**ISSN 2499-4553** 



Italian Journal of Computational Linguistics

Rivista Italiana di Linguistica Computazionale

> Volume 6, Number 1 june 2020



editors in chief

Roberto Basili Università degli Studi di Roma Tor Vergata Simonetta Montemagni Istituto di Linguistica Computazionale "Antonio Zampolli" - CNR

advisory board

**Giuseppe** Attardi Università degli Studi di Pisa (Italy) Nicoletta Calzolari Istituto di Linguistica Computazionale "Antonio Zampolli" - CNR (Italy) Nick Campbell Trinity College Dublin (Ireland) Piero Cosi Istituto di Scienze e Tecnologie della Cognizione - CNR (Italy) Giacomo Ferrari Università degli Studi del Piemonte Orientale (Italy) Eduard Hovy Carnegie Mellon University (USA) Paola Merlo Université de Genève (Switzerland) John Nerbonne University of Groningen (The Netherlands) **Joakim Nivre** Uppsala University (Sweden) Maria Teresa Pazienza Università degli Studi di Roma Tor Vergata (Italy) Hinrich Schütze University of Munich (Germany) Marc Steedman University of Edinburgh (United Kingdom) **Oliviero Stock** Fondazione Bruno Kessler, Trento (Italy) Jun-ichi Tsujii Artificial Intelligence Research Center, Tokyo (Japan)

#### editorial board

Cristina Bosco Università degli Studi di Torino (Italy) Franco Cutuqno Università degli Studi di Napoli (Italy) Felice Dell'Orletta Istituto di Linguistica Computazionale "Antonio Zampolli" - CNR (Italy) **Rodolfo Delmonte** Università degli Studi di Venezia (Italy) Marcello Federico Fondazione Bruno Kessler, Trento (Italy) Alessandro Lenci Università degli Studi di Pisa (Italy) **Bernardo Magnini** Fondazione Bruno Kessler, Trento (Italy) **Johanna Monti** Università degli Studi di Sassari (Italy) Alessandro Moschitti Università degli Studi di Trento (Italy) Roberto Navigli Università degli Studi di Roma "La Sapienza" (Italy) Malvina Nissim University of Groningen (The Netherlands) **Roberto Pieraccini** Jibo, Inc., Redwood City, CA, and Boston, MA (USA) Vito Pirrelli Istituto di Linguistica Computazionale "Antonio Zampolli" - CNR (Italy) **Giorgio Satta** Università degli Studi di Padova (Italy) **Gianni Semeraro** Università degli Studi di Bari (Italy) **Carlo Strapparava** Fondazione Bruno Kessler, Trento (Italy) Fabio Tamburini Università degli Studi di Bologna (Italy) Paola Velardi Università degli Studi di Roma "La Sapienza" (Italy) Guido Vetere Centro Studi Avanzati IBM Italia (Italy) Fabio Massimo Zanzotto Università degli Studi di Roma Tor Vergata (Italy)

editorial office **Danilo Croce** Università degli Studi di Roma Tor Vergata **Sara Goggi** Istituto di Linguistica Computazionale "Antonio Zampolli" - CNR **Manuela Speranza** Fondazione Bruno Kessler, Trento Registrazione presso il Tribunale di Trento n. 14/16 del 6 luglio 2016

Rivista Semestrale dell'Associazione Italiana di Linguistica Computazionale (AILC) © 2020 Associazione Italiana di Linguistica Computazionale (AILC)



Associazione Italiana di Linguistica Computazionale

direttore responsabile Michele Arnese

isbn PDF 9791280136404

Accademia University Press via Carlo Alberto 55 I-10123 Torino info@aAccademia.it www.aAccademia.it/IJCoL\_6\_1



# **IJCoL**

Volume 6, Number 1 june 2020

CONTENTS	
Editorial Note Roberto Basili, Simonetta Montemagni	7
Biodiversity in NLP: modelling lexical meaning with the Fruit Fly Algorithm <i>Simon Preissner, Aurélie Herbelot</i>	11
Building and Comparing Lemma Embeddings for Latin. Classical Latin versus Thomas Aquinas Rachele Sprugnoli, Giovanni Moretti, Marco Passarotti	29
Lost in Text: A Cross-Genre Analysis of Linguistic Phenomena within Text <i>Chiara Buongiovanni, Francesco Gracci, Dominique Brunato, Felice Dell'Orletta</i>	47
Towards Automatic Subtitling: Assessing the Quality of Old and New Resources <i>Alina Karakanta, Matteo Negri, Marco Turchi</i>	63
"Contro L'Odio": A Platform for Detecting, Monitoring and Visualizing Hate Speech against Immigrants in Italian Social Media Arthur T. E. Capozzi, Mirko Lai, Valerio Basile, Fabio Poletto, Manuela Sanguinetti, Cri stina Bosco, Viviana Patti, Giancarlo Ruffo, Cataldo Musto, Marco Polignano, Giovanni Semeraro, Marco Stranisci	i- 77

# Lost in Text: A Cross-Genre Analysis of Linguistic Phenomena within Text

Chiara Buongiovanni\* Università di Pisa

Dominique Brunato<sup>†</sup> Istituto di Linguistica Computazionale "Antonio Zampolli" (ILC–CNR) ItaliaNLP Lab Francesco Gracci\*\* Università di Pisa

Felice Dell'Orletta<sup>‡</sup> Istituto di Linguistica Computazionale "Antonio Zampolli" (ILC–CNR) ItaliaNLP Lab

Moving from the assumption that formal, rather than content features, can be used to detect differences and similarities among textual genres and registers, this paper presents a new approach to linguistic profiling – a well-established methodological framework to study language variation – which is applied to detect significant variations within the internal structure of a text. We test this approach on the Italian language using a wide spectrum of linguistic features automatically extracted from parsed corpora representative of four main genres and two levels of complexity for each, and we show that it is possible to model the degree of stylistic variance within texts according to genre and language complexity.

### 1. Introduction

The combination of corpus-driven and Natural Language Processing (NLP)-based approaches to study language variation and use from a stylistic and sociolinguistic perspective has become an established line of research. The heart of this research is the so-called 'linguistic profiling', a framework of analysis in which a large number of counts of linguistic features extracted from linguistically annotated corpora are used as a text profile and can then be compared to average profiles of texts (or groups of texts) to identify those that are similar, at least similar in terms of the profiled features (van Halteren 2004; Montemagni 2013). Although it has been originally developed for authorship verification and recognition, i.e. to select one author from a set of known authors or to confirm or deny authorship by a single known author, this methodology has been successfully applied to research on genre and register variation, as well as in a variety of scenarios focused on modeling the 'form' rather than the content of texts: from the identification of developmental patterns in typical (Lu 2009; Lubetich and Sagae 2014) and atypical language acquisition (Prud'hommeaux et al. 2011; Rouhizadeh,

<sup>\*</sup> Dipartimento di Filologia Letteratura Linguistica Piazza Torricelli, 2, 56126, Pisa. E-mail: c.buongiovanni@studenti.unipi.it

<sup>\*\*</sup> Dipartimento di Filologia Letteratura Linguistica Piazza Torricelli, 2, 56126 Pisa. E-mail: f.graccil@studenti.unipi.it

<sup>†</sup> ItaliaNLP Lab (www.italianlp.it), ILC-CNR - Via Moruzzi, 1 - 56124, Pisa, Italy. E-mail: dominique.brunato@ilc.cnr.it

<sup>‡</sup> ItaliaNLP Lab (www.italianlp.it), ILC-CNR - Via Moruzzi, 1 - 56124, Pisa, Italy. E-mail: felice.dellorletta@ilc.cnr.it

Sproat, and van Santen 2015) and in school learners' writing (Barbagli et al. 2015), to the detection of linguistic markers of adult cognitive impairments (Roark, Mitchell, and Hollingshead 2007); besides, from a computational sociolinguistics perspective, for studying variations related to the social dimension of language (Nguyen et al. 2016) or for modeling stylometric characteristics of authors or author groups (Daelemans 2013).

