

# Beyond Canonicity: Modeling Canon/Archive Literary Change in French Fiction

Jean Barré<sup>1,2,\*</sup>, Thierry Poibeau<sup>1,2</sup>

<sup>1</sup>École normale supérieure - Université PSL, 45 rue d'Ulm, Paris, 75005, France

<sup>2</sup>Lattice (Langues, Textes, Traitements informatiques, Cognition), 1 rue Maurice Arnoux, Montrouge, 92049, France

## Abstract

This study offers a fresh perspective on the Canon/Archive problem in literature through computational analysis. Following Tynianov's understanding of literature, we adopt a dynamic approach to literature by proposing a model of literary variability using the Kullback-Leibler divergence. We retrieve key authors and works that shape the broad outlines of literary change. Our aim is to evaluate the importance of canonical authors on literary variability. We opt for a cohort-driven setup to analyze the variability contributed by a given text, focusing on specific formal and semantic aspects of texts such as topics, lexicon, characterization, and chronotope. The findings reveal that canonical authors tend to contribute slightly more to literary change than those from the archive.

## Keywords

literary history, computational literary studies, distant reading, literary variability, canon/archive, cohort-driven model, cultural analytics, natural language processing,

## 1. Introduction

The Canon/Archive problem is a well-known issue in the field of the Computational Literary Studies (CLS). It has been and continues to be a fundamental aspect of the CLS field, as computational methods allow researchers to expand their investigations beyond the limited study of the Canon and its restricted number of texts. With the ability to process vast amounts of digitized texts in a matter of hours, researchers can now engage in distant reading, as proposed by Moretti [23], and conduct experiments on the textual content of literary works. This approach enables scholars to zoom in and out from the literary past, leading to a better understanding of general trends describing literary evolution.

This introduction of new perspectives and alternative modes of inquiry raises a fundamental question: "Do we understand the outlines of literary history?". Underwood [31] eloquently poses this question, contemplating whether the texts preserved thus far adequately represent the entire spectrum of literary production, or if the discipline of literary has been constrained by narrow perspectives throughout its existence.

---

CHR 2023: Computational Humanities Research Conference, December 6 – 8, 2023, Paris, France


\*Corresponding author.

✉ jean.barre@ens.psl.eu (J. Barré); thierry.poibeau@ens.psl.eu (T. Poibeau)

🌐 <https://crazyjeannot.github.io/> (J. Barré); <https://www.lattice.cnrs.fr/en/members/direction/thierry-poibeau/> (T. Poibeau)

🆔 0000-0002-1579-0610 (J. Barré); 0000-0003-3669-4051 (T. Poibeau)

© 2023 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

This line of investigation is not entirely new, as Iouri Tynianov expressed similar concerns in 1927 when he stated that “The theory of value in literary scholarship fueled the temptation to study major (but also isolated) phenomena and has turned literary history into a ”history of generals””[30]. However, Tynianov offered a way out by suggesting that the value of a given literary phenomenon should be understood in terms of its “significance and evolutionary qualities”. According to Tynianov [30], literary recognition is a dynamic process, and analyzing it requires studying “literary variability”. This refers to the diversity and range of formal elements present in literary works. It encompasses the different ways authors employ language, style, themes, narrative structures, characterization, settings, and other literary elements to create unique and distinct pieces of literature. This perspective sees literary history not as a linear chronology but considers literature within a dynamic and indivisible process that is constantly evolving. Every written text available in a library has the potential to influence the process of writing. Accounting for this perpetual movement necessarily requires an understanding of how each new text seeks to formally distinguish itself from its predecessors while still being shaped, for example, by a specific, conscious or unconscious, generic intertextuality.

This study aims at evaluating to what extent canonical works are reliable witnesses in terms of literary variability. Previous research uncovered disparities in the textual content between what is considered canonical and non-canonical across various corpora and cultural backgrounds ([1], [33], [14], [9]). These previous studies showed that canonical sets share to some extent (at least for specific timespans) an intrinsic norm. As outlined by Barré, Camps, and Poibeau [9], this research identifies what Altieri [2] terms a “cultural grammar”, suggesting that canonical literary works function as foundational texts shaping the norms, values, and conventions within a specific cultural tradition.

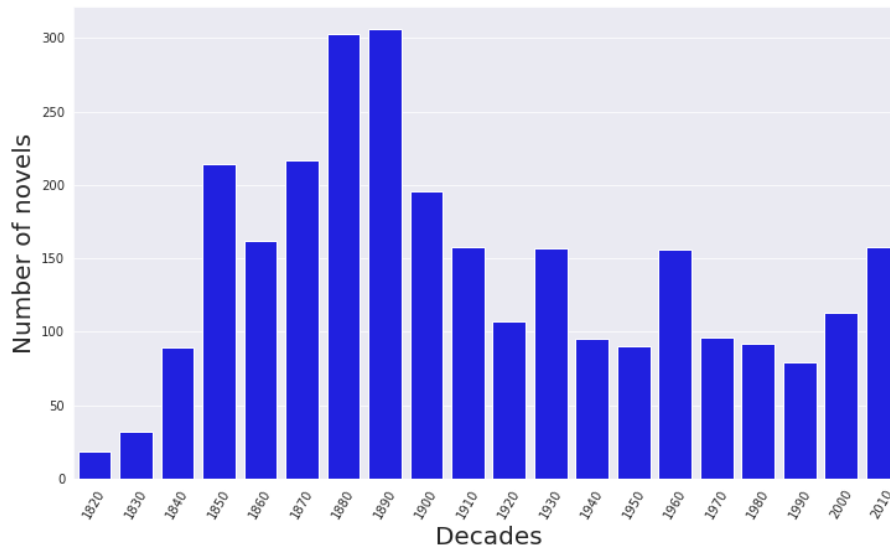
Hence, a pertinent question emerges: Does this specific norm comprehensively account for literary variability, or is it missing something? On the one hand, canonical novels, renowned for their significance and influence in literary traditions, can indeed act as pivotal benchmarks for both writers and readers [28]. Their impact on literary practices can manifest in various ways: inspiring new writing styles, introducing innovative themes, or encouraging formal experiments. Writers can be influenced by these canonical novels, either seeking to differentiate themselves or to align with their influence [25]. On the other hand, the concept of canonicity can be seen as biased and limited in capturing the evolution of formal practices. Indeed, the canon acts as a framework shaping not only present literary creation but also influencing how we retrospectively perceive the past, aligning it with contemporary norms. Therefore, the complex nature of the canonization process, influenced by external factors such as the school system and editorial policies [13], may hinder its ability to incorporate *avant-garde* literary changes.

