

Scale=MatchLowercase []Ligatures=TeX,Scale=1

texgyrepagella[Extension = .otf, UprightFont = *-regular, BoldFont = *-bold,
ItalicFont = *-italic, BoldItalicFont = *-bolditalic]

texgyrepagella-math[Extension = .otf,]

FiraSans[Extension = .otf, UprightFont = *-Regular, BoldFont = *-Bold, Ital-
icFont = *-Italic, BoldItalicFont = *-BoldItalic]

FiraCode[Extension = .ttf, Path = ./fonts/firacode/, Contextuals = Alternate,
UprightFont = *-Regular, BoldFont = *-Bold,]

Ubuntu[Extension = .ttf, Path = ./fonts/ubuntu/, UprightFont = *-R, Bold-
Font = *-B, ItalicFont = *-RI, BoldItalicFont = *-BI]

UbuntuMono[Extension = .ttf, Path = ./fonts/ubuntu/, UprightFont = *-R,
BoldFont = *-B, ItalicFont = *-RI, BoldItalicFont = *-BI]

Cláudia Freitas ¹

Diana Santos ²

1. Department of Letters, Pontifical Catholic University of Rio de Janeiro, Rio de Janeiro.

2. Department of Literature, University of Oslo, Oslo.

Keywords:

distant reading, annotation, Brazilian literature, Portuguese literature

Licenses:

This article is licensed under:



Abstract. In this paper we look at how masculine and feminine characters are described in literature in Portuguese, using a publicly available literary corpus: *Literateca*. We investigate the words used to characterise human beings, after classifying them in four broad categories, namely those related to the social, appearance, character and emotional axes. We study the influence of genre, literary school, author gender, and time, among others.

1. Introduction

1

The way people are described is a rich source of information about societies and cultures, revealing values and beliefs of those who describe. In addition to proper names, there are many other ways of human designation, such as the use of human general nouns like *man*, *woman*, *person*, *gentleman*, *lady*, and designation by traits or functions of the people mentioned (using places of origin, professions, family ties, etc, such as *Brazilian*, *doctor*, *mother*, *foreigner*).

2

3

4

5

6

7

In this paper, we look into how human beings are characterised in literature in Portuguese – also called lusophone literature – using a distant reading approach. In particular, we want to investigate the influence of features such as authorship, geographical origin, historical period and gender (both character gender, and authorial gender).

8

9

10

11

Inspired by Moretti and Sobchuk 2019’s warning, we try to go beyond simple visualisations by date or author, and add other ways to look at the data. Following their “dissecting table” analogy, our aim is to find which pieces are able to provide pertinent analysis, triggering meaningful readings. So, we search for “creative cuttings”, – such as the “volume” of speech verbs in Katsma 2018 – to give us new insight. Specifically, we add the class *human depiction* to our data; still, we aim for consensual and understandable categories, like “century” in history.

12

13

14

15

16

17

18

1.1. Gender in literature 19

The theme of gender roles in fiction texts has received increasing attention in the digital humanities community, as the following works testify. 20
21

Underwood, Bamman, and Lee 2018, looking at English literature (104 thousand works, 22
from 1703 to 2009), found that the gender difference between characters became less 23
pronounced from the middle of the nineteenth century to the present day: actions and 24
attributes of characters became less defined by gender categories. In other words, gender 25
roles tend to become more flexible. At the same time, they also found a decrease in the 26
number of feminine characters as the volume of fiction written by women from 1850 to 27
1950 drops by a half. 28

Exploring the *Black Drama* collection, which contains plays written between 1950-2006, 29
Argamon et al. 2009 reports poor results when trying to automatically distinguish the 30
gender of the author and/or character. However, they found differences in the way 31
masculine and feminine authors and characters use language. Feminine playwrights 32
allocate more than half (52.1%) of speeches to feminine characters, while 34.7% of 33
speeches in plays by masculine authors belong to feminine characters. 34

Working with present-day Dutch literary fiction (170 novels published in one sample 35
year), Smeets 2021 found the same imbalance between masculine and feminine characters. 36
However, the author questions what he describes as a “perhaps naive mimetic assumption” 37
according to which the relative absence of feminine characters is a result of their unequal 38
status in society. From the results of his investigation, feminine characters, although 39
fewer in number, occupy a relatively central position in their fictional social networks 40
– they display more relations, both more relations in general and more relations with 41
important characters. 42

Hoyle et al. 2019, using 3.5 million digitized books in English, analysed the lexical 43
choice (adjectives and verbs) associated to feminine gendered nouns and found that 44
positive adjectives used to describe women were more often related to their bodies than 45
adjectives used to describe men. Following the same trend, Schulz and Bahník 2019 46
explores the depiction of male and female characters using the Google Books Ngram 47
corpus, focusing on twentieth-century English-language fiction. The study analyses 48
adjective-noun bigrams associated with the words *man*, *woman*, *boy* and *girl*, and reports 49
that adjectives associated with *men* are more positive (“honest”, “wise”, “honorable”, 50
and “able”) than those associated with *women* (“vulgar”, “foolish”). As to preferences, 51
“charming”, “fashionable” and “warm” were relatively feminine words, while “lazy” and 52
“mean” were relatively masculine words. Men were described in decreasingly masculine 53
terms throughout the beginning and end of the 20th century; on the other hand, the 54
masculinity of adjectives used to describe women started to slightly increase from 1968 55
to 2000. 56

Weingart and Jorgensen 2013 performed a computational analysis of gendered bodies 57
in ca 200 European fairy tales (German, French and Italian folklore texts translated 58
into English). They show that feminine characters are described more than masculine 59

characters with appearance-evaluative words, suggesting that men are associated with the mind and women with the body.

Cermáková and Mahlberg 2022 explores linguistic descriptions of gendered body language and compare 19th-century British children’s literature (ChiLit Corpus) with contemporary fiction for children (the OCC2000+ corpus, a subcorpus of the Oxford Children’s Corpus). Using a corpus linguistic approach, the authors study sequences of 5 words which contain at least one body part noun and a marker of gender. They found fewer clusters for feminine characters in the 19th century. The contemporary data suggests, on the other hand, a trend for feminine and masculine clusters to become more similar, and an increasing range of options for the description of feminine characters and their interactional spaces. Using the same ChiLit corpus, Cermáková and Mahlberg 2021 focused on nouns – excluding proper names - frequently used to label people, and found that *Mothers* are the most frequent occurring feminine character in the corpus.

It is also worth noting the existence of studies such as Cao and Daumé 2021 and Lucy and Bamman 2021. The first one explores the consequences of gender bias for machine learning. The paper investigates how different aspects of linguistic notions of gender impact an annotator’s judgements of anaphora, and points out that a significant possible source of bias comes from the annotations themselves – from underspecified annotation guidelines and the human annotators. The authors emphasise that both humans and systems should not over-rely on cues such as names, semantically gendered nouns and terms of address, relying on “relatively safe” cues like syntax instead. At the other pole of the machine learning approach, the study conducted by Lucy and Bamman 2021 raises questions on how to avoid unintended social biases when using large language models for storytelling. Focusing on how GPT-3 may perceive a character’s gender based on textual features such as personal pronouns (*he/she/her* etc), the work finds that stories generated by GPT-3 place masculine and feminine characters in different topics and exhibit many gender stereotypes: for example, feminine characters are more associated with family and appearance than masculine characters.

In this paper, we also try to contribute to the investigation of gender roles, using works written in Portuguese. As a crossover between corpus linguistics and digital humanities, we use morpho-syntactic and semantic information automatically provided by the PALAVRAS parser Bick 2014, and we add extra semantic annotation, which will be described below.

With Larson 2017, we recognize that using gender as a variable in Natural Language Processing is an ethical issue, and that we need to explicitly explain what “gender” means along this work. As Larson 2017 points out, there are many views of how gender functions as a social construct. In this study, we treat gender as binary, since in the vast majority of works in our corpus gender was mainly constructed in terms of the binary distinction femininity/masculinity. But we acknowledge that the category “gender” can be more complex than this binary distinction, and that these kinds of studies, which describe the cultural apparatus around gender for an extended period of time do not in any way purpose to assert what gender is, but only how it was/is perceived. So they

should not be used for reinforcing gender stereotypes, as warned against by Mandell 2019. 102
103

1.2. Previous work for Portuguese 104

For distant reading of Portuguese, we are aware of some works dealing with characters in literature Santos and Freitas 2019, as well as of the DIP challenge for automatic character identification in Portuguese Santos, Willrich, et al. 2022, to which we come back later. 105
106
107
108

Our point of departure is the work by Freitas, Martins, and Biar 2022¹ – and later extended in Silva 2021's master thesis – who have suggested a fourfold classification for human characterisation. Human attributes were organised in social, appearance, character and emotional characteristics. 109
110
111
112

Using OBRas, a corpus of Brazilian literature in the public domain Santos, Freitas, and Bick 2018, they studied 223 works by 25 Brazilian authors, two of them women (authoring 3 novels altogether), and observed the following trends: 113
114
115

- Men were more frequently described than women (60%-40%), something which may be related to the fact that there were more masculine characters than feminine ones, roughly in the same proportion. 116
117
118
- The most frequent masculine characterising words were *bom* (good), *sério* (honest), *rico* (rich) and *alto* (tall), while *bonita* (beautiful) was by far the top characteristic for women 119
120
121
- Almost 50% of women depicting words were about beauty (namely *bonita* and *bela*) 122
123
- Character and social predication were most frequent for men; for women, social characterisation reduces to *married* and *rich*. 124
125
- Emotional characterisations like *feliz* ('happy') were (almost) exclusively used for women. 126
127

We wanted to check whether these observations held for a wider collection, including Portuguese literature as well. 128
129

1.3. A brief comparison with DIP 130

It is useful to compare and contrast our study with the recent DIP challenge for Portuguese, an evaluation contest for identifying literary characters and some information about them in Brazilian and Portuguese works Santos, Mota, et al. 2022; Santos, Mota, et al. 2023. By describing it and pointing out the differences, we throw some light on different ways to look at (roughly) the same data. 131
132
133
134
135

