

A Computational Authorship Attribution in the Children's Tales of Oscar and Constance Wilde: The Case of "The Selfish Giant"

Liviu P. Dinu^{1,2*} Alina Iacob^{1*} Cosmin Ciotlos³

¹Faculty of Mathematics and Computer Science, University of Bucharest

²Human Language Technologies Research Center, Bucharest, Romania

³Faculty of Letters, University of Bucharest

liviu.p.dinu@unibuc.ro

alinaiacob9517@gmail.com

cosmin.ciotlos@litere.unibuc.ro

Abstract

This study introduces and analyzes a novel authorship attribution case: the children's stories published by Oscar and Constance Wilde. We analyzed the corpus of stories with both supervised (SVM with string kernel) and unsupervised (Hierarchical Clustering via Rank Distance) methods and found stylistic resemblance between the story "The Selfish Giant" published by Oscar Wilde and the stylometric profile of Constance Wilde. Starting from this baseline, we also explored the capabilities of LLMs in authorship attribution via Perplexity. Our finding suggests that the story "The Selfish Giant" might be the result of a collaboration between Oscar and Constance Wilde. Moreover, our results pointed to the distinct stylistic fingerprints of the two authors with regards to the rest of the corpus, confirming that their respective styles are separable despite shared genre and period.

1 Introduction and Literature Review

Authorship attribution is a fundamental task in Computational Linguistics. The current paper applies three methods of authorship attribution to a literary case that has not yet been examined using computational methods.

Oscar Wilde was a prolific and intriguing figure of the 19th century, whose life and writings continue to fascinate and prompt scholarly interest. No less intriguing has been the figure of his wife, Constance Wilde, and their close collaboration in social and cultural pursuits is well-documented (Powell and Raby, 2013; Moyle, 2012; Amor, 1999). However, recent evidence suggests their collaboration may have extended beyond editorial work, including literary works as well. In 2008, the *Morgan Library & Museum* in New York City received a donation of manuscripts pertaining to Oscar Wilde, among which was a manuscript of the story "The

Selfish Giant" (published by Oscar Wilde in 1888 as part of the fairytales volume *The Happy Prince and Other Tales*) written in the hand of Constance Wilde with pencil annotations by Oscar Wilde. However, the final version of the publication exhibits substantial differences from the handwritten original. In the same year, Constance Wilde had published herself a volume of children stories titled *There was once: Grandma's stories*.

This makes *The Selfish Giant* both a compelling case study for stylometric analysis. Accordingly, the aims of the present study are to:

1. Explore the connection of the tale to both writers;
2. Determine whether the evidence points to either editing or collaboration;
3. Explore authorship attribution via supervised, unsupervised and LLM-based methods.

Since the problem at hand concerns a novel and previously unexplored authorship case, we employed authorship attribution models that have shown robust results in prior studies. Moreover, Constance Wilde's literary publications are limited to a short volume of fairy tales. Thus, the data availability constraint represents a main challenge in the current study. Consequently, we first selected two approaches described by Dinu et al. (2008): Classification based on Support Vector Machines (SVM) with a string kernel, and Clustering via Rank Distance - which have yielded positive results in authorship attribution tasks. This selection of complementary models allows us to cover both supervised and unsupervised modeling paradigms.

Furthermore, LLMs have increasingly been studied for their ability to capture complex linguistic and stylistic patterns, showing promise in the field of authorship attribution (Huang et al., 2024; Habib et al., 2025; Huang et al., 2025b,a). In this study,

*These authors contributed equally to this work.

we aim to further explore their potential and evaluate their performance on small, author-specific corpora by drawing from the methodology of Huang et al. (2025b) and exploring authorship attribution via LLM Perplexity.

At the core of stylometry lies the interplay between style, genre and stylistic fingerprint. While *style* is concerned with conscious aesthetic choices made by writers (Sotirova, 2015) and *genre* emerges from structural and rhetorical conventions of different types of texts (Biber and Conrad, 2019), it is the *stylistic fingerprint* of a writer that computational attribution techniques are most often concerned with - capturing the so-called "human stylome" (Van Halteren et al., 2005). So far, computational stylometric methods have been employed in the cases of a great number of authorship uncertainties, in both literary and non-literary fields (Mosteller and Wallace, 1964; Craig and Kinney, 2009; Labbé and Labbé, 2001; Juola, 2015; Dinu et al., 2012a; Tuzzi and Cortelazzo, 2018).

