# The Vocabulary of Medical English

# The Vocabulary of Medical English:

## *A Corpus-based Study*

By

Renáta Panocová

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# ACKNOWLEDGEMENTS

# CHAPTER ONE

# INTRODUCTION:
# DEFINING MEDICAL ENGLISH

This monograph explores the vocabulary of medical English from a corpus-based perspective. In investigating medical corpora, I will highlight the question of the characterization of medical vocabulary in English. One of the central issues I address is how to design a methodology appropriate for the purpose of description of medical English. This contrasts with the pedagogical perspective, which uses medical corpora for compiling word frequency lists used as a basis for developing teaching materials. The difference in aims of the two perspectives, pedagogical and descriptive, points to the need for a different methodology to be applied in research of the phenomenon referred to as *medical language* or *medical English* (ME).

Language is an important tool in professional communication in medicine. The history of medicine clearly points to Latin as a dominant language in the field throughout the middle ages and the early modern era, when it was the main international language not only in medicine, but also in religion and philosophy (Fischbach, 1993: 94). Even today, the influence of Latin in medical language should not be ignored. Several textbooks in medicine and medical terminological dictionaries in a number of different languages take Latin as a basis, for instance Vojteková's (2015a) trilingual dictionary of anatomical terms in Latin, Slovak, and Polish. A discussion of the importance of Latin in current medical terminology is given in Vojteková (2015b).

From the 17th century onwards a new tendency for the use of national languages such as German, French, and English in medical writings emerged (Ferguson, 2013: 282). The relatively equal status of German, French, and English changed in the second half of the 20th century, resulting in English taking over the most prominent role in medical texts. A piece of evidence comes from Maher (1986), who reports that in 1980 72.2% of the articles in the *Index Medicus* database were published in English. A similar tendency is observed by Giannoni (2008) who reports that more than 99% of medical journal papers by Italian authors are in

English. Gunnarsson's (2009) findings confirm this tendency in Scandinavian countries.

In this context it is interesting to report some results from a search in Google. A Google search for *medical English* gives 1 030 000 000 hits in half a second.[1] The top hits include medical English online exercises and games, medical English worksheets and teaching resources, exercises for doctors and their patients, reading and listening exercises for medical workers, offers for courses of medical English. Many of these webpages promise improvement of communication in a medical environment by mastering ME, not only for native speakers of English but also for doctors, nurses, and other health-care workers with a native language other than English. This also demonstrates the importance of ME at the international level.

The prominent role of English in medicine raises at least two important questions; the first addresses the nature of ME and its definition and the second concentrates on the position of ME in relation to general English. Two main perspectives have been adopted in determining the notion of ME. Firstly, ME can be defined in terms of the distinction from other language variants and common or general language, as in Lankamp (1989). The second perspective considers ME as a sublanguage.

A central hypothesis in Lankamp's (1989) research is that ME is so distinct from general language "that it would have to be acquired or learnt (two interchangeable terms in here) by language users with only general language knowledge" (1989: 14). Then he raises the question in what sense ME differs from general language or other language variants. In his study, Lankamp (1989: 14) focuses on "the investigation of the ways in which written English medical language differs on the various linguistic levels of analysis (discourse, syntax, semantics, lexicon and morphology) from other English written language variations". Following Hudson (1980) and Picht and Draskau (1985), Lankamp (1989: 20) views ME in terms of register and language for special purposes. On this basis Lankamp (1989: 20-22) distinguishes the five dimensions of variation given in (1).

(1)      a. medical specialism
         b. manner of transmission of medical language
         c. relations between participants in medical exchange
         d. communicative purpose
         e. national language