1. Although published in 2022, the work was conducted in 2018

For DIP the unit is the literary character, and so the challenge looked at their gender, 136
 their profession, occupation and/or social status, and their family relations with other 137
 characters. But the unit is the character. In addition, "literary character" in DIP 138
 does not include all people. In the present study, we try to look at all mentions of 139
 characterisation of people in the works, so our numbers are not per characters, but per 140
 mentions of people. 141

We will discuss and compare the findings about character gender in section 4.7. 142

1.4. The importance of studying literature in Portuguese 143

Portuguese has a rich literary tradition, but unfortunately the digitisation efforts are 144
 lagging behind other languages. This has for example been discussed in Schöch, Erjavec, 145
 et al. 2021. 146

Also, major actors in the big data landscape, no matter the high number of Portuguese 147
 speakers in the world, have not endowed Portuguese with the "current" tools that are 148
 available for other languages, even with much fewer speakers/readers/writers, like Hebrew 149
 or Italian: there is, for example, no Google Book N-grams² service for Portuguese. 150

Likewise, recent reviews of the computational literature landscape, because they do not 151
 find enough internationally published DH papers on Portuguese, have decided not to 152
 review or include them, therefore contributing actively to the lack of information on 153
 lusophone materials and studies. For example, Schöch, Fileva, and Dudar 2022, page 4 154
 state: 155

As several languages, however, were represented only with relatively low 156
 numbers of articles or papers, and in order not to misrepresent the research 157
 communities these publications stem from, we decided not to take the 158
 materials in several languages into account: (...) 159

This is one of the reasons why we are writing this paper for an international audi- 160
 ence. Maybe the results are not so different than the ones our English-speaking or 161
 English-studying colleagues obtained, but they are novel because they are obtained from 162
 completely different data. 163

2. The material 164

We provide here an overview of the data used, also with the purpose of making it known, 165
 and hopefully, useful, for other researchers. And not the least because it shows the 166
 methodological problems it invites. 167

Attempting to complement close readings of canonical authors with a wider material, 168
 following Moretti 2000; Moretti 2013 and Underwood 2019, we use as many books whose 169
 full text is currently publicly available in Portuguese to investigate properties of literary 170
 text which can be identified in an automatic way. 171

2. <https://books.google.com/ngrams/>

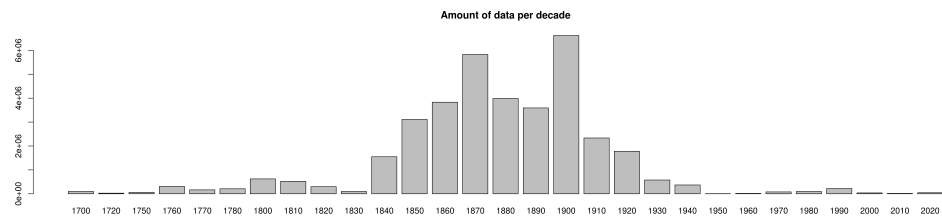


Figure 1: Distribution of words per decade

In order for these data to be shareable and studies replicable, we restrain our data (mostly³) to books in the public domain. We are aware that many more electronic texts exist in electronic form, but by using them we would incur either on law infringement, or at least we would risk creating materials only for our own study, not shareable by others.

Also, it is important to stress that we are referring to textual versions of the works, not simply images. Optical character recognition for Portuguese, especially for old books, is not good enough yet, so all books have been revised by humans, if not born digital.

2.1. Corpus

We used Literateca version 11.1, created on 26 May 2023, comprising ca 32 million tokens of (original) prose (excluding drama) from 1700 on.

A quantitative overview of the material is in Table 1.

Literature	no. of tokens	no. works	no. authors
Total	32,718,621	669	200
Portuguese	20,639,007	306	127
Brazilian	12,079,614	355	73

Table 1: Size of the material, prose from 1700 to the present

Figure 1 shows the distribution of the material in time, by size in words.

Literateca is the merge of several literary corpora written in Portuguese, and thus has some particularities:

- It includes literary works by canonical authors, but also other works by those canonical writers which are not usually or necessary deemed literary, such as newspapers chronics, letters, memoirs, and even scholarly works such as history books or ethnographic studies, and travel reports. For previous centuries, even sermons are included. However, these genres are only included for canonical writers.⁴

3. Exceptions are excerpts of books existing in parallel corpora, or texts whose authors gave us permission to use.

4. By this we mean that established authors who belong to the Portuguese and Brazilian canons were fully digitized, that is, everything they published is available. This is in strong contrast with the works

• It includes drama, poetry and prose.	193
• Some of the works have updated orthography, others keep the original orthography. Given that there have been several norms of Portuguese spelling across the centuries, this means that there can be a variety of forms for the same word.	194 195 196
• While some authors have all their works included, others have only a few, or just one. Especially for non-canonical writers, there is no claim to completeness.	197 198
• Given that the works have been digitised by different bodies and with different tools and for different purposes, there is no claim to homogeneity: works can come from the first or the last paper version, they may keep their prefaces or not, they have different ways of describing chapters, etc.	199 200 201 202
• All works are marked with author, author gender, date of publication, variety of Portuguese, genre, and whether they are original or translated. Some texts are also classified by the literary school they belong to.	203 204 205
We tried to use as much as this material as we could, but we removed poetry and drama. Poetry is probably a natural choice to be removed, because of syntactic idiosyncrasies – and therefore a worse parser performance –, and because we believe that poetry has not so many mentions of fictional characters. We removed drama, also in prose, because it was heavily unbalanced, given that most of the plays were from Portugal.	206 207 208 209 210
As to prose, we started to use everything published since 1700. It is, anyway, important to recognise that we do not have a balanced corpus, and the lion’s part is fiction. We then selected different subsets for different research questions.	211 212 213
• Just fiction, and just non-fiction, to see whether depiction was different across the fiction divide	214 215
• Just works published after 1840, to be able to compare Brazilian and Portuguese authors	216 217
• Just fiction published after 1840, to be able to compare Brazilian and Portuguese literature	218 219
See figures 2 and 3 for a bird’s eye view of the genre distributions in total and in fiction.	220
Only in Figure 3 do we include the variable author gender, since it is only in fiction that we have text written by women.	221 222
In Table 2 we give the numbers of words involved for the material published after 1840.	223
2.2. Gender attribution	224
We explore the influence of gender both in characters’ description and authorship. Masculine and feminine gender labels were manually ascribed to writers, for our corpus contains	225 226

of non canonical authors, which may have had some of their (mainly) novels digitized in the context of other projects.

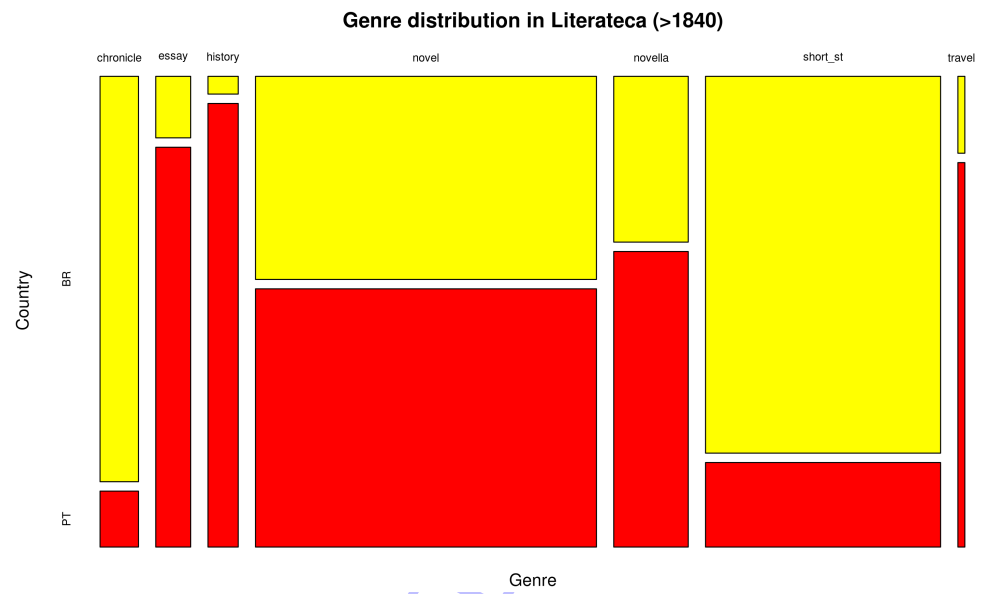


Figure 2: Genre in the full corpus. The unit is the work.

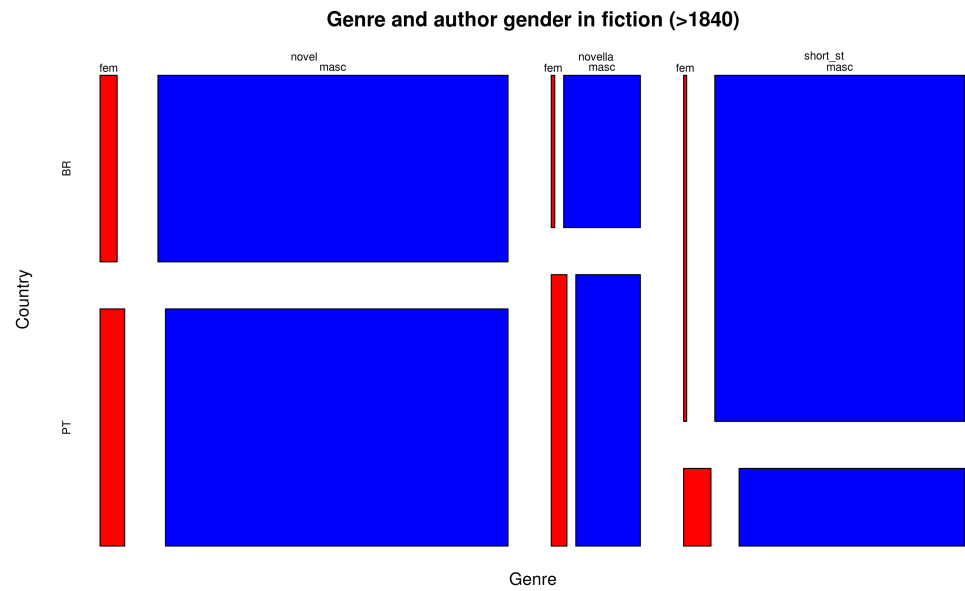


Figure 3: Genre in the fiction corpus. The unit is the work.