In this paper, we introduce an operational model aiming to define and measure formal variability in 19th and 20th century French novels. Our approach is firmly rooted in established canonical sets derived from prior research on contemporary reception in France ([9]). Our main objective is to explore whether the canonized works selected from contemporary reception accurately mirror the broader spectrum of change within the overall French novelist production. For this purpose, we try to identify the key works and key authors driving literary variability. By analyzing formal aspects and considering literature as a dynamic system, we seek to gain insights into the flow of literary variability.

## 2. Materials and Methods

### 2.1. Corpus

This study is based on the corpus collected in the framework of the “ANR Chapitres”<sup>1</sup>, a corpus of nearly 3000 French novels [3]. The goal of this project was to evaluate the pace of change in the length of chapters over two centuries. The corpus is structured in XML-TEI<sup>2</sup> (Text Encoding Initiative) encoding, to add metadata to the texts. The period concerned extends over two centuries of novel production, from the 19th to the 20th century, as can be seen in Figure 1.



**Figure 1:** Distribution of the number of novels over time

Each text in the corpus is enriched with metadata, including subgenre tags and authors' dates (birth and death). The latter are highly relevant for our work as we focus on the effect of cohorts on the pace of literary change.

### 2.2. Textual features

The concept of literary variability encompasses a broad spectrum of possibilities, and it can manifest itself in various ways within a text. By examining specific elements such as themes, characterization, vocabulary, and chronotope, we aim to understand how novels have evolved across different time periods. However, some notions (such as 'the plot') are hard to formalize and are thus not included in this study.

<sup>1</sup><https://chapitres.hypotheses.org/>

<sup>2</sup>TEI Consortium, eds. TEI P5: Guidelines for Electronic Text Encoding and Interchange. Version 1.0. TEI Consortium. <http://www.tei-c.org/Guidelines/P5/>.

We first implemented topic modeling methods to extract topics from the texts. The Python library Bertopic [16] was used in a guided setting. This refers to a set of techniques that influence the topic modeling process by providing predefined seed topics for the model to converge towards. These techniques enable users to specify a predetermined number of topic representations that are guaranteed to appear in the results. We constructed a list of 50 topics we found relevant for our study and retrieved their proportion within each novel.<sup>3</sup>

We also implemented a Bag-of-ngrams approach to retrieve the lexicon dimension of literary change. To do so, we rely on the 1000 most frequent lemmas and 1000 most frequent bigrams of lemmas. This may echo the paper by Cranenburgh and Koolen [15], which showed that using only unigrams and bigrams was sufficient to classify literary texts in terms of their literary quality. We did not remove stopwords, since they may reflect an unconscious and automatic structural way of writing [27], rather than less frequent words related to the content and themes of the text. Our hypothesis is that the structural way of writing novel changes over time and cohorts. Bag-of-words techniques work well for various experiments in the CLS field (stylometry and author attribution for example [18]), but they are quite controversial from a literary point of view, since they exclude a great deal of information, including word order and syntax. They are also limited in that they do not take into account the semantic drift of words over time. For instance the word “wild” does not refer to the same meaning when used in an adventure novel from the late 19th century or in a climate fiction from the late 20th century. Bearing this in mind, we assumed that bag-of-features still capture some dimensions of literary change, particularly as regards the very frequent structural elements.

One of the aims of this study was to capture *chronotope* information, which is a term coined by the Russian literary scholar Bakhtin [4]. In substance, the concept of *chronotope* explores how the relationship between time and space influences the portrayal of characters, the development of plotlines, and the themes conveyed within a literary text. In the Natural Language Processing (NLP) context, Kohlmeyer, Repke, and Krestel [19] demonstrated the limitations of traditional document embeddings (optimized for shorter texts) in capturing complex facets in novels (such as time, place, atmosphere, style, and plot). To address this problem, they propose to use multiple embeddings reflecting different facets, splitting the text semantically rather than sequentially. Inspired by these findings, we adapted their methodology. By using an NLP pipeline specifically tuned for novels, (fr-BookNLP, part of the multilingual BookNLP project [7, 5]), we extracted literary entities representing the *chronotope*, specifically focusing on FAC, TIME, LOC, and VEH.<sup>4</sup> The presence of chronotope elements in a novel is highly influenced by its subgenre categorization. We believe that this type of information is crucial for our task as it has the potential to capture significant aspects of literary variability. To obtain vector representations of the *chronotope* elements in novels, we trained a Paragraph Vectors model [20] (Doc2Vec) using a subset of our novel dataset. We then generated four vector embeddings from our four spans of entities. Each facet has a vector with 300 dimensions, resulting in a 1200 dimensions vector that captures the *chronotope* information for each novel.