---

[1] Retrieved 29 January, 2016

The dimensions in (1a-c) are in line with Picht & Draskau's (1985) discussion of terminology and specialized language, the remaining two are added by Lankamp (1989: 22). The dimension in (1a) highlights the importance of linguistic differences among different medical specializations or professional groups. These are similar to linguistic differences between ME and other language variations. In (1b) it is emphasized that important differences are expected between spoken medical language and written medical language. For instance, medical jargon or slang used in spoken medical context will not occur in research articles in respectable medical scientific journals. The label *tenor* is sometimes used to refer to the dimension in (1c). It indicates that linguistic properties may differ depending on the roles of participants, e.g. doctor-doctor conversation, doctor-patient conversation, doctor-nurse conversation, equipment manufacturer-doctor, etc. (Lankamp, 1989: 21). An additional dimension in (1d) highlights "the function of a language for special purposes-type register to communicate information of a specialist nature at any level of complexity in the most economic, precise and unambiguous terms possible, i.e. as efficiently as possible, especially in the expert-to-expert tenor" (Lankamp, 1989: 22; cf. also Sager et al. 1980: 290-291). In this dimension the role of terminology based on the need for precise, and preferably non-synonymous language items to label relevant concepts is crucial. The central point of (1e) is the fact that "medical language is differentiated according to specific national languages expressing international medical concepts. In this dimension, medical language is differentiated in medical Dutch, medical English, medical French etc." (Lankamp, 1989: 22).

Lankamp (1989: 23) suggests that defining ME as a type of register is fully compatible with the lexical competence of a language user in the psycholinguistic model he presents in his book. However, it should be noted that English for Specific Purposes (ESP) and register are not one and the same phenomenon. Biber and Conrad (2009: 3) in their book *Register, genre, and style* explain that "ESP focuses on description of the language used in registers/genres from a particular profession or academic discipline, e.g. biochemistry or physical therapy". It is important to emphasize that Biber and Conrad (2009: 2) use the terms *register*, *genre*, and *style* to refer to three different perspectives on text varieties.

It is understood that "the register perspective combines an analysis of linguistic characteristics that are common in a text variety with analysis of the situation of use of the variety" (Biber and Conrad, 2009: 2). This means that linguistic features such as pronouns and verbs are functional,

and, therefore their use depends on situational context and a type of communicative purpose.

The genre and style perspectives each concentrate on similarities and differences with respect to the register perspective. The property shared by the genre perspective and the register perspective is that the purposes and situational context of a text variety can be immediately identified. By contrast, linguistic analysis concentrates on the conventional structures used to build a text within the variety, for example, the conventional way in which a letter begins and ends (Biber and Conrad, 2009: 2).

Finally, the style perspective is similar to the register perspective in its linguistic focus. This means that linguistic features occurring in a particular variety are analysed. The crucial difference from the register perspective is "that the use of these features is not functionally motivated by the situational context; rather, style features reflect aesthetic preferences, associated with particular authors or historical periods" (Biber and Conrad, 2009: 2).

In the history of the development of ESP, Hutchinson and Waters (1987: 13) identify an initial stage in which, they say, "the analysis had been of the surface forms of the language" in the form of register analysis, that is, the study at the sentence level of the use of language in different communicative settings, such as the language used by nurses, airplane mechanics, and bank tellers. In this stage, the teaching of reading received minimal attention. It was at the next stage of development that reading pedagogy in ESP took major steps forward: "Whereas in the first stage of its development, ESP had focused on language at the sentence level, the second phase of development shifted attention to the level above the sentence, as ESP became closely involved with the emerging field of discourse or rhetorical analysis" (Hutchinson and Waters, 1987: 10). A more detailed discussion of the role of reading in ESP from various perspectives can be found in Hirvela (2013). To sum up, ME may be approached from the register perspective, but a register perspective leads to a much more fine-grained division. ME is not a register, but a range of registers. ME covers for instance, research papers, doctor-patient conversation, and Patient Information Leaflets with instructions on the use of pharmaceuticals. In line with Biber and Conrad (2009: 32) "there is no single "right" level for a register analysis". Obviously, even in the register of medical research articles more specific subregisters can be identified, e.g. there may be some degree of variation between research articles in the domain of psychiatry and cardiology. The same applies to doctor-patient communication in different situational contexts such as medical examination, surgery or disclosure that a patient suffers from a lethal

disease. As Biber and Conrad (2009: 33) emphasize, differences between registers can be viewed as a continuum of variation.

Ten Hacken and Panocová (2015: 2-3) note that it is not common to ask whether someone speaks medical English as opposed to the question whether someone speaks English. Similar to Lankamp (1989: 22) they point to the fact that medical language is language-specific and medical English differs from medical Slovak or medical Dutch. They also raise the question of the relationship between English medical language and general English, but also of the relationship with Dutch and Slovak medical language (ten Hacken and Panocová, 2015: 3).

