

TEXTUAL DIMENSIONS IN BIOLOGICAL SCIENCE CORPORA

Grisel María García Pérez¹

ggarciaperez@ouc.bc.ca

1. Introduction

Defining the distinctive features of the variety of English used in scientific contexts has been a special tendency in the past few years (Barber, 1962; Ewer and Latorre, 1969; Swales, 1984; Halliday, 1988; Reid, 1991, Harmer, 2001). This is so because English has emerged as the predominant medium of scientific discussion and progress, hence theoretical and practical applications on the teaching of English have become a powerful need in all parts of the world.

Genres introduce certain stability into a discourse community and are flexible enough to participate in social changes, so from this point of view they function as language itself. Because of this, they have become key points in some investigations carried out in English for Special Purposes.

Swales' analysis of genre (1990) has served as a reference for different studies on the teaching of ESP using a genre-based approach (Widdowson, 1983; Crookes, 1986; Marshal, 1991; Nwogu, 1991). Adopting strategies similar to those embodied in schema-theoretic models, Swales posits a four 'move' schema for article introduction in ESP courses and specifically for scientific discussions. His study demonstrates not only an attempt to chunk texts into identifiable knowledge structures, but a concern with characterizing the linguistic features of each move and the means by which information in the move is signalled.

¹ Modern Languages Department, Okanagan University College, BC. Canada

Although the previous studies are grounded on a firm basis and offer a rationale for genre analysis, Biber (1988) offers a model in which texts can be compared along dimensions of linguistic variation. This is one of the most sophisticated studies on genre differences that has been published so far. By computing factor scores, that is, by summing up the frequency of each of the linguistic features in a factor for each text, he was able to average the factor score for each text across all texts in a genre and compute a mean dimension score for the genre. He then used these mean dimension scores to compare and to specify the relations among genres.

The utility of genre analysis in the teaching of English as a Second Language has been revealed in our literature review (Widdowson, 1983; Martin, 1985; Crookes, 1986 Biber, 1988; Swales, 1990). However, few of these studies offer hard data on specific fields such as medicine, physics, biology and math. Because Biber's study provides a foundation for cross-linguistic research, the present paper aims at pointing out what linguistic features are shared specifically by texts in the field of biological science and compare the findings to the general science corpora that have been described in Biber's analyses.

Biber (1988) compares texts along 'dimensions' of linguistic variations. He states that researchers have found out that texts are related along particular situational or functional parameters, e.g. formal and informal, interactive and non-interactive, literary and colloquial, or restricted and elaborated. These parameters can be considered as dimensions because 'they define continuums of variations rather than discrete poles.

In his work he uses frequency counts of particular linguistic features as a means to give exact quantitative characterization of a text; however these counts do not identify linguistic dimensions. Linguistic dimensions are characterized on the basis of a consistent co-occurrence pattern among features; that is, the consistent co-occurrence of a cluster of features in texts define a linguistic dimension.

The approach used by Biber in his study completely differs from previous studies. Other studies began with a situational or functional distinction and afterwards identified linguistic features associated with that distinction; Biber identifies the clusters of features in terms of shared function, but without necessarily representing a linguistic

dimension. Biber uses quantitative techniques to identify the groups of features and then interprets them in functional terms. The linguistic rather than the functional dimension is given priority.

He bases this approach on the idea that "if certain features consistently co-occur, then it is reasonable to look for an underlying functional influence that encourages their use".

Once the linguistic co-occurrence patterns are identified, the resulting dimensions can be interpreted in functional terms. The approach moves from stating WHAT features co-occur to explaining the WHY of their co-occurrence.

After identifying and interpreting the linguistic dimensions, they can be used to specify 'textual relations'. Textual relations are defined by a simultaneous comparison of the texts with respect to all dimensions.

So far, researchers have investigated linguistic textual variations using either a microscopic or a macroscopic analysis or a combination of the two (Schiffrin, 1981; Besnier, 1983; Biber, 1988). Microscopic analysis identifies the linguistic features and genre distinctions to be included in a macro analysis, and provides a functional analysis of the features, so as to be able to interpret the textual dimensions in functional terms. Macroscopic analyses pinpoint the underlying textual dimensions in a set of texts, enable the general description of a general account of linguistic variations among texts, and provide a framework for the discussions of the similarities and differences among texts and genres.