Our study brings a contribution to the field by presenting a novel case of authorship attribution in which questions of collaborative writing and editorial influence can be explored in a setting of little data availability. We investigate the potential involvement of Constance Wilde alongside Oscar Wilde, examining the manuscript evidence in relation to the final published version. In addition, we introduce a dataset of short fairy tales, enabling reproducible computational analysis and offering a resource for future studies in authorship attribution. This approach not only sheds light on a previously unexplored literary case but also demonstrates the applicability of stylometric techniques. Finally, we explore whether LLMs can yield superior results to traditional feature-based methods.

In order to achieve this, the rest of the article is structured as follows: Section 2 describes the dataset; Section 3 presents the methodology; Section 4 discusses the results; Section 5 draws the conclusions and Section 6 discusses the Ethical Considerations.

2 Data

The dataset used in this study consists of short fairy tales authored by Oscar Wilde and Constance Wilde. Given the fact that Constance Wilde only authored one volume of tales, we limited the Oscar Wilde dataset to a similar size and genre. Oscar Wilde’s published collection, *The Happy Prince*

Table 1: Corpus summary with word counts

Author	Text	Words
Oscar Wilde	The Happy Prince	3473
Oscar Wilde	Nightingale and Rose	2328
Oscar Wilde	The Selfish Giant (published version)	1658
Oscar Wilde	The Selfish Giant (manuscript)	1652
Oscar Wilde	The Devoted Friend	4300
Oscar Wilde	The Remarkable Rocket	4384
Constance Wilde	Cinderella	1181
Constance Wilde	Jack the Giant Killer	1830
Constance Wilde	Little Red Riding Hood	837
Constance Wilde	Puss in Boots	860
Constance Wilde	The Three Bears	759

and *Other Tales* (1888), contains 5 stories including "The Selfish Giant". In addition, we include the manuscript of "The Selfish Giant" in the handwriting of Constance Wilde with pencil annotations by Oscar Wilde, held at the *Morgan Library & Museum*, New York City. Constance Wilde’s collection, *There Was Once: Grandma’s Stories* (1888), is also included. We excluded the tales in verses *Little Bo Beep*, *Old Mother Hubbard*, *Babes in the Wood* and *The Three Little Kittens*, leaving five tales from her collection.

Table 1 shows a complete description of the corpus. All texts were processed to remove non-story content such as introductions, illustrations, and page numbers. The dataset is available for research purposes, and all texts are in the public domain. Oscar Wilde’s stories were obtained from Project Gutenberg.

Finally, we computed the Jaccard Similarity in order to assess the rough similarity between the two versions of "The Selfish Giant":

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

where A and B are the sets of unique words in each text. We obtained a score of 0.89, indicating a high degree of overlap between the manuscript and the published version.

3 Methods

As mentioned above, we aim to investigate the authorship question using three distinct methods: Classification based on Support Vector Machines (SVM) with a string kernel, Clustering via Rank Distance, and Authorship Attribution based on LLM Perplexity.

3.1 Classification Methodology

We followed the methodology described by Dinu et al. (2008), modeling a 2-class classification problem using SVM with a string kernel.

String kernels are particularly suited to authorship attribution problems because they treat texts as sequences of characters rather than words, meaning no preprocessing is needed other than space normalization. The similarity between two texts is measured by counting shared substrings of length p (p-spectrum kernel). By choosing a small value for p (usually 2-5), the kernel will mostly capture a mixture a suffixes, function words and punctuation - all strong indicators of stylome. To capture stylistic patterns at multiple scales, we used a blended spectrum kernel, which sums similarities over all substring lengths from 1 to p :

$$k_1^p(s, t) = \sum_{q=1}^p k_q(s, t)$$

To account for differences in text lengths, the kernel was normalized:

$$\hat{k}_1^p(s, t) = \frac{k_1^p(s, t)}{\sqrt{k_1^p(s, s) \cdot k_1^p(t, t)}}$$

Given the limited size of the dataset, a nested Leave One-Out Cross Validation (LOOCV) scheme was adopted. In the outer loop, each fairy tale was iteratively used as the test instance while the remaining texts were used for training. Within each training set, an internal leave-one-out procedure was performed to estimate the model's stability and generalization capacity. The internal accuracy reported corresponds to this inner validation step, while the final prediction for each tale corresponds to the outer test fold. We also created two datasets, to account for the two versions of "The Selfish Giant": one including the published version, the other the manuscript version.

