



UvA-DARE (Digital Academic Repository)

Why the Daisy Sisters are different

A stylometric study on the oeuvre of Swedish author Henning Mankell and the Dutch translations of his work

Wijers, M.

DOI

[10.5281/zenodo.8093597](https://doi.org/10.5281/zenodo.8093597)

[10.48694/jcls.3585](https://doi.org/10.48694/jcls.3585)

Publication date

2023

Document Version

Final published version

Published in

Conference Reader: 2nd Annual Conference of Computational Literary Studies CCLS 2023 Würzburg June 2023

License

CC BY

[Link to publication](#)

Citation for published version (APA):

Wijers, M. (2023). Why the Daisy Sisters are different: A stylometric study on the oeuvre of Swedish author Henning Mankell and the Dutch translations of his work. In *Conference Reader: 2nd Annual Conference of Computational Literary Studies CCLS 2023 Würzburg June 2023* Journal of Computational Literary Studies.

<https://doi.org/10.5281/zenodo.8093597>, <https://doi.org/10.48694/jcls.3585>

General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

UvA-DARE is a service provided by the library of the University of Amsterdam (<https://dare.uva.nl>)



Conference Reader
2nd Annual Conference of
Computational Literary Studies
CCLS 2023 Würzburg
June 2023

Venue: University of Würzburg | Zentrum für Philologie und Digitalität
Campus Hubland Nord | Room 01.001 | Emil-Hilb-Weg 23 | D-97074 Würzburg

Local Organizer: Priority programme [SPP 2207 Computational Literary Studies](#)

Contact: spp2207@uni-wuerzburg.de

Hashtag: #CCLS2023

Conference Programme

Thursday | June 22, 2023

1:00 p.m. to 2:15 p.m. | Session 1

- Opening
- Human Depiction in Portuguese – Cláudia Freitas*, Diana Santos
- A Novel Approach for Identification and Linking of Short Quotations in Scholarly Texts and Literary Works – Frederik Arnold*, Robert Jäschke

2:45 p.m. to 4:15 p.m. | Session 2

- Automatic Topic-Guided Segmentation of Holocaust Survivor Testimonies – Eitan Wagner, Renana Keydar*, Amit Pinchevski, Omri Abend
- InvBERT: Reconstructing Text from Contextualized Word Embeddings by inverting the BERT pipeline – Kai Kugler*, Simon Münker, Johannes Höhmann, Achim Rettinger
- What do characters do? – Andrew Piper

4:45 p.m. to 6:15 p.m. Session 3

- Need a Good Book about Privacy? – Erik Ketzan*, Jennifer Edmond
- The Authorship of Stephen King's Books Written Under the Pseudonym "Richard Bachman": A Stylometric Analysis – Vincent Neyt, Mike Kestemont, Dorothy Henriette Modrall Sperling*
- Extracting Geographical References from Finnish Literature. Fully Automated Processing of Plain-Text Corpora – Harri Kiiskinen*, Asko Nivala, Jasmine Westerlund, Juhana Saarelainen

6:30 p.m. | Keynote

- Jan Rybicki: Reading too many books: first results on 10,005 Polish original and translated texts.

8:00 p.m. | Joint Dinner

Friday | June 23, 2023

9:15 a.m. to 10:45 a.m. | Session 4

- Stylistic History of the Hungarian Novel Based on Sentence Structures – Botond Szemes
- Why the Daisy sisters are different. a stylometric study on the oeuvre of Swedish author Henning Mankell and the Dutch translations of his work – Martje Wijers
- Translation-based connotation visualization for classical poetic Japanese vocabulary of the Kokin Wakashū ca. 905 – Xudong Chen*, Hilofumi Yamamoto, Bor Hodošček

11:15 a.m. to 12:15 p.m. | Session 5

- What's that Scary Sound? – Svenja Guhr*, Mark Algee-Hewitt
- Connecting the Dots – Leonard Konle*, Merten Kröncke, Simone Winko, Fotis Jannidis
- Computational approaches to opera libretti – Luca Giovannini*, Daniil Skorinkin

12:15 p.m.: Closing

conference version

Why the Daisy sisters are different

A stylometric study on the oeuvre of Swedish author Henning Mankell and the Dutch translations of his work

Martje Wijers¹ 

1. Amsterdam Center for Language and Communication, University of Amsterdam, Amsterdam, The Netherlands.

Citation

Martje Wijers (2023). "Why the Daisy sisters are different. A stylometric study on the oeuvre of Swedish author Henning Mankell and the Dutch translations of his work". In: *Journal of Computational Literary Studies* (conference reader 2023). tbc

Date published 2023-06-09

Date accepted 2023-04-21

Date received 2023-02-17

Keywords

stylometry, Cluster analysis, PCA, delta, zeta, Mankell, translation

License

CC BY 4.0 

Reviewers

tbc

Note

This paper has been submitted to the conference track of JCLS. It has been peer reviewed and accepted for presentation and discussion at the 2nd Annual Conference of Computational Literary Studies at Würzburg University in June 2023.

Abstract. In this paper, 32 books by the Swedish writer Henning Mankell were investigated using stylometric methods, to find out whether his style varies in different genres, if his style changed measurably over time, or if his books differ from each other stylistically for other reasons. The results show that the time of publication can play a role, but that other factors, such as dominant verb tense used and narrative perspective, as well as register, are more important in determining whether and how the style of novels differs. This study also gives more insight into frequently used methods in stylometry, such as cluster analysis and PCA, that give little information about the stylistic features that differ between texts. For this purpose, the original Swedish texts were also compared to the Dutch translations of the same texts to determine how translation and language influence the results of stylometric analyses.

1. Introduction

In a conversation at the university of Tulsa in 2011, Swedish author Henning Mankell told his colleague Michael Ondaatje:

I'm like the farmer, who knows, that the land shouldn't be used for the same crops many years in a row. I try to cultivate the land in my head in the same way...[] That's why I switch between styles and between novels, essays and theater. One of the decisive things for me is, when I have an idea for a story, to decide what kind of story it is. Is it a theater play? Is it a film script? A novel? A crime novel? (Jacobsen 2012, 31) ¹

Although Henning Mankell is most known for his detective series Wallander, he indeed wrote a variety of genres during his 42 years long career as a writer. He wrote literary novels, crime novels, non-fiction, theatre plays, film scripts and children's literature.

In this paper the whole oeuvre of Mankell is scrutinized using stylometric analyses to see if his style changed measurably over time, or if some books deviate stylistically from his other works for other reasons. In this study, style is used in the definition by Herrmann et al. (2015, 44): "Style is a property of texts constituted by an ensemble of formal features which can be observed quantitatively or qualitatively." In the current study style is analysed by quantitative features, as measured by word frequency patterns.

1. My translation

The original works by Mankell are also compared to their Dutch translations. The goal of this comparison is to investigate to what extent language and translation in general can influence the results of stylometric analyses. Apart from more insight into the styles in Mankell's oeuvre, this study will yield interesting observations about the selected methods, and it can give new insights about frequently used methods in computational literary studies, such as cluster analyses and principal components analyses. These methods are generally based on the Most Frequent Words (MFW) of a text, although little is known about which type of words are decisive and what factors should be taken into account in this type of analysis.

The current paper is inspired by the computational research project 'The Riddle of Literary Quality' (2012-2019). In this project, Karina van Dalen-Oskam and her colleagues at the Huygens Institute for the History of the Netherlands in collaboration with the Fryske Akademy and the Institute for Logic, Language and Computation at the University of Amsterdam investigated readers' perceptions of what (good) literature is and to what extent these perceptions can be linked to formal patterns in novels Dalen-Oskam 2021, 15-16² Five of the many novels that van Dalen-Oskam investigated were written by Mankell and particularly one of them, the literary novel *Daisy sisters*, stood out in several ways compared to the other novels written by translated male authors in the project. She looked into the novels by Mankell in more detail and the main conclusion was as follows:

It seems, therefore, that although Mankell published books in two different genres, Suspense and Literary novel, his style as reflected in his use of words, perhaps is not very different. The fact that *De Daisy Sisters* was an outlier does not disprove this because the original is much older (1982) and it is known that an author's writing style may develop over time just like languages and the conventions that apply to different genres [...]. Further research into Mankell's complete oeuvre would be needed to confirm this. Dalen-Oskam 2021, 76

So, Mankell's style did not differ very much between genres when looking at word use compared in a corpus including books by other writers, but the *Daisy sisters* deviated clearly from the other books, possibly because it was written much earlier. By looking at the broader oeuvre of Mankell, some of the questions that remained after the Riddle of Literary Quality was finished can be answered.

The corpus compiled for this study consists of 32 Swedish books written by Mankell in four genres: crime-fiction (N=15), literary novels (N=11), children's books (N=4) and non-fiction (N=2). For comparison purposes, ten books by the following best-selling Swedish writers were added to the corpus: Johannes Anyuru, Majgull Axelsson, Marianne Fredriksson, Lars Kepler, John Ajvide Lindqvist, Camilla Läckberg, and Håkan Nesser. Six of the books by other Swedish writers are literary novels and four are crime novels. The translation corpus contains 42 translations of all the above-mentioned works by Henning Mankell and other Swedish writers into Dutch.

2. An updated English version of this book will be published in English June 2023 under the title *The Riddle of Literary Quality: A Computational Approach.*, Amsterdam University Press.

2. Multi-faceted Henning Mankell 60

Arvas and Nestingen (2011, 1) state that Mankell is the top selling Swedish crime-fiction author who, according to them “has sold 25 million copies, even outperforming Harry Potter in the German-language market.” Mankell was certainly one of Sweden’s most popular and well-read crime-fiction writers, although Berglund Berglund (2013, 10) puts these numbers in perspective. He shows that Henning Mankell was not necessarily the number one best-selling author in Sweden in the period 2004-2010, but that he indeed was among the top-sellers. He was in fact in fourth position after Camilla Läckberg, Stieg Larsson and Liza Marklund on the top 40 best-selling crime-fiction authors in Sweden (Berglund 2013, 81). Interestingly, compared to the even more popular authors Stieg Larsson, Camilla Läckberg and Liza Marklund, the books by Henning Mankell were borrowed much more frequently at libraries (Berglund 2013, 100–101).

Henning Mankell is an interesting author to investigate for multiple reasons. He modernized the already existing Swedish police novel that included criticism on modern society (started by Maj Sjöwall & Per Wahlöö) and he was the first Swedish author of crime novels to be published in many languages abroad with wide circulation. Therefore, Mankell played an important role in the rise of the Nordic noir genre (Berglund 2013, 114).

