



Classifying communicative formats in CHET, CEChET and others

Isabel Moskowich^{1*} and Begoña Crespo

¹MuStE group, University of A Coruña, A Coruña, Spain.
imoskowich@udc.es, bcrespo@udc.es

Abstract

This paper describes one of the concerns of corpus compilers when gathering samples of texts. In particular, it explores how to classify such samples in wider categories in the case of the Corpus of English Chemistry Texts (CEChET), one of the subcorpus of the Coruña Corpus of English Scientific Writing. To this end, authors have revised the literature to find (and try to solve) the terminological mess that includes laves such as genre, text-type and textual category. These laves have been widely related either to the form or the function of the text. In this paper the idea of “communicative format” is used to bring together form and function as they are seen as intermingled in texts at all levels.

1 The terminological mess

The classification of text samples into larger categories has been a perennial problem for corpus compilers. One issue to consider in gathering samples is the classification of extracts into hierarchically higher units. Three terms have been used widely in the literature: genre, text-type and (textual) category. The first, genre, has been defined by Biber (1988, p. 170) as “categories determined on the basis of **external criteria** relating to the speaker's purpose and topic; they are assigned on the basis of **use** rather than on the basis of form”. Biber's identification between genre and use was subsequently adopted by Lee (2001).

Other authors, such as Swales (1990, p. 58), take the same approach, arguing that communicative events share some **communicative purposes** and that “these purposes (...) constitute the rationale for the genre”. Such a position has also been defended by some functionalist linguists, such as Martin (2000, p. 13) when he claimed that genre could be regarded as a “social and cultural artefact”. For him, genres can be further interpreted as “staged, goal-oriented social processes through which social

* The research here reported on has been funded by the Spanish Ministerio de Economía y Competitividad (MINECO), grant number FF12013-42215-P. This grant is hereby gratefully acknowledged.

subjects in a given culture live their lives". This orientation is not limited to any particular school of thought, since authors like Taavitsainen (2012, p. 94) have also considered that "Genres are important operational tools and constitute dynamic systems which undergo change and variation. Sociocultural needs change over time, and genres change accordingly: old genres become adapted to new functions, new genres are created, and genres that have lost their function cease to exist." This perspective includes not only the idea of function as a factor in defining genre, but also embraces the concept of language dynamism, that is, change over time. However, dynamism operates synchronically as well as diachronically, especially across disciplines and in terms of the addressees involved.

The second term often used by linguists here is text-type. According to Biber and Finegan (1989, p. 6): "genre distinctions do not adequately represent the underlying text types of English ...; linguistically distinct texts within a genre represent different text types; linguistically similar texts from different genres represent a single text type." A finer-grained description is proposed by Lee (2001, p. 38) in stating that "lexical or grammatical (co-)occurrence features (...) are the internal (linguistic) criteria forming the basis of text type categories." More recently, Alonso-Almeida (2008, pp. 10-17) emphasized the complementary character of these terms in describing the concepts they denote: genre is differentiated from text type in the sense that "genre is externally defined", whereas "text type is characterized according to internal linguistic criteria." Therefore, it seems reasonable to conclude that in talking about textual taxonomy authors use the term genre and text-type interchangeably.

However, there is at least one further element in this terminological mess: (textual) category. The vague reference implicit in this term makes it useful for authors who do not want to show their exact position or who do not feel capable of making a conceptual choice.

It is undoubtedly the case that texts are produced with a clear function, in that the main aim of human language is to achieve some kind of response on the part of the receiver. However, depending on the kind of response the sender/addressor envisages, that is, the function of the text, form will vary. Hence, there is no absolute independence of form and function, and texts adopt forms depending on the function they perform (telegram, advertisement, treatise...). This mutual dependence means that form and function can be seen as a whole, one which ultimately cannot be wholly divided. For this reason we believe that the symbiosis between the form and function of a given communicative act deserves a new name: communicative format, the term we will use henceforth.

In what follows we address two research questions using the material proposed in the following section:

- 1) What are the preferred communicative formats in late Modern English scientific writing?
- 2) Are there any constraints on communicative format selection?