	Fiction	Non fiction	Total
Brazil	10,547,327	1,532,287	12,079,614
Portugal	15,280,938	5,358,069	20,639,007
Total	25,828,265	6,890,356	32,718,621

Table 2: Size in words of the different materials, after 1840.

works written by canonical authors that have been widely discussed in literary studies. 227
 For the non-canonical authors, gender was attributed either based on adjective/inflected 228
 forms used in prefaces or based on their proper names. As to the characters, gender 229
 labels were automatically assigned by PALAVRAS parser, and then manually revised by 230
 linguists (Rocha, Freitas, and Santos 2019; Silva 2021). The linguistic clues followed on 231
 attributing and revising gender were syntactic agreement and morphological features. 232

Portuguese is a Romance language that forces the speakers to specify the gender of nouns 233
 (both common and proper nouns) and adjectives. The main formal clue to distinguish 234
 masculine and feminine forms is the word’s ending: masculine forms tend to end in 235
 -o, feminine ones tend to end in -a, and those ending in -e can be both feminine and 236
 masculine – *ponte* (‘bridge’) is feminine, and *pente* (‘comb’) is masculine. However, there 237
 is no perfect equivalence between the ending in -o or -a and the masculine or feminine 238
 gender, respectively – *planeta* (‘planet’) is masculine, and *tribo* (‘tribe’) is feminine. 239
 Therefore, to observe syntactic agreement between the head noun and its modifiers is 240
 the most reliable way to assign morphological gender. 241

When calculating the gender of depicting words, we take into account the gender of the 242
 nominal head (noun, proper noun or pronoun) being characterised, not the gender of 243
 the words (modifiers) associated with it. This choice is due to the fact that, although 244
 adjectives can be inflected for gender in most of the cases, the search patterns we used 245
 also retrieve nouns, which do not admit inflection. Thus, nouns like *anjo* (‘angel’) will 246
 always be masculine, even if the mentioned angels are feminine. When considering the 247
 gender of nominal heads, *anjo*, although a masculine common noun, is classified as a 248
 feminine classifier if it modifies a feminine character. 249

3. The process 250

We wanted to identify all cases where human beings were mentioned to find out how 251
 they were described, or depicted. We extended the search patterns used by Silva 2021⁵ 252
 in two ways: (i) we enriched the lexicon of general human nouns, including names of 253
 professions as targets, and (ii) having extended the amount of works analysed to include 254
 works written by Portuguese authors, we broadened the lexicon of characterising words⁶, 255
 based on prose of the eighteenth, nineteenth and twentieth centuries in Literateca. Along 256

5. Which, in turn, are an improvement of the patterns used in Freitas, Martins, and Biar 2022.

6. The list comprises not only adjectives and nouns, but also verbs (for past participles), given that it is a feature of PALAVRAS that most participles are analysed as verbs even though in an adjectival context.

the process of data analysis, we were forced to discuss previous classification, which lead 257
to precise classification guidelines and to a reclassification of a few words. 258

We start from the idea that specific linguistic patterns indicate certain (semantic) 259
relationships. So, we have used a set of patterns – relying on the automatic morpho- 260
syntactic annotation – to search the material for instances of describing human beings. 261
Some examples of what the patterns yielded follow (the patterns are publicly available). 262

- (1) – Ouviste? – perguntou **ela inquieta**. [– Did you hear? **she** asked **restlessly**.] 263
- (2) ...acudiu logo o **padre**, muito **arisco**. [... came the **priest**, very **skittish**.] 264
- (3) Uma **mulher honesta** não tem segredos para seu marido! [A **honest woman** 265
has no secrets from her husband!] 266
- (4) **D. Joana Tecla** era **idiota**. [–Mrs. **Joan Tecla** was an **idiot**.] 267
- (5) Em todo o caso era uma bela **mulher**, **alta** e forte sem ser gorda... [In any case, 268
she was a beautiful **woman**, **tall** and strong without being fat...] 269
- (6) ...calado como a tarde triste, um **homem**, ainda **moço**, vestido como os essênios 270
taciturnos, caminhava... [...silent as the sad afternoon, a **man**, still **young**, 271
dressed like..] 272

Then we proceeded to classify each word of the aforementioned list – which are the 273
words associated with human beings in the examples –, in four (non-mutually exclusive) 274
classes, according to type of characterisation: social, emotional, physical (appearance) 275
and character. In order to group these idiosyncratic data and provide a better view 276
from afar, we analysed the most frequent words and came up with the four classes. We 277
also used the class **other** if none of the four could hold, and one or more of the four 278
otherwise. As to the assignment of the categories proper, follows their scope and the 279
major decisions associated: 280

social In addition to professions, occupations and social status, we also included 281
absence of profession like *mendigo* ('beggar'), nationality, civil status, family 282
relations, political opinions like *monárquico* ('monarchist'), and cases which are a 283
consequence of social intercourse, like *ignorante* ('ignorant') or *educado* ('civil' or 284
'knowledgeable'). 285

appearance Physical appearance, including clothing or lack of it, as well as those 286
features associated with time, as *jovem* ('young') or *velho* ('old'). 287

emotional Feelings, emotions and emotional tendencies. 288

character Personality traits, also including cognitive properties, such as intelligence or 289
lack of it. It also includes evaluations according to social conduct, such as *honesto* 290
(‘honest’), *malcriado* ('rude') or *pretensioso* ('snob'). 291

It is important to mention that each category works as a label, which in turn encodes 4 perspectives on people: 'appearance' refers to what is visible; 'social' refers to the various roles someone can play in society; 'character' refers to internal/cognitive characteristics; and 'emotion' refers to emotional traits. We could also, and more broadly, consider two large classes: internal characteristics ('character' + 'emotion') and external characteristics ('appearance' + 'social'). We note that the words classified can often refer to non-human entities, as is the case of the next example (7). But if they could modify a human person, they were classified accordingly. However, the results presented in the next sections refer only to those cases where the characterisation was assigned to human beings, such as example (8), since only they are retrieved by the patterns applied.

(7) – Que **triste** pensamento!... [What a **sad** thought!] 302

(8) – Mas a **triste** senhora continuava a choramingar. [But the **sad** woman kept weeping.] 303
304

We classified the retrieved words out of context, except in those rare cases where we had to check whether the adjective had been used as characterising at all in the corpus⁷. For example, initially we wanted to discard the words *granítico* ('made of granite') and *triumfal* ('of triumph'), but we checked the corpus and there were instances where both were applied to human characters, so they were retained in our list.

(9) – Sim, o velho Afonso é **granítico**... [– Yes, old Afonso is **granitic**...] 310

(10) Nunca as mulheres **triumfais** me fizeram bater o coração... [**Triumphal** women never made my heart beat...] 311
312

The classification was done manually by the authors of this paper, and divergences were heartily discussed. We dismissed mistakes, either because (i) they were not characterisation words, (ii) they resulted from wrong parsing, or (iii) we decided they were not relevant to our goals. As to the exclusion:

- We did not take into account "complex adjectives" in the sense of having more than one word, like *bem intencionado* ('having good intentions'), *mal intencionado* ('having bad intentions'), *bem educado*⁸ ('polite'), etc. 317
318
319
- We did not classify relational adjectives, such as *partidário (de...)* ('partisan'), *apologista (de...)* ('in favour of'), *comparável (a ...)* ('comparable to'), *emparelhado com* ('pairing with'), *semelhante a* ('similar to'), since a precise characterization would require a close reading of each sentence. 320
321
322
323

7. Actually, there was one case where we consistently considered the context: in Portuguese, the word *grande* can mean either *big* or *great*. Since each meaning corresponds, in general, to a different syntactic position – *grande homem* ('great man'); *homem grande* ('big man'), we used this information to correctly classify each of the occurrences: *character* or *appearance*, respectively.

8. But note that *educado* and *bem-educado*, as words of size one, were included.

- We threw away misspellings, except for lack of diacritics.⁹ Our rationale being that, in future improved versions of the corpus, the corrected words would be correctly annotated.

Following the annotation approach adopted in the AC/DC project Santos 2014, underlying Literateca, we used multiple classification when two or more categories/senses could be assigned to a characterising word (vague or ambiguous words). References to madness, for instance, were considered both social and character. The same for habits like *madrugador* ('early riser') and *bêbado* ('drunkard' or 'drunk'), which can be either due to biology or social bringing up. The word *acanhado* (shy), can be interpreted as a not-expansive person (character) or as someone fearful (emotion), and the same applies to *impaciente* (impatient), which can be interpreted as anxious (emotion) or restless (character).

Finally, cases such as *maravilhoso* ('wonderful'), *incomparável* ('incomparable'), *ideal* ('ideal') or *horrível* ('horrible'), where it is not clear to which axis they apply out of context, were classified as referring simultaneously to 'character', 'social' and 'appearance'.

To verify the degree of reliability of the classifications and the adequacy of the classes, Silva 2021 carried out a study on the inter-annotator agreement of 15 people in the classification of occurrences considered especially difficult. The degree of agreement was 80%. We have not carried any further studies on this matter.

After this classification, we ended up with a list of 4481 words¹⁰ which might be employed in depicting human beings, see Table 3. Due to the vagaries of the parser, we list the lemmas which can be verb infinitives for past participle forms, because we use the lemmas in our patterns.

type	size
social	1391
appearance	672
emotional	514
character	1578
other	326
total	4481

Table 3: Depicting words, by category. Recall that words can belong to more than one category.

In order to provide a richer description of this list, we show in Table 4 how often depicting words are vague or ambiguous.

We then annotated the corpus with this new classification¹¹ and computed how often and when the words were used to describe human beings.

9. That is, we considered missing accents something that could be present in the original paper edition, but not OCR mistakes.

10. Available from <https://www.linguateca.pt/Gramateca/ListaPredicadoresClassificados.txt>.

11. The classification is encoded in the following tags `pred:carater`, `pred:aparencia`, `pred:social` and `pred:emo`. To find them in Literateca, search for `[sema=".*pred:social.*"]`, etc.

type	size
appearance character	88
appearance emo	12
appearance social	8
appearance character social	10
character emo	107
character emo social	1
character social	80
emo social	9
total with more than one class	315

Table 4: Words belonging to several categories.