We also considered that characterization was a significant element in our task, as we be-

---

<sup>3</sup>For the topic modeling process, see appendix A.1

<sup>4</sup>Respectively Facilities, Time, Location, Vehicle - see [6] for more information on the NER labels, and appendix A.2 for the evaluation of Fr-BookNLP

lieved that changes in literature could influence how characters are portrayed to readers. We thus focused on identifying key verbs that drive the actions of the main characters and the adjectives used to describe them. In line with Woloch [34]’s concept of the character space as “the encounter between an individual human personality and a determined space and position within the narrative as a whole”, we used coreference resolution techniques, specifically those offered by fr-BookNLP<sup>5</sup>, to automatically detect and analyze the distribution of character mentions throughout the narrative [8]. We used the Spacy parser to extract the verbs and adjectives associated with each character mention. By analyzing the syntactic structure of the text, we identified the verbs that represented the actions of the characters and the adjectives that characterized them. For each novel, we selected the top five main characters and generated two vector embeddings of 300 dimensions: one representing the adjectives associated with the characters and the other representing the verbs. These embeddings capture the semantic information related to the characters’ traits and actions, providing a compact representation of their characteristics within the narrative.

By incorporating these various aspects, each novel can be represented as a concatenated multidimensional vector with 3850 dimensions, 50 for the topics, 2000 for the bag-of-ngrams, 1200 for chronotope elements and 600 for the characterization. Therefore, our vector representation provides a comprehensive formalization of the novel, enabling further analysis and comparisons.

### 2.3. Measuring Literary Variability

Basing our work on the aspects presented above (topics, bag-of-ngrams, chronotope, and characterization), we had to find a way to grasp literary variability. We decided to implement a commonly used measure, the Kullback-Leibler divergence (KLD). It is a type of statistical metric that makes it possible to quantify the dissimilarity between two probability distributions: the target distribution  $P$  and a reference distribution  $Q$ . Within this framework, we assessed the variability of a text by measuring the surprise or deviation of that text from a set of other texts. Specifically, the KLD from  $Q$  to  $P$  is defined as follows:

$$KL(P||Q) = \sum_x P(x) \log\left(\frac{P(x)}{Q(x)}\right)$$

where  $P$  represents our formal features for a text, normalized as a probability distribution, and  $Q$  stands for the average of all the texts we wish to compare  $P$  with. This measure, derived from information theory, finds application in various fields, including assessing sample diversity in ecology and examining elements of linguistic evolution [12]. Barron, Huang, Spang, and DeDeo [10] applied KLD to a corpus of debates in the French revolution’s first parliament, assessing both the novelty of a particular speech compared to prior speeches and its transience compared to future ones.

In the realm of literature, Algee-Hewitt, Allison, Gemma, Heuser, Walser, and Moretti [1] proposed this method to determine the informational content of texts by evaluating the predictability of word-to-word transitions, taking into account the range of possible transitions.

<sup>5</sup>A discussion on the evaluation of fr-BookNLP and its limitations is provided in the appendix A.2.

Liddle [21] also discussed the possibility that mathematical information theory may be relevant to literary analysis by showing statistically significant correlations between national histories of the novel and information-theoretical pressures.

Previous research showed that literary change throughout an author's life was powerful enough to predict the publication date of a given text [29]. However, other studies demonstrated that literary change brought about by an author throughout their life remains limited compared to the cohort effect [32]. In other words, literary change appears to be driven by cohort renewal, which is indeed relevant since events that shape an author's life, and that are likely to have an impact on their writing style, also influence all authors within the same generation.

Measuring the variability of a text in relation to a set of other texts immediately places us in a dual configuration: we can measure the variability of a given text with the works that precede it and with those that follow it. Studying the circulation, selection, and propagation of literary patterns in a group of texts, we can understand the dynamics of literary change and the extent to which an author's language patterns influence and are adopted by others. From a literary perspective, measuring change between two successive texts may not make much sense, as numerous factors can come into play, such as affiliation with a particular subgenre, a specific period, a literary school, or even the author themselves. Therefore, a specific framework is necessary to conduct our experiments, which revolves around the notion of generation, assuming that texts produced within the same generation share certain characteristics and influences.

Béhar defines a generation as a concept that aims to "understand the succession of aesthetic productions based on a community of upbringing, interests, and ideas specific to the same age group of writers, following a periodicity of approximately 30 years linked to historical and political cycles"[11]. This generation-based approach allows us to examine the changes and innovations introduced by authors within their respective cohorts, while also considering the broader historical and literary context in which these works emerge. It provides a more nuanced understanding of how literary variability is shaped and influenced by various factors, contributing to a richer analysis of the dynamic nature of literature.

To further support Béhar's argument, Moretti [24] also evaluated the regularity of the replacement of literary subgenres that he examines. He suggested that "a sort of generational mechanism seems to be the best way to account for the regularity of the cycle of novelistic production". Moretti's analysis focused on the cycle of change, considering a timeframe of 25 to 30 years. These studies are complemented by Underwood, Kiley, Shang, and Vaisey [32] seminal research, who showed the significance of cohorts on literary change. Their findings indicated that cohorts have such a substantial impact that they account for more than half of the amount of change in literature.