Mankell belongs to a group of authors that were already established writers of fiction before they started to write crime (in the late 90s) when there was a boom in crime-fiction in Sweden. His debut in the crime genre was in 1991 with *Mördare utan ansikte (Faceless killers)*, but his debut as a writer of fiction was much earlier: in 1973 with *Bergsprängaren (the Rock Blaster)*.

The fact that he has a broad oeuvre covering four genres over a time span of 42 years (1973-2015) also makes Henning Mankell useful for a computational study. Furthermore, his novels are widely translated into other languages. Almost his entire oeuvre is translated into Dutch. There is a limited number of Dutch translators from Swedish which makes it possible to compare translations by different translators. As mentioned earlier, these translations were sometimes published much later in Dutch than the original. It is important to bear in mind that a writer’s style can change over time, and so do ideas about translation (Can and Patton 2004; David L. Hoover 2020; Ríos-Toledo et al. 2022).

2.1 Mankell and the Riddle of Literary Quality 92

In the Riddle of Literary Quality, van Dalen-Oskam and her colleagues investigated if literary quality is measurable using stylometric methods. They selected 401 contemporary novels in Dutch published between 2007-2012 based on sales numbers and library borrowings in the three years prior to the survey (2009-2012) (Dalen-Oskam 2021, 44). The works included both novels originally written in Dutch as well as translated novels. These novels were rated for their literary quality on a scale from one (not literary at all) to seven (very literary) by almost 14000 readers in ‘Het Nationale Lezersonderzoek’ (the National Reader Survey) in 2013 (Dalen-Oskam 2021, 40–43). The ratings were then linked to the formal aspects of the books, such as vocabulary and sentence

length or contextual information, such as whether the author is male or female (van Dalen-Oskam).

One of the findings in *The Riddle of Literary Quality* was that many readers seemed to be somewhat more critical towards translated literary fiction compared to literary novels originally written in Dutch. In other genres, such as crime novels, the bias was just the opposite: on average, translated works received higher scores on literary quality than original Dutch crime novels (Dalen-Oskam 2021, 104–105).

However, there was a clear difference between books translated from English and books translated from other languages. Translated books from other languages than English scored higher on literary quality than works translated from English and books originally written in Dutch in the category literary novels as well as the category crime novels (Dalen-Oskam 2021, 105). Dalen-Oskam (2021, 112) suggests that readers are more critical toward translations from languages they know than from languages they are unfamiliar or much less familiar with.

In total there were 249 translated books in the survey. English was by far the language in which most of these books were written, namely 180. After English, the second most recurring original language was, somewhat surprisingly, Swedish (Dalen-Oskam 2021, 102). One Swedish author is represented with five books in the corpus: Henning Mankell. Three of these books are in the category crime. Remarkably, these three books end up relatively high in the ranking of literary quality among literary novels (Dalen-Oskam 2021, 178). The two literary novels, on the other hand, ended up among the lowest scoring literary novels, although the scores were still somewhat higher than his crime novels (Dalen-Oskam 2021, 193). The literary novel *Daisy sisters* turned out to have different frequency patterns of MFWs compared to other translated novels written by male authors. However, the frequency patterns of this book were remarkably close to the frequency patterns of one of the highest scoring translations: *Norwegian wood* by Haruki Murakami (Dalen-Oskam 2021, 190). Dalen-Oskam (2021, 189) wonders whether this could have something to do with the fact that both works were translated into Dutch much later than they were published in the original languages Swedish and Japanese (both in the eighties).

However, she did not have enough data in her corpus to investigate this assumption further. The corpus in the study I report on in this contribution, consisting of 32 books written by Mankell during his entire career, can confirm or reject this hypothesis. The following section reports on the results of the studies, and looks at genre differences, possible change over time and other factors that influence the clustering of texts.

3. Genre and style differences

When a book gets translated the genre classification chosen by the publisher could, at least theoretically be different in the source language. However, in the translations of the books by Henning Mankell into Dutch this is not the case. Squires (2007, 71–72) states that genre is a necessary part of book publishing. It is implemented in the whole publishing process, from cover design to advertising and what literary prizes the book qualifies for. The genre also determines on what shelf the book ends up in the bookstore

or library. Because of this, Squires (2007, 71–72) concludes that genre classification is not so much a literary boundary, but rather a marketing tool. Although this might be true to some extent, multiple studies in computational literary studies have shown that it is possible to distinguish genres based on style, measured by high frequency words (e.g. Dalen-Oskam 2021; Jautze 2014; Jautze et al. 2013; Jockers 2013).

Jockers (2013, 68–70) showed that genre and style are closely linked. Jockers and his colleagues looked at various subgenres in nineteenth century English novels. They divided the text into samples of 1000 words and performed an unsupervised clustering using the most frequent words (MFW). The high-frequency words turned out to not only be highly successful in distinguishing samples from the same author and novel, but also placed text samples that belonged to the same genre closely together. Jockers concluded that (sub)genres have a stylistic fingerprint that can be detected by looking at high-frequency words.

Jautze (2014) investigated whether the MFWs can distinguish chick lit from literary novels. She performed a stylometric analysis using the R package *stylo* (Eder et al. 2016) and found that chick lit was stylistically different from high literature. High literature turned out to have a more descriptive style, whereas chick lit seemed to be more informal.

In an earlier study, Jautze et al. (2013) compared high literature and chick lit syntactically and found that novels that are classified as high literature contain more complex sentences than chick lit. High literature was also found to be richer in prepositional phrases than chick lit.

To my knowledge, there are no studies that compare the style of high literature and crime novels, the genre that Henning Mankell is most known for. The genre is sometimes even referred to as literary crime novel, indicating that it has a higher literary quality than regular crime novels or thrillers. One might expect that it is harder to distinguish between high literature and ‘literary’ crime novels, especially if they are written by the same author.

To find out if an analysis of the MFWs can make this distinction, I performed a stylometric analysis on the Mankell corpus using the R package *Stylo* (Eder et al. 2016). The *Stylo* package automatically compiles a list of MFWs in the entire corpus and can check which words occur relatively frequently in the various texts, based on the Delta procedure for authorship attribution (Burrows 2002). Burrow’s delta looks at texts as a collection of data or ‘bag of words’ and disregards the context of sentences. The frequency of each word in the corpus is counted and the separate texts are compared to each other based on their frequency lists (MFWs). For this comparison, the relative, normalized z-scores are used, so differences in text length or the high impact of a small number of high-frequency words on the total outcome are ruled out (Eder et al. 2016). The distances between texts can, for instance, be visualized in a dendrogram representing the results of a cluster analysis, grouping texts that are similar to each other.

A cluster analysis was first performed on the Swedish corpus, to see if there are clear stylistic differences between various genres. The analysis is based on the 1000 most frequent words in the books. The results are visualized in Figure 1. As illustrated in figure 1, most books are neatly clustered by genre, where L stands for literary novel; C

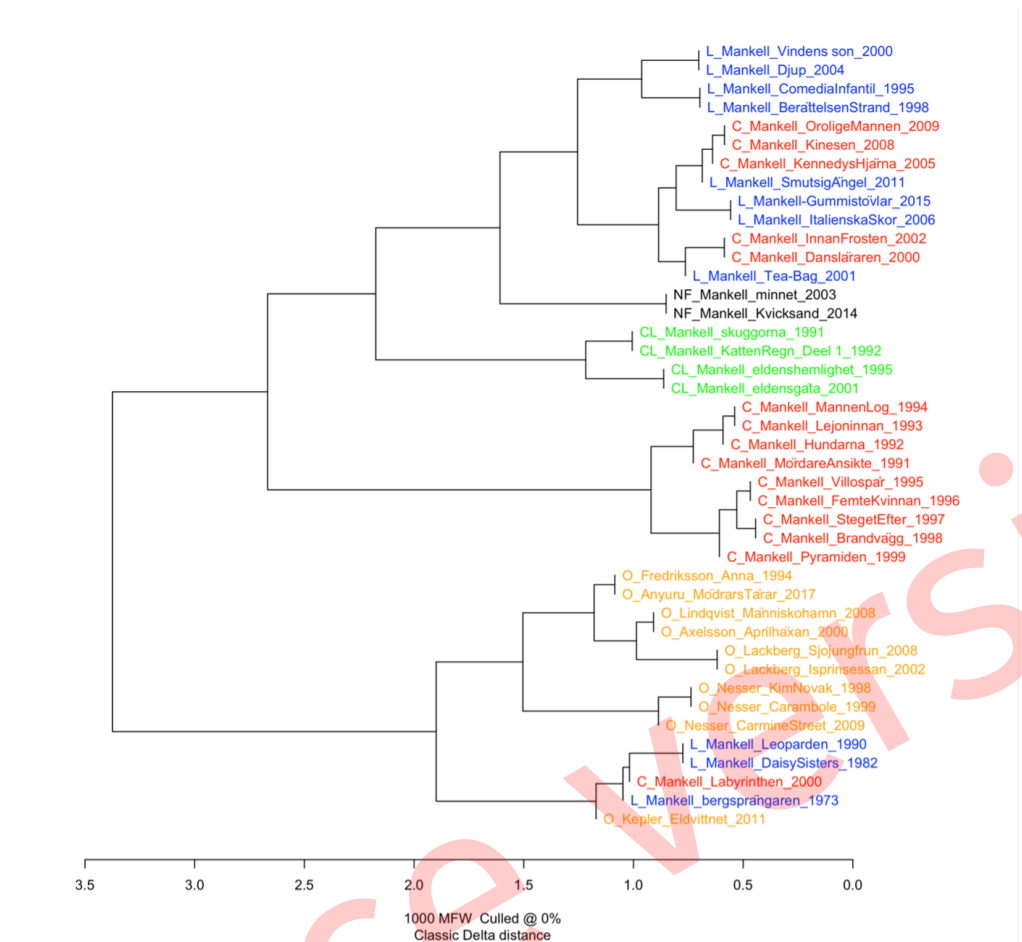


Figure 1: Cluster analysis of the Swedish books in the corpus based on the 1000 most frequent words (culling o, classic delta)

for Crime novel; NF for Non-Fiction and CL for Children's literature, although some 188
crime novels appear among literary novels or vice versa. This seems to be the case for 189
crime novels written from the year 2000 onwards. 190