The assumptions of linguistic profiling lay in the framework of the Multi-Dimensional Analysis (MDA) pioneered by Douglas Biber, a text-linguistic approach to characterize language use across social and communicative contexts through the quantitative and functional analysis of co-occurrence patterns of linguistic features and underlying dimensions of language use (Biber, Conrad, and Reppen 1998). The important hallmark of MDA is that "linguistic features from all levels function together as underlying dimensions of variation, with each dimension defining a different set of linguistic relations among registers" (Biber 1993). This set of features should capture properties of text that are 'signatures' of its style rather than the topic it deals with, since they are more appropriate to investigate differences and similarities between texts from a functional perspective. Since its foundation, MDA has also aimed to be "computerbased in that it depends on automated analyses of linguistic features in texts. This characteristic enables distributional analysis of many linguistic features across many texts and text varieties". However, as Argamon observes in its recent survey, traditional research in computational register analysis has exploited a limited number of stylistic features, such as the relative frequencies of function words, taken as indicative of different grammatical choices rather than topic ones, or of character n-grams assumed to model linguistic variation in lexical, grammatical, and orthographic preferences (Argamon 2019). This is easily explainable as these features are simpler to extract from text than more sophisticated syntactic and discourse features, yet effective enough to capture dimension of variations. In fact, if function words require language-specific lists of a few hundred words (such as pronouns, prepositions, auxiliary and modal verbs, conjunctions, and determiners), character n-grams do not need any form of languagespecific preprocessing, not even tokenization, which instead might be necessary for word n-grams.

More recently, the development of robust and fairly accurate NLP pipelines together with the increased computational power to process large volumes of data has allowed to automatize the process of feature extraction from large-scale corpora, enhancing the potential contribution of linguistic profiling for studying language variation. By modeling the 'form' of a text through large sets of features spanning across distinct levels of language description, it has been possible not only to improve automatic classification of genres (Stamatatos, Fakotakis, and Kokkinakis 2001), but also to get a better understanding of the impact of those features in classifying genres and text varieties (Cimino et al. 2017; Finn and Kushmerick 2006). This paper adopts this framework but, quite differently from much previous research, it presents a new application of linguistic profiling, in which the unit of analysis is not the document as a whole entity, but the internal parts in which it is articulated. A close perspective has been pursued by (Crossley, Dempsey, and McNamara 2011) but with a different aim, i.e. the automatic discrimination of paragraphs with a specific rhetorical purpose in English students' essays by using a feature engineering approach. In our contribution, we focus on the Italian language and we broaden the scope of the analysis to four traditional textual genres and two levels of linguistic complexity for each. Briefly, this paper intends to answer the following research questions:

*i*) to what extent is NLP–based linguistic profiling a viable approach to characterize the internal structure of a text?

*ii)* does the variance across different parts of a text change according to genre and level of complexity?

In what follows we first present the corpora on which the study was carried out and then move to the presentation of the approach, together with the description of the linguistic features used for analysis. In Section 5, we discuss the main findings we obtained and outline the primary conclusions.

# 2. Corpora

Our investigation was carried out on four traditional genres: Journalism, Educational material, Scientific prose and Narrative. For each genre we selected the two corpora described in (Brunato and Dell'Orletta 2017), which represent a 'complex' and a 'simple' language variety for that genre, where the level of complexity was established according to the expected reader. Specifically, the journalistic genre comprises a corpus of articles published between 2000 and 2005 on the general newspaper la Republica and a corpus of easy-to-read articles from Due Parole, a monthly magazine written in a controlled language for readers with basic literacy skills or mild intellectual disabilities (Piemontese 1996). The corpus belonging to the Educational genre is articulated into two collections targeting high school (AduEdu) vs. primary school (ChiEdu) students. For the scientific prose, the 'complex' variety is represented by a corpus of 84 scientific articles on different topics, while the 'simple' one by a corpus of 293 Wikipedia articles, extracted from the Italian Portal 'Ecology and Environment'. For the Narrative genre, we relied on a dataset specifically developed for research on automatic text simplification. It consists of 56 short novels for children and pieces of narrative writing for L2 high school students arranged in a parallel fashion, i.e. for each original text a manually simplified version is available. For our investigation, the original texts and the corresponding simplified versions were chosen as representative of the complex variety and the simple variety of Narrative, respectively.

All these corpora contain documents which are very different in terms of length: for instance, scientific articles are on average longer than others (215 sentences per document) and this reflects the fact that the body part is more dense and possibly articulated into more middle paragraphs. On the contrary, the easy-to-read newspaper articles (i.e. *Due Parole*) are made of  $\sim$ 12 sentences. Thus, for each document, we also defined an internal subdivision into six parts intended to enable the linguistic profiling investigation on the internal structure of text. The rationale of the splitting approach is outlined in the following section. As a result, we ended up with six sections per document, where each section is composed by a number of sentences that depends on the average document length, ranging from two sentences per section for the shortest documents, to  $\sim$ 35 for the longest ones. Moreover, according to this approach, documents shorter than six sentences were discarded. As a result of the whole process, the final corpus is made of 1,168 documents. Table 1 reports general statistics about the final dataset in terms of: number of documents in each corpus, corpus size in number of tokens, average number of sentences per document and average number of sentences per section in each corpus.

### 3. Methodology

As a first step we focus on inspecting differences and similarities across all corpora considering the document as the unit of analysis. This is a more traditional perspective

Genre	Corpus	N. doc	Tokens	Sent/doc	Sent/sec
Educational	c High-school texts (AduEdu)	69	47.854	23.4	3.9
Educational	s Primary school texts ( <i>ChilEdu</i> )	52	22.382	21.0	3.5
Journalism	c la Repubblica articles ( <i>Rep</i> )	304	230.789	30.5	5.1
	s Due Parole articles (2Par)	303	71.228	12.7	2.1
Namating	c Terence&Teacher original (TT orig	) 53	25.931	25.3	4.2
Nallauve	s Terence&Teacher simplified (TT simp)	54	23.866	25.5	4.3
Scientific prose	c Scientific articles (ScientArt)	84	471.883	215.6	35.9
Scientific prose	s Wikipedia articles (WikiArt)	249	200.681	29.1	4.9

#### Table 1

Corpus statistics ("c" and "s" stand for 'complex' and 'simple' variety for each genre.)

to address linguistic profiling and it is useful to provide quantitative and quality data on the average variation across the examined genres and internal level of complexity.

To carry out this analysis all corpora were firstly automatically tagged by the partof-speech tagger described in (Dell'Orletta 2009) and dependency parsed by the DeSR parser (Attardi et al. 2009). DeSR, trained on the ISST-TANL treebank (Montemagni et al. 2003) consisting of articles from newspapers and periodicals, achieves a performance of 83.38% and 87.71% in terms of Labeled (LAS) and Unlabeled Attachment Score (UAS) respectively, when tested on texts of the same type. However, it is widely acknowledged that even state-of-the art parsers have a drop of accuracy when tested against corpora differing from the typology of texts on which they were trained (Gildea 2001). Therefore, we can assume that the performance of DeSR would be possibly lower when parsing texts of a different textual genre, such as narrative or scientific writing. Despite this fact, we also expect that the distributions of parsing errors will be almost similar, at least when analysing texts of the same domain and language variety, thus allowing us to carry out an internal comparison with respect to examined linguistic parameters. Besides, the effect of genre variation on the performance of general-purpose parsers is likely to be less strong since all genres here considered contain standard texts, i.e. texts linguistically similar to the ones used in training.

