# Journal Pre-proof

Design, Development and Validation of a System for Automatic Help to Medical Text Understanding

Marco Alfano, Biagio Lenzitti, Giosuè Lo Bosco, Cinzia Muriana, Tommaso Piazza, Giovanni Vizzini

Please cite this article as: Alfano M, Lenzitti B, Bosco GL, Muriana C, Piazza T, Vizzini G, Design, Development and Validation of a System for Automatic Help to Medical Text Understanding, *International Journal of Medical Informatics* (2020), doi: https://doi.org/10.1016/j.ijmedinf.2020.104109

# Design, Development and Validation of a System
# for Automatic Help to Medical Text Understanding

## Marco Alfano[1,8] , Biagio Lenzitti[2], Giosuè Lo Bosco[2,3], Cinzia Muriana[4,6], Tommaso Piazza[7], Giovanni Vizzini[5,6]

1  Lero, Dublin City University, Dublin, Ireland

2 Dipartimento di Matematica e Informatica, Università di Palermo, Palermo, Italy

3 Dipartimento di Scienze per l'Innovazione tecnologica, Istituto Euro-Mediterraneo di Scienza e Tecnologia, Palermo, Italy

4 IRCCS-ISMETT (Istituto Mediterraneo per i Trapianti e Terapie ad Alta Specializzazione), Palermo, Italy

5 Department for the Treatment and Study of Abdominal Diseases and Abdominal Transplantation, IRCCS-ISMETT (Istituto Mediterraneo per i Trapianti e Terapie ad Alta Specializzazione), Palermo, Italy

6 University of Pittsburgh Medical Center (UPMC) Italy, Palermo, Italy

7 Università Ca' Foscari, Venezia, Italy

8 Anghelos Centro Studi sulla Comunicazione, Palermo, Italy


## Corresponding author
Marco Alfano
Lero, Dublin City University, Dublin, Ireland
marco.alfano@lero.ie
marco.alfano@anghelos.org
+39 3386366188

Higlights

- SIMPLE used for patient empowerment and for health literacy improvement

- SIMPLE helps medical text comprehension

- SIMPLE identifies medical terms, translates into consumer terms and adds explanations

- SIMPLE works with different languages. English and Italian already implemented.

**ABSTRACT**

**Objective**: The paper presents a web-based application, SIMPLE, that facilitates medical text comprehension by identifying the health-related terms of a medical text and providing the corresponding consumer terms and explanations.

**Background**: The comprehension of a medical text is often a difficult task for laypeople because it requires semantic abilities that can differ from a person to another, depending on his/her health-literacy level. Some systems have been developed for facilitating the comprehension of medical texts through text simplification, either syntactical or lexical. The ones dealing with lexical simplification usually replace the original text and do not provide additional information. We have developed a system that provides the consumer terms alongside the original medical terms and also adds consumer explanations. Moreover, differently from other solutions, our system works with multiple languages.

**Methods**: We have developed the SIMPLE application that is able to automatically: 1) identify medical terms in a medical text by using medical vocabularies; 2) translate the medical terms into consumer terms through medical-consumer thesauri; 3) provide term explanations by using health-consumer dictionaries. SIMPLE can be used as a standalone web application or can it be embedded into common health platforms for real time identification and explanation of medical terms. At present, it works with English and Italian texts but it can be easily extended to other languages. We have run subjective tests with both medical experts and non-experts as well as objective tests to verify the effectiveness of SIMPLE and its simplicity of use.

**Results**: Non-experts found SIMPLE easy to use and responsive. The big majority of respondents confirmed they were helped by SIMPLE in understanding medical texts and

declared their willingness to continue using SIMPLE and to recommend it to other people. The subjective tests, conducted with medical experts on a set of Italian radiology reports, showed an agreement between SIMPLE and the experts, on the highlighted medical terms, that ranges between 74.05% and 81.16% as well as an agreement of around 60% on the consumer term translation. The objective tests showed that the consumer terms, provided by SIMPLE, are, on average, eighteen times more familiar than the relative medical terms so proving, once more, the effectiveness of SIMPLE in simplifying the medical terms.

**Conclusions**: The performed tests demonstrate the effectiveness of SIMPLE, its simplicity of use and the willingness of people in continuing with its use. SIMPLE provides, with a good agreement level, the same information that medical experts would provide. Finally, the consumer terms are 'objectively' more familiar than the related technical terms and as a consequence, much easier to understand.

**Keywords**: e-health, patient empowerment, lexical simplification, consumer health vocabulary, term familiarity, infobutton.