2 Material

For the present paper we have used four of the subcorpora from the *Coruña Corpus of English Scientific Writing*: the Corpus of English Texts on Astronomy (CETA) (2012), the Corpus of English Philosophy Texts (CEPhiT) (2016), the Corpus of History English Texts (CHET) (in preparation), and the Corpus of English Chemistry Texts (CEChET) (in preparation). All these contain samples of texts published between 1700 and 1900. They represent disciplines that can be grouped into two larger blocks, CETA and CEChET representing the so-called hard sciences, and CEPhiT and CHET the soft ones. The total number of samples is 163, as set out in Table 1:

Sub-corpora	Samples	Science
CETA	42	Hard
CEPhiT	40	Soft
CHET	40	Soft
CEChET	41	Hard
TOTAL	163	

Table 1: samples in the study

Most samples contain around 10,000 words, and thus the total word number of words we will be using is around 1,600,000. However, within the samples there are three exceptional cases containing between 5,000-6,000 words each, two of these extracts having been used when no 10,000 word-samples were available.

The four sub-corpora represent two broad kinds of disciplines: CETA and CEChET, with text samples drawn from the disciplines of astronomy and chemistry respectively, both belong to the hard sciences; meanwhile, CEPHiT and CHET, with samples from philosophy and history respectively, are from the soft sciences.

3 Methodology

The compilation of the Coruña Corpus, in terms of classifying samples into different communicative formats, follows several steps relating to both the form and the function of texts.

1. Existing classifications: Görlach's textual typology (2004) is used as the initial source for establishing taxonomies.
2. *OED* definitions: The previous typology is complemented with those corresponding definitions recorded in the *Oxford English Dictionary*, looking specifically at the descriptions for the 18th and 19th centuries.
3. Titles of works: In many cases the title is self-explanatory, and in such cases has been taken as a starting point for classification.
4. Prefatory material/audience: The aim here is to find any stated opinions by the author about his/her work. Authors may also make reference to their readership, thus providing information as to the function of the text.
5. The sample itself: Books published in the eighteenth century often contain a variety of works in different formats. Since we do not compile whole texts, it is the specific sample which needs to be assessed carefully before ascribing any particular communicative format. In this way we overcome the philologist's dilemma (Rissanen, 1989) and ensure accuracy in our classification.

Although the steps in this process have been consecutively listed, not all of them have proved necessary in sample classification, or, at least, not in the same order, as we will see in the three examples below.

Our first example (Figure 1), a work written by Sir Benjamin Brodie in 1880, shows many clues which make classification easier.

In the first place, the title, *Ideal Chemistry. A Lecture*, contains the word *lecture* which is indicative of the sort of sample it may contain. When resorting to previous definitions, we find that both Görlach and the *OED* define *lecture* as "A discourse given before an audience upon a given

subject, usually for the purpose of instruction. (The regular name for discourses or instruction given to a class by a professor or teacher at a college or University.)”

This suggestion as to the text’s classification is further reinforced by the prefatory material shown in Figure 2. The preface itself includes the description “The following lecture was delivered before the Chemical Society on June 6, 1867, after the presentation of the Royal Society of my first memoir on the Calculus of Chemical Operations. The lecture, however...” (Brodie, 1880, p. iii)

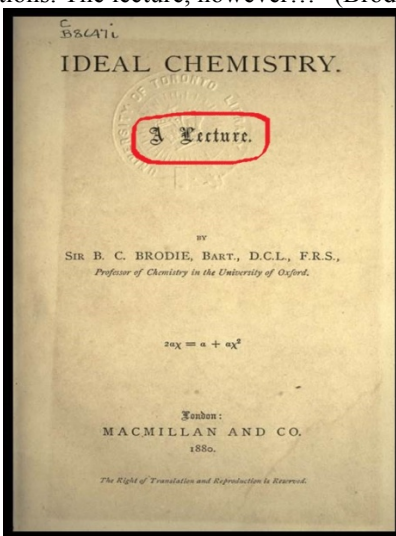


Figure 1: Benjamin Brodie's *Ideal Chemistry*. First page.