Based on the output of the different levels of linguistic annotation, we automatically extract from text a wide set of linguistic features, thus creating a feature-vector representation of each document, where each dimension of the vector corresponds to the average value of a given linguistic feature in the document. The set of features used in our study, and the motivation underlying their choice, are described in the following section.

For what concerns the second analysis, which is the most innovative perspective of this study, we looked at the internal structure of documents and we investigate how the same linguistic features vary across different sections of text. To allow this investigation we firstly defined the 'new' unit of analysis as follows. All documents were split into a fixed number of sections, where each section is composed by a certain number of paragraphs, roughly corresponding to the three main parts of the rhetorical structure of a text (i.e. introductory, body and concluding paragraphs). According to previous work, for some genres such as academic writing, the distinction into paragraphs is quite rigid and follows the so-called 'five-paragraphs' format (Crossley, Dempsey, and McNamara 2011) which adheres to the rhetorical goals of the document, i.e. the first and the last paragraph correspond respectively to the introduction and the conclusion, and the three middle ones to the body part. However, based on a preliminary investigation of our corpora showing that the average document length is highly variable, we preferred to define a six-section subdivision in order to avoid flattening too much the distinctions across genres.

Finally, we performed the feature extraction process, this time representing each section of a document as a vector of features, whose values correspond to the average value that each linguistic feature has in all sentences included in the section.

To understand whether and to what extent the different sections of a text represent distinctive varieties with a peculiar linguistic structure, we carried out two statistical analyses. The first one aimed to assess the significance of variation between the same features extracted from different sections. Specifically, we performed a pairwise comparison between each section and the following one (i.e. 1/2, 2/3, 3/4 etc), as well as between the first and the last section (i.e. 1/6). The latter was aimed at verifying whether our set of features alone is able to distinguish between the introductory and the closing part of a document, the two most distant sections of a text which are supposed to have a more codified structure. Secondly, we evaluated whether there is a correlation between the values of features in the two sections under comparison. For both analyses, all data were calculated across and within genre. The cross-genre analysis was focused on genre only, thus collapsing the internal distinction in terms of complexity and considering the two corpora as a unique one for each genre. In the within-genre condition, the two corpora were kept distinct thus allowing us to observe whether there is an effect of genre that is preserved despite changes in linguistic complexity.

### 4. Linguistic Features

The set of features used for our analysis models a variety of phenomena related to the sentence structure, with a particular focus on morpho–syntactic and syntactic properties. These features were selected since they proved to be effective predictors of systematic variations in automatic genre classification (Cimino et al. 2017), as well as in other tasks in which the 'form' of the text is more relevant than the content, such as the prediction of perceived sentence complexity by humans (Brunato et al. 2018), the assessment of text readability (Collins-Thompson 2014) or the identification of the native language of speakers from their productions in a second language (Malmasi et al. 2017).

The features can be distinguished into three different categories, according to the level of annotation from which they derive.

**Raw Text Features**: they include the average word and sentence length (*word\_length* and *sent\_length* in Tables 2 and Tables 3), calculated as the number of characters per token and of tokens per sentence, respectively.

**Morpho-syntactic Features**: i.e. distribution of unigrams of part-of-speech distinguished into 14 coarse-grained pos tags (*cpos\_*) and the 37 fine-grained tags (*pos\_*) according to the ISST-TANL tagset.

**Syntactic Features**: features modeling syntactic phenomena of different types, i.e.: - the *distribution of syntactic dependency types* (*dep\_*), e.g. subject, direct object, modifiers, calculated as the distribution of each typed out of the total dependency relations according to the ISST-TANL dependency tagset;

- the *length of dependency links*, i.e. the average length of all dependency links (*avg\_links\_l*) and of the longest link (*max\_links\_l*). For each link, the distance is calculated as the number of words occurring between the syntactic head and the dependent;



Figure 1



- the order of constituents with respect to the relative lexical head: this feature works as a proxy of canonicity effects and it is calculated for the main elements of the clause in terms of: the distribution of pre-verbal and post-verbal subjects (*pre\_subj*, *post\_subj*), of pre-verbal and post-verbal object (*pre\_obj*, *post\_obj*), of pre-nominal and post-nominal adjectives (*pre\_adj*, *post\_adj*), and of pre-verbal and post-verbal adverbs (*pre\_adv*, *post\_adv*);

- the *parse tree structure*, in terms of features corresponding to: the average depth of the whole parse tree (*parse\_depth*) (i.e. the longest path from the root of the dependency tree to a leaf); the average width of the parse tree (*parse\_width*), where the width is measured as the average number of nodes placed on the same level; the average number of dependents for all heads in the sentence (*avg\_dep\_all*), for the verbal heads and for the nominal heads (*avg\_dep\_verb, avg\_dep\_noun*);

- *subordination features*: a thorough analysis was devoted to investigate the use of subordination by computing: the average distribution of subordinate clauses with respect to the main clause (*sub\_main*) and of embedded subordinate clauses (i.e. subordinate clauses dependent on other subordinate clauses) out of all subordinate clauses (*sub\_minor*). In addition, both for the 'superordinate' subordinate clause and the embedded ones, it is calculated the relative order with respect to the clause on which they depend (*pre\_sub, post\_sub*), as well as the average depth (*subord\_depth*) and the average width (*subord\_width*) of the parse tree generated by the subordinate clause.

To exemplify some of the above-described features, we refer to the following sentence extracted from the Narrative corpus, whose graphical visualization is reported in Figure 1.

(1) Come girò pagina restò senza fiato perché c'era una vecchia fotografia della nave che aveva visto la sera prima.

The sentence has a length of 21 tokens (punctuation included) and the average word length is 4.62 characters. For what concerns the POS distribution, it presents e.g. an equal distribution of verbs and nouns, each one corresponding to the 23.81% of the whole POS. With respect to syntax–related features, the average length of dependency links is 1.84 and the maximum link (excluded punctuation) is 5–token long (corresponding to the *mod\_rel* dependency going from the nominal head 'fotografia' to the embedded verb of the relative clause 'visto'). The maximum parse tree depth is 6, corresponding to the longer path, in terms of number of intervening nodes, from the root of the sentence ('restò') to the more distant leaf, the definite article 'la'.

# 5. Data Analysis

### 5.1 Differences and similarities across genres and complexity

Table 2 reports the average value and standard deviation of all examined linguistic features in the four corpora considered as a whole, as well as considering the 'simple' and 'complex' sub-corpora as distinct varieties. Along with the mean and standard deviation value, we also compute the coefficient of variation, which corresponds to the ratio between the standard deviation and the mean value of the feature. This measure allows us to evaluate the dispersion of values around the average in a standardized way, that is comparing the stability of features pertaining to data measured on different scales. The assumption in exploiting this metric is that the more stable a feature in a given corpus, the more meaningful it is for characterizing it. Of course, we expect that features with the lowest coefficient of variation identify general tendencies of a language, in this case of Italian. However, the effect of genre and language variety can be inferred by the different positions that these 'most stable' features have in the ranking of each corpus. With this respect, in Table 2 we also indicate the first ten features ranked according to the coefficient of variation for each corpus.