This paper presents the results of both micro and macro analysis:

- a) A macroscopic outlook to analyze the co-occurrence patterns among ten linguistic features in 24 texts, identifying two textual dimensions; and
- b) A microscopic analysis to identify the features and to interpret the dimensions in functional terms.

Biber identifies six textual dimensions in his study:

- a) Involved versus Informational Production,
- b) Narrative versus Non-Narrative Concerns,
- c) Explicit versus Situation-Dependent Reference,

- d) Overt Expression of Persuasion,
- e) Abstract versus Non-Abstract Information, and
- f) On-line Informational Elaboration.

From these six textual dimensions, relevant salient loadings were reported in the analysis of the general science corpora in 'Explicit versus situation dependent' and 'Abstract versus Non-abstract information'. These two textual dimensions will be key points for main objective of this paper which is to compare the results of the analysis between our biological science corpora and the scientific texts included in Biber's study.

2. Method

2.1 Corpus

This study is based on a corpus of 24 texts from the field of biology selected from the book *Reading Selections for Biological Science Students*. This book is used by professors at the University of Havana to conduct reading classes in English to second year university students in the Faculty of Biology.

All texts in the book were published between 1986 and 1989. The professors in charge of editing the articles included in the book (García-Pérez et.al., 1991,) took into consideration the fact that the biological sciences are divided into three main branches in that faculty: Microbiology, Biology and Biochemistry. So they had to ensure balance among the branches when selecting the articles representing each branch in the book.

Another issue relating to the character of texts is text length. They should be long enough to represent reliably the linguistic characteristics of the full text, but not so long as to add unnecessary information not to be used in the analysis. Texts in this study are identified as 'continuous segments of naturally occurring discourse' (Biber and Finegan, 1991) Few empirical investigations of variation within texts and optimal text sample length propose the analyses of the distribution of linguistic features across 1000-word texts samples extracted from larger texts (Biber, 1988).

The study includes 700-word texts samples² extracted from larger texts, inasmuch as previous studies have indicated that such shorter extracts do reliably represent at least certain linguistic characteristic of a text (Biber, 1988). To analyze all these texts without the aid of a computer would require several years, but the use of a computerized corpus in this study enabled automatic inclusion of the texts in readable codes for the computer with the use of the scanner, automatic counting of words, and automatic identification of linguistic features in a collection of texts. The automatic identification of linguistic features was done with the use of AnyText™, a Hypercard® based program that allows one to do fast word searches on any text-only files.

2.2 Features

For the purpose of this study, Biber's research was surveyed to identify the relevant features characteristic of the scientific genre. Among the 67 linguistic features for all genres that Biber identifies, ten of those features that co-occurred the most in scientific writing were selected and grouped into six major grammatical categories:

a) Passives

1. agentless
2. by-passives

b) Pronouns

1. third person pronouns
2. it
3. demonstrative

c) Modals

1. possibility
2. predictive

d) Nominalizations (-tion, -ment, -ness, -ity including the plural forms)

² Because balance had to be ensure among the three sub-genres when selecting the texts from the book *Reading for Biological Science Students*, there were some articles which did not have 700 words. The distribution of words per article is as follows: Fifteen articles have 700 words and 9 have between 407 and 689 words. As all the counts were normalized to a text length of 1000, the difference between text length does not constitute a problem.

- e) Perfect tense
- f) Conditionals

2.3 Frequency counts

The frequency counts of linguistic features were normalized to a text length of 1000 words³. Normalizing text length is mandatory for any comparison of frequency counts across texts because, even though a text length may not be very relevant in relation to another, the fact that the amount of words differs, may lead to an inaccurate assessment of the frequency distribution in texts

The frequency of occurrence of the linguistic features analyzed in the study are given in five different values.

- a) the mean frequency,
- b) the maximum frequency,
- c) the minimum frequency,
- d) the range
- e) the standard deviation

2.4 Factors

Factors represent an area of high-shared variance in the data, a grouping of linguistic features that co-occur with a high frequency. Factors are defined by correlations among the frequency counts of linguistic features; that is, when several linguistic features are highly correlated, then a factor is defined.