The two key hyperparameters of the model were chosen empirically. The substring length $p=4$ was selected to balance sensitivity to local stylistic patterns and computational feasibility. Similarly, $\nu=0.05$ was chosen empirically after preliminary experiments indicated that this value achieved stable performance while limiting margin violations to at most 5% of the training samples. Both parameters were thus fixed for the final cross-validation and leave-one-out experiments to ensure consistency across all classification runs.

3.2 Clustering Methodology

Hierarchical Agglomerative Clustering is a bottom-up approach that iteratively merges clusters (initially represented by each individual data point)

a	been	had	its	one	that	was
all	but	has	may	only	the	were
also	by	have	more	or	their	what
an	can	her	must	our	then	when
and	do	his	my	shall	there	which
any	down	if	no	should	things	who
are	even	in	not	so	this	will
as	every	into	now	some	to	with
at	for	is	of	such	up	would
be	from	it	on	than	upon	your

Table 2: List of functional words used in the clustering experiments

by calculating the distance between every pair of clusters. The final resulting dendrogram shows the hierarchy of clusters and the relationships between them. Thus, when working with this distance-based clustering technique, the two main aspects to be considered are the selected features and the distance measure.

Since our aim is to capture the stylistic fingerprints of the two authors, the selected features need to unambiguous, quantifiable and used in an unconscious way.

Function words, first introduced by [Mosteller and Wallace \(1964\)](#), have been proposed as linguistic features that fit the criteria. It has subsequently been proven that they are indeed used unconsciously and carry part of the stylistic fingerprint of individuals ([Chung and Pennebaker, 2011](#)). We therefore based our experiments on the list of functional words introduced by [Mosteller and Wallace \(1964\)](#), which can be seen in Table 2.

Previous studies in authorship attribution have shown that Rank Distance is a metric that yields consistent results in clustering together works pertaining to an author ([Dinu et al., 2008](#); [Popescu and Dinu, 2008](#); [Dinu et al., 2012b](#)). Rank distance ([Dinu, 2003](#)) is an ordinal metric related to Spearman's footrule. To calculate the Rank Distance, the following steps were followed:

1. Calculate the raw frequencies of the selected functional words in each text
2. Replace the raw frequencies with the rank each functional word occupies in each text. The word with the highest frequency will be assigned the highest rank, while the words that are present on the function words list but absent in the text will be assigned the rank 0.

Tied objects are assigned the average rank of the positions they share.

3. Compute the pairwise rank distance between each two texts according to L1 norm, where $\sigma(i)$ represents the rank of the object:

$$D(\sigma_1, \sigma_2) = \sum_{i=1}^n |\sigma_1(i) - \sigma_2(i)|$$

4. Obtain a square similarity matrix, which serves as input for the Hierarchical Clustering algorithm.

For our experiments we used Hierarchical Agglomerative Clustering with average linkage.

3.3 LLM Perplexity Methodology

While traditional authorship attribution methods are based on hand-crafted features, attribution via LLMs leverages deep contextual representations learned from vast corpora of text. Attribution via LLMs is still a relatively new area of research, and since the current study concerns an authorship case where the ground truth is not known, we opted for a simple approach to examine whether LLM-based analysis can provide additional insights. We followed the approach of Huang et al. (2025b), which proved competitive with state-of-the-art methods in authorship attribution. In this framework, a separate language model is fine-tuned on the texts of each candidate author. The resulting model is intended to capture the stylistic and lexical patterns characteristic of a specific writer, including preferences in vocabulary, phrasing, and contextual word usage. The fine-tuned models are then used to evaluate how well each one predicts an unknown text by computing its perplexity. For a causal language model m and a text tokenized as a sequence $T = \{x_1, x_2, \dots, x_t\}$, perplexity is defined as:

$$\text{PPL}(m, T) = \exp \left(-\frac{1}{t} \sum_{i=1}^t \log p_m(x_i | x_{<i}) \right) \quad (1)$$

where $p_m(x_i | x_{<i})$ denotes the probability assigned by the model m to the token x_i given all preceding tokens $x_{<i}$, and t denotes the length of the token sequence T . Intuitively, perplexity measures how "surprised" a model is by a text: a model that has been fine-tuned on an author's writing will assign higher probabilities to new texts written in a similar style. The candidate author whose model

yields the lowest perplexity is considered the most likely author of the questioned document.