The earlier crime novels: from the 1991 *Mördare utan ansikte* (*Faceless Killers*) until 191
Pyramiden (*The Pyramid*) from 1999, all belonging to the Wallander series, are in a 192
separate cluster. This cluster has two subclusters: one for the books written in the first 193
half of the 1990s (1991-1994), and one for the Wallander books published in the second 194
half of the 1990s (1995-1999). Remarkably, Mankell's last Wallander book *Den oroliga* 195
mannen (*The troubled man*), that was published in 2009, falls outside of this cluster. This 196
could again be explained by the fact that this last book of the series was written ten 197
years after the previous Wallander book, and that Mankell's style changed over time. 198
The crime novel *Innan frosten* (*Before the frost*) from 2002, which is written from the 199
perspective of detective Wallander's daughter, does not belong to the Wallander series 200
either. However, this book is closer in time to the other Wallander books, indicating that 201
there are other factors that weigh in. 202

The books by other authors than Mankell are also clearly different from the books by 203
Mankell. The crime novel *Carambole*, from the Van Veeteren series by Håkan Nesser, for 204
instance, is very comparable genre-wise to Mankell's Wallander series. However, author 205
seems to be a stronger factor in the clustering of the text than genre, because *Carambole* 206

ends up in a separate cluster and clusters with other literary novels by Nesser. Genre, 207
 in its turn, plays a more important role than time overall. If we for instance look at the 208
 two non-fiction books by Mankell, they clearly form a cluster, even though they were 209
 published eleven years apart. 210

However, something remarkable is going on with the three oldest books by Mankell 211
 in the corpus. These three literary novels: Mankell's debut, *Bergsprängaren* (*The Rock* 212
Blaster) from 1973, *Daisy Sisters* from 1982 and *Leopardens öga* (*The Eye of the Leopard*) from 213
 1990 appear closer to other authors and even cluster with the 2011 crime novel *Eldvittnet* 214
 (*The Fire Witness*) by Lars Kepler. The same is true for the crime novel *Labyrinthen* (*The* 215
Labyrinth) from a much later period (2000), that clusters with Mankell's early literary 216
 novels. *Labyrinthen* is different from other works, because it was originally written as a 217
 film script and later turned into a novel. This might have influenced the style of this 218
 particular novel. 219

The fact that the three early Mankell novels are stylistically different from his later works 220
 seems to indicate that Mankell's writing style and word choice indeed has changed over 221
 time and confirms the findings by Van Dalen-Oskam that *Daisy Sisters* is different from 222
 other novels by Mankell. However, the texts and their MFWs have to be investigated in 223
 more detail to see how his style has changed and to ensure there are no other factors at 224
 play. 225

The Dutch translation corpus was analyzed using the same procedure as shown for the 226
 Swedish corpus to see if the texts appear in different clusters when they are translated. 227
 The Dutch corpus consists of the same books by Mankell and by the ten books by 228
 the aforementioned Swedish authors (Johannes Anyuru, Majgull Axelsson, Marianne 229
 Fredriksson, Lars Kepler, John Ajvide Lindqvist, Camilla Läckberg, and Håkan Nesser). 230
 This comparison can give important information about what type of MFWs influence 231
 the clustering of texts in stylometric analyses. The results are shown in Figure 2. The 232
 different titles were all labeled by genre first (L for literary novel; C for Crime novel; NF 233
 for Non-Fiction, CL for Children's literature and O for different author than Mankell). 234
 The second tag is the translator's initials, followed by the author's last name and two 235
 years, the first one indicates the year the original novel was published, the second one 236
 stands for the year the translation was published. 237

Overall, the results are similar to the results of the cluster analysis of the Swedish corpus. 238
 However, the genre differences seem to be slightly bigger in the translated works. Unlike 239
 the results in the Swedish corpus, all the non-Wallander crime novels end up in one 240
 cluster together. Two novels stand out in particular: the literary novel *Tea bag* from 2001, 241
 that appears close to Mankell's later crime novels and *Labyrinthen*, which just like in the 242
 Swedish novels clusters with the three early literary novels by Mankell. 243

Another noticeable difference between the Swedish and the Dutch cluster analysis, is 244
 that unlike in the Swedish originals, the early literary novels in the translations are more 245
 similar to other works by Mankell than to the other Swedish authors, although Lars 246
 Kepler's *Eldvittnet* shows up in this cluster again. 247

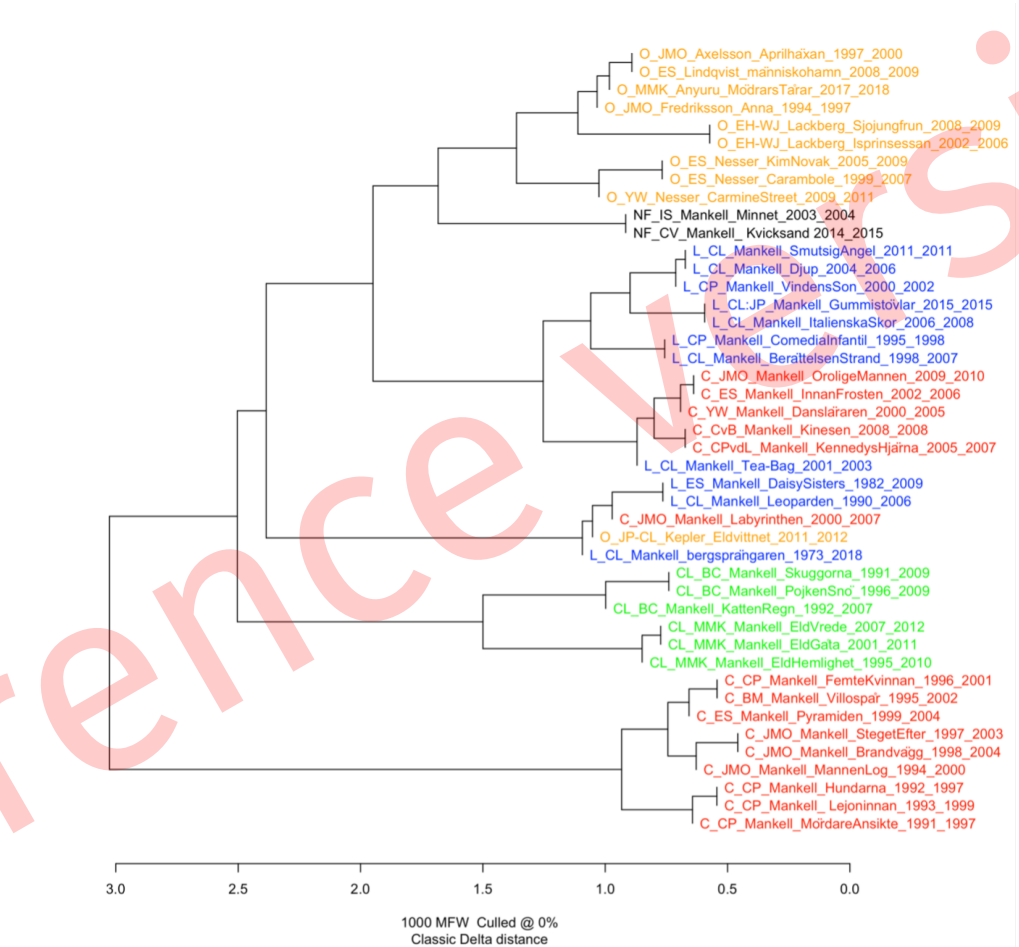


Figure 2: Cluster analysis of the Dutch translation corpus based on the 1000 most frequent words (culling o, classic delta)

4. Network analysis of Mankell's oeuvre 248

As pointed out by Eder 2015, there are a couple of problems with cluster analysis pertaining to the distance, linkage and number of features (MFWs) used for analysis. Outcomes can differ depending on the MFWs used and there is no real consensus about what the optimal number of MFWs is. These problems can partially be overcome by using a bootstrap consensus tree, because it repeats measurements for multiple numbers of MFWs, and looks for the most robust groupings across different measurements.

However, Eder 2015, 55–56 notes there is still some arbitrariness involved in the production of a consensus tree, such as how many times the analysis should be repeated, for how many words in total are considered and the underlying algorithm used for linkage. A bigger caveat for the current study, however, is that a consensus tree only looks for the closest ranking text, which means it mainly looks at the strongest similarities. In most cases, this is the authorial fingerprint.

In this paper, the central question is rather why some works within one oeuvre deviate from the majority of works and what other factor beside the author are decisive for clustering of texts. These weaker patterns might better be detected by producing a network analysis as proposed by Eder 2015. In a network analysis, not only the closest text in rank is taken into account, but also the second and third closest neighbours. These links are visualized in a network in which close similarities are shown with thicker lines and weaker links with thinner lines.

I performed a bootstrap consensus tree in Stylo and used the CSV output to create a network analysis in the open-source tool GEPHI Gephi using the ForceAtlas2 algorithm. I ran a Modularity Analysis (resolution 0.6) in GEPHI which detects communities in the network, helping to distinguish closely related topological subgroups of nodes from each other and to make clusters more visible in the network. Finally, I applied eigenvector centrality, to measure the influence of nodes in the network. Ranking the function size of nodes indicates the centrality of a work for the cluster it is in.

The results of the network visualizations are shown in Figures 3 and 4. A short description of the clusters is given in the titles in red for ease of interpretation.

In general, the results shown earlier in the cluster analyses are confirmed by the consensus networks. Works cluster mainly by author and genre, although there is some overlap between crime novels and literary novels. There is a separate cluster for Mankell's Wallander series and the older literary novels are in a separate cluster.

However, there are some remarkable differences between the Swedish consensus network and the translated Dutch one. In the Swedish network the older Mankell novels cluster with two literary novels and a crime novel by Nesser as well as a crime novel by Kepler. In the Dutch network, they are only grouped together with the crime novel by Kepler only. The consensus network of the original Swedish corpus distinguishes six clusters, whereas the consensus network of the translated corpus contains eight. Unlike the cluster analysis, where the texts were clustered more clearly by genre in the translated corpus, the network looks somewhat more messy in the translated corpus with smaller and less clearly defined clusters.

