāṣya</i>
ChUBh	<i>Chāndogyopaniṣadbhāṣya</i>
DCS	Digital Corpus of Sanskrit. <a href="http://www.sanskrit-linguistics.org/dcs/">http://www.sanskrit-linguistics.org/dcs/</a> (accessed 18 January 2024).
GauBh	<i>Gauḍapādīyabhāṣya</i>
GI	General imposters
GRETEL	Göttingen Register of Electronic Texts in Indian Languages. <a href="http://gretel.sub.uni-goettingen.de/gretel.html">http://gretel.sub.uni-goettingen.de/gretel.html</a> (accessed 18 January 2024).
HastBh	<i>Hastāmalakastotrabhāṣya</i>
Hetu	<i>Hetubindu</i>
IAST	International Alphabet of Sanskrit Transliteration
ĪsvV	<i>Īśvarapratyabhijñāvimarśinī</i>
ĪUBh	<i>Īśopaniṣadbhāṣya</i>
KaUBh	<i>Kaṭhopaniṣadbhāṣya</i>
KeUBh	<i>Kenopaniṣadbhāṣya</i>
MadhyBh	<i>Madhyāntavibhāgabhāṣya</i>
MadhyŚ	<i>Madhyamakaśāstra</i>
MīmŚV	<i>Mīmāṃsāslokaṣṭakavārttika</i>
MMK	(Nāgārjuna's) <i>Mūlamadhyamakakārikā</i>
MSBh	<i>Mahāyānasūtrālamkārahāṣya</i>
MuUBh	<i>Muṇḍakopaniṣadbhāṣya</i>
NaiṣS	<i>Naiṣkarmyasiddhi</i>
NṛsTBh	<i>Nṛsimha-(pūrva)-tāpanīyopaniṣadbhāṣya</i>
NyāyB	<i>Nyāyabindu</i>
NyāyK	<i>Nyāyakusumāñjali</i>
NyāyVṬ	<i>Nyāyavārttikatātparyapariśuddhi</i>

ParS	<i>Paramārthasāra</i>
PātŚVi	<i>Pātañjalayogaśāstravivarāṇa</i>
PramV	<i>Pramāṇavārttika</i>
Pras	<i>Prasannapadā</i>
PraśUBh	<i>Praśnopaniṣadbhāṣya</i>
SanatBh	<i>Sanatsujātīyabhāṣya</i>
ŚDV	<i>Śāṅkaradigvijaya</i>
SiddhB	<i>Siddhāntabindu</i>
ŚūnyS	<i>Śūnyatāsaptati</i>
ŚvUBh	<i>Śvetāśvataropaniṣadbhāṣya</i>
SWS	Sanskrit word splitting
TaittUBh	<i>Taittirīyopaniṣadbhāṣya</i>
TaittUBhV	(Sureśvara's) <i>Taittirīyopaniṣadbhāṣyavārtika</i>
Tan	<i>Tantrāloka</i>
TanV	<i>Tantravārttika</i>
TSS 41	<i>The Adhyātmapaṭala of the Āpastambadharmā with Vivaraṇa of Śrī Śāṅkara Bhagavatpāda</i> . Ed. by T. Gaṇapati Śāstrī. Trivandrum Sanskrit Series XLI. Trivandrum: Travancore Government Press, 1915.
Upad	<i>Upadeśasāhasrī</i>
Vād	<i>Vādanyāya</i>
VedS	<i>Vedārthasaṃgraha</i>
VibhV	<i>Vibhramaviveka</i>
VigrV	<i>Vigrahavyāvartanī</i>
ViṃS	<i>Viṃśatikasiddhi</i>
VVP 10	<i>The Works of Sri Sankaracharya</i> . Vol. 10: <i>Bṛhadāraṇyakopaniṣadbhāṣya, Chapters 5 and 6, and Nṛsimhapūrvaṭāpanīyabhāṣya</i> . Srirangam: Sri Vani Vilas Press, 1910.
VVP 13	<i>The Works of Sri Sankaracharya</i> . Vol. 13: <i>Viṣṇusahasranāmabhāṣya and Sanatsujātīyabhāṣya</i> . Srirangam: Sri Vani Vilas Press, 1910.
VVP 18	<i>The Works of Sri Sankaracharya</i> . Vol. 18: <i>Stotras and Lalitā-triśatisotrabhāṣya</i> . Srirangam: Sri Vani Vilas Press, 1910.