In contrast to Huang et al. (2025b), who perform full fine-tuning of GPT-2, we fine-tune the GPT-2 model via LoRA (Hu et al., 2022), a parameter-efficient method that updates only a small set of low-rank adapter matrices while keeping the base model weights frozen. This reduces computational cost while preserving the base model's language understanding.

We trained the Oscar model on a corpus of 8 short stories and a novel authored by Oscar Wilde, and the Constance model on three fairy tales attributed to her, leaving one for evaluation. Below are the configurations used for training the two models, chosen empirically:

Parameter	Value
LoRA rank (r)	16
LoRA alpha	32
Target modules	c_attn, c_proj
LoRA dropout	0.05
Bias	none
Task type	CAUSAL_LM
Chunk size	128 tokens
Gradient accumulation steps	2
Number of epochs	20
Learning rate	1×10^{-4}
Optimizer	AdamW (PyTorch)

Table 3: Training configuration for LoRA fine-tuning of the Constance-finetuned model

Parameter	Value
LoRA rank (r)	16
LoRA alpha	32
Target modules	c_attn, c_proj
LoRA dropout	0.05
Bias	none
Task type	CAUSAL_LM
Chunk size	512 tokens
Tokenizer pad token	EOS token
Gradient accumulation steps	2
Number of epochs	20
Learning rate	3×10^{-4}
Optimizer	AdamW (PyTorch)

Table 4: LoRA fine-tuning and training configuration for the Oscar-finetuned model

4 Results

The following section presents the results and analysis of our experiments.

4.1 Classification Results

We ran the classification experiments on two datasets: the first one including the manuscript version of the story "The Selfish Giant" and the second one including the published version of the text.

The results with the first dataset are presented in Table 5, and the second set of results is presented in Table 6. The results are highly similar: when the tale "The Selfish Giant" is included in the inner LOOCV training loop, the training accuracy of the model is low: and average of 0.55 for the first dataset and 0.58 for the second one. The presence of the story "The Selfish Giant" labeled as pertaining to Oscar Wilde in the training set obfuscates the model's ability to learn meaningful patterns and differentiate between the two authors, especially with the training set being small (just 9 entries). However, when the story "The Selfish Giant" is held out as a test document, the model reaches perfect accuracy and high confidence in both datasets (highlighted entry in Tables 5 and 6).

This fact that including the story "The Selfish Giant" in the training set degrades both training accuracy and confidence, inhibiting the model's ability to learn meaningful patterns, stands as evidence for the story's pronounced stylistic dissimilarity with the rest of the Oscar Wilde texts.

We reiterate the fact that the changes made to the manuscript version were most likely carried out by Oscar Wilde (based on the fact that the pencil notes on the manuscript are in Oscar Wilde's handwriting). The fact that both versions of the text are attributed to Constance Wilde with relatively high confidence (second highest in the first dataset, third highest in the second) points to the fact that Constance Wilde's input regarding "The Selfish Giant" goes deeper than mere editing. Moreover, Oscar Wilde's changes do not influence or improve the model's ability to differentiate between the two authors.

The results of this set of experiments bring initial insights into the case of "The Selfish Giant", showing evidence towards the hypothesis of at least a tight collaboration between Constance and Oscar Wilde in writing this tale.

4.2 Clustering Results

We ran the Clustering experiments on four datasets:

1. Full Corpus, including both versions of "The Selfish Giant" (Figure 1a)

2. A corpus including only the published version of "The Selfish Giant" (Figure 1b)
3. A corpus including only the manuscript version of "The Selfish Giant" (Figure 1c)
4. A corpus excluding all versions of "The Selfish Giant" (Figure 1d)

For ease of interpretation, in all of the plots the texts published by Constance Wilde are written in blue and those published by Oscar Wilde in red.

As is evident from the first three figures, reducing the texts to only functional words renders the two versions of "The Selfish Giant" are virtually indistinguishable. The two hierarchies in Figures 1b and 1c are identical: there is a clear separation between the cluster of texts published by Constance Wilde and the cluster of texts published by Oscar Wilde, with the exception of "The Selfish Giant" which is well integrated among Constance Wilde's writings.