The first step in a factor analysis is choosing a method for extracting the factors. In linguistics, the use of factor analysis is generally exploratory. Although there are several options of factor analysis available, the study will include the most widely used known as 'common factor analysis' (Biber, 1988).

³ Biber also normalized his corpora to a text length of 1000 words. In order to compare the results of this study to those of Biber's, the texts have to be normalized to the same amount of words.

Common factor analysis extracts the minimum amount of shared linguistics features. So the first factor extracts the maximum amounts of shared linguistic features; that is the first factor would correspond to the largest group of co-occurrence in the data (passive-nominalizations, for example); the second would then extract the maximum amounts of shared linguistic features from the tokens left over after the first factor has been analyzed, and so on.

3. Results and Discussion

After obtaining the raw number of all the occurrences of the linguistic features in each text, the counts were normalized to a text length of 1000 words. Table 1.0 presents descriptive statistics for the frequencies of the linguistic features in the entire corpus of texts used in the study.

This table does not include the characterization of particular sub-genres, but provides an assessment of the overall distribution of particular features in biological science texts. Some features occur very frequently, for example nominalizations with a mean of 20.2 per 1000 words; other features occur very infrequently, for example, *by-passives* with a mean of 2.3 per 1000 words.

The variability in the frequency of features also differs from one feature to the next; some show a small difference of distribution across the corpus, such as conditional clauses. They have a maximum frequency of 8.5 per 1000 words and a minimum of 0.0 per 1000 words; other features show large differences, for example predictive modals occurred 42 times in some texts but not at all in other texts.

Table 1: Descriptive Statistics for the corpus of biological science texts as a whole

Linguistic feature	mean	minimum value	maximum value	range	standard deviation
agentless passives	12.8	5.7	30.0	24.3	10.6
<i>by</i> -passives	2.3	0.0	7.3	7.3	2.3
3rd person pronouns	11.6	3.6	27.1	23.5	6.6
pronoun <i>it</i>	7.2	1.4	28.5	27.1	9.8
demonstrative pr.	11.1	2.8	24.2	21.4	5.9
possibility modals	9.2	0.0	24.5	24.5	6.1
predictive modals	4.2	0.0	42.0	42.0	8.3
nominalizations	20.2	0.0	38.5	38.5	10.2
perfect aspect	9.6	0.0	27.1	27.1	5.9
conditionals	2.4	0.0	8.5	8.5	2.4

The distribution of the mean frequency of features that highly co-occur within each sub-genre and across sub-genres compared to the general corpus can be seen in Table 2.0 A cut was made in the features having salient loadings of 9.6 and over. Table 2.1 presents the results of the co-occurrence of those features having loadings of less than 9.6. As just ten linguistic features were analyzed in the study, and as the ten constitute the main reason for comparison, no exclusions of linguistic features were made in spite of the fact that some of the features had very low loadings. The abbreviations used in the tables stand for the following:

- | | |
|--------------------------------------|---------------------------------|
| 1. A-P.: agentless passive | 6. Poss. M.: Possibility modals |
| 2. <i>By</i> -P.: <i>by</i> -passive | 7. Pred. M.: Predictive modals |
| 3. 3rd P.P.: 3rd person pronouns | 8. N.: Nominalizations |
| 4. P. <i>it</i> : pronoun <i>it</i> | 9. P. A.: Perfect Aspect |
| 5. D. Demonstrative pronouns | 10. Cond.: Conditionals |

Table 2.0: Distribution of the linguistic features that had a co-occurrence of 9.6 and over within each subgenre and across sub-genres compared to the general corpus

Biochemistry	Microbiology	Biology	General
N.	N.	N	N.
A-P	A-P	A-P	A-P
D	D.	D.	D.
3rd P.P.	Poss. M.	3rd P.P.	3rd P.P.
P.A.		P.A.	
Poss. M.			

Table 2.1: Description of the distribution of the linguistic features that had a co-occurrence of less than 9.6 within each subgenre and across sub-genres compared to the general corpus.