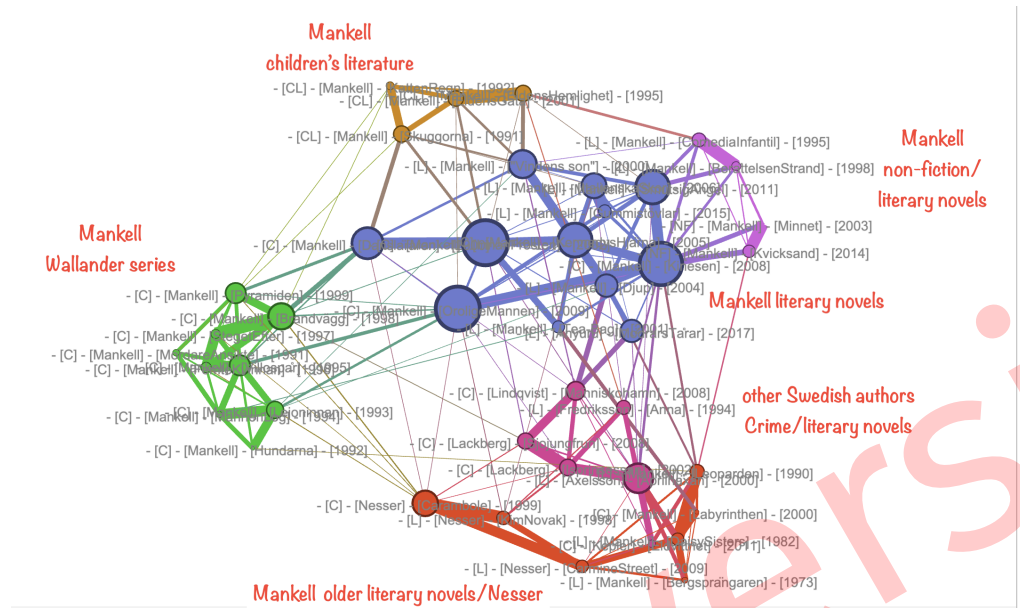


Figure 3: Consensus network of the Swedish corpus: classic Delta distance, 100–1000 MFWs, modularity 0.6

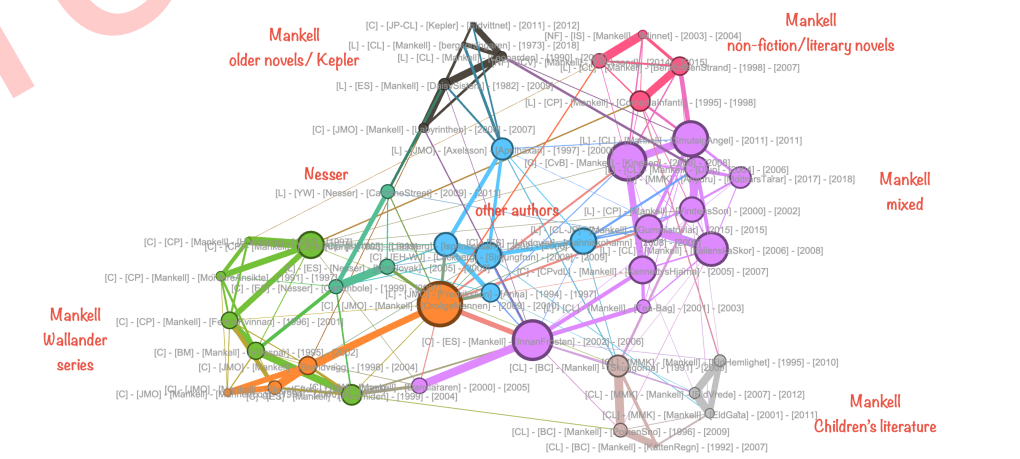


Figure 4: Consensus network of the translated Dutch corpus: classic Delta distance, 100–1000 MFWs, modularity 0.6

Most importantly for the current study, the same four novels that showed up as outliers in the cluster analyses (*Bergsprängaren* (*The Rock Blaster*) from 1973, *Daisy Sisters* from 1982, *Leopardens öga* (*The Eye of the Leopard*) from 1990 *Labyrinthen* (*The Labyrinth*)), again form a separate cluster. In the following section, the MFWs associated with these works are analysed to see why these particular novels stand out from the rest of Mankell's novels.

5. A closer look into the MFWs

To look at more dimensions in the data, a Principal Components Analysis (PCA) was performed on the Swedish corpus. Like a cluster analysis, a PCA also analyzes the MFWs in the dataset, but they are visualized in a scatterplot instead of a dendrogram. In a PCA multiple features are combined in an artificial variable, the so-called principal component that explains the largest proportion of the variance in the data (Jockers 2013, 65–67). On the x-axis, the first principal component is shown. The first principal component is often related to the author (David L. Hoover 2020) The y-axis shows the second principal component. The second principal component is less obvious to interpret, it could be explained by variables like chronology or genre (David L. Hoover 2020). These two principal components are unrelated.

I performed a classic PCA on the Swedish data in Stylo, with the Classic Delta and the correlation option, analyzing the 1000 MFWs. The results of the PCA of the Swedish corpus are presented in Figure 5 below. The x-axis, showing the first principal component, explaining 12.2% of the variance in the data, can clearly be linked to author and is in line with the findings in the cluster analysis. The same four books that were mentioned earlier are deviant from Mankell's other works and more similar to the other Swedish writers in the corpus. The books in question are the crime novel *Labyrinthen* from 2000, the two oldest books in the corpus, namely *Bergsprängaren* (*The Rock Blaster*) from 1973 and *Daisy Sisters* from 1982 and *Leopardens öga* (*The Eye of the Leopard*) from 1990. Unlike in the cluster analysis, we can now see that Lars Kepler's *Eldvittnet* is further away on the x-axis and probably clustered with these books because of the variance in the data that is represented on the y-axis.

Figure 5 shows that author and genre are still the most important factors in distinguishing between texts. However, there are a few books by Mankell that clearly behave differently and that end up closer to books by other authors. What makes these four stand out from the rest of Mankell's works?

If we perform the same PCA again, but with the option 'loadings' in Stylo, showing which words occur significantly more frequently in the texts they are close to in the graph, we might get a first impression about an important difference between the four atypical books and the rest of Mankell's work. In Figure 6, the results of the PCA with the option loadings are shown. This analysis was performed on the 100 MFW, because a figure with 1000 words would become illegible. This also means that the distribution of the novels on the graph is somewhat different. For instance, Kepler's novel *Eldvittnet* is now closer to Mankell's *Labyrinthen*, whereas Nesser's *Maskarna på Carmine street* (2009) appears close to Mankell's older novels. Importantly, Mankell's four diverging novels still stand apart from his other novels. In Figure 6, they are shown a bit below the upper

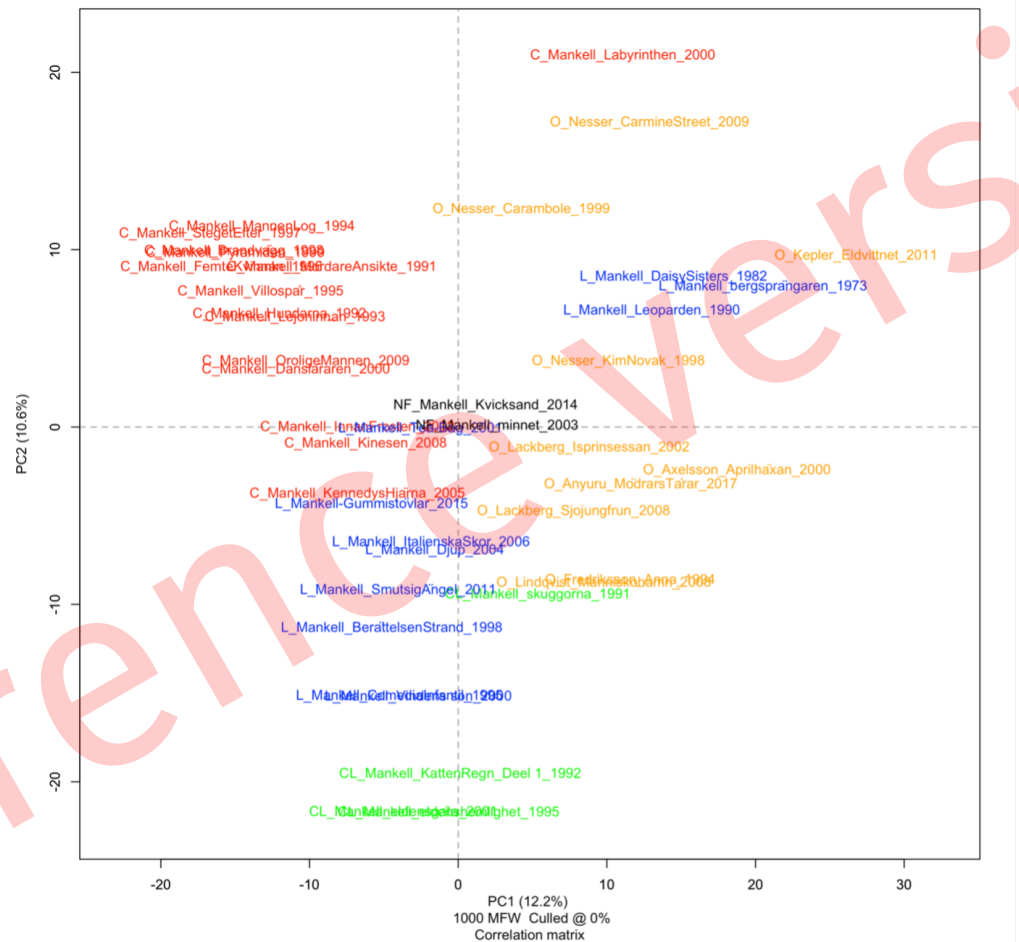


Figure 5: Principal component analysis of the Swedish corpus (1000 MFW, Classic Delta correlation, culling o)