Moreover, when we ran the Clustering experiment without including "The Selfish Giant" (Figure 1d), the algorithm produced a hierarchy that perfectly matches our dataset, the writings of the authors producing two clearly separated clusters. Since the clustering experiments are based on functional words only (the words that are unconsciously used by individuals), this result suggests that, if any literary collaboration occurred between the two spouses, it did not extend to a depth that would cause their writings to become stylistically "contaminated" by one another. In other words, this result proves that, despite the similarity of genre and the fact that the stories were published and probably written in the same year in a period when the two were still living together, both of them produced original and stylistically independent works.

The only exception is, then, "The Selfish Giant". When included in the dataset, both versions of the text perfectly integrate within Constance Wilde's cluster. Furthermore, despite the fact that the changes made to the manuscript might seem, on close reading, "significant" (according to Moyle (2012)), when the text is analyzed only based on those features that occur independently of the author's design the two versions of the text become so close that they form their own subcluster within that of Constance Wilde (Figure 1a)..

The results of this set of experiments reinforce the findings of the classification analysis: there is a clear separation between the stylometric profiles of

Table 5: SVM classification results with each fairy tale iteratively used as the test text. Dataset including the published version of “The Selfish Giant”

Test Text	True Author	Predicted Author	Conf.	Int.	Acc.
The Happy Prince (Oscar)	Oscar	Oscar	58.71	0.78	
Nightingale and Rose (Oscar)	Oscar	Oscar	57.62	0.44	
The Selfish Giant published (Oscar)	Oscar	Constance	76.78	1.00	
The Devoted Friend (Oscar)	Oscar	Oscar	54.09	0.78	
The Remarkable Rocket (Oscar)	Oscar	Oscar	60.92	0.44	
Cinderella (Constance)	Constance	Oscar	50.00	0.44	
Jack the Giant Killer (Constance)	Constance	Oscar	54.62	0.56	
Little Red Riding Hood (Constance)	Constance	Constance	78.67	0.44	
Puss in Boots (Constance)	Constance	Constance	58.81	0.56	
The Three Bears (Constance)	Constance	Oscar	53.63	0.56	

Table 6: SVM classification results with each fairy tale iteratively used as the test text. Dataset including the manuscript version of “The Selfish Giant”

Test Text	True Author	Predicted Author	Conf.	Int.	Acc.
The Happy Prince (Oscar)	Oscar	Oscar	60.58	0.78	
Nightingale and Rose (Oscar)	Oscar	Oscar	63.70	0.56	
The Selfish Giant manuscript (Oscar)	Oscar	Constance	77.22	1.00	
The Devoted Friend (Oscar)	Oscar	Oscar	58.22	0.67	
The Remarkable Rocket (Oscar)	Oscar	Oscar	79.71	0.44	
Cinderella (Constance)	Constance	Oscar	50.00	0.56	
Jack the Giant Killer (Constance)	Constance	Oscar	61.11	0.56	
Little Red Riding Hood (Constance)	Constance	Constance	90.08	0.56	
Puss in Boots (Constance)	Constance	Constance	54.69	0.56	
The Three Bears (Constance)	Constance	Oscar	54.05	0.56	

the two authors, with "The Selfish Giant" emerging as the only outlier. The strong stylistic resemblance between this story and Constance Wilde’s writings suggests that it may have been the product of a collaboration, or even authored by her alone.

4.3 LLM Perplexity Results

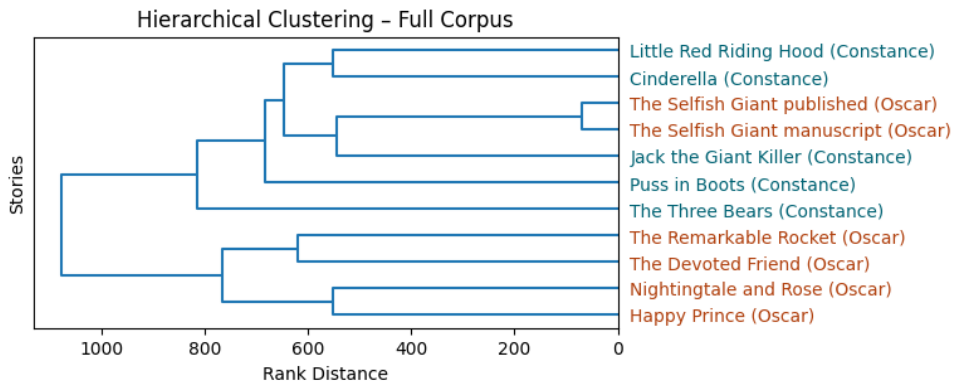
The results of the two authorial language models are presented in Tables 7 and 8. A positive difference indicates that the fine-tuned model finds the text more predictable than the base model, suggesting stylistic similarity to the training author, while a negative difference suggests the opposite.