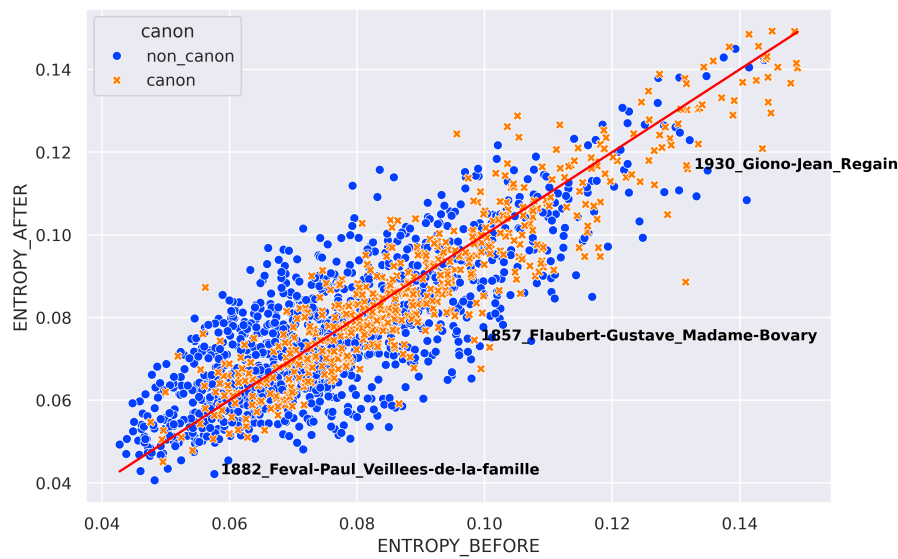
Building upon their conclusions, we computed KLD comparing successive cohorts in a timespan of 30 years. For instance, if analyzing a novel published in 1970 by an author born in 1930, we compare it with all the books written by authors born between 1870 to 1900 to evaluate the extent of change. This approach enables us to consider cohorts as a rolling phenomenon, since defining arbitrary cohorts would not be representative of the phenomenon of cohort succession. This methodology allows us to view literature as a continually evolving synchronic system, framed by cohorts.



### 3. Results

#### 3.1. Literary dynamics: between novelty and influence

Figure 2 represents the amount of variability for each text: The x-axis represents the entropy of each text, relative to the cohort preceding the text’s author, indicating the level of surprise or formal novelty in the text compared to its preceding cohort. A text with a high ”surprise” score on this axis would be considered formally innovative. The y-axis represents the entropy of each text, relative to the cohort following the text’s author, indicating the level of surprise or influence that the text has on the next cohort. A text with high surprise on this axis indicates that it has little influence on what follows. A text that is both highly influential and innovative would receive a high value on the x-axis and a low value on the y-axis.



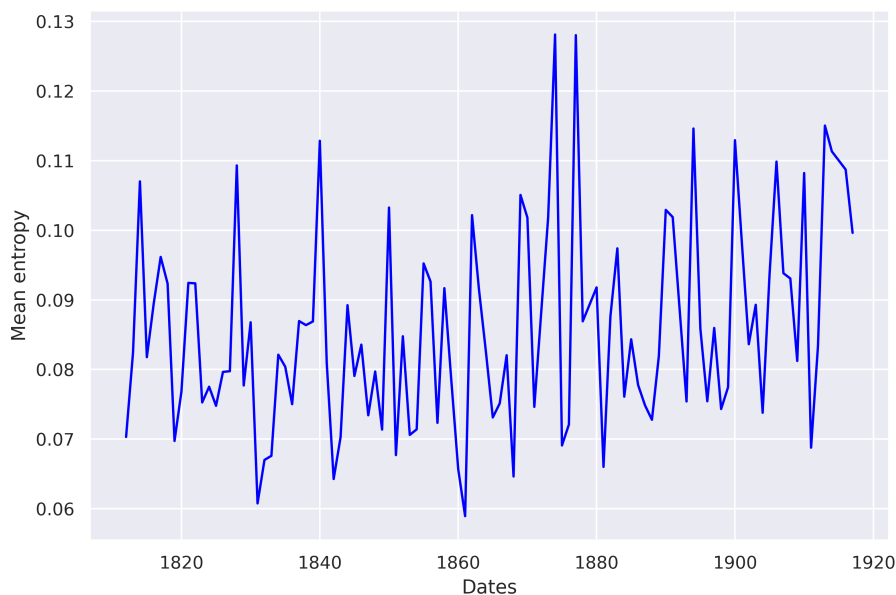
**Figure 2:** Amount of variability comparing cohort before and cohort after

As the graph is complex, three novels are depicted in order to make it easier to understand how the graph works. Gustave Flaubert’s *Madame Bovary* stands out with high novelty and high influence scores, thanks to its groundbreaking narrative style, character development, and enduring impact on literature. Jean Giono’s *Regain* receives a high novelty score for its exploration of resilience and human connection to nature, but its low influence suggests that it did not gain immediate widespread recognition. Émile Souvestre’s *Un Philosophe sous les Toits* addresses social issues and garnered significant influence despite not being part of the French literary canon, making it historically significant. Canonical texts are highlighted in orange, following the canonical sets at the author scale from previous work [9]. It is notable that most of these texts are positioned below the  $x=y$  line. This suggests that canonical works exhibit greater variability compared to the preceding cohort but slightly less variability compared to the following one. This may indicate that canonical works have a higher level of innovation and influence.