**1. Introduction**
The interest of patients towards health information is increasing as new information systems, like electronic health records (EHRs) or personal health records (PHRs), involve patients in the healthcare process. The possibility of electronically storing healthcare events allows them to recover information wherever and whenever needed so to participate as 'empowered' users. In approaching healthcare information, they employ their knowledge base characterized by

informal terms rather than medical jargon and, thus, they often find medical texts difficult to understand. In fact, the comprehension of a medical text requires semantic abilities that can differ from a person to another based on his/her health-literacy level and considering that health documents often contain abbreviations, neologisms and words of Latin and Greek origin. Studies report that even brochures written for patients or web sites dealing with medical issues can be difficult to understand [1-2]. Patients have a particular difficulty with test results, radiology reports, and medication lists, mainly for what concerns medical terminology and abbreviations [3]. Therefore, they either surf the Internet for terms explanation or ask for help to physicians [4] but, often, they do not find a prompt response to their needs.  As a consequence, [5] stressed the importance of improving patient health literacy and, at the same time, of writing easy-to-understand medical texts.

Text Simplification (TS) "aims to rewrite sentences so to reduce their syntactic or lexical complexity while preserving their meaning" [6]. In particular, lexical simplification (LS) identifies and replaces complex terms with simpler ones, while syntactic simplification identifies grammatical complexities in a text and rewrites them by using simpler structures [7]. There are many papers in the literature related to generic TS automation [7-11]. Others specifically address TS in a medical context [12-16]. There are TS tools that accept a document, with an original reading level, and convert it to a target reading level by using Natural Language Processing (NLP) techniques [8]. Other tools identify difficult terms and replace them with simpler terms by using the Unified Medical Language System (UMLS) and the Open-Access Collaborative Consumer Health Vocabulary (OAC CHV) [17]. There are ubiquitous web-based text simplification systems that use a term frequency criterion to replace technical terms with common ones and others that use a lexical density criterion to restrict the use of synonyms (for

Spanish language) [18]. The literature also presents semi-automatic methods for simplification of English medical texts [15]; tools that rely on MeSH and SNOMED-CT vocabularies for word replacement (for Swedish language) [19]; tools that are equipped with training, simplification and ranking modules [20]; and tools that use substitution selection and substitution ranking approaches [21]. Furthermore, some researchers have investigated alternative methods for synonym replacement in automatic TS, based on word length, frequency and level of synonymy [22] and others have used supervised learning methods for lexical simplification and dependency-based word embedding for replacing words with similar ones (for Japanese language) [23]. None of the listed works have considered the importance of adding information in consumer's language for facilitating understanding and improving the health literacy of patients.

We believe that, when dealing with medical texts, the greatest difficulty laypeople encounter is the understanding of medical terms rather than grasping the general meaning of the sentences, even the complex ones. For this reason, we have decided to focus on the lexical aspect of text simplification rather than on the syntactic one. Moreover, we believe that it is important for a non-expert to get both a translation of the technical term in its simple counterpart and additional information on the term (such as its explanation) that allows him/her to really understand the meaning of the term (also when the simplified term is not available). Finally, we consider of paramount importance that the meaning of a medical term is univocally determined so to eliminate any ambiguity and to use the same methodology and system with different languages.

This paper presents an automatic tool, called SIMPLE, targeted to create text simplification of health/medical documents. It recognizes health-related terms and provides lexical simplification and additional information based on medical vocabularies, thesauri and dictionaries. SIMPLE

has mainly been designed to 'empower' patients or, in general, laypeople but it can also be useful for people with different levels of expertise. It uses English medical vocabularies/thesauri/dictionaries, as knowledge base, since it is the most diffused and complete, but it has been built to operate with other languages, such as Italian, for which the knowledge base is more limited.

SIMPLE has been implemented through a web application that accepts, as input, any medical text. The output presents the same text, with the technical terms highlighted and one infobutton next to each technical term. By clicking on the infobutton, the term translation in common language (consumer term) appears together with the term explanation. Our system, to the best of our knowledge, constitutes one of the first solutions in this field. In fact, infobuttons have been used for years mainly to support clinicians' decisions and only recently, they have been used to bring information to patients [24-25]. Moreover, there are a couple of systems that present some similarities with ours [13, 17], but they replace medical terms with the consumer ones and add further information within the text, altering the original text. Our system provides the consumer information alongside the original text leaving it intact. Moreover, our system provides term explanations in consumer language whereas the other two systems do not. Finally, our system works with different languages. English and Italian are already implemented but other languages can be easily added.