The Constance-finetuned model assigns the largest positive difference to Cinderella by a substantial margin (9.45), far exceeding all other texts. Oscar Wilde’s fairy tales all receive negative values (ranging from -1.54 to -2.55), indicating that the model does not recognise them as stylistically similar to Constance’s writing. Both version of *The Selfish Giant* receive small positive values (0.53 and 0.10 respectively), suggesting that after fine-tuning, the model assess them as being slightly

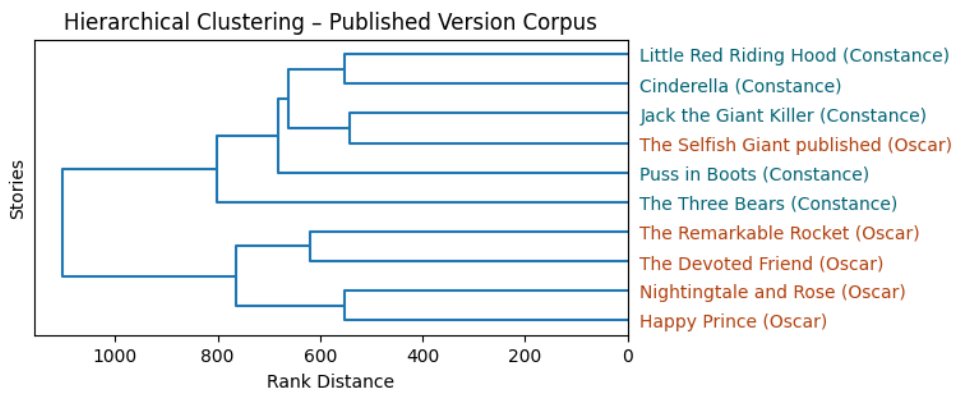
closer to Constance’s style.

The Oscar-finetuned model shows a complementary pattern. Oscar Wilde’s four fairy tales receive the highest positive differences, confirming that the model successfully captured his stylistic fingerprint. All of Constance’s fairytale receive negative perplexity scores, suggesting the model distinguishes between the style of the two authors. Interestingly, the published version of *The Selfish Giant* receives a small positive score, while the manuscript version receives a negative score, suggesting the model finds the latter more aligned with Constance’s style.

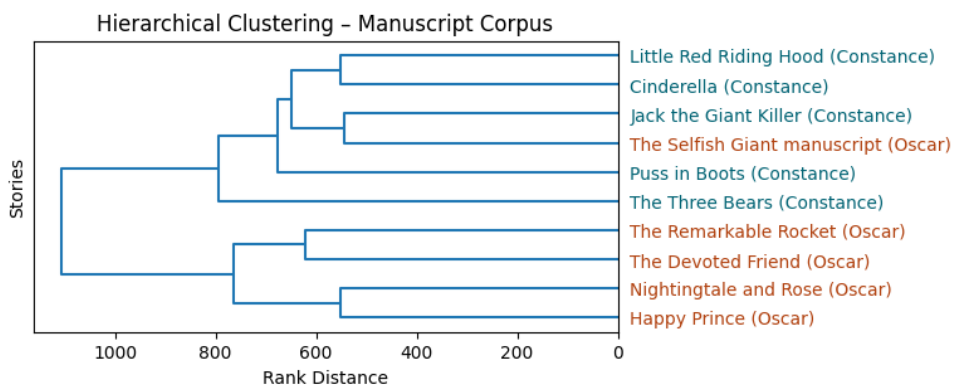
Both models indicate towards a shared stylistic profile between the story *The Selfish Giant* and Constance’s writings. There is a compelling difference in the case of the manuscript under Oscar-finetuned model, where the manuscript is clustered closely with Constance Wilde, while the published version is closer to the Oscar texts. We also note that, due to data availability, the Oscar model was fine-tuned on a greater amount on text than the Constance model. However, our results prove that fine-tuning on even a small amount of data is sufficient for