A first comparison across corpora based on the average feature value shows some expected tendencies, such as the highest sentence and word length characterizing scientific prose, also in the simple variety. Both these features are considered as raw indicators of sentence complexity and they are used by traditional indexes to evaluate the readability of a text (Collins-Thompson 2014). Focusing on the relative frequencies of core parts of speech, which is traditionally used as a marker for text genres (Biber, Conrad, and Reppen 1998), we can also note that scientific and journalistic texts use more nouns and less verbs in comparison to other genres, a marker of the nominal style featuring these texts. The opposite holds true for the narrative prose, which indeed has the lowest noun/verb ratio. Note that with the only exception of sentence length, all these parameters occur among the ten most stable ones in almost all genres and varieties, with the average word length always in the top two positions. The distribution of nouns is much prone to variation, especially in Narrative.

If we consider features modeling the syntactic profile of text, we can observe that Scientific articles again display values typically reported in readability studies to describe difficult-to-read texts, such as deeper parse tree. If this is particularly evident in the complex variety of this genre (where the average parse tree depth is 8.71), we can observe that also the simple one (i.e. Wikipedia articles) shows values that are much more higher than those reported by the simple varieties of all other genres. However, if scientific prose appears as the most difficult genre when sentence length and parse tree height are taken into account, it presents e.g. less complex verbal predicates in terms of number of dependents for verbal head (*avg\_dep\_verb*). With respect to this feature, Educational texts are those with the richest predicate structures, both in the easy and complex variety (2.16 and 2.25 respectively).

### 5.2 Studying the internal structure of text

We now focus on discussing the results of linguistic profiling carried out for the distinct sections of corpora, which were extracted according to the procedure described in Section 3. More specifically, to characterize the degree of variation within the different parts, and to evaluate how this variance derives from either genre-related or complexity-related features, we calculated two different statistical tests: *i*) the Wilcoxon

#### Table 2

An extract of linguistic features used for linguistic profiling across genres and complexity levels. For each feature it is reported the average value and (standard deviation) in the whole corpus (g), in the 'simple' (s) and in the 'complex' subcorpora (c). Values marked with exponent <sup>n</sup> mean that the corresponding feature is more stable for that corpus according to the coefficient of variation. The exponent index indicates the rank that these features have in the top ten positions of the ranking derived by the coefficient of variation.

faaturas	J	ournalisn	1	Scientific Prose					
leatures	g	s	с	g	s	с			
		Raw tex	t features						
a an t-lan a th	22.75	18.98	26.51	27.75	26.29	32.05			
sent_length	(11.39)	(7.65)	(13.13)	(13.63)	(13.96)	(11.62)			
word longth	$4.61^{2}$	$4.58^{2}$	$4.64^{2}$	$5.07^{1}$	$5.05^{1}$	5.13 <sup>1</sup>			
word_length	(0.5)	(0.52)	(0.47)	(0.62)	(0.64)	(0.57)			
	Mo	orpho-syn	tactic feat	ures					
	6.21	6.1	6.32	8.63	8.58	8.8			
cpos_ADJ	(4.59)	(4.77)	(4.4)	(5.07)	(5.59)	(3.09)			
	4.47	3.75	5.19	3.55	3.57	3.49			
cpos_ADV	(3.87)	(4.05)	(3.55)	(2.86)	(3.14)	(1.8)			
cros CONI	3.49	3.52	3.46	3.21	3.22	3.19			
cpos_conj	(2.93)	(3.32)	(2.49)	(2.32)	(2.51)	(1.63)			
cros PREP	15.17	14.83	15.5	$15.2^{9}$	$14.64^{9}$	16.06			
	(5.56)	(5.94)	(5.13)	(5.63)	(5.92)	(4.52)			
nos PROPN	5.44	5.49	5.39	3.99	4.48	2.54			
	(6.12)	(6.39)	(5.19)	(8.07)	(9.04)	(3.59)			
pos NOUN	$22.53^{7}$	$23.51^{5}$	$21.56^{7}$	$26.04^4$	$26.07^{5}$	$25.94^2$			
	(6.3)	(6.32)	(6.14)	(7.3)	(7.9)	(4.93)			
pos VERB	10.93	11.29	10.57	8.08	7.98	8.37			
	(4.52)	(4.8)	(4.18)	(3.76)	(3.96)	(3.06)			
		Syntacti	c features						
den dobi	4.18	4.91	3.45	2.43	2.38	2.58			
	(3.33)	(3.87)	(2.48)	(1.94)	(2.14)	(1.2)			
den subi	5.19	5.96	4.41	3.51	3.61	3.22			
ucp_subj	(3.15)	(3.53)	(2.51)	(2.22)	(2.44)	(1.86)			
dep mod	16.7	16.51	16.9	18.81 <sup>8</sup>	18.359	$20.18^{7}$			
	(6.63)	(7.35)	(5.81)	(7.15)	(7.62)	(5.29)			
avg dep noun	1.195	$1.19^{8}$	1.195	1.25	$1.19^{4}$	1.245			
0= 1 I =	(0.32)	(0.33)	(0.3)	(0.34)	(0.35)	(0.3)			
avg_dep_verb	$2.02^{\circ}$	2.1'	1.93°	$1.79^{10}$	$1.78^{10}$	1.81			
	(0.61)	(0.58)	(0.62)	(0.71)	(0.75)	(0.61)			
avg_dep_all	$(0.92^{\circ})$	$(0.92^{\circ})$	$(0.91^{\circ})$	$(0.18)^{-1}$	$(0.85^{2})$	$0.86^{\circ}$			
	0.03)	(0.05)	(0.07)	(0.16)	(0.10)	(0.19)			
avg_links_l	(0.49)	(0.49)	(0.47)	$2.04^{\circ}$ (0.67)	(0.71)	(0.52)			
	6 87 <sup>10</sup>	6 26 <sup>6</sup>	7 510	7.837	7 566	8 71 <sup>8</sup>			
parse_depth	(2.24)	(1.71)	(2.51)	(2.6)	(2.6)	(2.41)			
	1 959	<u> </u>	5.08 <sup>9</sup>	5 27 <sup>9</sup>	5 18 <sup>8</sup>	5 56			
parse_width	(1.58)	(1.45)	(1.68)	(2.02)	(2.14)	(1.57)			
	36.396	37 519	35 296	26.77	27.12	25 74 <sup>9</sup>			
link_pre	(10.04)	(10.46)	(9.48)	(11.29)	(12.35)	(7.19)			
1.1.	62.82 <sup>3</sup>	62.2 <sup>3</sup>	63.44 <sup>3</sup>	65.95 <sup>3</sup>	65.92 <sup>3</sup>	66.06 <sup>3</sup>			
link_post	(10.17)	(10.49)	(9.8)	(16.09)	(17.01)	(12.99)			
	2.86	2.33	3.41	3.4	3.21	4.03			
subord_depth	(2.22)	(2.13)	(2.18)	(2.14)	(2.25)	(1.63)			
and and 110	1.59	1.3	1.89	1.86	1.79	2.11			
subora_width	(1.29)	(1.28)	(1.23)	(1.24)	(1.33)	(0.9)			
auh main	47.29	45.5	49.5	47	45.89	50.47			
sub_main	(32.76)	(37.25)	(27.29)	(27.86)	(30.23)	(18.61)			