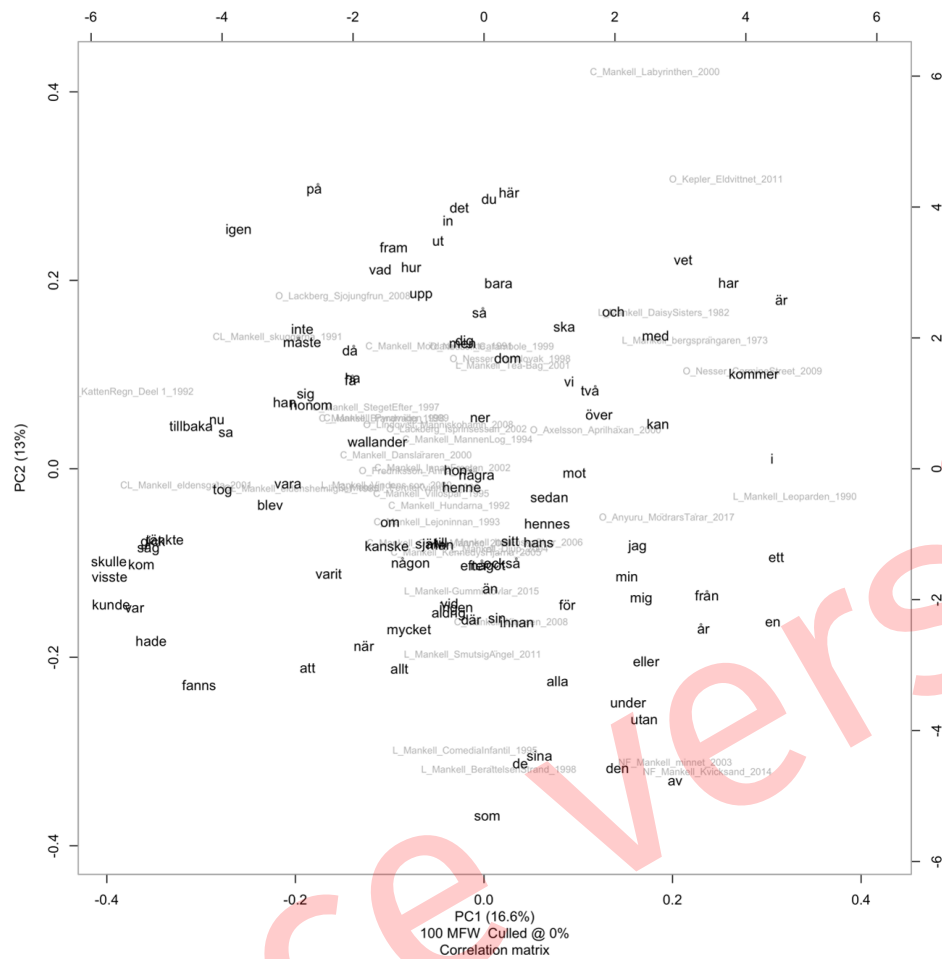


Figure 6: Principal Components Analysis showing the 100 MFW in the Swedish corpus

right corner. The words that are associated with these novels, and less with the other 333
 books, are: *kommer* ‘come’, *vet* ‘know’, *har* ‘have’, *är* ‘is/are’, *och* ‘and’ and *med* ‘with’. 334
 The first four are verbs in the present tense, whereas the verbs associated with other works 335
 are all in past tense or past participles. 336

This indicates that rather than the chronology, the tense primarily used in the narrative, 337
 established by verb tense, might be a decisive factor in why the four mentioned books 338
 are different from other Mankell books. On closer inspection, these books as well as 339
Eldvittnet by Lars Kepler are primarily written in the present tense, whereas the other 340
 works by Mankell are primarily written in the past tense. Of course, this may be related 341
 to a chronological development: over time a writer can also change their preference for 342
 which tense to narrate a story in. 343

The same procedure was followed for the translated Dutch corpus. The results are 344
 shown in Figure 7 and 8. In Figure 7 the PCA for the translated Dutch corpus is shown, 345
 which is in many ways comparable to the results of the Swedish PCA. One remarkable 346
 outcome is that Mankell’s Children’s books are quite different on the x-axis, where 347
 this was not the case at all in the Swedish results. Another remarkable finding is that 348
 some novels by other writers in the corpus, namely Håkan Nesser, Camilla Läckberg 349
 and Marianne Fredriksson, end up very close to the literary novels by Mankell and in 350
 between books by Mankell in different genres. 351

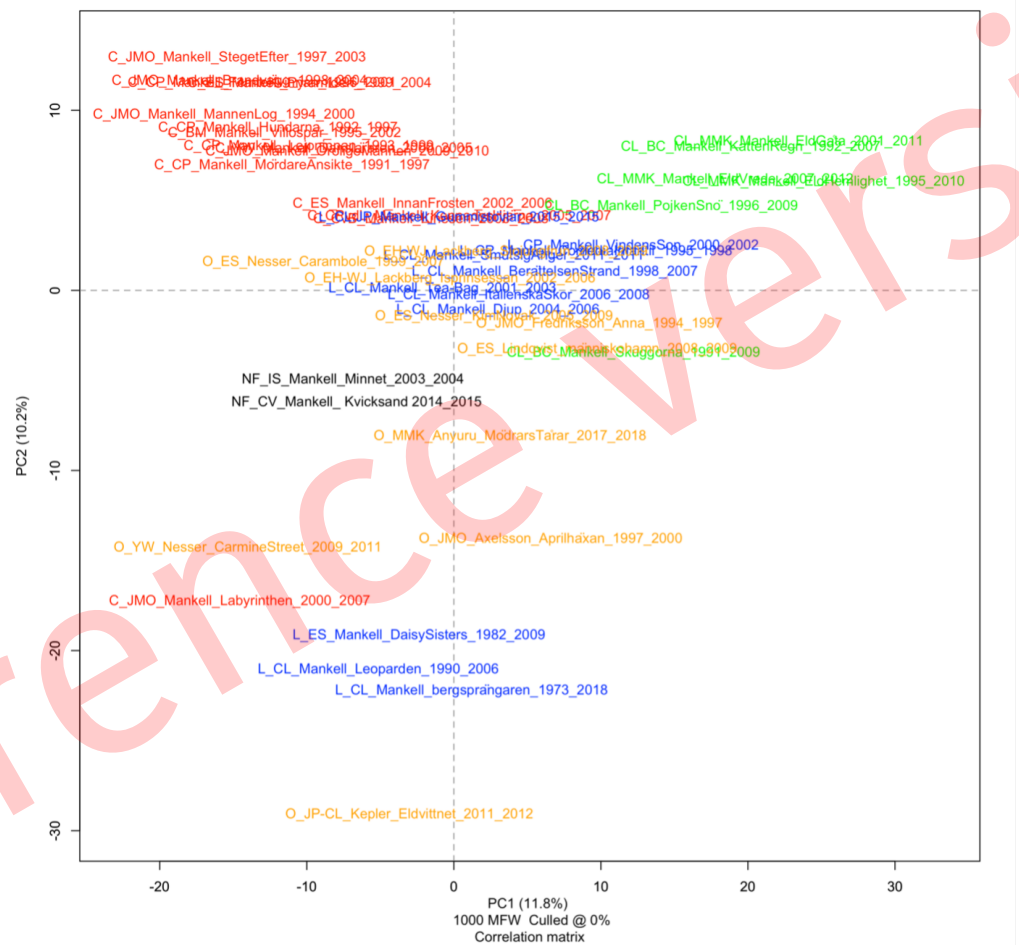


Figure 7: PCA of the translated Dutch corpus based on the 1000 MFWs

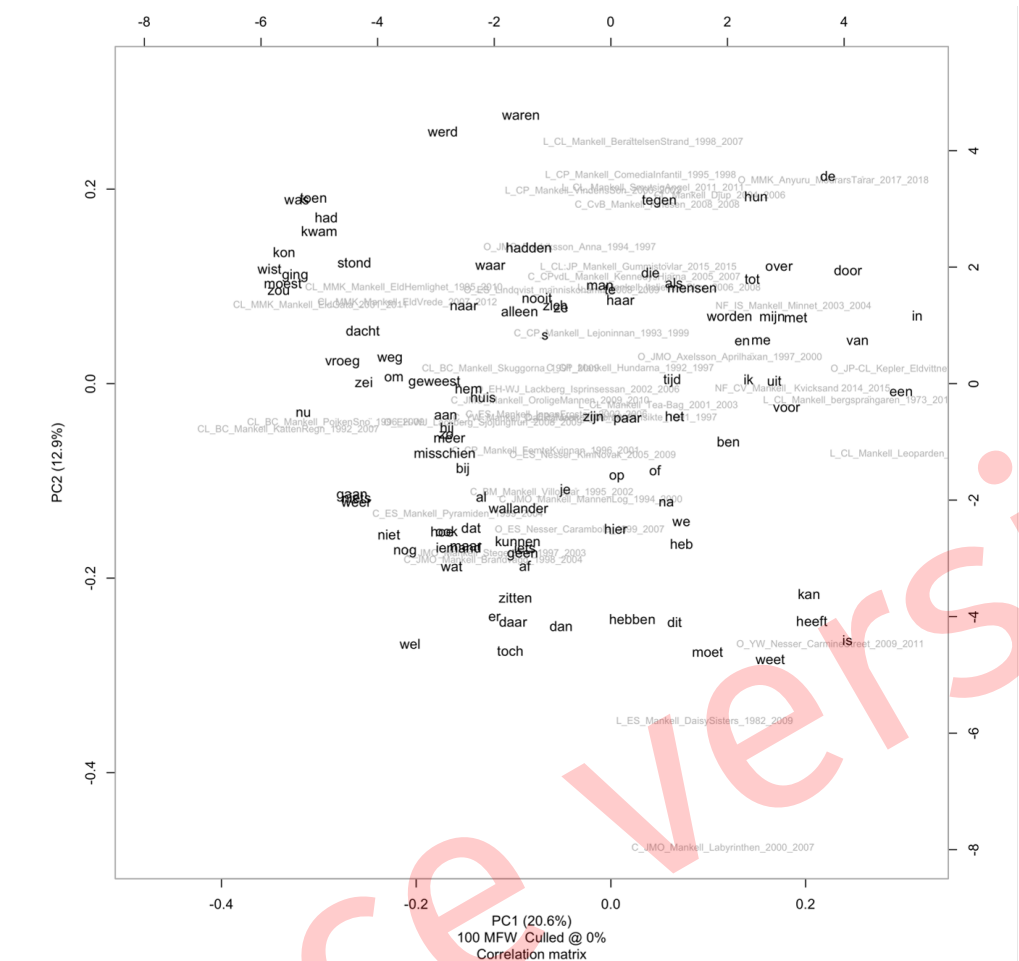


Figure 8: PCA (loadings) of the translated Dutch corpus based on the 100 MFWs

Otherwise, the same four books (*Leopardens öga*, *Bergsprängaren*, *Labyrinthen*, and *Daisy Sisters*) diverge in the translation corpus. The PCA with the loadings function (Figure 8) clearly shows that this is likely caused by the narrative tense again. Words that occur more frequently in these books are: *moet* ‘has to’, *kan* ‘can’, *is* ‘is’, *heeft* ‘have’ and *weet* ‘know’ whereas past tense verbs occur more frequently in other works. An important difference between Swedish and Dutch is that Swedish only has tense marking on verbs whereas Dutch has tense and person marking. This also means that Swedish verbs probably tend to end up higher in the list of MFWs because there are fewer possible forms compared to Dutch where the same verb is spread out over more possible forms.

5.1 Verb tense and perspective

In the remaining part of this section, I will elaborate on the results of the study so far, to get more insight into the stylometric methods used and the studied texts. It is important to look beyond the analyses of Delta distances to see what is behind the measurements and which words are decisive in the clustering of texts.