### 3.2. Signal of literary variability

These initial findings were highly intriguing, leading us to pursue a complementary approach to gain a deeper understanding. We focused on the x-axis, which we deemed more comprehensible from a literary standpoint. The change introduced by a text (or an author) in relation to the previous generation is intuitively grasped as each text finds a way to differentiate itself from the broader literary production of a given period. By associating the entropy value obtained for each text with its author's birth date, we were able to represent in figure 3 the signal of change over time.

Through visual representation, we showcased the patterns and fluctuations observed in the analysis of KL divergence or entropy. This approach allowed us to capture the dynamic nature of literary change and observe how it manifests itself over different historical periods. The resulting graph provides a visual narrative of the evolving literary landscape, shedding light on the formal shifts in the realm of literature. It offers a compelling visual representation of the signal of literary change, enabling a more nuanced understanding of the complex processes at play in cultural and artistic evolution.



**Figure 3:** Signal of change comparing previous cohort

Each peak corresponds to a significant variability introduced by a specific author (or a group of authors sharing the same birth date). This approach allows us to identify the names of key works and key authors that drive literary variability. It should be mentioned that the last peak should be ignored due to the lack of authors born around 1910 and later in our corpus.

Jules Verne, born in 1828, is the author who mainly explains the second peak in variability. His works, particularly *Vingt-mille lieues sous les mers*, published in 1870 (with 0.149 KLD), and



*L'Île Mystérieuse*, published in 1875 (with 0.232 KLD) exemplify his innovative approach to literature. In the former one, Verne introduced readers to Captain Nemo's underwater vessel, the Nautilus, which travels beneath the seas and explores uncharted depths. This visionary depiction of a futuristic submarine, powered by electricity and equipped with advanced technology, set the stage for the emergence of the science fiction genre. In *L'Île Mystérieuse*, Verne combined elements of adventure and survival on a remote island with the exploration of technology and engineering. The novel tells the story of a group of castaways who use their knowledge and resourcefulness to survive and thrive on the island.

Verne's innovative storytelling can be seen as a response to the social fascination with progress and exploration. During the late 19th century, the world was witnessing rapid advancements in science and technology, driven by the Industrial Revolution and scientific discoveries. This era of progress and innovation deeply influenced the literary landscape, as writers such as Verne sought to capture the spirit of exploration and curiosity prevalent in society. This way, Verne laid the foundation for a mixture of adventure and science fiction subgenres.

The peak from 1877 is led by Raymond Roussel, a lesser-known but highly innovative writer, particularly with his works *Impressions d'Afrique*, published in 1910 (with 0.145 KLD) and *Locus Solus*, published in 1914 (with 0.21 KLD). The former one is a novel that defies traditional narrative conventions and follows a dreamlike, non-linear structure. The story revolves around a group of travelers who embark on a journey through Africa, encountering strange and surrealist occurrences along the way. The second narrative is also characterized by its intricacy and complexity, as it contains multiple layers of storytelling that takes the reader on a tour of the estate of a scientist named Martial Canterel, where he showcases a series of bizarre and macabre inventions. Roussel's experimental and imaginative storytelling style set him apart as a pioneer in avant-garde literature, and his works can be seen as early surrealism.

René Crevel and Nathalie Sarraute lead the peak in 1900. René Crevel's work *Le roman cassé*, published in 1935 (with 0.17 KLD) and Nathalie Sarraute's *Tropismes*, published in 1939 (with 0.159 KLD) both exemplify their innovative approaches to literature, as they challenged conventional narrative structures and scrutinized the inner workings of human consciousness.

Crevel's novel breaks away from traditional linear storytelling and embraces a fragmented and non-linear narrative style. The author's exploration of the subconscious mind and his use of stream-of-consciousness writing make the novel a precursor to the surrealist and modernist movements. The novel centers around the mental states of its characters, delving into their thoughts, dreams, and desires. The title *Le roman cassé* itself, which translates to *The Broken Novel*, is indicative of Crevel's intention to dismantle traditional narrative conventions and explore new modes of expression. His work can be seen as an early example of the deconstruction of the novel form, where the focus shifts from external events and plot-driven storytelling to an exploration of the characters' inner lives and psychological states.

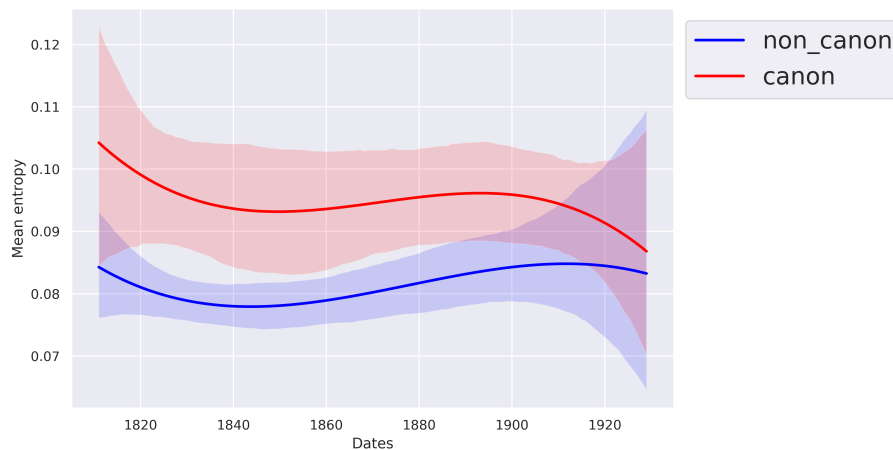
Nathalie Sarraute's *Tropismes* is a collection of interconnected short prose pieces that explore the subtle and fleeting movements of the characters' inner thoughts and feelings. Sarraute's writing style is characterized by its precision and attention to the nuances of human behavior. She coined the term "tropismes" to describe these brief and involuntary movements of the characters' consciousness. Her innovative use of language and her focus on the psychological subtleties of her characters set her apart as a pioneer of the *Nouveau Roman* movement.