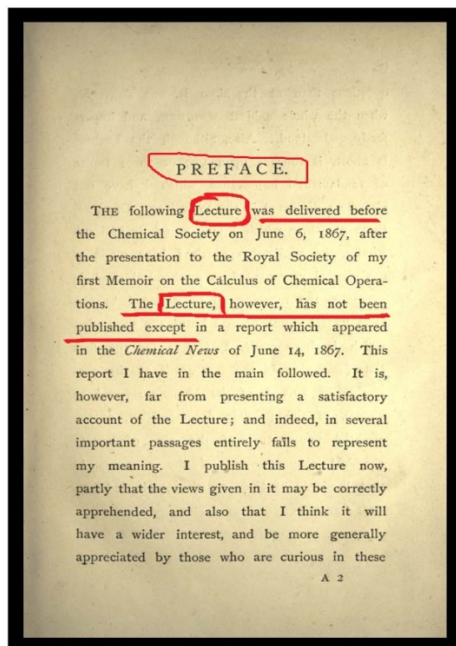


Figure 2: Benjamin Brodie's *Ideal Chemistry*. Prefatory material

Although this seems to be a clear case of Lecture, stated as such in the title of the work and also by the author, the sample must be examined in order to arrive at a complete picture of form and function. On inspection it quickly becomes evident that the form is of a text written to be orally delivered, in that it contains a direct address to the audience (“Mr President”, “the Chemical Society”), deictic elements referring to the context of delivery (“this evening”, “in the brief space of one hour”) and the object of the message itself (“an account of an abstruse and difficult subject”). Indeed, all this information is revealed in the first four lines reproduced in Figure 3.

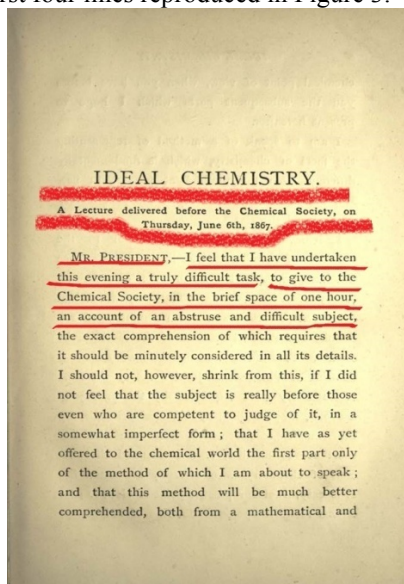


Figure 3: Benjamin Brodie's *Ideal Chemistry*. First Page

Our second example illustrates a slightly more complex text, as well as a different procedure regarding the order of the steps taken in classification. The sample is taken from a (1708) work by Christopher Packe, the title of which is extremely long (typical of many works at that time): *Medela Chymica: or, an Account of the Vertues and Uses of a select Number of Chymical Medicines Adapted to the Cure of the most Chronick and Rebelious Diseases To which is Subjoyned a Brief History of Cures Effectuated by Them. As also, An essay upon the Acetum Acerrimum Philosophorum, or Vinegar of Antimony, with some Experiments made therewith*. In terms of its communicative format, the title mentions “account”, which is certainly vague and not identifiable with any particular format, and “essay”; hence we must turn to previous definitions. Görlach defines essay as a “Short prose composition, first draft” (2004, p. 88), and in definition number 8 for the same entry in the OED it says:

8. A composition of moderate length on any particular subject, or branch of a subject; originally implying want of finish, ‘an irregular undigested piece’ (Johnson), but now said of a composition more or less elaborate in style, though limited in range. The use in this sense is app. taken from Montaigne, whose *Essais* were first published in 1580.

In both definitions, the length of the work (that is, pertaining to form) seems to play a part. However, since this step has not provided information leading to a definite classification, the prefatory material must be read in order to see whether the author himself expressed any notion as to the nature of his work.

Figures Figure 4 and Figure 5 below illustrate Packe's words in two different parts of his dedication. As can be seen, the references he makes in this prefatory text are contradictory, since he refers to his work as both a "treatise" and an "essay".

to me. So many Cures of Venereal Patients have been made by the *Astrum Mercurii* and *Sal Solutivum cum Sulph. Veneris*, that I have inserted none of them here directly, both because they would be sufficient of themselves for a large **Treatise** upon that Subject, but chiefly because I thought it an Undecency to offer them to many People who may perhaps vouchsafe this Treatise a Reading. The

Figure 4: Packe's *Medela Chymica*. Fragment.