The results so far, show that verb tense is an important factor for the outcome in stylometric methods based on the MFWs, because there are many verbs (with tense marking) among the MFWs. In similar lines, narrative perspective might also play an important role, because pronouns are very frequent words. In order to get a good indication of

the predominant narrative perspective in the books in the current corpus, I applied Van Rossum's I-index to the data (Van Rossum et al. 2020). Van Rossum et al. (2020) applied both a machine learning and a narratology-based approach in which they computed the ratio of pronouns. Both methods turned out successful in determining the narrative perspective of texts, although the second approach was slightly more robust and yielded a perfect 1.00 score. This perfect score was possible, because Van Rossum et al. (2020) cleaned the data from dialogue. The narrative perspective was already known, so the predictions could be tested for their accuracy. For now, this is not possible in the Mankell corpus, but the ratio of pronouns can still give a good indication of a book's narrative perspective.

Van Rossum's I-index (Van Rossum et al. 2020) is focused on the first person narrative perspective but can be applied to other perspectives as well. I computed the I-index and the he-index, she-index and (singular) you-index (du-index) for both the Swedish originals and the Dutch translations. For the he-index (han-index), for instance, I did this by adding the relative frequency scores of *han* 'he', *honom* 'him' and *hans* 'his' as calculated in Stylo and divided this number by 1 + the relative frequency scores for all the pronouns in the text. The reflexive possessive pronouns *sin*, *sitt* and *sina* were left out of the equation, because they are used to refer to both male and female antecedents. For the she-index (hon-index) I did the same but with *hon* 'she', *henne* 'her' (object form) and *hennes* 'her' (possessive). Finally for the singular you-index (du-index) I divided the sum of the relative frequencies of *du* 'you' (singular), *dig* and *dej* 'you' (object form in two spelling variants), *din/ditt/dina* 'your' (singular in three inflection forms) by the relative frequencies of all pronouns combined.

Figure 9 shows the results of the indexes in a graph. The ratio of pronouns gives a good indication of the narrative perspective(s) in the texts. Only the results of the Swedish corpus are shown here, because I observed no big differences between the Swedish and the Dutch ratios. Mankell's texts are ordered chronologically from oldest to most recent. The other authors are in random order. The first part of the bars on the bottom left side shows the I-index. In most texts, this index is between 0,10 and 0,30. Clear peaks in the I-index can be detected for the two non-fiction books *Jag dörr, men minnet lever* (2003) and *Kvicksand* (2014), which indeed are mainly written from first person perspective.

Peaks in the I-index can also be observed for *Italienska skor* (2006) and *Svenska gummistövlar* (2015) which are both literary novels with the same main character Fredrik Welin written from an I perspective. Two books by Håkan Nesser: *Maskarna på Carmine street* (2009) and *Kim Novak badade aldrig i Genesarets sjö* (1998) also score high on the I-index. Recall that Nesser's books also appeared close to Mankell's outliers in the PCA.

The second part of the bars in Figure 9 show the han-index (he-index). Most books by Mankell score high on the han-index (he-index) and have indeed a male main character. The third part of the bar show the she-index (hon-index) and it is most interesting to compare these two indexes directly. *Daisy sisters* (1982) is one of the books that score very high on the she-index (hon-index), which makes sense, because it is a novel about three generations of women. The combination of a deviant verb tense (present tense) and a female perspective could very well explain why this particular book appears to be an outlier in the cluster analyses and the PCAs. Two of the children's books (*Eldens gåta* and *Eldens hemlighet*) also score relatively high on the she-index. Some books have a

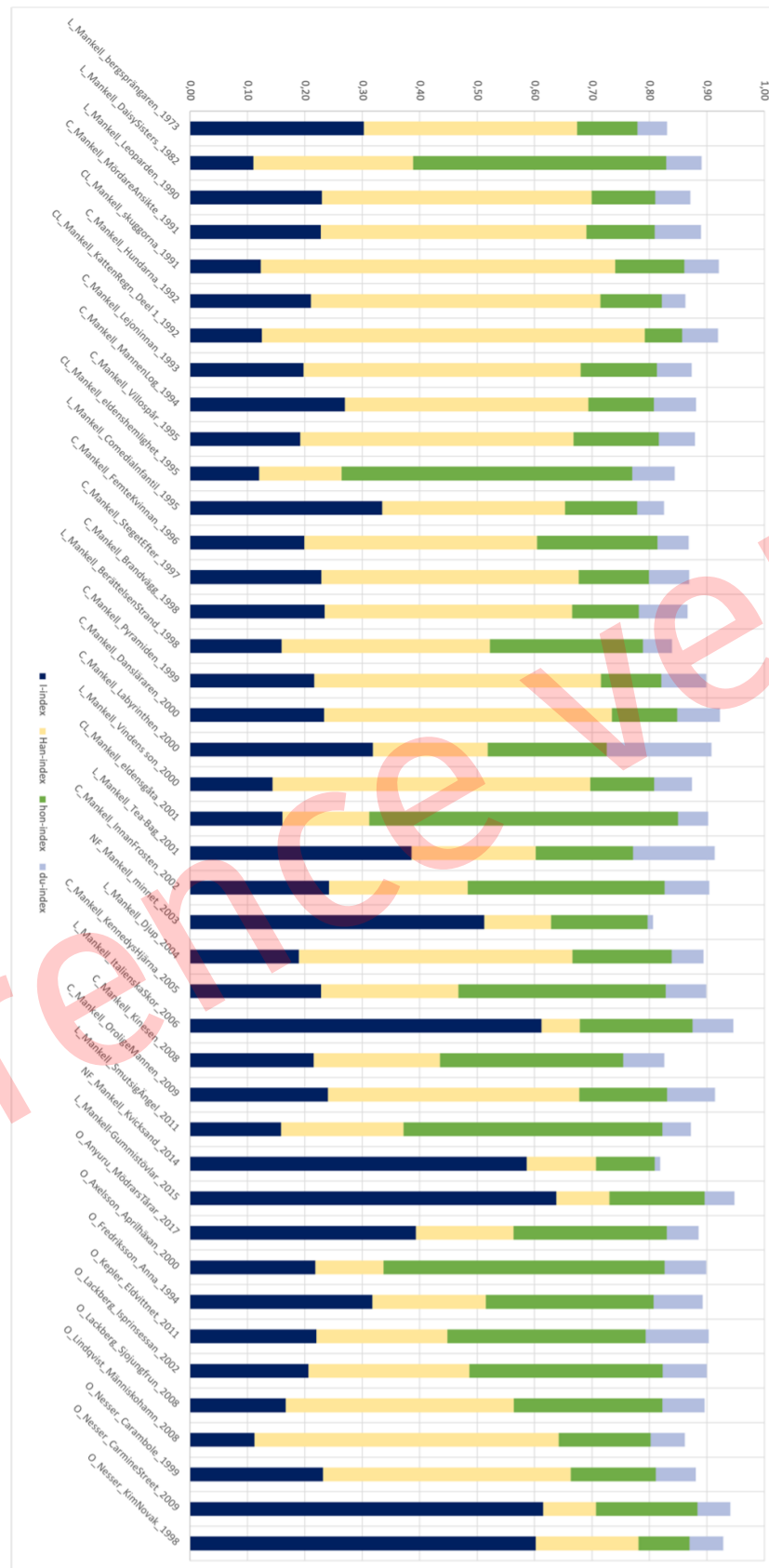


Figure 9: Indication of narrative perspective in the books in the Swedish corpus, measured by I-index (dark blue), Han-index (yellow), Hon-index (green) and du-index (light blue)

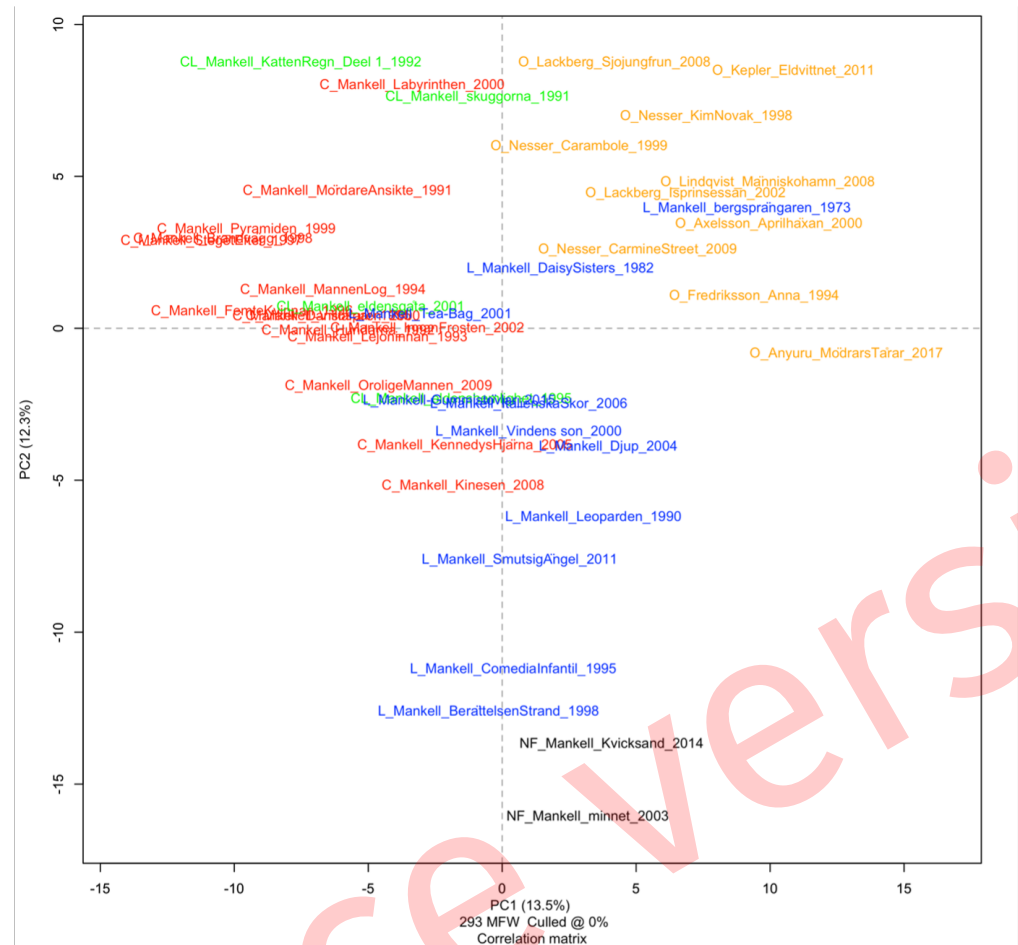


Figure 10: PCA (classic) of the Swedish corpus excluding tense-marked verbs and personal pronouns based on the 1000 MFWS

more evenly divided ratio between pronouns. This is especially the case in *Labyrinthen* (2000) which also was one of the clear outliers in the PCA and cluster analysis together with the earliest Mankell novels. Narrative perspective thus seems to be an important explanatory factor. The analysis of narrative perspective and narration tense leads to useful new observations about what can influence MFW scores for Swedish and Dutch and shows how a novel like *Daisy sisters* differs from other novels by Henning Mankell.