Thus, our novelty signal highlights works and authors who have significantly distanced

themselves from the dominant formal rules of the previous generation. We identified authors who have contributed to the creation of new sub-genres or avant-garde writers with varying degrees of recognition, such as Raymond Roussel, an author from the *Archive*. These peaks might represent pivotal moments in literary history where new ideas, styles, or narrative techniques emerge, leading to a distinct shift in the literary landscape. Nevertheless, it is worth noting that any work that differs formally from the majority of other novels stands out. For instance, children’s novels such as *Le petit prince* (Antoine de Saint-Exupéry) also emerge prominently in the scores.

### 3.3. Canonical novels, the drivers of literary variability?

When examining the list of authors who stand out in their contribution to literary change, many of them are well-known and directly associated with the literary canon. To assess the amount of change among canonical works compared to non-canonical works from the archive, we project in figure 4 two distinct curves onto the graph based on their canonicity labels, at the author level (considering all works by an author as canonical). By considering the canonicity distinction, we gain a deeper understanding of how these different subsets of texts contribute to the overall landscape of literary variability.



**Figure 4:** Mean Novelty per birth year, broken out by canonicity tag

Thus, we observe two distinct curves on the graph: the red curve representing canonical authors and the blue one representing non-canonical authors. The margin of error for canonical authors is larger due to their smaller number compared to the archival authors. The gap between the two sets remains relatively stable over a century of authors’ birth dates. The clear conclusion from the graph is that canonical authors tend to introduce more variability in their novels compared to non-canonical authors.

The smaller difference observed towards the end of the period suggests a few possibilities. It implies that the overall corpus becomes more limited towards the end, with a smaller pool of texts available for analysis. Furthermore, it indicates that the criterion of canonicity might be less relevant for the last generations of authors.

## 4. Limitations

Our approach is subject to the inherent accuracy limitations of many NLP algorithms used, including Fr-BookNLP<sup>6</sup>, Spacy, and Bertopic, all of which are prone to error.

The choice of a 30-year time frame for cohort succession in our experiments is somewhat subjective. Although it is reasonable to assume that significant changes occur within this window compared to shorter intervals such as 5 or 10 years, the selection remains debatable. Furthermore, our comparisons are limited to successive cohorts, neglecting the potential influence of earlier literary works. Authors are likely to have been influenced by canonical texts published decades or even centuries before their own works, which warrants further consideration.

We faced the challenge of conducting close readings. While we have identified distinctive authors and texts, we have not provided textual evidence of the observed changes. Given the large-scale nature of our study, this is understandable, but future research should strive to incorporate detailed textual analysis to support our findings. Future work will be dedicated understanding which features contribute to what extent to the observed literary variability.

## 5. Conclusion

In conclusion, this study has provided valuable insights into the dynamics of literary variability and the role of canonical works in the French literary landscape. Through our operationalization of formal variability and our cohort-driven model, we succeeded in identifying the names of key works and key authors that drive literary variability. Analyzing formal aspects such as topics, styles, chronotopes, and characterization in a large corpus of novels, the study aims to uncover patterns of literary change and explore the relationships between texts, authors and cohorts. By organizing texts into generations, we established a temporal and contextual framework that allows us to capture and analyze the evolving literary dynamics over time. This approach acknowledges that texts produced within the same generation share certain characteristics and influences, providing a meaningful basis for measuring and understanding literary variability.

We then investigated how accurately the canon reflects the overall degree of change in literature. Surprisingly, our findings indicate that canonical authors contribute more variability than non-canonical authors. At first glance, this might seem counter-intuitive, given that the canonization process historically favors well-established, conventional works over avant-garde and experimental ones. This tendency could lead one to expect that the canon might display less variation in its literary characteristics. However, our results demonstrate that the canon is far from a monolithic entity. It is not a rigid, uniform collection that uniformly represents a particular literary style or period. This suggests that within the canon, there exists a spectrum of works, spanning from those aligning with existing norms to those that truly challenge boundaries and introduce novel literary elements. This implies that canonization is not an entirely conservative process.

One possible explanation could also be related to cultural and economic factors. Writers whose works are part of the archive often aim for a widespread readership and commercial

---

<sup>6</sup>see appendix A.2 for the evaluation of Fr-BookNLP

success, especially in subgenres associated with mass literature. In such subgenres, the “horizon of expectations” [17] of the audience might induce the authors to adhere to certain expected norms and styles. This emphasis on reaching a larger audience and meeting certain expectations might impact the level of experimentation and deviation from established norms, leading to a perceived lower mean variability in the archive compared to the canon.

Further analyses are needed to conclude on this issue. We plan to delve into a more granular examination of texts, shifting scale towards close reading. The intention is to meticulously analyze pivotal passages that significantly contribute to literary variability. The goal here is to discern the relevance of various facets and assess their textual manifestation. To achieve this, a more detailed exploration is planned, honing in on a specific subgenre. This narrower scope will facilitate a more intricate analysis and interpretation of the texts, enabling a deeper understanding of their distinctive literary attributes.