## 2. Materials and Methods

### 2.1 Description of SIMPLE

Laypeople reading medical texts often need some external help to understand the technical terms. It usually comes in the form of different resources (online or not) such as vocabularies, dictionaries and thesauri. In particular:

- A 'medical vocabulary' is a selective list of words and phrases used in the medical field; it can be used to find the technical terms in a medical text;

- A 'medical-consumer thesaurus' contains synonyms and antonyms of medical terms; it can be used to find consumer synonyms of the technical terms;

- A 'health-consumer dictionary' gives information about the meaning of the words; it can be used to find additional information on the technical terms in a simple language.

In some cases, a single resource can have multiple functionalities, e.g., it can contain both definitions and synonyms. Of course, there are numerous resources of each type.

SIMPLE is designed in order to automatically find the medical/technical terms (words or combination of words) in an online medical document, translate them in simple or consumer terms and provide additional information in simple language. The architecture of SIMPLE is shown in Fig. 1 and presents the following three main modules:

1. The *HIGHLIGHT* module takes an arbitrary text as input, uses a medical vocabulary to find the technical terms and highlights them when a consumer term and/or a consumer explanation exists. Moreover, an infobutton with an information icon (i) is put next to each highlighted item. When clicked, it will provide the consumer translation and/or the explanation in a tooltip next to the item.

2. The *MAP* module connects the technical terms found by the HIGHLIGHT module to the equivalent consumer terms by using a medical thesaurus.

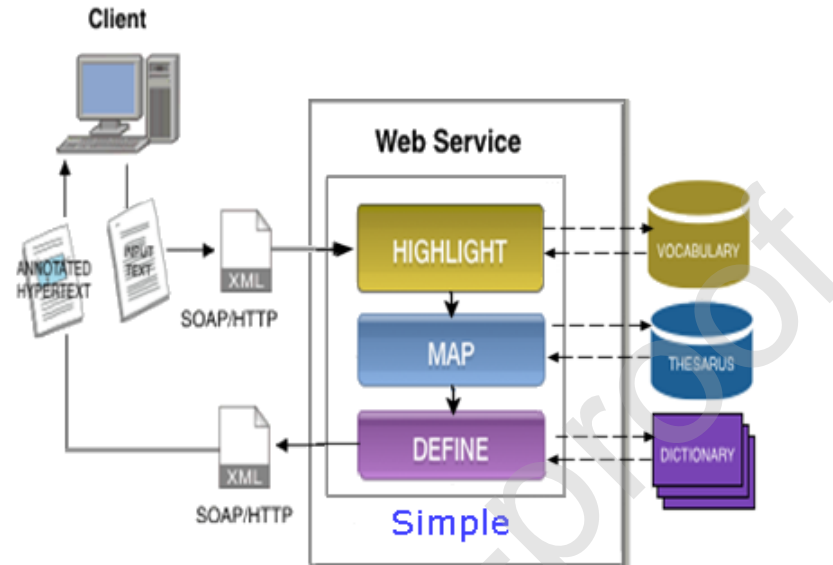3.  The *DEFINE* module provides a simple explanation of the term by using a consumer medical dictionary.



Figure 1. SIMPLE architecture.

The implementation of SIMPLE has been carried out by using a SOAP/HTTP approach. In particular, a client makes a request to a web service — written in php and javascript and using a MySql database — which receives the text to be processed and performs the work by means of the HIGHLIGHT, MAP and DEFINE modules. We now provide some details on the implementation of the three modules. Further information on the SIMPLE architecture and implementation can be found in [26-28].

**2.1.1 The HIGHLIGHT Module**

For the HIGHLIGHT module, we mainly use the 'Unified Medical Language System' (UMLS)[1] created and maintained by the US National Library of Medicine. It is a large collection of multilingual vocabularies that contains information about biomedical and health-related

---

[1] https://www.nlm.nih.gov/research/umls/

concepts. It uses a 'Concept Unique Identifier – CUI' to create a unique identifier for each concept and a mapping among vocabularies, thus allowing translation among the various terminology systems. Since our system has been designed to work with medical texts written in English and/or Italian, the HIGHLIGHT module uses the vocabularies of UMLS that have both English and Italian versions, namely the 'Medical Subject Headings' (MeSH), the 'International Classification of Primary Care' (ICPC), the 'Medical Dictionary for Regulatory Activities Terminology' (MedDRA) and the 'Metathesaurus Version of Minimal Standard Terminology Digestive Endoscopy' (MTHSMS). Notice that the CUI identifier allows immediate translation between English and Italian terms.