An alternative is to approach medical language from the perspective of sublanguages (Harris, 1968; Kittredge, 1987; Lehrberger, 2014). This means that English medical language is taken to be a sublanguage of English. In the first perspective, ME is like a register of English. In the second one, it is considered as a subset. In either perspective, the vocabulary of a particular area of study or professional use, for instance medicine, is an example of specialized vocabulary.

Harris (1968) introduced the notion of *sublanguage* in linguistics "by analogy with *subsystem* used in mathematics" (Lehrberger, 2014: 20). A sublanguage is viewed as a theoretical construct. Lehrberger (2014: 20) points out that whereas in mathematics "a subsystem can be readily defined in terms of restrictions on the sets and operations of the system of which it is a part", in a natural language and its sublanguage "the relation between part and whole is not so clear-cut". This is in line with Kittredge's (1982: 110) observation that "[i]n considering which samples of specialized language can be regarded as representing "genuine" sublanguages we are immediately faced with the lack of an empirically adequate definition of the term" and there is a need for more precise delimitation criteria. Kittredge (1982: 110) claims that "the closure property proposed by Harris (1968) is not in itself sufficient to resolve this question" and the main reason is that "[i]f a sublanguage can be *any* subset of sentences which is closed under the transformational operations, this definition could identify a very large number of linguistic subsets as sublanguages". In mathematics, *closure* refers to the property that if a particular operation is applied to members of a set, the result will always again be a member of that set. Kittredge (1982: 110) understands the closure property as a necessary condition. This means that even if we have a set of sentences which can be considered as a sublanguage, "we must include in it all sentences generated from the candidate set by means of transformational operations of negation, question formation, clefting, conjunction, etc." (Kittredge, 1982: 110). It is important to add that Harris (1968) points out that a

linguistic subsystem can be closed only under some, but not all of the operations.

Another type of closure is vocabulary closure. This has been investigated by McEnery and Wilson (2001) and Temnikova (2013). They examine relationships between types and tokens in a corpus of the genre. If a genre shows closure properties, the number of types stops growing after some number of tokens has been processed. On the other hand, if it does not exhibit closure, then the number of types will continue to rise continually as the number of tokens increases (Temnikova, 2013: 72).

Kittredge (1983: 49) repeatedly emphasizes that although sublanguages have been investigated in a number of ways and perspectives, "there is no widely accepted definition of the term", but there is an agreement about certain factors that are usually present in a subset of a particular natural language and are essential for semantic processing. The factors mentioned by Kittredge (1983: 49) are presented in (2).

(2)       a. restricted domain of reference
          b. restricted purpose and orientation
          c. restricted mode of communication
          d. community of participants sharing specialized knowledge

In (2a) the main point is that linguistic expressions refer to a set of objects and relations and their number is relatively small. In (2b) it is emphasized that there are clearly identifiable relationships among participants. The same applies to the goals of the exchange. The factor in (2c) indicates that there are differences not only between spoken and written communication, but "there may be constraints on the form because of "bandwidth" limitations (e.g. telegraphic style)" Kittredge (1983: 49). In medicine one could think of prescriptions. The last factor in (2d) suggests it is easier to determine properties of a sublanguage if a community of users who share it can be identified. This is important in order to determine characteristic patterns of usage which contribute to a complete characterization of a sublanguage as a linguistic system.

If we compare the lists of features in (1) and (2) we can see that some of the features are shared by the register perspective and sublanguage perspectives. It is obvious that taken together they must overlap to a certain extent. A striking difference is that (1e) is not explicitly mapped in (2). What the two perspectives share is that a specialized language is central in a more in-depth research. The main difference between the two perspectives is that register is about the use of competence whereas sublanguage concentrates on a subset of a language.