fasturas		Narrative		Educational							
leatures	g	s	С	g	s	с					
Raw text features											
cont lon oth	18.99	17.62	20.4	28.02	22.79	31.96					
sent_length	(7.69)	(6.2)	(8.92)	(13.25)	(10.52)	(13.73)					
word longth	$4.27^{1}$	$4.22^{1}$	4.33 <sup>2</sup>	$4.55^{2}$	$4.33^{2}$	$4.72^{2}$					
word_length	(0.43)	(0.42)	(0.44)	(0.53)	(0.5)	(0.5)					
Morpho-syntactic features											
	6.2	6.1	6.29	8.01	6.83	8.91					
cpos_ADJ	(4.7)	(4.76)	(4.66)	(4.29)	(4.08)	(4.24)					
cros ADV	6.52	6.5	6.54	6.04	6.03	6.05					
	(4.23)	(4.26)	(4.21)	(3.8)	(4.21)	(3.47)					
cpos_CONJ	5.24	5.31	5.16	4.56	4.17	4.86					
	(2.97)	(2.96)	(2.98)	(2.61)	(2.76)	(2.45)					
cpos_PREP	12.06	11.86	12.28 (4.74)	(4.7)	(5.3)	(4.1)					
	/ 99	5.28	(4.74)	1 78	2.3)	1.16					
pos_PROPN	(4.99)	(5.13)	(4.84)	(2.76)	(3.11)	(2.41)					
	$17.93^{9}$	$\frac{(0.10)}{18.07^9}$	17 789	21.335	20.857	$\frac{(2.11)}{21.69^4}$					
pos_NOUN	(5.93)	(5.93)	(5.94)	(5.21)	(6.1)	(4.4)					
1/EDD	14.37 <sup>10</sup>	14.510	14.24	12.119	13.189	11.31					
POS_VERB	(4.89)	(4.8)	(4.99)	(3.89)	(4.26)	(3.37)					
		Syntacti	c features								
dan dahi	4.53	4.66	4.39	3.78	3.93	3.67					
aep_aobj	(2.95)	(3.07)	(2.83)	(2.6)	(3.1)	(2.16)					
don subi	6.17	6.42	5.92	5.24	5.64	4.94					
	(3.22)	(3.26)	(3.17)	(2.66)	(2.9)	(2.42)					
dep mod	15.05	14.76	15.35	16.46	15.23	$17.4^{10}$					
uop_mou	(6.23)	(6.11)	(6.34)	(5.6)	(6.18)	(4.94)					
avg_dep_noun	1.09	1.06	$1.13^{10}$	$1.29^{8}$	$1.2^{10}$	1.35°					
0= 1=	(0.38)	(0.36)	(0.39)	(0.37)	(0.43)	(0.3)					
avg_dep_verb	2.1°	$2.12^{\circ}$	$2.09^{\circ}$	(0.54)	$2.16^{\circ}$	(0.52)					
	0.802	(0.30)	0.01	0.021	0.011	(0.32)					
avg_dep_all	(0.09)	(0.09)	(0.9)	(0.93)	(0.91)	(0.94)					
	$2.01^4$	$\frac{(0.07)}{1.97^4}$	$\frac{(0.00)}{2.06^4}$	$\frac{(0.00)}{2.22^4}$	$\frac{(0.00)}{2.08^6}$	2 335					
avg_links_l	(0.48)	(0.47)	(0.49)	(0.53)	(0.55)	(0.48)					
1 (1	6.287	6.025	6.5410	7.618	6.448	8.45					
parse_depth	(1.86)	(1.64)	(2.03)	2.56	(1.89)	(2.7)					
marca width	4.58	4.337	$4.68^{8}$	5.62	4.98	6.08					
parse_width	(1.37)	(1.22)	(1.49)	(2.02)	(1.89)	(2.02)					
link pre	$38.65^{6}$	38.92 <sup>8</sup>	38.39 <sup>6</sup>	36.56 <sup>7</sup>	$39.27^4$	$34.52^{8}$					
	(11.43)	(11.23)	(11.63)	(9.24)	(9.75)	(8.27)					
link post	59.81 <sup>3</sup>	59.29 <sup>3</sup>	$60.33^3$	$62.52^3$	$59.5^{3}$	$64.8^{3}$					
I	(12.11)	(12.04)	(12.17)	(9.53)	(10.14)	(8.36)					
subord_depth	2.79	2.56	3.03	3.57	2.81	4.2					
	1 54	(1.04)	(1.73)	2.19)	(1.00)	2 54					
subord_width	(1.05)	(1)	1.7	(1 4)	(1.04)	(1.41)					
	50.76	48.68	52.89	54 53	52.07	56.43					
sub_main	(28.92)	(29.41)	(28.31)	(28.6)	(31.87)	(25.22)					
	8.92	7	1.12	12.74	8.94	16.27					
sub_minor	(13.5)	(11.1)	(7.37)	(15.52)	(13.78)	(16.6)					

rank-sum test and *ii*) the Spearman's correlation test. As previously stated, each metric was calculated between the average value of features extracted from two consecutive sections and from the first and the last section of each document. Results of this comparison are reported in Table 3, which displays all features that turned out to have a statistically significant variation in at least one of the six pairwise comparisons, or a correlation score > 0.3 according to the Spearman's correlation coefficient.

As a first remark, we clearly notice that a higher number of features varying in a statistically significant way occurs in the journalistic and scientific genre, both considered as whole (i.e. row *g* for each feature) and with respect to the language complexity variety (rows *s* and *c*). The opposite trend is reported for texts of the Educational domain, which is probably due to the heterogeneous nature of this genre that includes documents of different textual typologies (course books, pieces of literature etc.).

If journalism and scientific prose are the two genres with the highest internal variance, the comparison between sections allows us to get a better understanding of this data. Specifically, for both genres, the majority of significant variations are observed between the first and the second section and between the first and the last one. This suggests that the introduction is a stylistic unit with a peculiar linguistic structure with respect to the body and the conclusion. For instance, it is characterized by shorter sentences (Figure 2), likely due to the presence of the title in both newspaper and scientific articles, and by a distinctive distribution of Parts-of-speech (Figure 3). With this respect, this data is consistent with other studies in the literature, e.g. (Voghera 2005), and also with what we observed in the global analysis of our corpora reported in paragraph 5.1, showing that scientific prose and newswire texts rely more on the nominal style. However, with the proposed approach, we were able to go further in this analysis, highlighting that noun/verb ratio is always higher in the first section than all other ones. Besides, at least for newspaper articles, this feature appears as a genre marker which is not affected by language complexity, since the same tendency is observed when the 'simple' and the 'complex' corpus are analyzed independently. The same does not hold for other features related to syntactic structure and, in particular, to the use of subordination. In this case, the 'shift' between the introduction and the subsequent part of texts yields significant variations only for articles of la Repubblica. Specifically, the first section contains less embedded sentences (*parse\_depth*: 1<sup>st</sup> section: 5.55; 2<sup>nd</sup> section: 7.76), and a lower presence of subordinate clauses, which appear as structurally simpler e.g. in terms of depth (subord\_depth: 1st section: 1.67; 2nd section: 3.5) and width (subord width: 1<sup>st</sup> section: 0.94; 2<sup>nd</sup> section: 1.97). Conversely, for the simple variant of this genre (i.e. the articles of the easy-to-read newspaper Due Parole), we do not observe significant changes affecting these features: this is not particularly surprising since subordination is always less represented in this corpus with respect to all the other ones (as reported in Table 2, the distribution of subordinate clauses with respect to the main clause (*sub\_main*) is 45.5.)

Leaving aside the similar tendencies characterizing the introduction, Journalistic and Scientific prose show a different behavior when we focus on the internal structure of text. While in this case much fewer features vary in a significant way, the majority occurs in the journalistic genre only, especially between the second and the third section. Again, they concern a different distribution of morpho-syntactic categories but also some syntactic features related to the hierarchical structure of the parse tree (e.g. *parse\_width, avg\_dep\_all*) and to the presence and level of embedding of subordinate clauses (*subord\_depth, subord\_width*). According to these data, we can conclude that the journalistic genre has a more rigorous structure and that it is possible to capture the





Average sentence length in the 6 sections across genres.