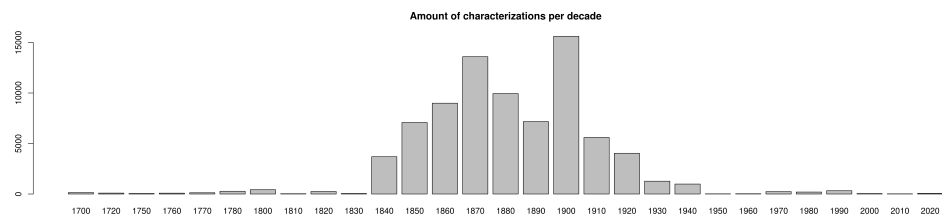


Figure 4: Characterised people in the corpus per decade.

We start by providing a picture of the distribution of human characters in time, in [Figure 4](#), as well as how many depicting events we were able to identify, in [Figure 5](#).

A comment is in place: the decade of 1830 is a clear outlier, because it contains one short text only, of 19,334 words, a political pamphlet by Alexandre Herculano, in the whole decade. The same happens with 1950, which in the material is only represented by 4,777 words of Jorge Amado's *Gabriela, Cravo e Canela*.

4. Analysis

The first thing we report is the proportion of these subclasses in our material. [Table 5](#) shows the raw numbers, and also those referring to masculine and to feminine

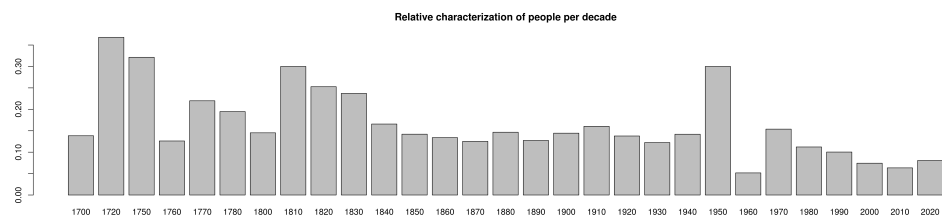


Figure 5: Relative characterisation per person, per decade

characters.¹² Figure 6 displays the overall distribution of characterisation words. 360

	Total	Masc. characters	Fem. characters
People	578,815	352,851	173,370
Characterised people	80,415	52,252	24,664
Social	11793	7813	3534
Appearance	15394	9099	5862
Emotion	9670	5562	3895
Character	23880	16542	6394

Table 5: Different depiction classes, in general and per gender of the characterised person, using the subject's gender.

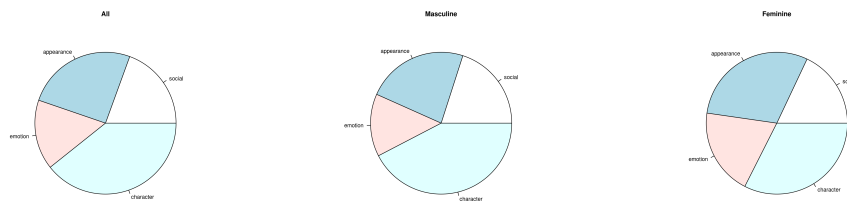


Figure 6: Distribution of characterisation words among the 4 classes, for all, masculine and feminine depictions

The first observation is that there are way more masculine than feminine characters in 361
the material (ca. twice as many). Feminine characters are, however, almost as often 362
characterised as the masculine ones: 14.2% against 14.8%. 363

The second remark is that the by far most frequent subclass deals with character (most 364
frequent words: *bom* ('good'), *grande* ('great'), *honrado* ('honourable, honest'), *simples* 365
(('simple'), *digno* ('with dignity'), *excelente* ('excellent')), followed by appearance (most 366
frequent words: *velho* ('old'), *novo* ('young')¹³, *antigo* ('old-fashioned'), *jovem* ('young'), 367
belo ('beautiful'), *formoso* ('beautiful'), *bonito* ('pretty')). 368

Social characterisation comes third (most frequent words: *rico* ('rich'), *ilustre* ('illus- 369
trious'), *nobre* ('noble'), *casado* ('married'), *célebre* ('famous'), *pobre* ('poor'), *livre* 370
(('free'), *famoso* ('famous')). while emotional characterisation is the least frequent (*pobre* 371
(('poor'), *infeliz* ('unhappy'), *valente* ('brave'), *feliz* ('happy'), *triste* ('sad'), *desgraçado* 372
(('miserable), *alegre* ('joyful'), *humilde* ('humble')). 373

Thirdly, feminine characters have a higher chance of being characterised by their 374
appearance compared to masculine ones (23.8% vs. 17.4%), something that corroborates 375
the findings in previous studies, and to which we return in subsection 4.2. 376

12. It should be noted that the numbers do not add up because in some cases the parser is not able to assign a morphological gender, marking them as M/F. Also, remember that by "character" here we mean mentions to people, not distinct characters.

13. It may seem surprising at first to include age as appearance, but it is something that we visually assess.

4.1. Does textual genre matter? 377

Does it make more sense to look only at literary texts, removing travel writing, essays, history and political writings? 378
379

On the one hand, we had left all material because we wanted to look at the way people described people in Portuguese, but then it is also conceivable that the kinds of information about people are rather different when you write the history of the Inquisition, an essay about your fellow writers, or a report of how you crossing Africa, compared with a narrative where you introduce fictional characters. 380
381
382
383
384

So, we reproduced our queries removing all texts not classified as novels, novellas or short stories, and the new numbers are in Table 6. 385
386

	Total	Masculine	Feminine
Words	25,828,265		
People	490,892	291,403	159,216
Characterised people	47,450	30,036	16,620
Social	8968	5720	2979
Appearance	12,951	7401	5226
Emotion	8767	4922	3665
Character	19,002	12587	5773

Table 6: Different depiction classes, in general and per gender of the characterised person, using the subject's gender, only in novels, novellas and short stories.

It is interesting to see that removing the non-fictional prose genres does not change the relative order of the subcategories, but increases the percentage of feminine characters, that raises from 30.0% to 32.4/%, and characterised feminine characters, from 33.2% to 35.0/%. 387
388
389
390

As to the characterisation of masculine and feminine characters, we have trends similar to the ones we present for the full material, as shown in Figure 7: masculine targets are characterised, by far, by their character, while feminine targets are (almost) equally characterised by their appearance and character. 391
392
393
394

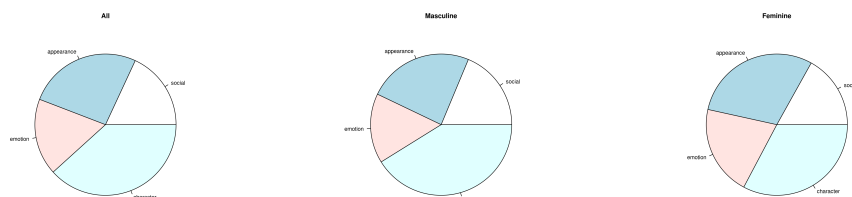


Figure 7: Relative characterisation per gender in novels, novellas and short stories.

Let us see now for the non-fiction part, whether the picture is different. In Table 7 we describe the masculine and feminine characterisations in the (considerably smaller) non-fiction part. 395
396
397

	Total	Masculine	Feminine
Words	6,890,356		
People	87,923	61,448	14,154
Characterised people	10,537	8033	1899
Social	2825	2093	555
Appearance	2443	1698	636
Emotion	966	688	245
Character	4878	3955	621

Table 7: Different depiction classes, in general and per gender of the characterised person, using the subject's gender, only in non fiction.

The percentage of feminine characters, and feminine characterisations shrunk considerably: 16% and 18%, confirming that women are even less important in the public sphere. 398 399 400

Now social characteristics are – globally – more frequent than appearance. Character remains the most important form of describing people, and emotion the least. 401 402

In Figure 8 we present the distribution of the four kinds of features, and see that the few women that are mentioned have a large proportion of appearance descriptions, even more in non-fiction than in fiction. 403 404 405

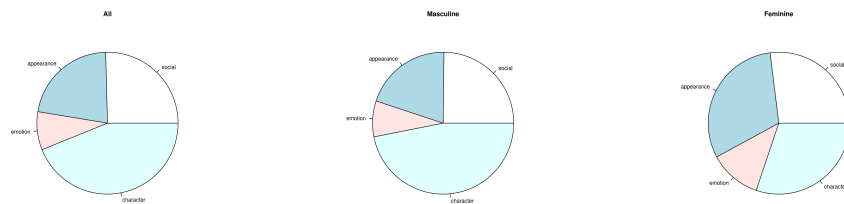


Figure 8: Relative characterisation per gender in non-fiction.

4.2. Differences when describing masculine and feminine characters 406

The previous figures have shown that appearance is more frequent when describing feminine characters. This can also be appreciated in the barplot of figure 9 407 408

However, this is just the tip of the iceberg. The analysis of depictive words preferentially used with masculine and feminine characters can be more revealing than the general analysis we presented in figure 9, which takes into account the whole bunch of depictive words. In order to be evaluated as 'preferred', a word must (i) be used for masculine targets at least for 80% of the occurrences, or for feminine targets more than 60% of the occurrences; (ii) have a total frequency of 4 or more. 409 410 411 412 413 414

In cases where different lexical items correspond to gendered male/female pairs (*mãe/pai* ('mother/father'); *rainha/rei* ('queen/king'); *namorada/namorado* ('girlfriend/boyfriend') etc), we manually grouped the elements of the pair as if they shared the same lemma so they could be included in the preference count. 415 416 417 418

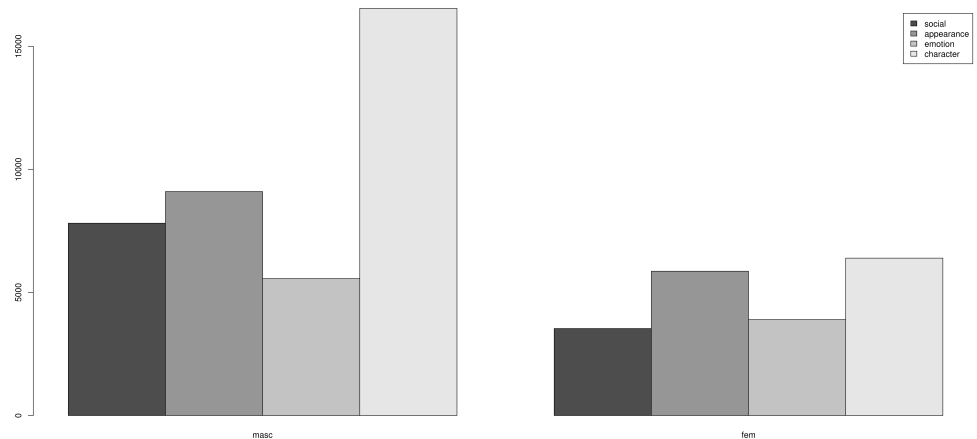


Figure 9: Relative characterisation per gender in the whole material, as a barplot.