The DEDICATION.
The **Essay** upon the *Vinegar of Antimony* I have Written purely for the sakes of those Ingenious Souls who believe that the Writings of the Adept Philosophers are not altogether Fabulous, or Published with no other intent but to amuse or deceive the Credulous. The Art of Transmuting or

Figure 5: Packe's *Medela Chymica*. Fragment

Hence, an inspection of this work reveals just how careful compilers must be in classifying text samples based on the characteristics of a selected excerpt. Indeed, the second part of Packe's work includes text in the form of a letter (Figure 6), which is itself mentioned as such by the author.

THE
SECOND PART
Containing a short History of Cures, &c.

Epileptick Fits Cured.

ON the 30th. of November, 1696. the following Letter was brought to me, with a desire to send what I should think fit, in Answer to it.

Loving Brother,

HERE is a Direction Enclos'd which *Mathew Bearder* desires you or my Brother *John* to go to the Dr. with, for his Son *Mathew* hath strange Fits, either *Convulsions* or *Falling Sickness*, he hath had 6 or 7. They take him so suddenly, that they give him no warning of their coming, and he can remember nothing of them when they are gone. *Robert Ryley's* Daughter had just such

Figure 6: Packe's *Medela Chymica*. Fragment

Other parts of the volume by Packe employ more colloquial language and are of a less technical nature, intended for practitioners and professionals. In those pages, the language is direct and conveys clear instructions in the form of what could be understood as a catalogue (see Figure 7).

Arcanum Universale.

THIS is a Medicine very small in Dose, but powerful in Operation, penetrating the whole Body, resolving all tenacious stubborn diseasly Matter it meets with, and expelling it by the nearest Emunctories, as either by Vomit, Stool, Urine or Sweat. It is therefore administred with great benefit in all Diseases where there is Peccant Matter to be expell'd, or Nature languishing under a certain Stupor or Sluggishness, wants to be excited and stirr'd up to Action.

Figure 7: Packe's *Medela Chymica*. Fragment

Taking all these factors into consideration, the final decision to include a 10,000-word sample from this work involved the compilers evaluating the whole extract again, independently of the rest of the work. As already announced by the author, it is in fact an essay, containing the rhetorical features expected in such a format (see Figure 8).

Of the Acetum Acerrimum Philo-
sophorum; or Vinegar of An-
timony.

MANY have been the attempts of those who have exercised themselves in the more secret Chymia to prepare the Vinegar of Antimony, being perswaded from the Writings of Arcepinus, Riply, Paracelsus, Basilus, Isaacus Hollandus and others, that it is a true Philosphick Menstruum, by which Metals and Minerals may be radically Dissolv'd, and thence excellent Medicines for the Cure of Diseases in Humane Bodies (at least) may be prepar'd.

Figure 8: Packe's *Medela Chymica*. Fragment

Finally, this example is illustrative of those cases in which explicit mention of an assumed or expected audience by the author gives compilers a firm indication of the probable classification. Such is the case with the (1870) work by James M. Crafts, *A Short Course in Qualitative Analysis, with the New Notation*. It is not only the title containing the word *course*, but also the first lines of the preface that give clues as to the readership and, therefore, the format chosen, in this case to convey chemical knowledge. Indeed, the author affirms: “This little work was written for the use of a class of students in the Cornell University, who take a year's course of chemistry, including four hours a week of laboratory practice”.

Previous definitions by Görlach (2004) and by the OED have the following definitions respectively:

“Book used as a standard reference work.”

And

“A book used as a standard work for the study of a particular subject; now usually one written specially for this purpose; a manual of instruction in any science or branch of study, esp. a work recognized as an authority.”

Also, when looking at the sample to resolve the philologist’s dilemma, the layout of the material confirms our suspicions that this may be a “textbook” (see Figure 9), as suggested by the use of bold font-types, summaries, etc.