YuK	<i>Yuktiṣaṣṭikakārika</i>

## References

- ANDRIJANIĆ, Ivan 2019. “Śaṅkara and the authorship of Śvetāśvataropaniṣad-Bhāṣya”. *The Journal of Hindu Studies* 12(3): 273–291. <https://doi.org/10.1093/jhs/hiy014>
- ANDRIJANIĆ, Ivan 2020a. “Śaṅkara and the Authorship of the *Īsopaniṣadbhāṣya* and the *Kaṭhōpaniṣadbhāṣya*”. *International Journal of Hindu Studies* 24: 257–282. <https://doi.org/10.1007/s11407-020-09279-z>
- ANDRIJANIĆ, Ivan 2020b. “The authorship of the *Chāndogyopaniṣad-Bhāṣya*: A stylometric approach”. [In:] Michalak-Pikulska, Barbara, Marek Piela and Tomasz Majtczak, eds, *Oriental Languages and Civilizations*. Cracow: Jagiellonian University Press, pp. 103–116.
- BADER, Jonathan 2000. *Conquest of the Four Quarters. Traditional Accounts of the Life of Śaṅkara*. New Delhi: Aditya Prakashan.
- BELVALKAR, S. K. 1929. *Shree Gopal Basu Mallik Lectures on Vedānta Philosophy*. Part one: *Lectures 1–6*. Poona: Bilvakuṅja Publishing House.
- BELVALKAR, S. K. 1930. “An Authentic but Unpublished Work of Śaṅkarācārya”. *Journal of the Bombay Branch of the Royal Asiatic Society* 6: 241–246.
- BURROWS John 2002. “‘Delta’: a Measure of Stylistic Difference and a Guide to Likely Authorship”. *Literary and Linguistic Computing* 17(3): 267–287. <https://doi.org/10.1093/lc/17.3.267>
- DAELEMANS, Walter 2013. “Explanation in Computational Stylometry”. [In:] Gelbukh, Alexander, ed., *Proceedings of the 14th International Conference on Computational Linguistics and Intelligent Text Processing*. Vol. 2. Berlin, Heidelberg: Springer, pp. 451–462. [https://doi.org/10.1007/978-3-642-37256-8\\_37](https://doi.org/10.1007/978-3-642-37256-8_37)
- DAVID, Hugo 2017. “Towards a Critical Edition of Śaṅkara’s ‘Longer’ Aitareyopaniṣadbhāṣya: a Preliminary Report Based on Two Cambridge Manuscripts”. [In:] Vergiani, Vincenzo, Daniele Cuneo and Camillo A. Formigatti, eds, *Indic Manuscript Cultures through the Ages: Material, Textual, and Historical Investigations*. Berlin/Boston: De Gruyter, pp. 727–754. <https://doi.org/10.1515/9783110543100-022>
- EDER, Maciej 2015. “Does size matter? Authorship attribution, small samples, big problem”. *Digital Scholarship in the Humanities* 30(2): 167–82. <https://doi.org/10.1093/lc/fqt066>
- EDER, Maciej 2018. “Authorship verification with the package *stylo*”. <https://computationalstylistics.github.io/docs/imposters> (accessed 18 January 2024).

- EDER, Maciej, Jan RYBICKI and Mike KESTEMONT 2016. “Stylometry with R: a package for computational text analysis”. *R Journal* 8(1): 107–21. <https://doi.org/10.32614/RJ-2016-007>
- EVERT, Stefan, Thomas PROISL, Fotis JANNIDIS, Isabella REGER, Steffen PIELSTRÖM, Christof SCHÖCH and Thorsten VITT 2017. “Understanding and explaining Delta measures for authorship attribution”. *Digital Scholarship in the Humanities* 32 (suppl. 2), ii4–ii16. <https://doi.org/10.1093/llc/fqx023>
- HACKER, Paul 1950. “Eigentümlichkeiten der Lehre und Terminologie Śāṅkaras: Avidyā, Nāmarūpa, Māyā, Īśvara”. *Zeitschrift der Deutschen Morgenländischen Gesellschaft* 100: 246–286. Reprinted in *Kleine Schriften*, 1978, Wiesbaden: Franz Steiner Verlag, pp. 69–109.
- HACKER, Paul 1951. *Untersuchungen über Texte des frühen Advaitavāda. I. Die Schüler Śāṅkara's*. Wiesbaden: Verlag der Akademie der Wissenschaften und der Literatur in Mainz: Franz Steiner Verlag.