The new data are in figure 10, which shows a slightly different picture, in which (i) words 419 of the *emotional* axis are almost not seen at all, and almost disappear in the feminine 420 characters, (ii) the balance between *appearance* and *character* in feminine depiction 421 gives way to a characterisation based mainly on *appearance*, which accounts for half 422 of all preferred feminine characterisations, and (iii) and *appearance*, the second most 423 frequent characterisation (of both masculine and feminine characters), drops to the third 424 position when associated with masculine characters, and up to the first position, when 425 associated with feminine characters. 426

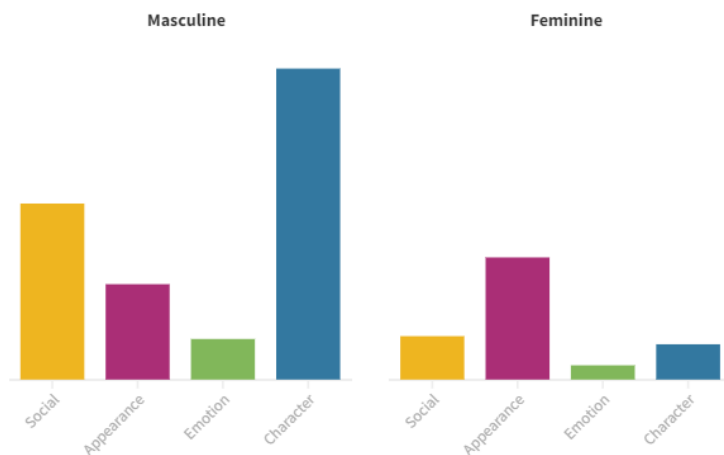


Figure 10: Preferred characterisation per gender

The appearance axis is the second most common for both genders, but figures 11 and 427 12, complementary to figure 10, provide a few details that enrich the analyses. 428

14. In figures 11 and 12, words such as *beautiful_1* and *pretty_2* relate to different Portuguese words that could be translated into the same English word, such as *bonita e formosa*, which could be both translated as *pretty*.

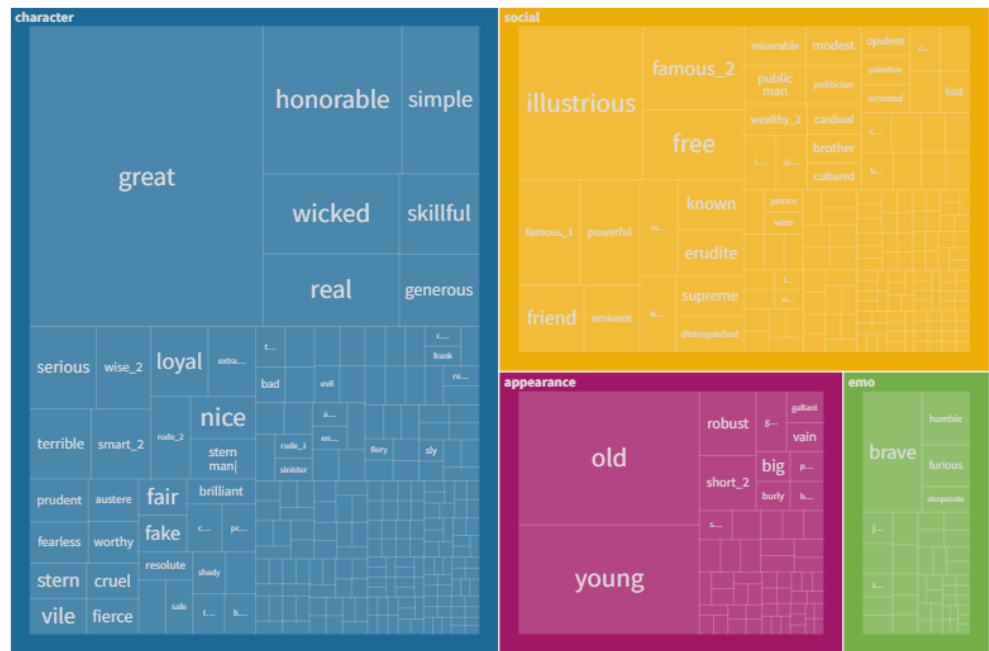


Figure 11: Preferred characterisation of masculine characters.

As noted in previous studies, typically feminine *social* characterisations relate to the 429
familiar environment (*mãe* ('mother'), *prima* ('cousin')), but mentions to marital status 430
are the highlight (*casada* ('married') and *viúva* ('widow') are the most frequent words, 431
but *adúltera* ('adulteress') is frequent as well). Marital status, in turn, is absent as typical 432
masculine social characterisation. These are related to (positive) social recognition such 433
as *ilustre* ('illustrious'), *célebre* ('famous'), *famoso* (another word for 'famous') and 434
poderoso ('powerful'). 435

On the feminine emotional axis, words associated with love and sweetness (*adored* 436
and *sweet*) stand out, but also words associated with sadness and insecurity (*poor*, 437
tearful, *jealous*, *offended*) and fear (*terrified*). On the other hand, bravery is the 438
masculine highlight: *valente* ('brave'), is, by far, the most frequent word, and *atrevido* 439
(('cheeky/audacious') is in the 6th. Anxiety also appears (*desesperado* ('desperate') is 440
the third most frequent). 441

Finally, masculine characters seem to be taken by surprise more often than feminine 442
ones, being *maravilhado* ('marveled'), *assombrado* ('haunted') and *surpreso* ('surprised'), 443
which might be due to their roles in narrative events. 444

Appearance, although highly typical for feminine targets, varies relatively little as to the 445
most mentioned attributes: beauty (*bonita*; *formosa*, *bela*, *linda*, Portuguese words for 446
beautiful) or the lack of it (*feia* ('ugly')) is the most frequent feature. In the masculine 447
appearance axis, age and size, instead of beauty, are the most mentioned attributes: 448
velho ('old') and *jovem* ('young'); *robust*, *big* and *short*. 449

In the character axis, typically masculine, stands out *grande* ('great'), *honrado* ('hon- 450

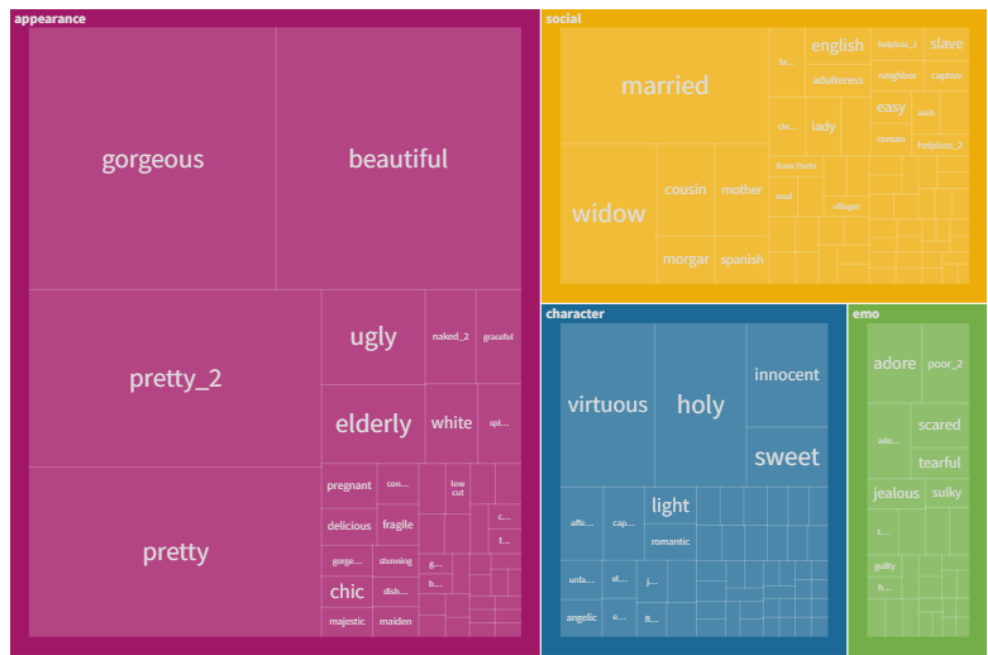


Figure 12: Preferred characterisation of feminine characters.

ourable’). Other highly mentioned positive traits are *generoso* (‘generous’), *habilidoso* 451
 (‘skillful’), *real*, *sério* (‘serious’) and *leal* (‘loyal’). For the feminine targets, the highlights 452
 are *virtuosa* (‘virtuous’), *inocente* (‘innocent’) and *meiga* (‘sweet’). 453

4.3. Does the gender of the author matter? 454

Do these findings change according to the author’s gender? In our material, see Table 8, 455
 feminine authors use more appearance descriptions than masculine ones, as shown in 456
 Figure 13. 457

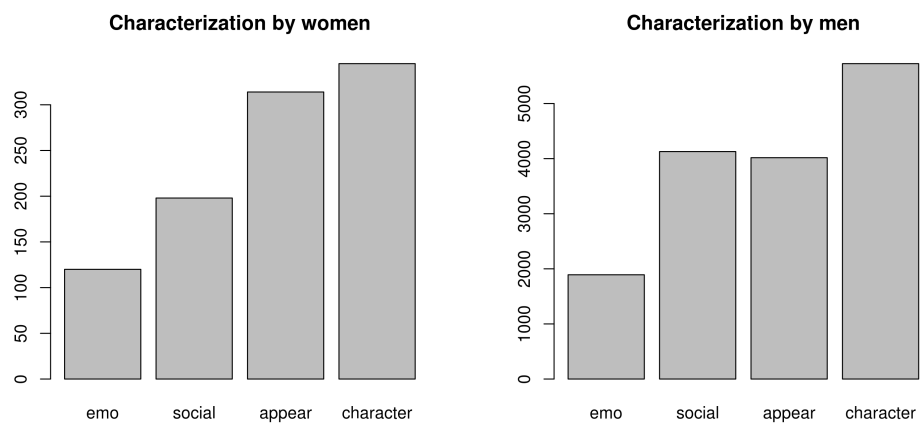


Figure 13: Characterisation by masculine and feminine authors. E PRECISO REFAZER

However, there is a huge difference in the size of the material compared: there are only 1.2 million words written by women compared to almost 32 million words by men. In fact, this is an inescapable problem, given the reduced number of writings by women in our corpus: only 19 authors¹⁵ who wrote 33 works in prose.