Moving forward, future research could also focus on investigating the role of specific subgenres in driving literary change. By examining whether the emergence or growth of certain subgenres corresponds to peaks of change, we can gain a deeper understanding of how different literary trends influence the overall dynamics of literature.

## **Acknowledgments**

Jean Barré’s PhD is supported by the EUR (Ecole Universitaire de Recherche) Translitteræ (programme “Investissements d’avenir” ANR-10- IDEX-0001-02 PSL and ANR-17-EURE-0025). This work was also funded in part by the French government under management of Agence Nationale de la Recherche as part of the “Investissements d’avenir” program, reference ANR-19-P3IA0001 (PRAIRIE 3IA Institute). The authors also wish to thank the anonymous reviewers whose comments have helped us to substantially improve this paper.

## References

- [1] M. Algee-Hewitt, S. Allison, M. Gemma, R. Heuser, H. Walser, and F. Moretti. “Canon/Archive. Large-scale Dynamics in the Literary Field”. In: *Pamphlets of the Stanford Literary Lab*. Pamphlets of the Stanford Literary Lab 11 (2016). URL: <https://litlab.stanford.edu/LiteraryLabPamphlet11.pdf>.
- [2] C. Altieri. “An Idea and Ideal of a Literary Canon”. In: *Critical Inquiry* 1 (Sept. 1983), pp. 37–60. DOI: 10.1086/448236.
- [3] ANRChapitres. *Corpus Chapitres*. Version v1.0.0. 2022. DOI: 10.5281/zenodo.7446728.
- [4] M. Bakhtin. *The dialogic imagination: four essays*. Slavic series 1. Austin, Tex: University of Texas Press, 2011. 443 pp.
- [5] D. Bamman. *BookNLP*. 2021. URL: <https://github.com/booknlp/booknlp>.
- [6] D. Bamman, S. Popat, and S. Shen. “An annotated dataset of literary entities”. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Minneapolis, Minnesota: Association for Computational Linguistics, 2019, pp. 2138–2144. DOI: 10.18653/v1/N19-1220.
- [7] D. Bamman, T. Underwood, and N. A. Smith. “A Bayesian Mixed Effects Model of Literary Character”. In: *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics*. Acl 2014. Baltimore, Maryland: Association for Computational Linguistics, 2014, pp. 370–379. DOI: 10.3115/v1/P14-1035.
- [8] J. Barré, P. Cabrera Ramírez, F. Mélanie, and I. Galleron. “Pour une détection automatique de l’espace textuel des personnages romanesques”. In: *Humanistica 2023*. Corpus. Association francophone des humanités numériques. Genève, Switzerland, 2023, pp. 56–61. URL: <https://hal.science/hal-04105537>.
- [9] J. Barré, J.-B. Camps, and T. Poibeau. “Operationalizing Canonicity. A Quantitative Study of French 19th and 20th Century Literature”. In: *Journal of Cultural Analytics* (2023). DOI: 10.22148/001c.88113.
- [10] A. T. J. Barron, J. Huang, R. L. Spang, and S. DeDeo. “Individuals, institutions, and innovation in the debates of the French Revolution”. In: *Proceedings of the National Academy of Sciences* 115.18 (2018), pp. 4607–4612. DOI: 10.1073/pnas.1717729115.
- [11] H. Béhar. *La littérature et son golem*. Vol. 1: Travaux de linguistique quantitative 58. Paris: H. Champion, 1996. 2 pp.
- [12] C. Bentz, D. Alikaniotis, M. Cysouw, and R. Ferrer-i-Cancho. “The Entropy of Words–Learnability and Expressivity across More than 1000 Languages”. In: *Entropy* 19.6 (2017), p. 275. DOI: 10.3390/e19060275.
- [13] P. Bourdieu. *Les règles de l’art. genèse et structure du champ littéraire*. Paris: Éditions du Seuil, 1992.

- [14] J. Brottrager, A. Stahl, A. Arslan, U. Brandes, and T. Weitin. “Modeling and Predicting Literary Reception”. In: *Journal of Computational Literary Studies* (2022). DOI: 10.48694/jcls.95.
- [15] A. van Cranenburgh and C. Koolen. “Identifying Literary Texts with Bigrams”. In: *Proceedings of the Fourth Workshop on Computational Linguistics for Literature*. Proceedings of the Fourth Workshop on Computational Linguistics for Literature. Denver, Colorado, USA: Association for Computational Linguistics, 2015, pp. 58–67. DOI: 10.3115/v1/W15-0707.
- [16] M. Grootendorst. *BERTopic: Neural topic modeling with a class-based TF-IDF procedure*. 2022. arXiv: 2203.05794.
- [17] H. R. Jauß. *Toward an aesthetic of reception*. Trans. by T. Bahti. Minneapolis, Minn: Univ. of Minnesota Press, 1982.
- [18] M. Kestemont. “Function Words in Authorship Attribution. From Black Magic to Theory?” In: *Proceedings of the 3rd Workshop on Computational Linguistics for Literature (CLFL)*. Proceedings of the 3rd Workshop on Computational Linguistics for Literature (CLFL). Gothenburg, Sweden: Association for Computational Linguistics, 2014, pp. 59–66. DOI: 10.3115/v1/W14-0908.
- [19] L. Kohlmeyer, T. Repke, and R. Krestel. “Novel Views on Novels: Embedding Multiple Facets of Long Texts”. In: *2021 Association for Computing Machinery*. (2021).
- [20] Q. V. Le and T. Mikolov. *Distributed Representations of Sentences and Documents*. 2014. arXiv: 1405.4053.
- [21] D. Liddle. “Could Fiction Have an Information History? Statistical Probability and the Rise of the Novel”. In: *Journal of Cultural Analytics* (2019). DOI: 10.22148/16.033.
- [22] L. Martin, B. Muller, P. J. O. Suárez, Y. Dupont, L. Romary, É. V. de la Clergerie, D. Seddah, and B. Sagot. “CamemBERT: a Tasty French Language Model”. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics* (2020).
- [23] F. Moretti. “Conjectures on world literature”. In: *New Left Review* (2000).
- [24] F. Moretti. *Graphs, maps, trees: abstract models for literary history*. London New York: Verso, 2007. 119 pp.
- [25] A. Mukherjee. *Canonicity*. 2017. DOI: 10.1093/obo/9780190221911-0054.
- [26] I. G. Nils Reimers. *Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks*. 2019. arXiv: 1908.10084.
- [27] J. W. Pennebaker. *The secret life of pronouns: what our words say about us*. New York: Bloomsbury Press, 2011. 352 pp.
- [28] G. Pollock. *Differencing the canon: feminist desire and the writing of art’s histories*. Revisions. London ; New York: Routledge, 1999. 345 pp.
- [29] O. Seminck, P. Gambette, D. Legallois, and T. Poibeau. “The Evolution of the Idiolect over the Lifetime: A Quantitative and Qualitative Study of French 19th Century Literature”. In: *Journal of Cultural Analytics* 7.3 (2022). DOI: 10.22148/001c.37588.