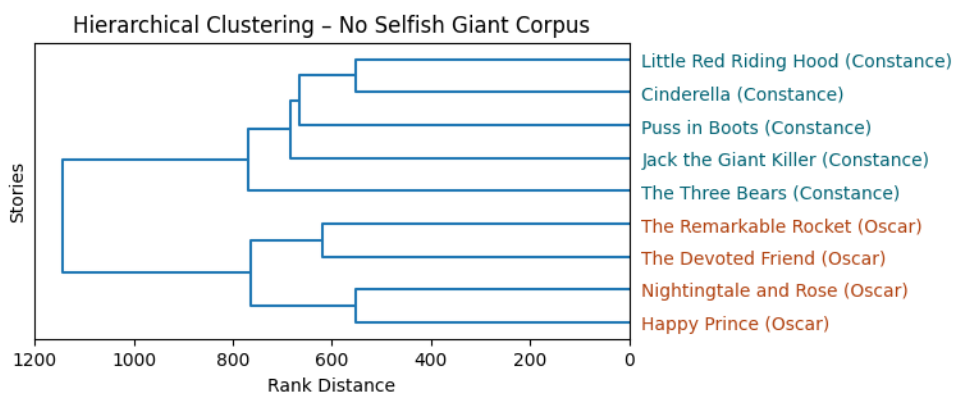
(a) Full Corpus



(b) Published Version Corpus



(c) Manuscript Version Corpus



(d) No "Selfish Giant" Corpus

Figure 1: Clustering Results for different versions of the corpus.

Table 7: Perplexity results for the Constance-finetuned model

Text	Base	Fine	Diff.
Cinderella (Constance Wilde)	50.57	41.12	9.45
The Selfish Giant	21.99	21.46	0.53
The Selfish Giant (manuscript)	18.58	18.48	0.10
The Star-Child (Oscar Wilde)	19.56	21.10	-1.54
The Fisherman and his Soul (Oscar Wilde)	20.97	22.92	-1.95
The Birthday of the Infanta (Oscar Wilde)	35.76	38.05	-2.29
The Young King (Oscar Wilde)	25.36	27.74	-2.37
A House of Pomegranates (Oscar Wilde)	29.79	32.34	-2.55

Table 8: Perplexity results for the Oscar-finetuned model

Text	Base	Fine	Diff.
The Birthday of the Infanta (Oscar Wilde)	35.76	30.89	4.87
A House of Pomegranates (Oscar Wilde)	29.79	26.14	3.65
The Fisherman and his Soul (Oscar Wilde)	20.97	18.80	2.17
The Star-Child (Oscar Wilde)	19.56	17.69	1.87
The Selfish Giant	21.99	20.27	1.72
Cinderella (Constance)	50.57	50.63	-0.06
The Selfish Giant (manuscript)	18.58	18.64	-0.06
Little Red Riding Hood (Constance)	25.08	25.27	-0.19
The Three Bears (Constance)	12.41	13.51	-1.10
Puss in Boots (Constance)	21.70	22.98	-1.29
Jack the Giant Killer (Constance)	21.60	24.06	-2.46

capturing nuances in stylistic profiles.

5 Conclusion

We introduced and analyzed a novel authorship attribution case - the children's stories authored by Constance and Oscar Wilde, which, to the best of our knowledge, have not been included in a computational authorship analysis up to date.

We analyzed the corpus using supervised (SVM classification with a string kernel) unsupervised (Hierarchical Clustering) and LLM-based methods (LLM-perplexity of fine-tuned models). Despite the small amount of available data, all three independent methods converged towards a shared conclusion: Constance Wilde might have authored or participated in the writing of one of the stories published in Oscar Wilde's collection.

Moreover, our results show that Constance Wilde was, at the time, a writer in her own right, with a style distinguishable from that of her husband.

Further studies could include collaborative authorship analyses through paragraph-level exploration in order to expand the collaboration hypothesis and lexical stylometric methods to provide insights into the similarities of "The Selfish Giant" and the rest of the texts authored by the two writers.