Figure 3

Distribution of lexical parts-of-speech in the four genres.

boundaries between different parts by using linguistic features that are not related to the content of the article.

#### Table 3

A set of linguistic features resulting as significant in at least one pairwise comparison.  $\checkmark \checkmark$  means highly statistically significant (p < 0.001),  $\checkmark$  statistically significant (p < 0.05), - no significance; \* correlation related to the Spearman's rank correlation coefficient (rho > 0.3), g=global corpus, s=simple variety of the corpus, c=complex variety of the corpus.

faaturas				Journa	alism				Sc	ientifi	c Pros	e	
leatures		1/2	2/3	3/4	4/5	5/6	1/6	1/2	2/3	3/4	4/5	5/6	1/6
Raw text features													
	g	$\checkmark\checkmark$	√ *	- *	- *	- *	$\checkmark\checkmark$	$\checkmark\checkmark$	-	-	- *	- *	$\checkmark\checkmark$
sent_length	s	$\checkmark\checkmark$	-	-	-	-	-	$\checkmark\checkmark$	-	-	-	-	$\checkmark\checkmark$
	с	$\checkmark\checkmark$	-	- *	- *	- *	$\checkmark\checkmark$	- *	- *	- *	- *	- *	- *
	g	- *	- *	- *	- *	- *	<ul> <li>✓</li> </ul>	✓	-	- *	- *	-	✓
word_length	s	- *	- *	- *	- *	-	$\checkmark\checkmark$	$\checkmark$	-	- *	-	-	$\checkmark\checkmark$
	с	-	- *	- *	- *	- *	-	- *	√ *	- *	- *	- *	- *
				Mor	pho-sy	ntacti	c featur	es					
	g	-	-	-	-	-	-	$\checkmark$	$\checkmark$	- *	-	-	*
cpos_ADJ	s	√	-	-	-	-	-	$\checkmark\checkmark$	-	- *	-	-	~
	С	<b>√</b>	-	-	-	-	<b>√</b>	-	- *	- *	- *	- *	- *
	g	<b>√√</b> *	√ *	-	-	-	~~	<b>√√</b> *	-	-	- *	-	<u> </u>
cpos_ADV	s	<u> </u>	-	-	-	-	<u> </u>	<b>√√</b> *	-	-	-	-	$\checkmark\checkmark$
	С	<u> </u>	- *	-	-	-	<u> </u>	- *	- *	- *	- *	- *	- *
601 H	g	<u> </u>	√	-	-	-	<u> </u>	<b>√</b> √	-	-	-	-	<u> </u>
cpos_CONJ	s	<u></u>	-	-	-	-	<u> </u>	<b>v v</b>	-	-	-	-	<u> </u>
	с	<u> </u>	-	-	-	-	<u> </u>	√ *	- *	- *	- *	- *	<u>√*</u>
NOUN	_ <u>g</u>	<b>√</b> √ *	√ *	- *	- *	- *	<u></u>	<b>√</b> √	-	-	- *	-	<u></u>
cpos_NOUN	s	<b>√</b> √ *	- *	-	-	- *	<u></u>	<b>v v</b>	-	-	-	-	~ ~
	С	V V *	<b>v v</b> *	- *	- *	- *	<b>V V</b>	- *	- *	- *	- *	- *	- *
THE PROPN	g	<b>√</b> √ *	- *	- *	- *	- *	<b>√</b> √ *	~ ~	- *	- *	√ *	- *	<u></u>
pos_proprin	s	<u>vv</u> *	- *	- *	- *	- *	<b>v v</b> *	<b>v v</b>	- *	- *	- *	- *	V V .
	c	••	~ ~	-	- *	-	••	-*	- *	- *	- *	- *	- *
mag VEPP	<u> </u>	~~	✓	-	-	-	<u></u>	V V	-	-	-	-	<u></u>
CPOS_VERD	<u>s</u>	<u> </u>	-	-	-	-	<u> </u>	V V	-	-	-	-	V V
	с а	<b>v v</b>	V V (	-	-*	- *	<u> </u>	-*	- *	- *	- *	- *	- *
DOS ALIX	<u> </u>	<b>v</b> *	<b>v</b> *	- *	-*	- *	v	<b>V V</b> *	-		-	-	<b>v v</b>
pos_AUX		- *	- *	- *	-*	- *	-	• •	-		- *	-	<u> </u>
		V V *	<b>v v</b> *	- *		- *	••	- *			- *		- *
					Syntac	ctic fea	tures						
dan dahi	_ <u>g</u>	V	<u> </u>	-	-	-	<u></u>	V V	-	-	-	-	<u></u>
dep_dobj	<u>s</u>	-	✓	-	-	-	<u></u>	V V	-	-	-	-	<u> </u>
	<i>c</i>	~ ~	-	-	-	-	~ ~	-*	- *	- *	- *	- *	- *
don subi	<u> </u>	-	-	-	-	-	-	<b>v v</b>	-	-	-	-	<u> </u>
dep_subj	<u>s</u>	-	-	-	-	-	-	V V	-	-	-	-	V V
	- c	•		- (*	- *		••	- *	- *	- *	- *	- *	- *
may links l	<u> </u>	<b>V V</b>		• *	- *		••	<b>V V</b>			- *		<b>v v</b>
IIIdx_IIIIK5_I		• •		-	- *	-	<b>v</b>	-*	- *	- *	- *	- *	- *
	σ	<u> </u>	-	• •	-	-	<u> </u>	<u> </u>	-	-	-	-	<u> </u>
ave links l	<u> </u>	<u> </u>	-	-	-	-	-	<u> </u>	-	-	-	-	<u> </u>
uvg_inuo_i	<u>c</u>	· · ·	-	~	-	-	<b>√</b> √	- *	- *	- *	- *	- *	-
	g	<u> </u>	- *	- *	- *	- *	<u> </u>	<i>√ √</i>	-	-	-	<b>√</b> ∗	<u></u>
parse depth	<u>-</u> 8	-	-	-	~	-			-	-	$\checkmark$	- *	~~
L moo-mol m	c	$\checkmark$	-	- *	- *	- *		- *	- *	- *	- *	- *	- *
	g	~~	$\checkmark$	-	-	-	~~	$\checkmark$	-	-	-	-	<b>\</b>
parse width	s	~~	-	-	-	-	-	$\sqrt{}$	-	-	-	-	~~
1	c	~~	-	$\checkmark$	- *	-	<b>√</b> √	- *	- *	- *	- *	- *	- *
	g	$\checkmark\checkmark$	$\checkmark$	- *	- *	- *	~~	$\checkmark\checkmark$	- *	- *	- *	√ *	$\checkmark\checkmark$
avg_dep_all	s	~	-	-	-	- *	-	$\checkmark\checkmark$	-	-	-	-	~~
0- 1-	с	<b>√</b> √	$\checkmark$	- *	- *	- *	~~	- *	- *	- *	- *	- *	- *
				Sul	bordir	nation	features						
	g	$\checkmark\checkmark$	$\checkmark$	-	-	-	<ul> <li></li> <li></li> </ul>	$\checkmark\checkmark$	-	-	-	$\checkmark$	$\checkmark\checkmark$
subord depth	<u> </u>	-	-	-	-	-		$\sqrt{}$	-	-	-	-	$\overline{\sqrt{}}$
ru	c	~~	-	-	-	-	~~	- *	- *	- *	- *	- *	- *
	g	~~	$\checkmark$	-	-	-	~~	$\checkmark\checkmark$	-	-	-	$\checkmark$	<b>√</b> √
subord width	s	-	-	-	-	-	√ 	$\checkmark$	-	-	-	-	V V
	c	~~	-	-	-	-	<b>√</b> √	- *	- *	-	- *	- *	- *
	g	~~	$\checkmark$	-	-	-	~~	$\checkmark\checkmark$	-	-	-	-	<b>√</b> √
sub_main	s	-	-	-	-	-	~~	$\checkmark\checkmark$	-	-	-	-	$\checkmark\checkmark$
-	с	$\checkmark\checkmark$	-	-	-	-	<b>√</b> √	- *	- *	- *	- *	- *	- *
	g	$\checkmark\checkmark$	- *	-	- *	-	$\checkmark\checkmark$	-	-	-	-	-	$\checkmark$
sub_minor	s	-	-	-	-	-	$\checkmark$	-	-	-	-	-	$\checkmark$
	с	$\checkmark\checkmark$	-	-	-	-	$\checkmark\checkmark$	- *	- *	-	- *	-	-