- HACKER, Paul 1968–1969. “Śāṅkara der Yogin und Śāṅkara der Advaitin. Einige Beobachtungen”. *Wiener Zeitschrift für die Kunde Süd- und Ostasiens* 12–13 (Festschrift für Erich Frauwallner): 119–148.
- HACKER, Paul 1972. “Notes on the Māṇḍūkyaopaniṣad and Śāṅkara's Āgamaśāstravivarāṇa”. [In:] Ensink, J. and P. Gaeffke, eds, *India Maior: Congratulation Volume Presented to J. Gonda*. Leiden: Brill, pp. 115–132.
- HACKER, Paul 1978. “Śāṅkarācārya and Śāṅkarabhagavatpāda. Preliminary remarks concerning the authorship problem” (Korrigierte Neufassung). [In:] Hacker, Paul, *Kleine Schriften*. Wiesbaden: Franz Steiner Verlag, pp. 41–59. Originally published in *New Indian Antiquary* 9 (1947): 175–186.
- HARIMOTO, Kengo 2006. “The date of Śāṅkara: Between the Cāḷukyas and the Rāṣṭrakūtas”. *Journal of Indological Studies* 18: 85–111.
- HARIMOTO, Kengo 2014. *God, Reason, and Yoga: A Critical Edition and Translation of the Commentary Ascribed to Śāṅkara on Pātañjalayogaśāstra 1.23–28*. Hamburg: Department of Indian and Tibetan Studies, University of Hamburg.
- HELLWIG, Oliver 2016. “Detecting sentence boundaries in Sanskrit texts”. [In:] Matsumoto, Yuji and Rashmi Prasad, eds, *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*. Osaka: The COLING 2016 Organizing Committee, pp. 288–297. <http://aclweb.org/anthology/C16-1028> (accessed 18 January 2024).
- HELLWIG, Oliver and Sebastian NEHRDICH 2018. “Sanskrit Word Segmentation Using Character-level Recurrent and Convolutional Neural Networks”. [In:] Riloff, Ellen, David Chiang, Julia Hockenmaier and Jun'ichi Tsujii, eds, *Proceedings of the 2018 Conference on Empirical Methods in Natural*

- Language Processing*. Brussels: Association for Computational Linguistics, pp. 2754–2763. <https://doi.org/10.18653/v1/D18-1295>
- HOLMES, David I. 1994. “Authorship Attribution”. *Computers and Humanities* 28 (2): 87–106. <https://doi.org/10.1007/BF01830689>
- HOUVARDAS, John and Efsthios STAMATATOS 2006. “N-Gram Feature Selection for Authorship Identification”. [In:] Euzenat, J. and J. Domingue, eds, *Proceedings of Artificial Intelligence: Methodologies, Systems, and Applications (AIMSA 2006)*. Springer, pp. 77–86. [https://doi.org/10.1007/11861461\\_10](https://doi.org/10.1007/11861461_10)
- JACOB, George A. 1886. “The Nṛsiṃhatāpanīya-Upaniṣad”. *The Indian Antiquary: A Journal of Oriental Research* 15: 69–74.
- JANNIDIS, Fotis, Steffen PIELSTRÖM, Christof SCHÖCH and Thorsten VITT 2015. “Improving Burrow’s Delta – An empirical evaluation of text distance measures”. [In:] *Digital Humanities Conference*. Sydney: University of Western Sydney.