6 Ethical Considerations

This study makes use of literary works that are no longer under copyright protection. The research does not involve human participants, experiments, or any data that would raise ethical or privacy concerns.

7 Acknowledgements

This research was supported by the Ministry of Education and Research, CNCS-UEFISCDI, project SIROLA, number PN-IV-P1- PCE-2023- 1701, within PNCDI IV.

References

- Anne Clark Amor. 1999. Constantly undervalued: A centenary appreciation of constance wilde. *The Wildean*, (14):8–25.
- Douglas Biber and Susan Conrad. 2019. *Register, genre, and style*. Cambridge University Press.
- Cindy Chung and James Pennebaker. 2011. The psychological functions of function words. In *Social communication*, pages 343–359. Psychology Press.
- Hugh Craig and Arthur F Kinney. 2009. *Shakespeare, computers, and the mystery of authorship*. Cambridge University Press.
- Anca Dinu, Liviu P Dinu, Alina Resceanu, and Ion Resceanu. 2012a. Some issues on the authorship identification in the apostles' epistles. In *LRE-Rel Workshop on Language Resources and Evaluation for*

- Religious Texts, co-located with LREC 2012, Istanbul*, pages 18–24.
- Liviu P Dinu. 2003. On the classification and aggregation of hierarchies with different constitutive elements. *Fundamenta Informaticae*, 55(1):39–50.
- Liviu P Dinu, Vlad Niculae, and Octavia-Maria Şulea. 2012b. Pastiche detection based on stopword rankings. exposing impersonators of a romanian writer. In *Proceedings of the Workshop on Computational Approaches to Deception Detection*, pages 72–77.
- Liviu Petrisor Dinu, Marius Popescu, and Anca Dinu. 2008. **Authorship identification of romanian texts with controversial paternity**. In *Proceedings of the International Conference on Language Resources and Evaluation, LREC 2008, 26 May - 1 June 2008, Marrakech, Morocco*. European Language Resources Association.
- Nudrat Habib, Tosin Adewumi, Marcus Liwicki, and Elisa Barney. 2025. Trends and challenges in authorship analysis: A review of ml, dl, and llm approaches. *arXiv preprint arXiv:2505.15422*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Liang Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. *Iclr*, 1(2):3.
- Baixiang Huang, Canyu Chen, and Kai Shu. 2024. Can large language models identify authorship? In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 445–460.
- Baixiang Huang, Canyu Chen, and Kai Shu. 2025a. Authorship attribution in the era of llms: Problems, methodologies, and challenges. *ACM SIGKDD Explorations Newsletter*, 26(2):21–43.
- Weihang Huang, Akira Murakami, and Jack Grieve. 2025b. Attributing authorship via the perplexity of authorial language models. *PloS one*, 20(7):e0327081.
- Patrick Juola. 2015. The rowling case: a proposed standard analytic protocol for authorship questions. *Digital Scholarship in the Humanities*, 30(suppl_1):i100–i113.
- Cyril Labbé and Dominique Labbé. 2001. Inter-textual distance and authorship attribution corneille and molière. *Journal of Quantitative Linguistics*, 8(3):213–231.
- Frederick Mosteller and David Wallace. 1964. *Inference and disputed authorship: The Federalist*. Reading, MA: Addison-Wesley.
- Franny Moyle. 2012. *Constance: The Tragic and Scandalous Life of Mrs. Oscar Wilde*. Open Road Media.
- Marius Popescu and Liviu P Dinu. 2008. Rank distance as a stylistic similarity. In *Coling 2008: Companion volume: Posters*, pages 91–94.
- Kerry Powell and Peter Raby. 2013. *Oscar Wilde in Context*. Cambridge University Press.
- Violeta Sotirova. 2015. *The Bloomsbury Companion to Stylistics*, 1st edition. Bloomsbury Companions. Bloomsbury Academic, London.
- Arjuna Tuzzi and Michele A Cortelazzo. 2018. *Drawing Elena Ferrante’s Profile. Workshop Proceedings, Padova, 7 September 2017*. Padova University Press.
- Hans Van Halteren, Harald Baayen, Fiona Tweedie, Marco Haverkort, and Anneke Neijt. 2005. New machine learning methods demonstrate the existence of a human stylome. *Journal of Quantitative Linguistics*, 12(1):65–77.