	Total	Feminine author	Masculine author
Words	25,828,265	1,206,744	24,621,521
People	490,892	24,271	466,621
Characterised people	57,680	2235	55,445
Social	8968	355	8613
Appearance	12951	595	12356
Emotion	8704	533	8171
Character	19002	887	18115

Table 8: Different depiction classes, for masculine and feminine authors, in novels, novellas and short stories.

Even though the material is heavily unbalanced, we tried to discern any interesting trend in works written by women as far as whose appearance was more described – could it be that they would emphasise or concentrate more on the appearance of masculine characters?

We get 265 appearance descriptions of feminine characters, and 319 of masculine characters, in 985 characterisations of feminine characters and 1195 characterisations of masculine characters. In other words, 26.9% of feminine characterisations and 26.7% of masculine characterisations involve their appearance. But we acknowledge that numbers are too small to be conclusive. In any case it is conspicuous that both genders have roughly the same characterisation frequency in literature written by women.

Despite the imbalanced data, figure 14 shows preferential characterisation regarding gender of both characters and writers. In what follows we sketch some differences between human depiction in works written by men and women. The main difference is the increase of appearance in masculine characterisation in works written by woman.

Beginning with feminine characters and focusing on women writers only, we found that *married* is no longer among the most frequent social depictions, but *widow* and *single* remain. Despite still being frequent, less space is devoted to beauty in works written by women. On the other hand, age is more present: *young* and *old*. As to emotional characterisation, *happy* and *adorable* are the highlights, and none of the preferred emotional words relate to sadness. As to character, the highlights of feminine depiction words are *honest*, *infamous*, *crazy*, *refined* and *dangerous*. In the social axe, masculine characters are mainly *married* and *noble*. Positive emotions are present for masculine characters as well (*happy/pleased*, *enthusiastic*), but bravery (*brave*) has only

15. Júlia Lopes de Almeida, Virgínia de Castro e Almeida, Ana Plácido, Teresa Margarida da Silva e Orta, Maria Amália Vaz de Carvalho, Maria O'Neill, Maria Firmina dos Reis, Florbela Espanca, M.M.S.A. e Vasconcelos, Cláudia Campos, Maurícia C. de Figueiredo, Maria Luísa Marques da Silva, Matilde Isabel de Santana e Vasconcelos Moniz Bettencourt, Ana de Castro Osório, Alice Moderno, Maria Peregrina de Sousa, Paulina Filadélfia, Clarice Lispector and Sónia Coutinho, by decreasing number of words in the corpus

one occurrence. Masculine appearance follows the general trend, and masculine character are mainly *kind* and *honourable*.

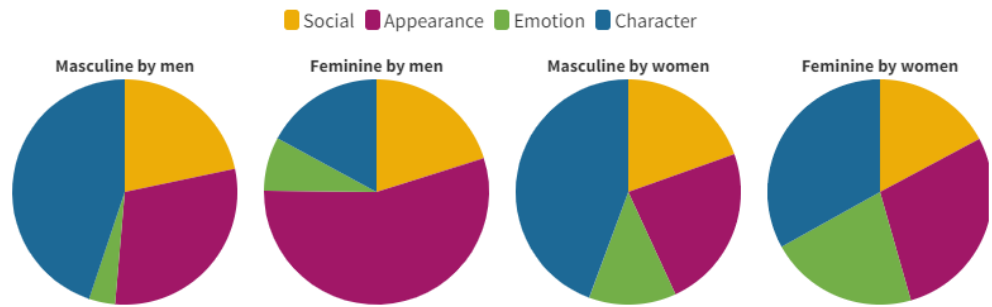


Figure 14: Preferred characterisation by masculine and feminine authors.

4.4. Difference between Brazil and Portugal

Are there differences between the two countries as regards people’s characterisation?

We compared the works from 1840 to the present (Brazil became independent in 1822, and, as already mentioned, for the 1830 decade we only have one work by a Portuguese author).

We decided to compare only novels, novellas and short stories between the two countries, because the non-fiction parts differ widely: While we have a large body of texts on history in the Portuguese side, we have almost only short essays in newspapers on the Brazilian side. The results are presented in Table 9.

	Total	Brazil	Portugal
People	486,575	209,283	277,292
Characterised people	46,704	19,642	27,062
Social	8887	3545	5342
Appearance	12877	6199	6678
Emotion	8704	4874	3650
Character	18782	7649	11133

Table 9: Different depiction classes in novels, novellas and short stories, in general and per author nationality, after 1840.

We see that *character* and *social* characterisation are somewhat higher in Portuguese literature, while the other categories – especially emotion – are more pronounced in Brazilian literature. One may wonder if this is due to a more socially rigid society in Portugal, or whether the cause lies with the historical novels (almost absent in the Brazilian material, and quite frequent in the Portuguese material).

We also investigated whether the differences among genders are more obvious in the Brazilian material, or different from the ones in the Portuguese material. For this, we created Table 10, where we can see that Brazilian literature has a higher proportion of mentions of feminine characters (36.5%) than the Portuguese (29.7%). This may again be due to the historical novels, but will have to be investigated closer.

	Br total	Br fem.	Br masc.	Pt total	Pt fem.	Pt masc.
People	202,829	74,020	118,088	275,301	81,847	165,796
Characterised	17453	6381	10591	24548	8452	15372
Social	3545	1216	2217	5342	1753	3434
Appearance	6199	2579	3472	6678	2618	3885
Emotion	3474	1444	1949	5230	2206	2925
Character	7649	2446	4955	11133	3292	7452

Table 10: Different depiction classes in novels, novellas and short stories after 1840, per author nationality and per gender of the characterised.

Here we see that the social status of male characters is more important in Portuguese literature. 506 507

If we now compare the distribution by country by gender, presented in 15, masculine characters seem to be similarly depicted. But for feminine characters, there are significantly relatively less mentions of their social status and more mentions of their appearance in Brazilian authored works. 508 509 510 511

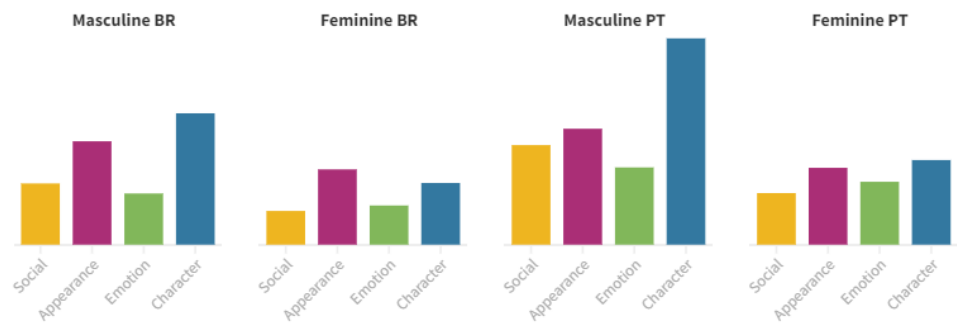


Figure 15: Characterisation by country

4.5. Differences among authors 512

In table 11, we show the distribution of the kinds of characterisation for 12 canonical authors, 6 Brazilian and 6 Portuguese. 513 514

We see there are some differences among these authors. They agree in that none emphasises an explicitly emotional description, and several authors follow the "general" pattern in fiction: first *character*, then *appearance*, then *social*, and last *emotion*: Machado de Assis, Eça de Queirós, Aluísio de Azevedo, José de Alencar, Júlio Dinis, Teófilo Braga and Alexandre Herculano. 515 516 517 518 519

But in José Manuel de Macedo, Coelho Neto and Raul Brandão *appearance* is the most frequent characterisation, and *character* is the second most frequent. 520 521

As to the relative order of *character* and *social* characterisation, Humberto de Campos is the only one who reverts the "canonical" order, using more *social* characterisations than those reflecting *character*, while Camilo Castelo Branco (incidentally the author 522 523 524

Author	land	nr	total	char	soc	app	emo	mfreq
Camilo Castelo Branco	PT	42	4045	1781	938	845	481	pobre
Machado de Assis	BR	140	1864	793	219	643	209	bom
Eça de Queirós	PT	16	2487	1019	420	923	125	bom
JM de Macedo	BR	7	1325	411	232	515	167	velho
Aluísio Azevedo	BR	13	1307	513	191	374	229	pobre
José d'Alencar	BR	15	887	331	154	370	32	velho
Coelho Neto	BR	17	966	369	81	440	76	velho
Humberto de Campos	BR	6	766	169	193	368	36	velho
Júlio Dinis	PT	9	1038	430	127	302	179	pobre
Teófilo Braga	PT	4	419	144	82	112	81	pobre
Alexandre Herculano	PT	8	809	321	201	228	59	velho
Raul Brandão	PT	5	206	73	24	102	7	grande

Table 11: Different depiction classes per authors, ordered by number of characterisations. "nr" shows the number of different fiction works by that author in Literateca, and "mfreq" the most frequent characterising word.

with more works in Literateca) is the only that describes more *social* than *appearance* 525

In any case, there are also differences in the amount of characterisation provided by 526
each author: Figure 16 illustrates how much each author depicts, i.e. how many 527
characterisations he uses per number of words. 528

Relative characterization per author

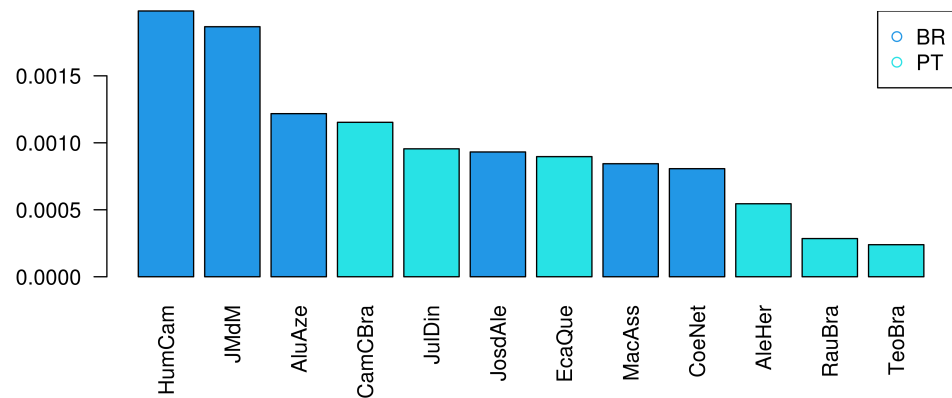


Figure 16: Characterisation by author.