Generally, it is well-known that language users on advanced and proficient levels must have implicit knowledge about register, word meaning, and lexical and grammatical patterns, because otherwise it would not be possible to write and speak appropriately (Nesi, 2013: 451). It is also obvious that language users have intuitions about language, but it is questionable whether and when they are reliable or misleading. According to Sinclair (1991: 4) "human intuition about language is highly specific, and not at all a good guide to what actually happens when the same people actually use [it]". This may be seen as a sufficiently strong argument for the use of corpora and especially specialized corpora in the research into specialized languages, e.g. medical language. At present, large corpora for English are available for online search of a number of linguistic features, for instance the Corpus of Contemporary American English (COCA), the British National Corpus (BNC), etc. There is a strong tendency for other languages to compile similar national corpora, e.g. the Russian National Corpus, the Slovak National Corpus, the National Corpus of Polish.

Tognini-Bonelli (2001) differentiates between *corpus-based* research and *corpus-driven* research. Corpus-based research makes it possible to verify intuitions a researcher has about language use, whereas corpus-driven research uses corpus data to formulate relevant observations and generalizations about language use. Although it is always necessary to start from a theoretically informed research question, the present research assigns corpora a much more central role than as only a tool to test intuitions. Therefore it should be considered as corpus-driven in Tognini-Bonelli's sense.

Against this background, it is possible to raise the most relevant questions that guide the research in the remaining chapters:

- How can the structure of medical vocabulary in English be determined on the basis of a specialized corpus?
- How does the choice of a particular perspective (pedagogical versus characterizing/descriptive) influence the methodology of corpus-based research?
- How does the text type influence the structure of medical vocabulary?
- How does the choice of a corpus influence the results?

The monograph is divided into four chapters and a conclusion. After this introduction, chapter 2 presents an overview of previous efforts to characterize and determine medical English. Defining the key notions of *lemma/lexeme*, *word family*, *specialized vocabulary*, and *terminology* and

the relationships between them are central issues which are addressed. Then, an overview of methods relevant for identifying ESP vocabulary with an emphasis on medical English is given. The role of corpora and specialized corpora in determining the vocabulary of medical English is discussed in detail. Word lists of academic vocabulary by Coxhead (2000) and Gardner and Davies (2013) and of medical vocabulary by Wang et al. (2008) are described. All these word lists are based on specialized corpora of academic texts and medical research journal papers. The methodologies applied in these word lists are compared and critically evaluated.

In chapter 3 I argue that the methodology used in a pedagogical approach, which results in medical word lists, is neither sufficient nor adequate if the main aim is to characterize or describe medical vocabulary and modifications of methodology are suggested. First, the chapter explains why it is reasonable to use The Corpus of Contemporary American English (COCA) to find answers to the above mentioned research questions. Then, the chapter concentrates on the description of the medical subcorpus ACAD: Medicine in COCA. It discusses how the structure of the medical corpus influences the characterization of medical vocabulary. Finally, an overview of the procedure applied to arrive at the characterization of medical vocabulary is presented. It explains why it is better to approach medical vocabulary from the perspective of a cline or continuum based on two dimensions: absolute and relative frequency. Determining the threshold values for each of these dimensions is a crucial decision. It also demonstrates what effect the different threshold values might have on the structure or description of medical vocabulary in English. The chapter concludes by presenting a model of medical vocabulary as a two-dimensional continuum based on the interaction of absolute frequency and relative frequency.

Chapter 4 compares the results based on the subcorpus of medicine in COCA with an alternative corpus of medical texts, a specially compiled corpus for illustrative purposes. This medical corpus is based on the Wikipedia corpus, which was made available as a supplement to COCA in 2015. The Wikipedia corpus is based on the full text of the English version of the Wikipedia at a particular point in time and it contains 4.4 million Wikipedia articles with 1.9 billion words. The Wikipedia articles on medical topics represent a different type of medical text to medical journal articles. The chapter compares the results based on the two specialized corpora and evaluates their usefulness with respect to the characterization of medical vocabulary in English.

Finally, the conclusion summarizes the most relevant findings of the previous chapters and indicates their significance. On one hand, it is

argued that the perspective to ME adopted here contributes to a better understanding of language use in medical communication. On the other, lines of further research are outlined.