To get more insight into how much of the outcome was influenced by narration tense and narrative perspective, we should only look at the words that are not clearly linked to verb tense and narrative perspective. Stylo has the option to analyze the corpus using an 'existing word list' which enables the researcher to look at specific sets of words. I excluded all verbs marked for tense and all personal pronouns the influence of verb tense was left out of the analysis to better determine how big their influence is on the analyses. I then ran another PCA with the 'loadings' function. The resulting PCA without personal pronouns and verbs indicating tense is shown in Figure 10. This figure clearly shows that now three of the four deviant books are much closer to Mankell's other works, at least on the x-axis, and they no longer form a separate cluster. *Leopardens öga* is also closer to other books in the same genre, but especially *Bergsprängaren* and *Labyrinthen*, and to a lesser extent also *Daisy Sisters*, are still more distant from other works by Mankell.

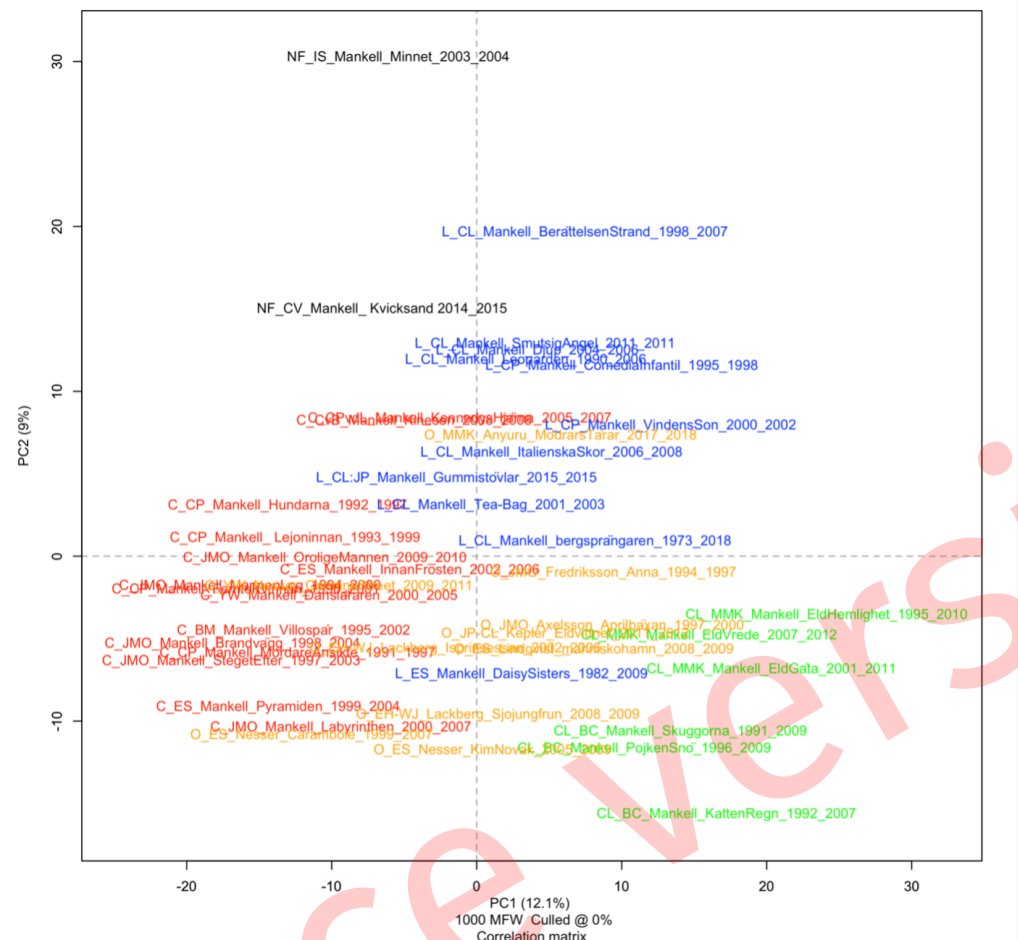


Figure 11: PCA (classic) of the translated Dutch corpus excluding tense-marked verbs and personal pronouns based on the 1000 MFWs

Figure 11 shows the Dutch PCA excluding tense-marked verbs and personal pronouns. In this graph it becomes clear that Daisy sisters (together with Labyrinth) diverges more from other Mankell novels on the Y-axis than Bergsprängaren. So, the translations and the original Swedish texts are different in this perspective. The word frequency patterns of the translations of the other Swedish authors are also much harder to distinguish from the frequency patterns in Mankell's books compared to the results of the analysis of the Swedish texts. This implies that some aspects of style get lost in translation.

5.2 The influence of register

Due to space limitations, I will exclusively focus on one of the earliest Mankell novels *Daisy sisters* in this section. We have seen that this novel is deviant in style, partly because of the use of the present tense and because it is one of the relatively few books by Mankell written from a female third person perspective. A third reason for why *Daisy sisters* has deviant word frequency patterns compared to Mankell novels that were published later, can be detected if we look at Zeta scores.

Zeta was initially introduced by Burrows (2007), and later on improved by Hugh Craig (Craig and Kinney 2009). Burrows's Delta, the method used in this study so far and

the method most often used in stylometry, relies on high frequency words (MFWs). Burrows's Zeta and Craig's Zeta, on the other hand, analyze the middle frequency words and measure distinctiveness or keyness of keywords in a corpus relative to a reference corpus (David L Hoover 2010). Middle frequency words are usually more meaningful than high frequency words, because high frequency words generally are function words (Rybicki 2016, 751). In a Zeta analysis, the texts to be analyzed are first divided into equal segments, then the dispersion of each word in the two separate corpora is registered, by counting how many segments it occurs in at least once (Craig and Kinney 2009).

Stylo can generate wordlists containing the most distinctive keywords in two opposing texts or corpora (Eder et al. 2016). I compiled a primary corpus, consisting of *Daisy Sisters* and a secondary, reference corpus containing all other books in the initial corpus. I did this for both the Swedish corpus and the Dutch translation corpus. I then performed the command `oppose()` in Stylo to analyze *Daisy Sisters* and the reference corpus using Craig's Zeta. Zeta is the sum of the proportions of sections from *Daisy sisters* in which each word occurs and the sections of other works in the corpus in which it does not (David L Hoover 2010). This can point out stylistically interesting characteristics of a text or a corpus. I used samples of 3000 words.

This generates two word frequency lists: one list with words that are preferred (or relatively more frequent compared to the other books in the initial corpus) in *Daisy Sisters* and a list with words that are avoided (or relatively less frequent compared to the other books in the initial corpus). From these lists I selected the twenty most distinctive words, excluding names and verbs, because I was interested in whether there were other stylistic differences besides tense and narrative perspective. The results for the Swedish data are shown in Table 1.

preferred		avoided	
mej	'me'	mig	'me'
dej	'you' (object)	dig	'you' (object)
jo	'yes' (after negation)	genast	'immediately'
fan	'damn'	polis	'police'
herregud	'lord'	mina	'my' (plural)
sej	(reflexive pronoun)	oss	'us'
vadå	'what'	min	'my' (singular)
såjer	'say'	skäl	'reason'
lust	'desire'	polishuset	'police station'
visst	'certainly'	våra	'our' (plural)
jävla	'fucking'	samtalet	'the conversation'
ju	'of course'	telefonen	'the phone'
världen	'the world'	papper	'paper'
omedelbart	'immediately'	nånting	'something'
värre	'worse'	mannen	'the man'
lov	'holidays'/'permission'	sedan	'then'
tåget	'the train'	frågor	'questions'
knappt	'hardly'	din	'your'
morsan	'mom'	rummet	'the room'
full	'drunk'	bland	'among'

Table 1: 20 most distinctive keywords based on Craig's Zeta in *Daisy sisters* compared to other books in the Swedish corpus, excluding verbs and names

Firstly, the results of the Zeta analysis show some clear genre differences. Crime-related words, such as *polis* 'police', *polishuset* 'police station', and possibly words like *skäl* 'reason(s)', *samtalet* 'the conversation' and *frågor* 'questions' occur clearly less frequently in *Daisy sisters*. Words related to murder were also on the list. If the books in the reference corpus had only included literary novels, these types of words had probably not been included.

Another obvious difference has to do with spelling conventions and register. *Mej* and *mig* are spelling variants of the same word: 'me'. The variant *mej*, occurring relatively more frequent, in *Daisy sisters*, is the less formal variant which is closer to speech, whereas *mig* is the official variant. The same is true for *dig* and *dej* and *sej* and *sig*. This pattern could also be detected in the spelling of certain verbs, like *säga* 'write', which occurred relatively more frequent in the alternative, informal spelling variant *säja* in *Daisy sisters*. There are other words on the keyness list that confirm the idea that *Daisy sister* is written in a more speech-like, colloquial style. Examples are *jo* 'yes' used after a negation and *ju*, a discourse particle that is especially frequent in spoken language. Similarly, *morsan* is a colloquial form for 'mother' and *vadå* a colloquial form for *vad* 'what'. The keyness list also contains swear words and curse words, which are clearly associated with everyday, informal language, *jävla* 'fucking', *herregud* 'lord' and *fan* 'damn'.

The following example from *Daisy sisters* contains three keywords from the keyness list:

Men vad spelar det för roll att **morsan** är här och säger att hon skäms? Hon kan **ju** inte veta något. Mer än... Ja, **vadå**? Så minns hon allt blod och förstår att det var därför hon måste gå till sjukhuset.

The English translation of this passage is as follows: ³

What does it matter that **mom** is here saying she's ashamed? She can't know anything, **right**? More than.. well **what**? Then she remembers all the blood and realizes that's why she had to go to the hospital.

In both the English translation and the official Dutch translation of this passage the colloquial style is at least partially lost:

Maar wat maakt het uit dat **haar moeder** hier is en zegt dat het een schande is? Ze weet **toch** nergens van. Alleen dat ... Ja, **wat**? Dan herinnert ze zich al het bloed en ze begrijpt dat ze daarom naar het ziekenhuis moest.