In Figure 17, we represent each author in a plane formed by internal and external 529
characteristics. 530

4.6. The influence of literary school 531

For a subset of the works of Literateca we have metadata about the literary school they 532
belong, as has been described in Santos, Pires, et al. 2020. 533

We selected all works marked as romantic in one group (11,850,395 words, 175 books), 534

Authors by relative characterization

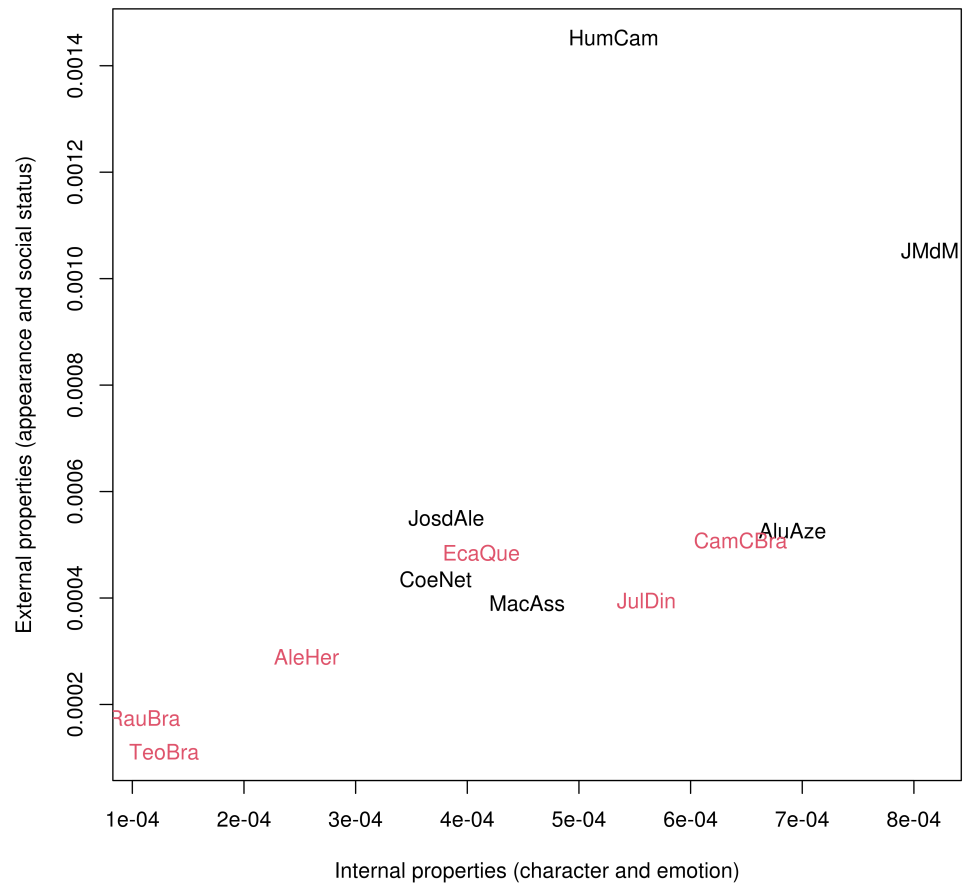


Figure 17: Characterisation by author as far as kind and relative weight of characterisation.

and those marked as realist or naturalistic (7,616,384 words, 121 different books) in [535](#) another group,¹⁶ to see whether one could identify differences as to people's depictions, [536](#) just based on this fourfold sub-classification, and also according to the gender of who [537](#) gets characterised. The results are presented in Table [12](#) and in Figure [18](#). [538](#)

The first interesting remark is that there are (relatively) more mentions of feminine [539](#) characters in realist works than in romanticism. However 10.9% of the feminine occurrences [540](#) are characterized in romantic books (and 9.9% of masculine occurrences), but only 9.8% [541](#) in realist ones (compared to 9.5% for masculine). [542](#)

We see that in romanticism there are far more *character* characterisations of masculine [543](#) characters than in realism, where the relationship across all kinds of characterizations is [544](#) stable across genres. In addition, realism describes physical appearance of both genders, [545](#) while romanticism prefers feminine appearance. [546](#)

16. Note that the groups are not mutually exclusive: there are a few books classified as both romantic and realist, which correspond to the transition between the two schools.

	Romantic	fem	masc	Realist	fem	masc
People	238,338	74,991	142,245	149,699	52,771	86,861
Characterised	22,733	8140	14041	13834	5187	8244
Social	4629	1510	3002	2516	946	1501
Appearance	5573	2279	3179	3944	1678	2147
Emotion	4370	1932	2350	2635	1112	1464
Character	9389	2899	6237	5649	1819	3650

Table 12: Different depiction classes in novels, novellas and short stories, per literary school and per gender of the characterised.

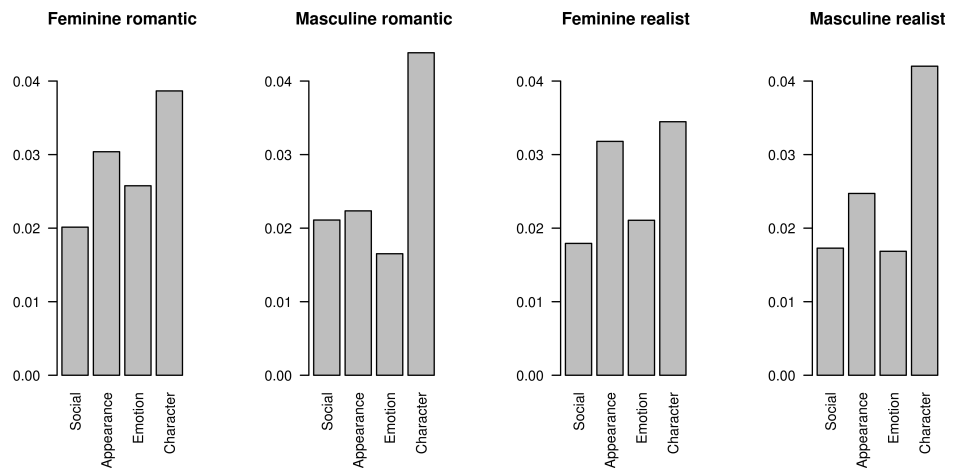


Figure 18: Characterisation per literary school and per gender

4.7. Going back to DIP

547

DIP has clearly demonstrated that there are fewer feminine characters in lusophone literature.

548

549

But in this study we see that those feminine characters are relatively more characterised, at least for appearance, than the masculine ones.

550

551

Ideally, and for the near future, we would like to connect the two studies/activities/forms of distantly looking at literature and provide, for each literary work, not only their description in terms of characters (as DIP does), but also how each character is characterised, using the present work and some form of anaphoric resolution of the non-proper name depictions and of those cases where human subjects (whether or not proper names) are omitted (Freitas and Souza 2021 found omitted subjects in 41% of clauses in Brazilian literature material).

552

553

554

555

556

557

558

We might therefore link kinds of characters with particular clusters of properties, like the beautiful rich woman and the poor honest lad and the evil old priest.

559

560

5. Concluding remarks 561

In this paper, we offered some insights into human depiction based on distant reading 562 literature in Portuguese. We can summarise our results as follows: human depiction 563 seems to obey the pattern *character, social, appearance* and *emotion* for masculine 564 characters, and *character* and *appearance, social* and *emotion* for feminine characters. If 565 we consider only preferred depiction words, differences between feminine and masculine 566 characters become more pronounced, and changing lens – from distant to close reading – 567 reveals that features associated with characters are related to their genders. The results 568 also suggest an impact of the author’s gender in the types of characterisation used, but 569 the limited number of works written by women hinders a more definite conclusion. 570

We acknowledge that the material we used (works and words) is smaller than those used 571 in other studies conducted under the umbrella of Digital Humanities. However, our 572 findings show that an advantage of annotated data is the opportunity to see trends and 573 patterns even in moderate-sized collections. On the other hand, we stress that another 574 intention with this work is to convince (the Portuguese-speaking community, mainly) to 575 enlarge literary Portuguese-language collections with machine readable texts. 576

In the near future, we would like to assess the precision of each rule used, and to correct 577 the detected mistakes, as well as to widen the scope of characterisation. We are aware 578 that human depiction is not restricted to the lexical-syntactic patterns we used, and to 579 detect other ways Portuguese language manifests characterisation is, therefore, a natural 580 route to continue the investigation. 581

We are also aware that our study reflects mainly the vision of male authors of nineteenth 582 and early twentieth century, and that therefore it is by no means an unbiased description 583 of gender. 584

Other studies that we may still do on this material is to add an evaluation view: of 585 these ways of depicting, which ones are positive, negative, or neutral? This is more 586 straightforward for character and emotional words, but also possible for appearance and 587 even social descriptions. 588