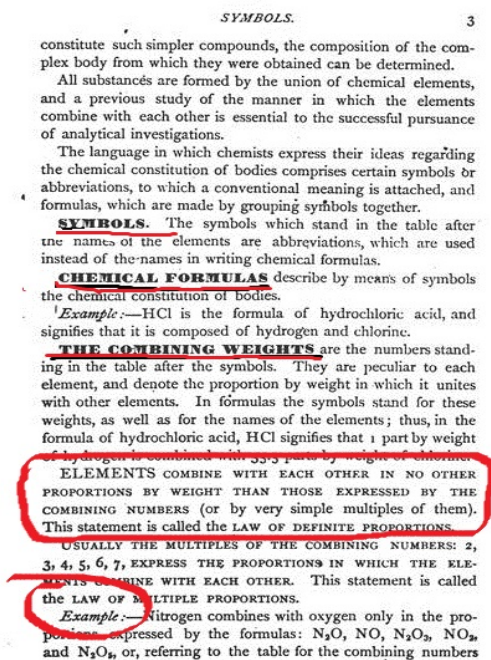


Figure 9: James M. Crafts, *A Short Course in Qualitative Analysis*

All the samples in the four sub-corpora under analysis were classified following this procedure. Hence, the resulting formats can be grouped and analysed in light of the socio-external factors that led late modern English authors to choose particular formats to convey their knowledge in a world where there was a rich array of textual possibilities open to them.

4 Analysis: Communicative formats

All 163 samples compiled in the sub-corpora here can be classified into ten different communicative formats: letter, manual, dictionary, dialogue, article, travelogue, lecture, textbook, essay, treatise.

Each of these terms is defined in the OED as follows:

Letter: Two possible senses in the OED:

Sense II.4.a.= A written communication addressed to a person, organization, or other body, esp. one sent by post or messenger; an epistle. Recorded in the Ancrene Riwe ca. 1200 with this sense.

Sense II 4. c.= An article or report describing the social, political, or cultural aspects of a particular situation or place, esp. one by a journalist or correspondent in another country. Chiefly with from, and in titles and headings. (In early use not clearly different from sense 4.a)

This second sense is in fact exemplified in the OED with the title of a book from 1782 by J. H. St. J. de Crèvecoeur, *Letters from an American Farmer; describing certain provincial situations, manners, and customs, not generally known; and conveying some idea of the late and present interior circumstances of the British colonies in North America*. This is the use most often found in our corpus.

Manual

According to the OED, from 1475 onwards this refers to “A handbook or textbook, esp. a small or compendious one; a concise treatise, an abridgement. Also in extended use”.

Dictionary

“In extended use: a book of information or reference on any subject in which the entries are arranged alphabetically; an alphabetical encyclopedia”. This meaning is recorded as early as 1576 by the OED.

Dialogue

“A literary work in the form of a conversation between two or more persons, in which opposing or contrasting views are imputed to the participants.” With this sense the term can be found as early as the Anglo-Saxon Charters.

Article

“A non-fictional piece of writing forming part of a journal, encyclopedia, or other publication, and treating a specific topic independently and distinctly.” First used in 1701, according to the OED

Travelogue

Defined as “An (illustrated) lecture about places and experiences encountered in the course of travel; hence a film, broadcast, book, etc., about travel; a travel documentary.” This term is originally from American English, and was coined in the early 20th century.

Lecture

The OED, in its sense 4. a. of the term, defines a lecture as “A discourse given before an audience upon a given subject, usually for the purpose of instruction. (The regular name for discourses or instruction given to a class by a professor or teacher at a college or University.)” Recorded for the first time in the year 1536, in Act 27 of Henry XII, Acts of Parliament.

Textbook

“A book used as a standard work for the study of a particular subject; now usually one written specially for this purpose; a manual of instruction in any science or branch of study, esp. a work recognized as an authority”. OED dates the first use to 1779.

Essay

“A composition of moderate length on any particular subject, or branch of a subject; originally implying want of finish, ‘an irregular undigested piece’ (Johnson), but now said of a composition more or less elaborate in style, though limited in range. The use in this sense is app. taken from Montaigne, whose *Essais* were first published in 1580.”

Treatise

Of the senses given in the OED, the one that can best be applied here is “A book or writing which treats of some particular subject; commonly (in mod. use always), one containing a formal or methodical discussion or exposition of the principles of the subject; formerly more widely used for a literary work in general”. However, there is a more general meaning, now obsolete: “A descriptive treatment, description, account (of something)”. Given the period under survey, some of the authors in the Coruña Corpus may have had this sense in mind when naming and describing their works.