Discourse particles in general are very hard to translate, because they can have various meanings depending on context (Aijmer 2008). Here *ju* is translated, but there is no Dutch or English equivalent that is equally frequent and associated with speech as much as the Swedish word. The two other colloquial words in this short passage are translated into standard Dutch, which leads to a loss of this style feature.

A final result from the Zeta analysis is that different synonyms are used in the primary corpus and the reference corpus. In the list of distinctive words, *omedelbart* is preferred in *Daisy sisters* and *genast* is avoided. These words are synonyms and both mean 'immediately' with no difference in register.

3. My translation.

Table 2 shows the list with distinctive words based on Craig's zeta in the Dutch translation corpus. The genre differences are even more obvious in this list compared to the original Swedish list: words like *onderzoek* 'investigation', *politiebureau* and *bureau* 'police station' *vermoord* 'murdered', *waarheid* 'truth' and *lichaam* 'body' are clearly linked to the crime genre. This also confirms the earlier findings in the cluster analyses that genre differences seem to be magnified in the translations.

preferred		avoided	
<i>immers</i>	'after all'	<i>onze</i>	'our'
<i>want</i>	'because'	<i>onderzoek</i>	'investigation'
<i>nou</i>	'well'	<i>ineens</i>	'suddenly'
<i>gewoon</i>	'just'	<i>politiebureau</i>	'police station'
<i>ja</i>	'yes'	<i>dood</i>	'dead'/'death'
<i>opeens</i>	'suddenly'	<i>bureau</i>	'police station'/'desk'
<i>verdomme</i>	'damn'	<i>politie</i>	'police man'
<i>minder</i>	'less'	<i>vermoord</i>	'murdered'
<i>aardig</i>	'kind/rather'	<i>zee</i>	'sea'
<i>voorbij</i>	'(all) over'/'past'	<i>zeer</i>	'very'
<i>vieze</i>	'dirty'	<i>water</i>	'water'
<i>niks</i>	'nothing'	<i>telefoon</i>	'phone'
<i>kennelijk</i>	'apparently'	<i>vlak</i>	'right'/'flat'
<i>hemel</i>	'heaven'	<i> bezig</i>	'in process'
<i>baan</i>	'job'	<i>aantal</i>	'number'/'amount'
<i>hoekje</i>	'corner'	<i>waarheid</i>	'truth'
<i>raar</i>	'strange'	<i>papieren</i>	'paper'
<i>geluk</i>	'luck'	<i>vervolgens</i>	'then'
<i>nergens</i>	'nowhere'	<i>lichaam</i>	'body'
<i>zij</i>	'she'/'they'	<i>reden</i>	'reason'

Table 2: 20 most distinctive keywords in *Daisy sisters* based on Craig's Zeta, compared to other books in the Dutch translation corpus, excluding verbs and names

The register difference, on the other hand, is not as obvious as in the Swedish list, although *nou* 'well' *gewoon* 'just' *ja* 'yes' *verdomme* 'damn' do point in the direction of register and speech-like language. *Immers*, which is on top of the list of distinctive words in the Dutch translated corpus, is a good example of translationese. It is the translation of the previously mentioned discourse particle *ju*. In terms of meaning, this translation is accurate, but *immers* does not at all belong to the same register. While *ju* is associated with spoken language, *immers* is almost exclusively used in written language and has a somewhat archaic connotation. Again, this indicates that the speech-like, informal style gets partially lost in the Dutch translation. In the Dutch list with distinctive keywords there are also two synonyms both meaning 'suddenly': *opeens* is preferred in *Daisy sisters* whereas *ineens* is preferred in the other books in the corpus. This can likely be explained by the individual preference of the translator. However, more research about the influence of the translator on style is necessary to confirm this.

6. Conclusion

In this paper, 32 books by the Swedish writer Henning Mankell were investigated using stylistometric methods, to find out whether his style changed measurably over time, or if some of his books deviate stylistically from his other works for other reasons. 10 books

by other Swedish authors were added to the corpus as a reference. The study also gives more insight into the methods that are frequently used in stylometry, such as cluster analysis and PCA, that basically are black boxes, because they give little information about the stylistic features that differ between texts. For this purpose, the original Swedish texts were also compared to the Dutch translations of the same 42 texts to determine how translation and language influence the results of stylometric analyses.

Cluster analyses and PCAs of the data showed that works were clustered by author in the first place and secondly by genre, although there were a few exceptions. The division into genre was somewhat stronger in the translated corpus. The analyses also seemed to indicate that the factor time explains part of the variance. However, on closer inspection, verb tense rather than year of publication turned out to be the decisive factor: the most deviant books in the corpus were primarily written in the present tense, whereas most other books were predominantly written in the past tense. Moreover, narrative perspective also influenced the results noticeably. An analysis of the pronoun ratios in the works in the corpora indicated that the majority of the novels in the corpus had a dominant third person male perspective. Books that mainly had a first person perspective tended to cluster together, just like books with a third person female perspective. Leaving out pronouns and verbs marked for tense showed a very different picture.

Finally, an analysis of the data based on Craig's Zeta (Craig and Kinney 2009) showed that words most distinctively used in the original Swedish *Daisy sisters* were often colloquial words with a speech-like connotation. Words that were avoided in the novel were associated with the crime genre. However, the most distinctive words in the Zeta analysis for the translated Dutch corpus, were not as clearly related to register. The genre differences seemed magnified in the Zeta results of the translation corpus compared to the list of Swedish keywords. This confirmed the findings in the cluster analyses and the PCAs that books were more clearly clustered by genre in the translated texts. This can be due to the different language and language specific features or due to inherent characteristics of translated texts in general. More research on different languages and translations would be useful to get a better understanding of this process. In a follow-up study I intend to investigate the style differences between translators and how they can be detected and measured.

This study has shown that Zeta analysis and a closer look at word lists in stylometric studies can give useful insights into the specific style features that make texts different from each other instead of only focusing on the fact that they differ.

7. Data Availability

Data can be found here: data.example.edu/data

8. Software Availability

Software can be found here: github.com/something

9. Acknowledgements 579

I would like to thank Karina van Dalen-Oskam for our fruitful conversations and for her valuable feedback on earlier versions of this article. 580
581

10. Author Contributions 582

Martje Wijers: Conceptualization, Writing – original draft, Formal analysis, Investigation 583
584

References 585

- Aijmer, Karin (2008). "Translating discourse particles: A case of complex translation". 586
In: *Incorporating Corpora. The Linguist and the Translator*, 95–116. 587
- Arvas, Paula and Andrew Nestingen (2011). *Scandinavian Crime Fiction*. University of 588
Wales Press. 589
- Berglund, Karl (2013). *Deckarboomen under lupp: Statistiska perspektiv på svensk kriminallit-* 590
teratur 1977–2010. [10.1353/scd.2013.0008](https://doi.org/10.1353/scd.2013.0008). 591
- Burrows, John F. (2002). "'Delta': a measure of stylistic difference and a guide to likely 592
authorship". In: *Literary and linguistic computing* 17 (3), 267–287. 593
- (2007). "All the way through: testing for authorship in different frequency strata". 594
In: *Literary and Linguistic Computing* 22.1, 27–47. 595
- Can, Fazli and Jon M Patton (2004). "Change of writing style with time". In: *Computers 596*
and the Humanities 38, 61–82. 597
- Craig, Hugh and Arthur F Kinney (2009). *Shakespeare, computers, and the mystery of 598*
authorship. Cambridge University Press. 599
- Dalen-Oskam, Karina van (2021). *Het raadsel literatuur: Is literaire kwaliteit meetbaar?* 600
Amsterdam University Press. 601
- Eder, Maciej (Dec. 2015). "Visualization in stylometry: Cluster analysis using networks". 602
In: *Digital Scholarship in the Humanities* 32.1, 50–64. ISSN: 2055-7671. [10.1093/llc/fqv](https://doi.org/10.1093/llc/fqv061) 603
[061](https://doi.org/10.1093/llc/fqv061). eprint: [https://academic.oup.com/dsh/article-pdf/32/1/50/11046630/fqv](https://academic.oup.com/dsh/article-pdf/32/1/50/11046630/fqv061.pdf) 604
[061.pdf](https://doi.org/10.1093/llc/fqv061). <https://doi.org/10.1093/llc/fqv061>. 605
- Eder, Maciej, Jan Rybicki, and Mike Kestemont (2016). "Stylometry with R: A Package 606
for Computational Text Analysis". In: *The R journal* 8 (1). [https://github.com/com](https://github.com/computationalstylistics/stylo) 607
[putationalstylistics/stylo](https://github.com/computationalstylistics/stylo). 608
- Herrmann, J. Berenike, Karina van Dalen-Oskam, and Christof Schöch (Mar. 2015). 609
"Revisiting Style, a Key Concept in Literary Studies". In: *Journal of Literary Theory* 9 610
(1). ISSN: 1862-5290. [10.1515/jlt-2015-0003](https://doi.org/10.1515/jlt-2015-0003). 611
- Hoover, David L (2010). "Teasing Out Authorship and Style with T-tests and Zeta." In: 612
DH, 168–170. 613
- (2020). *Modes of Composition and the Durability of Style in Literature*. Routledge. 614
- Jacobsen, Kirsten (2012). *Mankell om Mankell*. Leopard Förlag. 615
- Jautze, Kim (2014). "Measuring the style of chick lit and literature". In: *DH*. 616
- Jautze, Kim, Corina Koolen, Andreas van Cranenburgh, and Hayco de Jong (2013). 617
From high heels to weed attics: a syntactic investigation of chick lit and literature, 72–81. 618
<http://literaryquality.huygens.knaw.nl>. 619

- Jockers, Matthew L (2013). *Macroanalysis: Digital methods and literary history*. University of Illinois Press. 620
621
- Ríos-Toledo, Germán, Juan Pablo Francisco Posadas-Durán, Grigori Sidorov, and Noé Alejandro Castro-Sánchez (2022). "Detection of changes in literary writing style using N-grams as style markers and supervised machine learning". In: *Plos one* 17.7, e0267590. 622
623
624
625
- Rybicki, Jan (2016). "Vive la différence: Tracing the (authorial) gender signal by multivariate analysis of word frequencies". In: *Digital Scholarship in the Humanities* 31.4, 746–761. 626
627
628
- Squires, Claire (2007). *Marketing Literature: The Making of Contemporary Writing in Britain*. Palgrave Macmillan. 629
630
- Van Rossum, Lisanne, Joris J. van Zundert, and K.H. van Dalen-Oskam (2020). "I Catching: Computationally Operationalising Narrative Perspective for Stylometric Analysis". In: DH Benelux. 631
632
633