# CHAPTER TWO

# DETERMINING THE VOCABULARY
## OF MEDICAL ENGLISH

A central question in any subfield of English for Specific Purposes (ESP) is how it relates to the lexicon. Johns (2013: 23) points out that this issue has been discussed since the early years (1961-1982) of the history of the ESP research by the *Washington School*. The main representatives of this school, John Lackstrom, Larry Selinker, and Louis P. Trimble, made the relationships of the grammar and lexicon of English for science and technology with the authors' rhetorical purposes central in their research (see e.g. Lackstrom et al. 1972). Since then, it has continued to be in the focus of ESP research. Although in most ESP research, the main aim is pedagogical, with emphasis on an accurate definition of which vocabulary ESP learners need for their professional communication, the approach also triggers interesting theoretical questions about the lexicon. These are the main focus of this chapter.

The terminology that is used to refer to ESP vocabulary includes the terms *specialized*, *technical*, *sub-technical*, and *semi-technical vocabulary* (Coxhead, 2013: 141). The term *sub-technical* vocabulary is used by Cowan (1974). Farrell (1990) prefers the term *semi-technical vocabulary*. According to Coxhead "such terms usually refer to the vocabulary of a particular area of study or professional use" (2013: 141). This shows the importance of defining the key notions of *general vocabulary*, *specialized vocabulary*, and *terms in the narrow sense* and the relationships between them. The definition of these terms is addressed in section 2.1. Section 2.2 gives an overview of methods relevant for identifying ESP vocabulary with an emphasis on medical English. It discusses in detail the role of corpora in determining the vocabulary of medical English. In 2.3 the importance of corpus-based methods applied in identifying ESP vocabulary is emphasized and Coxhead's (2000) Academic Word List (AWL) is described in more detail. In section 2.4 I will present the Medical Academic Word List (MAWL) by Wang, Liang and Ge (2008). A

more recent contribution, the New Academic Vocabulary List (AVL) by
Gardner and Davies (2013), is discussed in 2.5.

## 2.1 Defining specialized vocabulary

The vocabulary of medical English clearly belongs to a specialized
professional area. This means that non-professionals might not have
knowledge of medical vocabulary or at least of the specialized senses of
vocabulary items relevant in a medical professional environment. On the
other hand, there is a certain degree of overlap with general vocabulary.
This raises crucial questions about specialized vocabulary and its
relationship to words and terms.

Coxhead (2013: 141) emphasizes that "the range of a word is important
in ESP. That is, a specialized word would have a narrow range of use
within a particular subject area". She distinguishes three types of
specialized words: words of Greek or Latin origin, highly technical words,
and words used also in general language (Coxhead, 2013: 141). These
three different types of specialized words are exemplified in (1).

(1)      a. malleus
          b. trocar
          c. jacket

The first type of specialized words based on Coxhead (2013) includes
words with Greek or Latin elements. It is exemplified in (1a). The
specialized word (1a) is of Latin origin and means hammer in the sense of
'the largest ossicle of the three auditory ossicles' (Stedman, 1997). The
word in (1b) is a highly technical word and it represents the second class
of specialized words distinguished by Coxhead (2013). The meaning of
(1b) is 'an instrument for withdrawing fluid from a cavity, or for use in
paracentesis' (Stedman, 1997). A corpus search confirms that (1b) is a
highly specialized technical word, the Corpus of Contemporary American
English (COCA) gives 12 occurrences, and the British National Corpus
(BNC) only 1.[1] This also means that only experts are likely to store the
meaning of (1b) in their mental lexicon. In (1c) we see the third type in
line with Coxhead (2013), a word which is used in much narrower senses
in medicine than in general English. For (1c), Stedman (1997) gives two
senses typical of medicine. In one sense it may mean 'a fixed bandage
applied around the body in order to immobilize the spine' (Stedman, 1997)

---

[1] Corpus results were retrieved on 30 July, 2015 from COCA.

whereas in dentistry, it means 'an artificial crown composed of fired porcelain or acrylic resin' (Stedman, 1997).

Especially the third type of specialized vocabulary, exemplified in (1c) has been the main object of a number of research studies such as those by Crawford Camiciottoli (2007), Nation (2008), etc. The top ten word list in business studies by Crawford Camiciottoli (2007) includes *price*, *work*, and *market*, which are frequently used also in common contexts of general language use. In medical vocabulary, Nation (2008) reports that *neck* and *by-pass* occur frequently. However, they are also frequent in general lexis but in different senses, e.g. a *city by-pass* or a *bottleneck*. According to Coxhead (2013: 151), the question of polysemy of ESP and its vocabulary is a challenging issue in a pedagogical perspective. In her view, "new technical meaning requires […] learners to build their knowledge of both the concept of a word and its meaning" (Coxhead, 2013: 151).