- JUOLA, Patrick 2006. “Authorship Attribution”. *Foundations and Trends in Information Retrieval* 1(3): 233–334. <https://doi.org/10.1561/1500000005>
- KESTEMONT, Mike 2014. “Function Words in Authorship Attribution. From Black Magic to Theory?” [In:] Feldman, Anna, Anna Kazantseva and Stan Szpakowicz, eds, *Proceedings of the 3rd Workshop on Computational Linguistics for Literature (CLFL)*. Gothenburg, Sweden: Association for Computational Linguistics, pp. 59–66. <https://doi.org/10.3115/v1/W14-0908>
- KESTEMONT, Mike, Walter DAELEMANS and Dominiek SANDRA 2012. “Robust Rhymes? The Stability of Authorial Style in Medieval Narratives”. *Journal of Quantitative Linguistics* 19(1): 54–76. <https://doi.org/10.1080/09296174.2012.638796>
- KESTEMONT, Mike, Justin STOVER, Moshe KOPPEL, Folgert KARSDORP and Walter DAELEMANS 2016. “Authenticating the Writings of Julius Caesar”. *Expert Systems With Applications* 63: 86–96. <https://doi.org/10.1016/j.eswa.2016.06.029>
- KITAGAWA, Yoshiaki and Mamoru KOMACHI 2017. “Long Short-Term Memory for Japanese Word Segmentation”. [In:] *32nd Pacific Asia Conference on Language, Information and Computation Hong Kong, 1–3 December 2018*. <https://doi.org/10.48550/arXiv.1709.08011>
- KOPPEL, Moshe, Jonathan SCHLER and Shlomo ARGAMON 2009. “Computational methods in authorship attribution”. *Journal of the American Society for Information Science and Technology* 60 (1): 9–26. <https://doi.org/10.1002/asi.20961>

- KOPPEL, Moshe and Yaron WINTER. 2014. “Determining if two documents are written by the same author”. *Journal of the Association for Information Science and Technology* 65(1): 178–187. <https://doi.org/10.1002/asi.22954>
- KRISHNA, Amrith, Bishal SANTRA, Pavankumar SATULURI, Sasi Prasanth BANDARU, Bhumi FALDU, Yajuvendra SINGH and Pawan GOYAL 2016. “Word Segmentation in Sanskrit Using Path Constrained Random Walks”. [In:] *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, Osaka: COLING 2016 Organizing Committee, pp. 494–504.
- LEGGET, Trevor 1978. *The Chapter of the Self*. London: Routledge and Kegan Paul.
- LOVE, Harold 2002. *Attributing Authorship: An Introduction*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511483165>
- LUTOSŁAWSKI, Wincenty 1898. “Principes de stylométrie appliqués à la chronologie des œuvres de Platon”. *Revue des Études Grecques* 11 (Fasc. 41): 61–81. <https://doi.org/10.3406/reg.1898.5847>
- LUYCKX, Kim and Walter DAELEMANS 2011. “The effect of author set size and data size in authorship attribution”. *Literary and Linguistic Computing* 26(1): 35–55. <https://doi.org/10.1093/lc/fqq013>
- MARSCHNER, Käthe 1933. *Zur Verfasserfrage des dem Śaṅkarācārya zugeschriebenen Brhadāranyakopaniṣad-Bhāṣya*. Berlin-Charlottenburg: Alfred Lindner Verlag.
- MAYEDA, Sengaku 1965a. “The Authenticity of the Bhagavadgītābhāṣya ascribed to Śaṅkara”. *Wiener Zeitschrift für die Kunde Südasiens und Archiv für indische Philosophie* 9: 155–197.
- MAYEDA, Sengaku 1965b. “The Authenticity of the Upadeśasāhasrī Ascribed to Śaṅkara”. *Journal of the American Oriental Society* 85(2): 178–196. <https://doi.org/10.2307/597989>
- MAYEDA, Sengaku 1967. “On Śaṅkara’s Authorship of the Kenopaniṣad-bhāṣya”. *Indo-Iranian Journal*, 10(1): 33–55. <https://doi.org/10.1163/000000068792937937>
- MAYEDA, Sengaku 1967–1968. “On the author of the Māṇḍūkyaopaniṣad and the Gauḍapādīyabhāṣya”. *Adyar Library Bulletin* 31–32: 73–94.
- MAYEDA, Sengaku, ed. and trans. 2006. *Śaṅkara’s Upadeśasāhasrī*. Vols I and II. Delhi: Motilal Banarsidass.
- MOISL, Hermann 2015. *Cluster Analysis for Corpus Linguistics*. Berlin and New York: De Gruyter Mouton. <https://doi.org/10.1515/9783110363814>

- MORROW, Daniel G. 1986. “Grammatical Morphemes and Conceptual Structure in Discourse Processing”. *Cognitive Science* 10(4): 423–455. [https://doi.org/10.1207/s15516709cog1004\\_2](https://doi.org/10.1207/s15516709cog1004_2)
- NAKAMURA, Hajime 1983. *A History of Early Vedānta Philosophy*. Part One. Delhi: Motilal Banarsidass.
- PALMER, David D. 2010. “Text Preprocessing”. [In:] Indurkha, Nitin and Fred J. Damerau, eds, *Handbook of Natural Language Processing*. 2nd edition. A Chapman & Hall Book.