fasturas				Nar	rative					Educat	tional		
reatures		1/2	2/3	3/4	4/5	5/6	1/6	1/2	2/3	3/4	4/5	5/6	1/6
					Raw	text fe	atures						
	g	-	-	-	-	-	√*	- *	- *	- *	- *	- *	- *
sent length	<u>s</u>	- *	-	-	-	-	~	- *	- *	- *	- *	<b>√</b> *	- *
- 0	с	-	-	-	-	-	$\checkmark$	- *	-	- *	-	-	-
	g	-	- *	- *	-	-	-	- *	- *	- *	- *	- *	- *
word length	<u>s</u>	-	- *	- *	- *	-	-	<b>√</b> *	- *	- *	- *	- *	- *
	c	-	- *	- *	-	-	-	- *	- *	- *	- *	- *	- *
	-			Mc	rnho	watac	tic footu	roc					-
	~			WIC	ipno-s	syntac	lic leatu	lles					
	<u> </u>	-	-	-	-	-	-	-	- *	- *	-	- *	- *
cpos_ADJ	- 5	-	-	-	-	-	-	~ ~	-	-	-	-	- *
	c	-	-	-	- *	- *	- *	-	- *	- *	-	- *	-
mag ADV	<u> </u>	-	-	-	-	-	<u> </u>	-	-	- *	<u> </u>	-	-
cpos_ADV	- 5	-	-	-	-	-	v	v	-	- *	v	-	-
	C	-	-	-	-	-	-	-	-	- *	-	-	-
	g	-	-	-	√	-	-	-	<b>√</b>	√	-	-	-
cpos_CONJ	s	-	-	-	-	-	-	-	√	-	-	-	-
	с	-	- *	-	-	- *	-	-	-	-	-	-	-
NOURI	g	<u> </u>	-	-	- *	-	<b>√</b> √ *	-	- *	- *	- *	- *	- *
cpos_NOUN	s	✓	-	-	- *	-	$\checkmark$	$\checkmark$	- *	- *	-	-	- *
	С	√ *	-	-	- *	-	√ *	-	-	- *	- *	- *	-
	g	✓	- *	- *	- *	- *	- *	- *	- *	- *	- *	- *	- *
pos_PROPN	s	-	- *	- *	- *	- *	-	√*	- *	- *	- *	- *	- *
	С	-	- *	- *	- *	- *	-	- *	- *	- *	- *	-	- *
	g	$\checkmark$	-	-	-	-	-	-	-	-	$\checkmark$	$\checkmark$	-
cpos_VERB	s	$\checkmark$	-	-	-	-	-	-	-	-	$\checkmark$	$\checkmark\checkmark$	-
	С	$\checkmark$	-	-	-	-	-	-	- *	- *	-	-	-
	g	-	-	-	-	-	-	-	-	-	-	$\checkmark$	-
pos_AUX	s	-	-	- *	-	-	-	-	-	-	-	-	- *
	с	-	-	-	- *	-	-	-	-	-	-	-	-
					Synta	actic fe	atures						
-	g	✓	~	<b>√</b>		-	~	-	~	-	-	-	$\checkmark$
dep dobi	<u>-</u> s			- *	-	-	· ·	-		-	- *	-	· ·
		-	-	-	-	-	-	-	-	-	-	-	-
	0	1	-	-	-	-	-	-	- *	-	-	-	-
den subi	<u> </u>	-	- *	-	-	-	-	- *	- *	-	- *	-	-
ucp_subj	- <u>c</u>	~	-	-	-	-	-	-	-	-	-	-	-
	σ	-	-	-	-	-	- *	-	- *	- *	- *	- *	- *
max links l	<u> </u>	-	-	-	-	-	-	-	-	- *	-	- *	- *
Intax_III1K5_I		-	-	-	-	-	- *	-	- *	*	-	-	-
	α σ	-	-	-	-	-	-	-	-	-	-	- *	
ava linka l	<u> </u>				-			_		-		- *	
avg_mks_i		-	-	-	-	-	-	-	-	- *	-	- *	-
	с —	-	-	- *	-	-	- *	-	- *	-	-	-	-
and Janth	_ <u>g</u>	- *	-	- *	- *	- *	-	- *	- *	- *	- *	- *	- *
sent_depth	S	- *	-	-	- *	- *	-	-	- *	- *	- *	-*	<b>√</b> *
	C	-	-	- *	- *	-	-	- *	- *	-	- *	- *	- *
	g	-	-	-	-	-	-	- *	- *	-	-	- *	-
sent_width	s	-	-	-	-	-	-	- *	- *	-	-	- *	- *
	С	-	- *	-	-	-	- *	-	- *	√	- *	-	-
	g	$\checkmark$	- *	- *	-	- *	✓	- *	- *	- *	- *	- *	- *
avg_dep_all	s	-	-	- *	-	- *	✓	-	- *	- *	- *	- *	- *
	с	-	- *	- *	-	- *	√	- *	- *	√	- *	- *	- *
				S	ubordi	natior	n featur	es					
	g	- *	-	- *	- *	- *	-	-	-	- *	- *	-	-
subord_depth	s	-	-	-	- *	- *	$\checkmark\checkmark$	-	-	- *	- *	- *	$\checkmark$
- 1	с	-	-	-	- *	-	-	-	~	-	-	-	-
	g	-	-	-	-	-	-	-	-	- *	-	-	-
subord width	s	-	~	-	- *	- *	~	-	~	- *	- *	- *	-
	c	-	-	-	- *	-	-	-	-	-	-	- *	-
	g	-	-	-	-	- *	11	-	-	-	-	-	-
sub main	<u></u>	-	-	-	-	- *	· •	-	√ *	<b>√</b> ./*	- *	<b>√</b> *	-
sub_mam	- <u>s</u>	-	-	-	-	÷	-	-	• *	• • • •	÷	• ^ -	
	с 0	· ·	-						• - *	-			
sub minor	<u> </u>	-	-	-	-	v - *		-	- *	• -/	-		v
sub_mmor	<u> </u>	-	-	-		- *	- *	-	-	v		- *	- *
	C	v	-	-	- *	-	-		-	-	- *	-	v

# 6. Conclusion

In this paper we have presented a novel approach to the study of language variation, which relies on the prerequisites of the linguistic profiling methodology but with the specific purpose of modeling the stylistic form of the different parts within a text. A comparative investigation on four traditional genres in Italian, and two levels of complexity for each, showed that morpho-syntactic and syntactic features are differently distributed across subsections of texts representative of a given genre and language variety. From a linguistic perspective, this suggests that the study of genre and register variation can benefit by inspecting corpora from this fine-grained perspective. In this respect, we intend to carry out further analyses to prove the validity of this approach and how scalable it is when applied to novel texts of the same genres here considered. A first step in this direction is surely to assess the accuracy of our distinction into six parts, which was mainly driven by corpus-based considerations in terms of average document length. In particular, we would like to compare whether the findings obtained applying this splitting methodology are in line with those deriving form a manuallybased annotation, in which the sections are identified according to their structural coherence as perceived by readers.