Many ESP researchers identified another essential problematic question related to specialized vocabulary. Specialized vocabulary is dynamic and develops rapidly. It is important that the fast progress in specialized vocabulary development is reflected in teaching material. This is the main reason why, for instance, Crawford Camiciottoli (2007) questions the correspondence of specialized vocabulary between professional texts and university level texts.

A different view of specialized words combining lexicographic and terminological perspectives can be found in ten Hacken (2008, 2010, 2015). He discusses the relationship between *general vocabulary*, *specialized vocabulary*, and *terms*. Ten Hacken's approach is based on a theory of prototypes (e.g. Labov, 1973) and preference rules formulated by Jackendoff (1983). Labov's experiment with the concept of *cup* is a classical demonstration of the fact that a judgement whether a particular object is a cup or not is prototype-based. The informants had a stronger tendency to reject the label cup for an object which was further removed from the prototype. Scalar conditions and preference rules determine the distance from the prototype. Ten Hacken (2010: 917) gives the height-width relation as an example of a scalar condition in the case of cup and the presence of a handle exemplifies a preference rule. It is obvious that preference rules interact with scalar conditions in the sense that if the object has a handle, "it can be further removed from the prototypical height-width relation and still be judged a cup" (ten Hacken, 2010: 917).

Let me now turn to the consequences of the assumption of prototypes for the distinction between general words, specialized words, and terms in ten Hacken's perspective. Both general words and specialized words are based on prototypes. The difference between the two is that the latter is a

label for expressions "used only in specialized language" (ten Hacken, 2010: 918). The example in (1c) is a case in point. The meaning of (1c) in general language is distinct from its specialized meaning in medicine. This also means that specialized words "are in the mental lexicon of a much smaller group of speakers" (ten Hacken, 2015: 6) as opposed to general lexis. It is reasonable to assume that (1c) in the sense of 'an outer garment for the upper part of the body' must be stored in the mental lexicon of most speakers of English. However, retrieving its specialized meaning of 'a fixed bandage applied around the body in order to immobilize the spine' requires a specialized context, familiar to a much smaller number of speakers.

Ten Hacken (2008, 2010) observes that two conditions can be used for defining terms, specialization and the precise delimitation of the extension. He emphasizes the different nature of these conditions. Specialization represents a scalar condition, whereas "having a precisely delimited extension produces a dichotomy" (ten Hacken, 2010: 917). To put it differently, with the former we can decide where the cutoff point is, but for the latter we have to select one or another. According to ten Hacken (2010: 917), only expressions with precisely delimited meanings can be labelled as terms in the narrow sense. A direct consequence of the fact that the conditions for specialized words and terms are independent is that the overlap of the two categories is possible without triggerring any problems.

The overlap of the two categories in certain contexts, specialized vocabulary items and terms, brings us to another crucial distinction ten Hacken (2010, 2015) makes; the difference between *terms* and a subset of terms he labels *terms in the narrow sense*. Only the latter can be made distinct from specialized vocabulary in the sense of terminological definition proper. This is illustrated in (2).

(2)    **trocar** – an instrument for withdrawing fluid from a cavity, or for use in paracentesis; it consists of a metal tube (cannula) into which fits an obturator with a sharp three-cornered tip, which is withdrawn after the instrument has been pushed into the cavity; the name t. is usually applied to the obturator alone, the entire instrument being designated t. and cannula. (Stedman, 1997)