As can be seen in Figure 10 below, the format treatise was recorded in 74 samples, that is, in 45.39% of the samples. Textbook is the second most common format, used in 27 samples compiled (16.56%), followed by essays (21 samples; 12.88%), lectures (17; 10.42%) and article (10; 6.13%).

This illustrates broad tendencies in the use of communicative formats within late modern English scientific discourse.

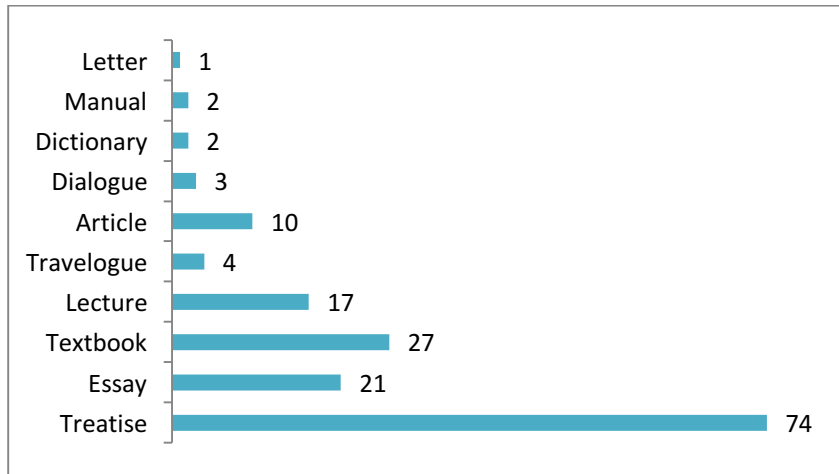


Figure 10: Communicative formats in CC, number of text samples

As we have noted, the samples in our analysis belong to four disciplines which can be grouped into the so-called soft (philosophy, history) and hard sciences (astronomy, chemistry). In the case of the former, and following the general tendency, treatise is the most common format across the two disciplines. Equally significant is the underrepresentation or absence of some other formats, such as dictionary and manual. Both the presence and absence of particular communicative formats will help us determine the kind of constraints underlying format selection, which will allow us to address research question number two.

Figure 11 below illustrates the formats used in philosophy and history:

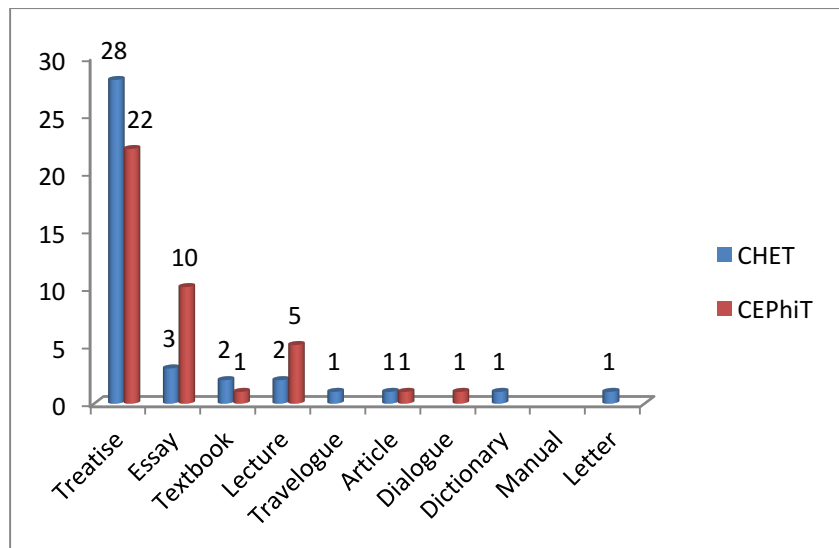


Figure 11: Communicative formats in the Soft Sciences

As for the hard sciences, different selection preferences were found. Textbook is the most frequently used format within the group of astronomy texts but, once more, treatise is equally represented in both astronomy and chemistry samples. The absent formats differ from those for the soft sciences: travelogue and letter. Figure 12 below illustrates the formats used in astronomy and chemistry:

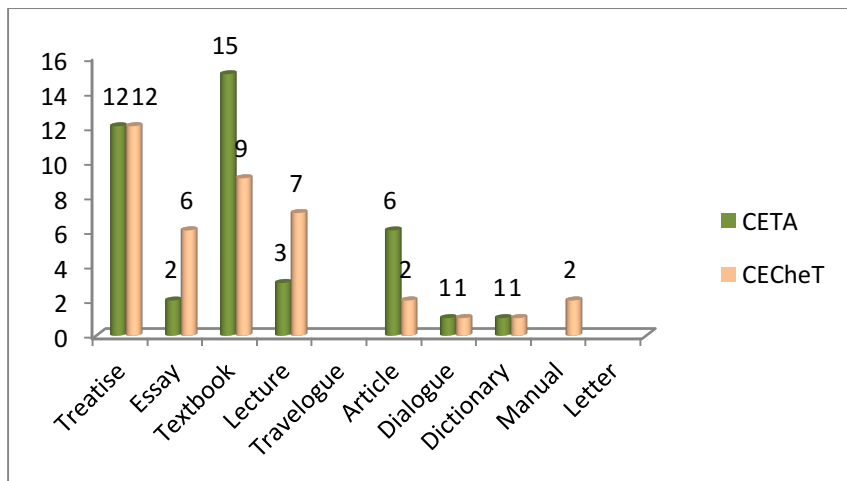


Figure 12: Communicative formats in the Hard Sciences

The frequent use of textbooks and manuals within the hard sciences may reflect a response to the growing social demand for knowledge which characterized post-empiricist times.

In the comparison between the Corpus of History English Texts (CHET) and the Corpus of English Chemistry Texts (CEChET), a total of 34.5 % of all samples are classified as treatises, although there is a remarkable difference in frequency between CHET (68.3%) and CEChET (30%). The information which history texts typically provide, representing here the field of humanities, seems to be conveyed mainly by a format which narrates previous facts or past events as a timeline or sequence; it evinces the hand of a distant third person narrator who seeks only to present straightforward facts through expository writing.

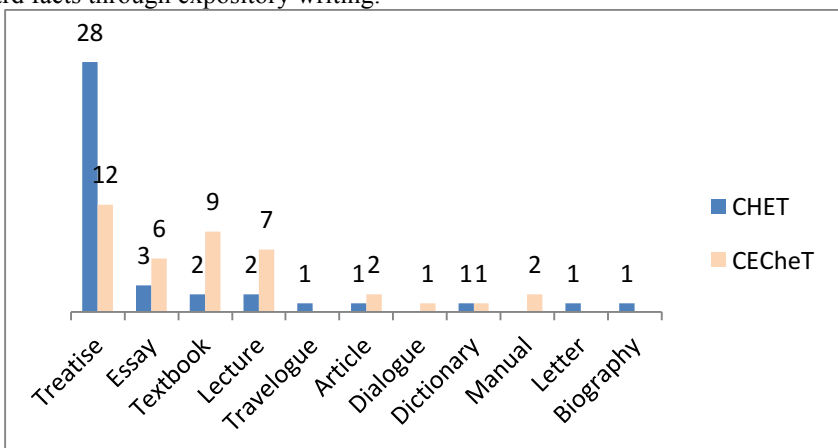


Figure 13: Communicative formats in CHET and CEChET

The second and third most common formats in history writing samples are travelogue, with four samples, and essay, with three. Travelogue, although formally similar to the expository nature of treatise, fulfills a different communicative function in that it describes different kinds of travel events. In the small discourse community of travellers, communicative practices and knowledge of the specific field or area were shared, thus rendering travelogues a characteristic communicative format in historical writing.

This is not the case with chemistry, where textbooks (9 samples) and lectures (7 samples) are the second and third most frequent formats. Works that could be used as standard references in the teaching/learning process in a society eager to advance and disseminate scientific knowledge are abundant in chemistry. This function determines the formal appearance of written texts, in which notions, formulae and experiments, among other ideas, are explained for the benefit of those seeking to learn. Similarly, lectures (oral to written mode) represent the same function in the oral medium, despite these having been subsequently transcribed. Lectures were conceived of as teaching-oriented oral channels for the communication of knowledge. There is a reciprocal relationship between the two participants in the communicative process, addressor and addressee, or, to be more specific, writer and reader, speaker and listener: Addressors present their findings, and addressees are amenable to the process and aim to assimilate this knowledge, independent of the medium of communication.