We believe that a better understanding and computational modeling of linguistic phenomena characterizing the introductory, middle and conclusive parts of a text can also serve to enhance automatic genre classification, as well as a number of NLP-based applications devoted to modeling style. For instance, in the educational domain, it could be used as a part of intelligent tutoring systems able to provide detailed feedback to students in writing courses and to support the development of automatic selfassessment tools. Also in the field of Natural Language Generation, such a rich featurebased description of the internal profile of texts can be useful to inform the automatic generation of texts, which are not only semantically meaningful and coherent, but also compliant with the stylistic fingerprints of a specific genre and level of complexity.

### References

- Argamon, Shlomo E. 2019. Computational register analysis and synthesis. *Register Studies*, 1(1):100–135.
- Attardi, Giuseppe, Felice Dell'Orletta, Maria Simi, and Joseph Turian. 2009. Accurate dependency parsing with a stacked multilayer perceptron. In *Proceedings of EVALITA 2009 Evaluation of NLP and Speech Tools for Italian 2009*, Reggio Emilia, Italy, December 2009.
- Barbagli, Alessia, Pietro Lucisano, Felice Dell'Orletta, Simonetta Montemagni, and Giulia Venturi. 2015. Il ruolo delle tecnologie del linguaggio nel monitoraggio dell'evoluzione delle abilità di scrittura: primi risultati. Italian Journal of Computational Linguistics (IJCoL), 1(1):99–117.
- Biber, Douglas. 1993. Using register-diversified corpora for general language studies. *Computational Linguistics*, 19(2):219–242.
- Biber, Douglas, Susan Conrad, and Randi Reppen. 1998. *Corpus linguistics: investigating language structure and use.* Cambridge University Press, Cambridge.
- Brunato, Dominique and Felice Dell'Orletta. 2017. On the order of words in italian: a study on genre vs complexity. In *International Conference on Dependency Linguistics (Depling 2017)*, Pisa, Italy, 18-20 September 2017.
- Brunato, Dominique, Lorenzo De Mattei, Felice Dell'Orletta, Benedetta Iavarone, and Giulia Venturi. 2018. Is this sentence difficult? do you agree? In *Proceedings of the 2018 Conference on Empirical Methods of Natural Language Processing*, Brussels, Belgium, October-November.
- Cimino, Andrea, Martijn Wieling, Felice Dell'Orletta, Simonetta Montemagni, and Giulia Venturi. 2017. Identifying predictive features for textual genre classification: the key role of syntax. In *Proceedings of 4th Italian Conference on Computational Linguistics (CLiC-it)*, Rome, 11-13 December 2017.

- Collins-Thompson, Kevyn. 2014. Computational assessment of text readability. *Recent Advances in Automatic Readability Assessment and Text Simplification. Special issue of International Journal of Applied Linguistics*, 165(2):97–135.
- Crossley, Scott A., Kyle B. Dempsey, and Danielle S. McNamara. 2011. Classifying paragraph types using linguistic features: Is paragraph positioning important? *Journal of Writing Research*.
- Daelemans, Walter. 2013. Explanation in computational stylometry. In *Proceedings of the International Conference on Computational Linguistics and Intelligent Text Processing*, pages 451–462. Springer Berlin Heidelberg, march.
- Dell'Orletta, Felice. 2009. Ensemble system for part-of-speech tagging. In *Proceedings of EVALITA* 2009 *Evaluation of NLP and Speech Tools for Italian* 2009, Reggio Emilia, Italy, December 2009.
- Finn, Aidan and Nicholas Kushmerick. 2006. Learning to classify documents according to genre. Journal of the Association for Information Science and Technology, 57:1506–1518.
- Gildea, Daniel. 2001. Corpus variation and parser performance. In *Proceedings of Empirical Methods in Natural Language Processing (EMNLP 2001)*, Pittsburgh, PA, USA, June.
- Lu, Xiaofei. 2009. Automatic measurement of syntactic complexity in child language acquisition. *International Journal of Corpus Linguistics*, 14(1):3–28.
- Lubetich, Shannon and Kenji Sagae. 2014. Data-driven measurement of child language development with simple syntactic templates. In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 2151–2160, Dublin, Ireland, August.
- Malmasi, Shervin, Evanini Keelan, Aoife Cahill, Joel Tetreault, Robert Pugh, Christopher Hamill, Diane Napolitano, and Yao Qian. 2017. A report on the 2017 native language identification shared task. In Proceedings of the 12th Workshop on Building Educational Applications Using NLP, held at EMNLP 2017, Copenhagen, Denmark, September.
- Montemagni, Simonetta. 2013. Tecnologie linguistico-computazionali e monitoraggio della lingua italiana. *Studi Italiani di Linguistica Teorica e Applicata (SILTA)*, XLII(1):145–172.
- Montemagni, Simonetta, Francesco Barsotti, Marco Battista, Nicoletta Calzolari, Ornella Corazzari, Alessando Lenci, Antonio Zampolli, Francesca Fanciulli, Maria Massetani, Remo Raffaelli, Roberto Basili, Maria Teresa Pazienza, Dario Saracino, Fabio Zanzotto, Nadia Mana, Fabio Pianesi, and Rofolfo Delmonte. 2003. Building the italian syntactic-semantic treebank. In A. Abeille, editor, *Building and Using syntactically annotated corpora*.
- Nguyen, Dong, A. Seza Doğruöz, Carolyn P. Rosé, and Franciska de Jong. 2016. Survey: Computational sociolinguistics: A Survey. *Computational Linguistics*, 42(3):537–593.
- Piemontese, Maria Emanuela. 1996. *Capire e farsi capire. Teorie e tecniche della scrittura controllata.* Tecnodid, Napoli.
- Prud'hommeaux, Emily T., Brian Roark, Lois M. Black, and Jan van Santen. 2011. Classification of atypical language in autism. In *Proceedings of the 2nd Workshop on Cognitive Modeling and Computational Linguistics, held at ACL HLT*, pages 88–96, Portland, Oregon, USA, June.
- Roark, Brian, Margaret Mitchell, and Kristy Hollingshead. 2007. Syntactic complexity measures for detecting mild cognitive impairment. In *Proceedings of the Workshop on BioNLP 2007: Biological, Translational, and Clinical Language Processing, held at ACL 2007, pages 1–8, Prague,* Czech Republic, June.
- Rouhizadeh, Masoud, Richard Sproat, and Jan van Santen. 2015. Similarity measures for quantifying restrictive and repetitive behavior in conversations of autistic children. In *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality, Association for Computational Linguistics*, pages 117–123, Denver, Colorado, June.
- Stamatatos, Efstathios, Nikos Fakotakis, and George Kokkinakis. 2001. Automatic text categorization in terms of genre and author. *Computational Linguistics*, 26:471–495.
- van Halteren, Hans. 2004. Linguistic profiling for author recognition and verification. In *Proceedings of the Association for Computational Linguistics (ACL04)*, pages 200–207, Barcelona, Spain, July.
- Voghera, Miriam. 2005. La misura delle categorie sintattiche. In *Chiari Isabella / De Mauro Tullio (eds.) Parole e numeri. Analisi quantitative dei fatti di lingua*. Aracne, Roma, pages 125–138.