First, the definition in (2) may seem as an example of a classical terminological definition. It classifies the instrument within a particular category, or a class of objects and specifies its typical properties. However, the definition in (2) does not determine precise boundaries consisting of necessary and sufficient conditions relevant for (1b). The

concept remains prototype-based and there is no need to impose clear-cut boundaries in order to determine whether a particular instrument is a trocar as in (1b) or not. Thus, it may be argued that (2) represents a well-formed lexicographic definition taken from a specialized medical dictionary. In (2), 'an instrument' fulfils the function of the hyperonym. A detailed description of the relevant parts specifies material, shape, and purpose. This information is an example of scalar conditions. *Usually* in the final part of the definition in (2) indicates a preference rule. It allows a user to select whether he wishes to refer to the instrument as a whole or only to a specific component. Although *trocar* is used only in specialized contexts, which makes it distinct from any general vocabulary item, similarly to natural concepts it is based on a prototype. a consequence is that *trocar* in (1b) is an example of a specialized vocabulary item and a term, but not of what ten Hacken (2010, 2015) labels as a term in the narrow sense.

Ten Hacken (2015: 7) demonstrates that "the distinction between terms (in the narrow sense) and specialized vocabulary is determined by the need to resolve conflicts. Unless there is such a need, we can continue to use prototypes, which correspond to the natural state of concepts." The usefulness of the difference between terms in the narrow sense and specialized vocabulary items arises in contexts where it is necessary to adopt clear boundaries of the concept in contrast to a continuum. Ten Hacken (2015) identifies two such contexts, legal disputes and scientific theories. The example in (3) illustrates the former.

(3)    The term "drug" means

(A) articles recognized in the official United States Pharmacopoeia, official Homoeopathic Pharmacopoeia of the United States, or official National Formulary, or any supplement to any of them; and

(B) articles intended for use in the diagnosis, cure, mitigation, treatment, or prevention of disease in man or other animals; and

(C) articles (other than food) intended to affect the structure or any function of the body of man or other animals; and

(D) articles intended for use as a component of any article specified in clause (A), (B), or (C). A food or dietary supplement for which a claim, subject to sections 343 (r)(1)(B) and 343 (r)(3) of this title or sections 343 (r)(1)(B) and 343 (r)(5)(D) of this title, is made in accordance with the requirements of section 343 (r) of this title is

not a drug solely because the label or the labeling contains such a claim. A food, dietary ingredient, or dietary supplement for which a truthful and not misleading statement is made in accordance with section 343 (r)(6) of this title is not a drug under clause (C) solely because the label or the labeling contains such a statement.

The source of the definition in (3) is the *Code of Laws of the United States of America*, commonly abbreviated to U.S. Code.[2] This document is the official compilation and codification of the general and permanent federal statutes of the United States.[3]

The definition in (3) is an example of a terminological definition. It delimits the precise boundaries of what is and what is not a *drug* by making necessary and sufficient conditons explicit. This means that *drug* as defined in (3) is a term in the narrow sense. On the basis of (3) it is possible to make a distinction between, for instance, a drug and a dietary supplement. Such a distinction is relevant in legal contexts and provides guidance in resolving legal disputes. An example of a relatively current issue is the case of Cholestin. The precise categorization of the product as a dietary supplement or drug was the main issue in a federal court case (cf. Havel, 1999; Heber et al., 1999).

Another type of context illustrating the need for a more precise terminological definition required for terms in the narrow sense is scientific theory. Ten Hacken (2010: 920) points out that terms in empirical science classify entities in the real world. This explains why theories in empirical science need a certain degree of precision for claims to be testable. In a medical context, this may be illustrated with different types of incisions. In (4), two definitions of *incision* are given.

(4)     a.  The action of cutting into something; esp. into some part of the body in surgery. (OED, 2015)

---

[2] Available at available at http://uscode.house.gov/, accessed on 2 August, 2015
[3] According to Wikipedia (United States Code, 2 August, 2015), the main edition is published every six years by the Office of the Law Revision Counsel of the House of Representatives, and cumulative supplements are published annually. The U.S. Code is organized in 52 titles, Food and Drugs can be found under title 21. The basic structure of the titles includes sections. The sections are numbered sequentially across the entire title without regard to the previously-mentioned divisions of titles. Frequently, the sections are further structured into subsections, paragraphs, subparagraphs, clauses, subclauses, items, and subitems. (3) is Title 21 › Chapter 9 › Subchapter II › § 321 (g).

   b. a cut; a surgical wound; a division of the soft parts made with   a
      knife. (Stedman, 1997)