- [30] Y. Tynianov. "On Literary Evolution (1927)". In: *Permanent Evolution*. Boston, USA: Academic Studies Press, 2019, pp. 267–282. DOI: 10.1515/9781644690635-015.
- [31] T. Underwood. *Distant horizons: digital evidence and literary change*. Chicago: The University of Chicago Press, 2019, pp. 1–33. 206 pp.
- [32] T. Underwood, K. Kiley, W. Shang, and S. Vaisey. "Cohort Succession Explains Most Change in Literary Culture". In: *Sociological Science* 9 (2022), pp. 184–205. DOI: 10.15195/v9.a8.
- [33] T. Underwood and J. Sellers. "The "Longue Durée" of Literary Prestige". In: *Modern Language Quarterly* 77.3 (2016), pp. 321–344. DOI: 10.1215/00267929-3570634.
- [34] A. Woloch. *The One vs. the Many*. Princeton University Press, 2003.



## A. Appendix

### A.1. Topic modeling: detailed approach

We provided Bertopic 50 specific topics with a list of 10 words associated with each topic. These topics served as seed topics to guide the model’s convergence during the analysis. Bertopic is an algorithm with several layers: for the embedding one we employed a CamemBERT base sentence vectorizer [26, 22] to create embeddings for the sentences. These embeddings capture the semantic meaning of the sentences and facilitate Bertopic analysis. Then we employed Principal Component Analysis which allowed us to transform the high-dimensional embeddings into a lower-dimensional space while preserving the essential information, making the subsequent steps more efficient. Then we clustered the sentences into distinct groups based on their semantic similarities, with the HDBSCAN algorithm. It is a density-based clustering method that identifies clusters of varying shapes and sizes, allowing us to group sentences that share similar topics or themes. Throughout the analysis, the model was capable of retrieving more than the initial 50 predefined topics. However, in order to maintain consistency and focus on the specific topics we had predefined, we chose to stick with the 50 seed topics provided by Bertopic. This allowed us to have a more targeted and interpretable analysis, focusing on the topics that were of particular interest for our research.

### A.2. Fr-BookNLP evaluation

#### A.2.1. NER

**Table 1**

NER evaluation of Fr-BookNLP on literary texts

	precision	recall	$F_1$
PER	85.0	92.1	88.4
LOC	59.4	54.3	56.8
FAC	73.4	66.0	69.5
TIME	75.3	36.4	49.1
VEH	68.9	63.6	66,1

When evaluating the performance of the model, having better precision than recall implies that when the model identifies literary entities, it is more likely to be accurate in its predictions. Precision measures the percentage of correctly predicted literary entities out of all the predicted entities. This is beneficial for the analysis as it ensures that the literary entities identified are more likely to be correct, even though some relevant entities may be missed (lower recall). In this context, prioritizing precision helps in reducing false positives and improving the reliability of the identified literary entities. One important thing to note is that literary entities are not exactly the same thing as NER in NLP. The specificities of literary texts make the detection of this kind of entity more complicated. Therefore the results obtained, even if they may seem far from NLP standards, are state-of-the-art for the specific processing of literary texts.

### A.2.2. Coreference resolution evaluation:

**Table 2**

Coreference resolution evaluation of Fr-BookNLP on literary texts

Metrics	$F_1$	
$MUC$	88,0	<i>Average 76.4</i>
$B^3$	69,2	
$CEAF_e$	71.8	

The issue of duplication arises when the model detects the same character multiple times within the analyzed text. In some cases, the top five literary entities identified by the model may contain instances where two or more main characters from a text are the same character in terms of name or attributes. While this duplication might seem problematic at first glance, it is essential to understand the context and purpose of the analysis. In this particular study, the primary objective was not to identify unique and distinct characters but rather to retrieve a proxy for characterization as a whole. We aimed to capture the prevalence and significance of certain characters across different texts and literary works. Therefore, the focus is more on character representation and the overall impact of these characters on the literary landscape, rather than identifying completely separate and non-repeating characters.