Biochemistry	Microbiology	Biology	General
Pred. M <i>P. it</i> Cond. <i>By-P.</i>	Pred. M <i>P. it</i> Cond. <i>By-P.</i> 3rd. P.P. P.A.	Pred. M <i>P. it</i> Cond. <i>By-P.</i> Poss. M	Pred. M <i>P. it</i> Cond. <i>By-P.</i> Poss. M

Given Biber's idea that 'if certain features consistently co-occur, then it is reasonable to look for an underlying functional influence that encourages their use', the study first compared the general results of the microscopic analysis done in the Biological Science Corpora (Table 1.0) to that of Biber's (Table 3.), then analysed if the features that highly co-occur in this study coincided with Biber's (Table 3.1). A comparison of the underlying functional dimensions in both studies from a macroscopic outlook follows.

Table 3: Descriptive Statistics for the Science Corpora presented in Biber's Study

Linguistic feature	mean	minimum value	maximum value	range	standard deviation
agentless passives	17.0	7.0	38.0	31.0	7.4
<i>by</i> -passives	2.0	0.0	8.0	8.0	1.7
3rd person pronouns	11.5	0.0	46.0	46.0	10.6
pronoun <i>it</i>	5.9	1.0	16.0	15.0	3.4
demonstrative pr.	2.5	0.0	9.0	9.0	1.9
possibility modals	5.6	0.0	14.0	14.0	3.1
predictive modals	3.7	0.0	14.0	14.0	3.4
nominalizations	35.8	11.0	71.0	60.0	13.3
perfect aspect	4.9	0.0	16.0	16.0	3.5
conditionals	2.1	0.0	9.0	9.0	2.1

Table 3.1: Comparison between the mean frequencies in the Biological Science Corpora and the Science Corpora presented in Biber's Study

Linguistic feature	mean Bio. Sc. Corpora	mean Biber's study
agentless passives	12.8	17.0
by-passives	2.3	2.0
3rd person pronouns	11.6	11.5
pronoun <i>it</i>	7.2	5.9
demonstrative pr.	11.1	2.5
possibility modals	9.2	5.6
predictive modals	4.2	3.7
nominalizations	20.2	35.8
perfect aspect	9.6	4.9
conditionals	2.4	2.1

A comparison between this study and Biber's brings out interesting results. *Nominalizations* have a high frequency of occurrence with a mean of 20.4 in this study and a mean of 35.8 in Biber's. This is the feature that most frequently occurred in the corpus analyzed. Although it is not our objective to analyze the occurrences of *nominalizations* separately; that is, reporting how many words ending in *-tion* occur in this text, and how many words ending in *-ment* occur in another; it is interesting to note that a high percentage of *nominalizations* fall into the *-tion* group.

The next feature with the second highest frequency of occurrence in both studies was *agentless passives*. In the Biological Science Corpora 12.8 *agentless passives* occurred per 1000 words, while in the general Science Corpus 17 *agentless passives* occurred per 1000 words. It was observed that whenever there was a high frequency of passives there were many nominalizations, a correlation that also exists in Biber's study. The *agentless passives* are used to present propositions with no emphasis on the agent. They are used to give prominence to the patient of the verb, the entity acted upon. Agentless passives are frequently used in procedural discourse where the agent is presupposed across several clauses and the specific agent of a clause is not important to the discourse purpose. This type of discourse is typically very technical in content and formal in style.

Keeping the order of features from those which co-occurred the most to those which co-occurred the least in both studies, the third position is shared by *third person*

pronouns. The word *shared* was used because the features shared the third position and almost the same mean frequencies. *Third person pronouns* have a mean frequency of 11.5 in Biber's study; and a mean frequency of 11.6 in this study. *Third person pronouns* mark reference to referents apart from the speaker and addressee. The results show that *agentless passives*, *nominalizations* and *third person pronouns* highly co-occur in both studies.

The fourth feature having a salient loading in the Biological Science Corpora was the *demonstrative pronouns*. This feature was not marked at all in Biber's study. It had a mean of 2.5. Demonstrative pronouns are highly used as referential elements, a device very much used in scientific texts. It was very interesting to observe that there was a correlation between the presence of *third person pronouns* and *demonstrative pronouns*. When one of them occurred frequently, the other one did not. This does not mean that the presence of one presupposed the absence of the other. Both occurred in texts and there was always one more marked than the other.; but they do co-occur with *passives* and *nominalizations*.

The fifth more marked feature was the *perfect aspect*. The markedness of this feature in this study compared to its unmarkedness in Biber's is amazing (the same with *demonstrative pronouns*). A mean of 4.9 was reported in Biber's work whereas in this work the mean is 9.6. *Perfect aspect* proved to be very much used in the general Biological Science Corpora as the feature describes past events.