At the other end of the scale of frequency we find that some formats are not used in certain disciplines. Thus, manual, dialogue and dictionary were absent from the history corpus, whereas travelogue and letter were not found in chemistry.

Our analysis of the relationship between communicative format and discipline in the four sub-corpus reveals that format selection is constrained by subject matter. The nature of what is being communicated in a text necessarily calls for a particular format, and vice versa, in the sense that a particular format also transmits a particular function. And this is something that can be deduced from the presence and/or absence of communicative formats in each of the disciplines. At the same time, this reinforces the idea that discipline or subject matter to some extent merges with the form, structure and organisation of a text composed by the author to elicit a particular response or provoke a certain reaction on the part of the reader. If this is so, the symbiosis between the form and function of a given communicative act is successful.

5 Concluding remarks

In order to account for the evident interrelationship between form and function, a new term is required: COMMUNICATIVE FORMAT.

The compilation process, involving gathering and classifying text samples, runs the risk of becoming subjective and clear methodological steps are needed to overcome this. An orderly and principled methodology leads to a classification which is maximally systematic. Formats are dynamic and depend on the nature of the relevant epistemic community and its audience, and it is precisely because of this dynamism that both change and variation can be attested. Although both participants in the communicative process are to be taken into account, it is undoubtedly the case that subject matter also plays its part.

References

- Alonso Almeida, F. (2008). The Middle English medical charm: Register, genre and text type variables. *Neuphilologische Mitteilungen*, 109/1, 9-38.
- Biber, D. (1988). *Variation across Speech and Writing*. Cambridge: Cambridge University Press.
- Biber, D., & Finegan, E. (1989). Drift and Evolution of English Style: A History of Three Genres. *Languages*, 65, 487-517.
- Brodie, B. (1880). *Ideal Chemistry. A Lecture*. London: Longmans, Green and Co.
- Crafts, J. (1870). *A Short Course in Qualitative Analysis, with the New Notation*. New York: John Wiley and son.
- Görlach, M. (2004). *Text Types and the History of English*. Berlin: Walter de Gruyter.
- Lee, D. (2001). Genres, registers, text types, domains, and styles: clarifying the concepts and navigating a path through the BNC jungle. *Language Learning & Technology* 5/3, 37-72.
- Martin, J. (2000). Analysing genre: functional parameters. In F. Christie, & J. Martin, *Genre and Institutions. Social Processes in the Workplace and School* (pp. 3-39). London/New York: Continuum.
- Moskowich, I. *et al.* (2012). *Corpus of English Texts on Astronomy (CETA)*. Amsterdam: John Benjamins.
- Moskowich, I. *et al.* (in preparation). *Corpus of English Chemistry Texts (CEChET)*.
- Moskowich, I. *et al.* (in preparation). *Corpus of History English Texts (CHET)*.
- Moskowich, I., Camiña-Rioboo, G., Lareo, I., & Crespo, B. (2016). *'The Conditioned and the Unconditioned' Late Modern English Texts on Philosophy*. Amsterdam/Philadelphia: John Benjamins.
- Oxford English Dictionary. <http://www.oed.com>. (2016, February 11).
- Packe, C. (1708). *Medela Chymica: or, an Account of the Vertues and Uses of a select Number of Chymical Medicines Adapted to the Cure of the most Chronick and Rebelious Diseases To which is Subjoyned a Brief History of Cures Effected by Them...* London: Printed for John Lawrence at the Angel in the Poultry.
- Rissanen, M. (1989). Three problems connected with the use of diachronic corpora. *ICAME Journal* 13, 16-19.
- Swales, J. (1990). *Genre analysis: English in Academic and Research Settings*. Cambridge: Cambridge University Press.
- Taavitsainen, I. (2012). Discourse Forms and Vernacularisation Processes in Genres of Medical Writing 1375–1550. In A. Aejmelaeus, & P. Pahta, *Studies across Disciplines in the Humanities and Social Sciences* 7. (pp. 91-112). Helsinki: Helsinki Collegium for Advanced Studies.