The definitions in (4) can be described as examples of lexicographic
definitions. In (4a), the OED gives a general definition of a vocabulary
item that is likely to be used in common everyday situations. There is also
a good reason to assume this item is stored in the mental lexicon of a large
number of individual speakers of English. The definition in (4b) is slightly
more specific, it indicates a degree of specialization of this vocabulary
item. A precise delimitation of the boundaries in the sense of necessary
and sufficient conditions is not given. Thus, the concepts described in (4a)
and (4b) are prototype-based. This implies that *incision* is another example
of the two overlaps. First, the general vocabulary item overlaps with the
specialized word, and second, the specialized word with the term.
However, neither (4a) nor (4b) can be considered as a definition of a term
in the narrow sense.

   In the medical theoretical literature, a number of different types of
incisions are distinguished and sets of conditions specifying precisely
when a particular incision is the best with respect to healing. The so-called
*MacFee incision* is a good example to discuss. It is a type of incision used
for neck dissection (Werner and Davies, 2004). Detailed descriptions and
discussions of the MacFee incision can be found in books on theoretical
medicine examining clinical judgement and reasoning. *Metastases in Head
and Neck Cancer* by Werner and Davies (2004) is an example of such a
book. It summarizes the types of health problems related to head and neck
cancer, explains and describes a range of methods for their treatment and
evaluates them with respect to a set of criteria. With the incisions used for
neck dissection, nine evaluation criteria for the selection of a particular
type are listed. These include the tendency of necrosis of the detached skin
parts, the planned extent of the tumor intervention, the primary defect
coverage in cases of more extended skin resections, the blood supply of
the flaps, the overview of the entire operation field, the additional
performance of tracheotomy, the possible excision of existing scars, the
potential for avoiding skin incision when mucosal incisions suffice, and
the possibility of an extension of the incision if additional cervical lymph
node regions must be dissected (Werner and Davies, 2004).

   Then, the illustrations of the incisions are presented. The illustration of
the MacFee incision shows the procedural details. Individual steps
represent necessary and sufficient conditions combined in the definition of
the concept. The choice of the name is less relevant. This means that in
contrast to *incision*, *MacFee incision* is an example of a scientific term or

in ten Hacken's terminology, a term in the narrow sense. The term *MacFee incision* has precisely delimited boundaries. The illustrations make it possible to determine the type of incision immediately and unambiguously.

The illustrations are followed by a summary of their advantages and in some cases by a comparison with competing incisions. a detailed description of the advantages of the MacFee incision is given in (5).

(5)     "The so-called MacFee incision probably has the best chance of
        healing because this type of incision addresses the blood supply of
        the neck […]. It leads to very good esthetic results as long as the
        incisions are performed along skin lines, especially in pre-formed
        creases. Furthermore, this type of incision protects the carotid
        artery. The operative procedure is more difficult to perform in
        patients with short necks. Additionally, exposure of the operative
        field is often impaired so that intensive retraction by the assistant is
        required. The MacFee incision is preferred for patients suffering
        from a peripheral vascular disease or for patients who have
        undergone prior radiotherapy […]. It is often used in young patients
        undergoing neck dissection for thyroid cancer." (Werner and
        Davies, 2004, references deleted)

The description of the incision in (5) shows at least two important facts. First, it delimits which factors for the selection of a particular type of incision are essential for the MacFee incision, e.g. blood supply and the overview of the operation field. Second, it demonstrates that the type of incision, i.e. the concept identified by this name, includes a number of benefits for a patient. This becomes a part of theorizing about the best practice in a particular context. The overview of the positive outcomes is based on a reasonable amount of empirical data, their collection, and evaluation. There may be discussion about advantages and disadvantages influenced by different empirical data, but this does not mean that the scientific term has fuzzy boundaries.

It is interesting to compare Coxhead's understanding of specialized vocabulary with ten Hacken's interpretation. An essential similarity is immediately obvious. Both linguists agree on the fact that specialized vocabulary requires a specialized context. Coxhead (2013) uses this label to cover general vocabulary items with a narrowed meaning, words of Greek and Latin origin, and highly technical words. These types were exemplified in (1). This typology is sufficient in the context of ESP for pedagogical purposes. The questions related to the representation of