What has been described so far is the procedure for constructing a factor. The first factor would then be the sum of the features that highly co-occurred in the study. That is:

$$20.2 \text{ (nominalizations)} + 12.8 \text{ (agentless passives)} + 11.6 \text{ (third person pronouns)} + 11.1 \text{ (demonstrative pronouns)} + 9.6 \text{ (perfect aspect)} = 65.3$$

So, 65.3 would be the first factor in the analysis. However if the dimension underlying this factor were to be analyzed, one would first have to think of the functions of the linguistic features:

- a) nominalizations: indicate a referentially explicit statement,
- b) agentless passives: present propositions with no emphasis on the agent,

- c) third person pronouns: mark reference to referent apart from the speaker and addressee,
- d) demonstrative pronouns: mark reference to referent apart from the speaker and addressee,
- e) perfect aspect: describes a past event that is psychologically relevant to the present.

When interpreting the functions of these features, and when comparing their co-occurrence underlying the dimensions presented in Biber's study, there is clear evidence that the features fall into four of Biber's dimensions:

- a) *nominalizations* are related to the 'Explicit versus Situation-Dependent' dimension,
- b) *agentless passives* (as well as *by-passives*) are related to the 'Abstract versus Non-Abstract Information' dimension,
- c) perfect aspect and third person pronouns are related to the 'Narrative versus Non-narrative Concerns' dimension, and
- d) demonstrative pronouns are related to the 'On-line Informational Elaboration' dimension.

The two dimensions in Biber study which were highly marked in the General Science Corpora were 'Explicit versus Situation Dependent' and 'Abstract versus Non-Abstract Information'. The two mostly marked features in our study also fall into these two dimensions; however, if we take into consideration that dimensions are characterized on the basis of a consistent co-occurrence pattern among features, we cannot take Biber's dimensions as point of departure for comparison in our study as the consistent co-occurrence of the cluster of features previously analyzed are scattered in different dimensions in Biber's study. So, if we were to name this dimension in our study we would call it '*Allusion to Experimental Versus Factual Information*'. The way this dimension is labelled in this study focuses much deeper on the general function of the features in the texts included in the corpus.

The sixth feature having a salient loading in the study is *possibility modals*. The mean is really high, 9.2, compared to Biber's study in which the mean is 5.6.

Possibility modals are pronouncements concerning the ability or possibility of certain events occurring, that they *can*, *may* or *might* occur. In the biological science world these possibilities are always present.

The pronoun *it*, in the seventh position, marks a reduced surface form that can be a noun or a phrase. In this study the mean (7.2) is close to that of Biber's (5.9).

Predictive modals fall in the eighth step. In Biber's study they have a mean of 3.7 not far from the mean in this study, 4.2. Predictive modals are direct pronouncements that certain events *will* occur (something always long for in natural sciences but not always achieved).

Numbers nine and ten are shared by *conditionals* (2.4) and *by-passives* (2.3) which (keeping the order) have a mean frequency in Biber's study of 2.1 and 2.

As the features unmarkedly co-occur they can be said to belong to the same factor. So the second factor in this analysis would be:

$$9.2 \text{ (possibility modals)} + 7.2 \text{ (pronoun it)} + 4.2 \text{ (predictive modals)} + 2.4 \text{ (conditionals)} + 2.3 \text{ (by-passives)} = 25.3$$

The function of each of these linguistic features is the following:

- a) possibility modals: pronounce that certain events can, may or might occur,
- b) pronoun *it*: marks a reduced surface form,
- c) predictive modals: pronounce that certain events will occur,
- d) conditionals: specify the conditions that are required in order for certain events to occur, and
- e) *by* passives: reduce the emphasis on the agent.

The functions underlying these linguistic features may be related to Biber's fourth dimension: 'Overt Expression of Persuasion', as possibility modals, predictive modals, and conditionals fall into this dimension. However 'Overt Expression of Persuasion' was not reported as a characteristic dimension of scientific texts in Biber's study. He explains that possibility modals, predictive modals and conditionals are often used to persuade; nevertheless, the analysis of these features in the corpus analyzed do not seem to indicate persuasion, but conceivable information. Hence, if this dimension

were to be named, a more comfortable expression for this analysis would be 'Conceivable Information'.