- PANDE, G. C. 1994. *Life and Thought of Śaṅkarācārya*. Delhi: Motilal Banarsidass.
- PEÑAS, Anselmo and Alvaro RODRIGO 2011. “A Simple Measure to Assess Non-response”. [In:] *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, vol. 1. Portland, Oregon: Association for Computational Linguistics, pp. 1415–1424.
- POTHA, Nektaria and Efstathio STAMATATOS 2017. “An Improved Impostors Method for Authorship Verification”. [In:] Jones, G. J. F., S. Lawless, J. Gonzalo, L. Kelly, L. Goeuriot, T. Mandl, L. Cappellato and N. Ferro, eds, *Experimental IR Meets Multilinguality, Multimodality, and Interaction. CLEF 2017*. Lecture Notes in Computer Science, vol. 10456. Springer, Cham, pp. 138–144. [https://link.springer.com/chapter/10.1007/978-3-319-65813-1\\_14](https://link.springer.com/chapter/10.1007/978-3-319-65813-1_14) (accessed 18 January 2024).
- REIGLE, David and Nancy REIGLE 2005. *Śaṅkarācārya’s Collected Works: An Annotated Bibliography of Published Editions in Sanskrit*. Cotopaxi, Colorado, U.S.A.: Eastern Tradition Research Institute.
- RUBIN, D. 1995. *Memory in Oral Traditions. The Cognitive Psychology of Epic, Ballads and Counting-out Rhymes*. Oxford: Oxford University Press. <https://doi.org/10.1093/oso/9780195082111.001.0001>
- RUŽIČKA, M. 1958. “Anwendung mathematisch-statistischer Methoden in der Geobotanik (synthetische Bearbeitung von Aufnahmen)”. *Biología* (Bratislava) 13: 647–661.
- RYBICKI, Jan and Maciej EDER. 2011. “Deeper Delta across genres and languages: do we really need the most frequent words?”. *Literary and Linguistic Computing* 26(3): 315–321. <https://doi.org/10.1093/lc/fqr031>
- SANDERSON, Alexis, ed. and trans. 2017. *The Smārta Śāktism of South India: Lalitātrīśatīstotra: The Hymn of the Three Hundred Epithets of the Goddess Lalitā, edited with a brief introduction, an annotated English translation, and an appendix containing the Nāmāvalī*. [https://www.academia.edu/34452056/The\\_Sm%C4%81rta\\_%C5%9A%C4%81ktism\\_of\\_South\\_India\\_Lalit](https://www.academia.edu/34452056/The_Sm%C4%81rta_%C5%9A%C4%81ktism_of_South_India_Lalit)

- [%C4%81tri%C5%9Bat%C4%ABstotra\\_edited\\_with\\_a\\_brief\\_introduction\\_an\\_annotated\\_English\\_translation\\_and\\_an\\_appendix\\_containing\\_the\\_N%C4%81m%C4%81val%C4%AB](#) (accessed 18 January 2024).
- SASTRY, R. Anantakrishna 1980. *Viṣṇusahasranāma with the Bhāṣya of Śrī Śaṅkarācārya*. Madras: The Adyar Library and Research Centre.
- SEIDMAN, Shachar 2014. “Authorship Verification Using the Impostors Method Notebook for PAN at CLEF 2013”. [In:] Forner, P., R. Navigli, D. Tufis and N. Ferro, eds, *CEUR Workshop Proceedings*. Vol. 1179: *Working Notes for CLEF 2013 Conference, Valencia, Spain, 23–26 September 2013*.
- STAMATATOS, Efstathios 2009. “A survey of modern authorship attribution methods”. *Journal of the American Society for Information Science and Technology* 60(3): 538–556. <https://doi.org/10.1002/asi.21001>
- TensorFlow. <https://www.tensorflow.org/> (accessed 18 January 2024).
- VETTER, Tillman 1968–1969. “Zur Bedeutung des Illusionismus bei Śaṅkara”. *Wiener Zeitschrift für die Kunde Süd- und Ostasiens* 12–13: 407–423.
- ZIPF, George Kingsley 1935. *The Psychobiology of Language: An Introduction to Dynamic Philology*. Boston, MA: Houghton Mifflin.
- ZIPF, George Kingsley 1949. *Human Behavior and the Principle of Least Effort: An Introduction to Human Ecology*. Cambridge, Mass.: Addison-Wesley.