We could also separate age from appearance, and check what this dimension may bring. 589

Anyway, all material is open for inspection, from the lists of the characterising words 590 to the patterns used, and the annotated works themselves, which allow interested 591 researchers to redo our searches and even refine them. 592

6. Software availability	593
Not relevant	594
7. Acknowledgements	595
Funding, Funding and thanks!	596
8. Author contributions	597
Cláudia Freitas: Conceptualization, Writing – original draft, review & editing	598
Diana Santos: Conceptualization, Writing – original draft, review & editing	599
References	600
Argamon, Shlomo et al. (2009). “Gender, Race, and Nationality in Black Drama, 1950-2006: Mining Differences in Language Use in Authors and their Characters”. In: <i>Digit. Humanit. Q.</i> 3.2. URL: http://www.digitalhumanities.org/dhq/vol/3/2/000043/000043.html .	601 602 603 604
Bick, Eckhard (2014). “PALAVRAS, a Constraint Grammar-based Parsing System for Portuguese”. In: <i>Working with Portuguese Corpora</i> . Ed. by Tony Berber Sardinha and Thelma de Lurdes São Bento Ferreira. Bloomsbury Academic, pp. 279–302.	605 606 607
Cao, Yang Trista and III Daumé Hal (Nov. 2021). “Toward Gender-Inclusive Coreference Resolution: An Analysis of Gender and Bias Throughout the Machine Learning Lifecycle*”. In: <i>Computational Linguistics</i> 47.3, pp. 615–661. ISSN: 0891-2017. DOI: 10.1162/coli_a_00413. eprint: https://direct.mit.edu/coli/article-pdf/47/3/615/1971880/coli_a_00413.pdf . URL: https://doi.org/10.1162/coli%5C_a%5C_00413 .	608 609 610 611 612 613
Cermáková, Anna and Michaela Mahlberg (2021). “The representation of mothers and the gendered social structure of nineteenth-century children’s literature”. In: <i>English Text Construction</i> 14.2, pp. 119–149. ISSN: 1874-8767. DOI: https://doi.org/10.1075/etc.00044.cer . URL: https://www.jbe-platform.com/content/journals/10.1075/etc.00044.cer .	614 615 616 617 618
— (2022). “Gendered body language in children’s literature over time”. In: <i>Language and Literature</i> 31.1, pp. 11–40. DOI: 10.1177/09639470211072154. eprint: https://doi.org/10.1177/09639470211072154 . URL: https://doi.org/10.1177/09639470211072154 .	619 620 621 622
Freitas, Cláudia, Flávia Martins, and Liana Biar (Feb. 2022). “Um ‘olhar discursivo’ sobre predicação e gênero: aproximações metodológicas entre corpus e discurso”. In: <i>Texto Livre</i> 15, e36213. DOI: 10.35699/1983-3652.2022.36213.	623 624 625

- Freitas, Cláudia and Elvis Souza (2021). “Sujeito oculto às claras: uma abordagem 626
 descritivo-computacional / Omitted subjects revealed: a quantitative-descriptive 627
 approach”. In: *Revista de Estudos da Linguagem* 29.2, pp. 1033–1058. ISSN: 2237- 628
 2083. DOI: [10.17851/2237-2083.29.2.1033-1058](https://doi.org/10.17851/2237-2083.29.2.1033-1058). URL: [http://www.periodicos 630](http://www.periodicos 629

 .letras.ufmg.br/index.php/relin/article/view/17439)
- Hoyle, Alexander Miserlis et al. (July 2019). “Unsupervised Discovery of Gendered 631
 Language through Latent-Variable Modeling”. In: *Proceedings of the 57th Annual 632
 Meeting of the Association for Computational Linguistics*. Florence, Italy: Association 633
 for Computational Linguistics, pp. 1706–1716. DOI: [10.18653/v1/P19-1167](https://doi.org/10.18653/v1/P19-1167). URL: 634
<https://aclanthology.org/P19-1167>. 635
- Katsma, Holst (2018). *Loudness in the novel*. workingpaper. URL: [https://litlab.sta 637](https://litlab.sta 636

 nford.edu/LiteraryLabPamphlet7.pdf)
- Larson, Brian (Apr. 2017). “Gender as a Variable in Natural-Language Processing: 638
 Ethical Considerations”. In: *Proceedings of the First ACL Workshop on Ethics 639
 in Natural Language Processing*. Valencia, Spain: Association for Computational 640
 Linguistics, pp. 1–11. DOI: [10.18653/v1/W17-1601](https://doi.org/10.18653/v1/W17-1601). URL: [https://aclanthology 642](https://aclanthology 641

 .org/W17-1601)
- Lucy, Li and David Bamman (June 2021). “Gender and Representation Bias in GPT-3 643
 Generated Stories”. In: *Proceedings of the Third Workshop on Narrative Understand- 644
 ing*. Virtual: Association for Computational Linguistics, pp. 48–55. DOI: [10.18653/v 646](https://doi.org/10.18653/v 645

 1/2021.nuse-1.5)
<https://aclanthology.org/2021.nuse-1.5>. 646
- Mandell, Laura (2019). “Gender and Cultural Analytics: Finding of Making Stereotypes?” 647
 In: *Debates in the Digital Humanities*. Ed. by Matthew K. Gold and Lauren F. Klein. 648
 Manifold Scholarship. 649
- Moretti, Franco (2000). “The slaughterhouse of literature”. In: *Modern Language Quar- 650
 terly* 61.1. 651
- (2013). *Distant Reading*. Verso. 652
- Moretti, Franco and Oleg Sobchuk (2019). “Hidden in Plain Sight: Data Visualization 653
 in the Humanities”. In: *New Left Review* 118, pp. 86–115. 654
- Rocha, Luísa, Cláudia Freitas, and Diana Santos (Oct. 2019). “Preparação para Leitura 655
 Distante em português: diálogos entre PLN e Humanidades Digitais”. In: *Anais do 656
 TILic 2019*. URL: [https://www.linguateca.pt/Diana/download/Rochaetal2019 658](https://www.linguateca.pt/Diana/download/Rochaetal2019 657

 .pdf)
- Santos, Diana (2014). “Corpora at Linguateca: Vision and roads taken”. In: *Working 659
 with Portuguese Corpora*. Ed. by Tony Berber Sardinha and Thelma de Lurdes São 660
 Bento Ferreira. Bloomsbury Academic, pp. 219–236. 661
- Santos, Diana and Cláudia Freitas (Oct. 2019). “Estudando personagens na literatura 662
 lusófona”. In: *STIL 2019 - Symposium in Information and Human Language Tech- 663
 nology and Collocates Events, October 15-18, 2019, Salvador, BA, Proceedings of 664
 conference*, pp. 48–52. 665
- Santos, Diana, Cláudia Freitas, and Eckhard Bick (Sept. 2018). “Obras: a fully annotated 666
 and partially human-revised corpus of Brazilian literary works in public domain”. In: 667
CorLex. URL: <http://www.linguateca.pt/Diana/download/CorLex.pdf>. 668

- Santos, Diana, Cristina Mota, et al. (2022). *Introduction to DIP: goal, setup, resources and results*. Encontro do DIP. URL: https://www.linguateca.pt/aval_conjunta/dip/apr_encontro/DIPpresentation.pdf.
- (2023). “DIP - Desafio de Identificação de Personagens: objetivo, organização, recursos e resultados”. In: *Linguamática*. to appear.
- Santos, Diana, Emanuel Pires, et al. (June 2020). “Periodização automática: Estudos linguístico-estatísticos de literatura lusófona”. In: *Linguamática* 12.1, pp. 81–95.
- Santos, Diana, Roberto Willrich, et al. (2022). “Identifying literary characters in Portuguese: Challenges of an international shared task”. In: *Computational processing of the Portuguese language, 15th International Conference, PROPOR 2022, Fortaleza, Brazil, March 21-23, 2022 Proceedings*. Ed. by Vlândia Pinheiro et al. Springer, pp. 413–419.
- Schöch, Christof, Tomaz Erjavec, et al. (2021). “Creating the European Literary Text Collection (ELTeC): Challenges and Perspectives”. In: *Modern Languages Open* 1, pp. 1–19. URL: <https://doi.org/10.3828/mlo.v0i0.364>.
- Schöch, Christof, Evgeniia Fileva, and Julia Dudar (Feb. 2022). *CLS INFRA D3.1 Baseline Methodological User Needs Analysis*. DOI: [10.5281/zenodo.6389333](https://doi.org/10.5281/zenodo.6389333). URL: <https://doi.org/10.5281/zenodo.6389333>.
- Schulz, Daniel and Štěpán Bahník (2019). “Gender associations in the twentieth-century English-language literature”. In: *Journal of Research in Personality* 81, pp. 88–97. ISSN: 0092-6566. DOI: <https://doi.org/10.1016/j.jrp.2019.05.010>. URL: <https://www.sciencedirect.com/science/article/pii/S0092656619300510>.
- Silva, Flávia Martins da Rosa Pereira da (2021). *Diferenciações de gênero na caracterização de personagens: uma proposta metodológica e primeiros resultados*.
- Smeets, Roel (2021). *Character Constellations: Representations of Social Groups in Present-Day Dutch Literary Fiction*. Leuven University Press. ISBN: 9789462702950. URL: <http://www.jstor.org/stable/j.ctv21wj5cb> (visited on 12/17/2022).
- Underwood, Ted (2019). *Distant Horizons: Digital Evidence and Literary Change*. University of Chicago Press.
- Underwood, Ted, David Bamman, and Sabrina Lee (2018). “The transformation of gender in English-language fiction”. In: *Journal of Cultural Analytics* 3.2, p. 11035.
- Weingart, Scott and Jeana Jorgensen (May 2013). “Computational analysis of the body in European fairy tales”. In: *Literary and Linguistic Computing* 28.3, pp. 404–416. ISSN: 0268-1145. DOI: [10.1093/llc/fqs015](https://doi.org/10.1093/llc/fqs015). eprint: <https://academic.oup.com/dsh/article-pdf/28/3/404/2784264/fqs015.pdf>. URL: <https://doi.org/10.1093/llc/fqs015>.