The range of variation of some linguistic features is very high in both studies:

Features	standard deviation in this study	standard deviation in Biber's study
passives	10.6	7.4
third p.p.	6.6	10.6
pron. it	9.8	3.4
predictive modals	8.3	3.4
nominalizations	10.27	13.3

5. Conclusions

The analysis presented here corroborates Biber's thesis that knowing what linguistic features co-occur in a text and across texts helps researchers understand why the linguistic features occur. The study offers a rationale for genre analysis. It provides support for both, the existence of genres and the importance in carrying out a study departing from simply the analysis of the markedness and/or unmarkedness of linguistic features in a typology of texts to the underlying function of their co-occurrence. This study uses a computerized text corpora and of a not-grammatically tagged computer program for the automatic identification of ten linguistic features, and it narrowed down the analysis from the General Science Corpora, to the Biological Science Corpora. Such an analysis provided accurate and valuable information for future comparative analysis.

We began this research by investigating the co-occurrence of features in the General Biological Science Corpora in relation to the Science Corpora presented in Biber's study, however, the biological sciences are divided into three main branches: Biology, Microbiology and Biochemistry. The data related to these subgenres will be provided in the near future.

Additional research is required to find out the relations among texts in other fields, such as the Social Sciences, the Natural Sciences, and the Exact Sciences. The present model of genre analysis should prove useful for such related studies in ESL and it is hoped that it will provide a foundation for research to identify the relevance of genre analysis in reading and writing.

REFERENCES

1. BARBER, C. L. (1962). "Some measurable characteristics of modern scientific prose",
2. BESNIER, Niko (1986). "Register as a sociolinguistic unit: defining formality", in *Social and Cognitive Perspectives on Language*, Los Angeles: University of Southern California.
3. BIBER D. and FINEGAN E. (1991). "On the exploitation of computerized corpora in variation studies", in *Corpus Linguistics* pp. 204-220, Longman.
4. BIBER, Douglas (1988). *Variation Across Speech and Writing*, Cambridge: Cambridge University Press.
5. COUTURE, Barbara (1986). *Functional approaches to writing: research perspectives*, Norwood, NJ: Ablex.
6. CRANDALL, J. A. (1984). "Adult ESL: The other ESP", *The ESP Journal*, No. 3, pp. 91-96.
7. CROOKES, Graham (1986). "Towards a validated analysis of scientific text structure", *Applied Linguistics*, No. 7, pp. 57-70.
8. EWER, J. R. and LATORRE, G. (1969). *A Course in Basic Scientific English*, Longman.
9. GARCÍA-PÉREZ et al. (1991). *Readings for Biological Science Students*, Havana: University of Havana Press.
10. HALLIDAY, M. A. K., McINTOSH, A. and STREVEN, P. (1964). *The Linguistic Sciences and Language Teaching*, Longman.
11. HARMER, Jeremy (2001). *The Practice of English Language Teaching* (Third Ed.). Longman. in *Contributions to English Syntax and Phonology*, Stockholm: Almqvist and Wiksill.
12. MARSHALL, Stewart (1991). "A genre-based approach to the teaching of report writing", *English for Specific Purposes*, Vol. 10, No. 1, pp. 3-13.

13. MARTIN, James (1985). "Process and Text: Two aspects of human semiosis", in *Systemic perspectives on discourse*, Vol. 1., Norwood, NJ: Ablex.
14. NWOGY, Kevin N. (1991). "Structure of Science Popularization: A genre-analysis approach to the schema of popularized medical texts", *English for Specific Purposes*, Vol. 10, No. 2, pp. 111-123.
15. REID, Joy (1991). "Responding to different topic types: a quantitative analysis from a contrastive rhetoric perspective", in *Second Language Writing: Research Insights for the Classroom*, Cambridge: Cambridge University Press.
16. SCHIFFRIN, Deborah (1981). "Tense variation in narrative". *Language* No. 57, pp. 45-62.
17. SWALES, John. M. (1990). *Genre Analysis. English in Academic and Research Settings*, Cambridge, Cambridge University Press.
18. WIDDOWSON, H. G. (1983). *Learning purpose and language use*. Oxford: Oxford University Press.