It is worth mentioning that while the English version of the UMLS provides a quite complete and updated list of medical terms, other versions, such as the Italian one, have a restricted subset of terms. For this reason, we additionally considered other English/Italian medical vocabularies so to increase the terms found by the HIGHLIGHT module but, at the same time, we avoided using too many vocabularies that would overload the system. Overall, we have built a metavocabulary with around 125,000 English/Italian entries but we have found, so far, consumer terms and/or explanations for around 80,000 terms that are the terms presently highlighted by SIMPLE. Nevertheless, we maintain the larger metavocabulary and we are in the process of adding further consumer terms and/or explanations.

Notice that, the modular nature of SIMPLE allows easy addition and removal of vocabularies, as well as of thesauri and dictionaries of different languages. This allows creating a tailored system for specific needs such as particular medical fields or languages.

**2.1.2 The MAP Module**

Consumer health vocabularies (CHVs) have been used/created for the MAP module in order to translate medical terms into their corresponding terms for consumers [29]. One of the best-known examples of CHV is the 'Open Access Collaboratory Consumer Health Vocabulary (OAC-CHV)', created and maintained by the Consumer Health Vocabulary Initiative [30]. It is a relationship file that links commonly used terms to associated medical terminology represented by the UMLS. The OAC-CHV focuses on expressions and concepts that are employed in health-related communications from or to consumers and contains around 160,000 rows (one for each term) and different fields among which:

- 'Term': The term as found in the text;

- 'Concept Unique Identifier' (CUI): The unique identifier of a concept as found in the UMLS;

- 'CHV Preferred Name': The preferred consumer term as defined in the Consumer Health Vocabulary.

The MAP module uses the OAC-CHV as a thesaurus. The mapping from technical to consumer terms is accomplished by means of the CUI when available [31]. This is the case of the technical terms that are found in the UMLS. For the terms without CUI, we have created a custom concept identifier.

Notice that the OAC-CHV is written in English and there is no Italian CHV available, beside the 'Italian Consumer-oriented Medical Vocabulary' (ICMV) that only contains a few items [32]. We have then translated the OAC-CHV (with its 160,000 entries) from English to Italian using English-Italian medical glossaries and automated translation tools.

**2.1.3 The DEFINE Module**

For what concerns the DEFINE module, we have analyzed many health-consumer web sites that often contain health and medical dictionaries specifically created for health consumers and then use a language that can easily be understood by them. We have then used the Italian dictionaries for health consumers and have got around 15,000 entries. Moreover, we have added some entries from some slightly more technical dictionaries in order to be able to provide as many definitions as possible, so creating a metadictionary of around 100,00 entries. For the English definitions, we have used the 28,000 entries of the WebMD online dictionary. Overall, the combination of the metavocabulary, metathesaurus and metadictionary presently presents, as seen above, around 80,000 technical terms that have consumer terms and/or explanations.

**2.1.4 SIMPLE Graphical Interface**

The English SIMPLE client can be accessed at the address http://www.math.unipa.it/simplehealth/simple2/ and the Italian client at the address http://www.math.unipa.it/simplehealth/ita/simple/. It is mainly used by a non-expert (e.g., a patient) with a medical text (e.g., medical report or lab result) that contains one or more technical terms he/she does not understand. The web interface presents a text area — where it is possible to insert any medical text — and a selection of preloaded medical reports found on the Internet (mainly used for testing purposes).

Fig. 2.a shows a bone density report, chosen among the preloaded reports. Fig. 2.b shows the same report after being processed by SIMPLE. Next to each technical term, we find an infobutton, i.e., an information-icon that, when clicked, shows the consumer term and explanation in a tooltip next to the term. In the figure, the 'DEXA' term is selected and the consumer translation (in this case the acronym meaning, i.e., Dual-energy X-ray absorptiometry) and an explanation ("a means of measuring bone mineral density – BMD") are shown into the

tooltip. This is a typical example of how SIMPLE facilitates understanding of an acronym that is

likely to be unknown by most non-expert users.



(a)                                                    (b)

Figure 2. (a) Original medical report and (b) processed report with highlighted medical terms.

Notice that, as said above, SIMPLE does not create any change in the original text (by replacing

words or inserting explanations into the text) because, in our opinion, this could disorient the

user. It only provides a translation and additional info (in a tooltip) on request, leaving the user

fully in charge of his/her navigation path through the text as it was originally created.

## 2.2 Evaluation of SIMPLE

In order to measure the effectiveness of SIMPLE, we ran different types of tests, i.e., subjective

tests with non-expert and expert users and objective tests. The next subsections present a

description of the evaluation process.

### 2.2.1 Subjective tests with non-expert users

After obtaining the ethical approval for our tests, we asked a number of Italian non-expert users

to use SIMPLE either by inputting any medical text or by using the preloaded medical reports, as

shown above. Moreover, we asked them to fill out a short questionnaire made of two parts, i.e.,

user information, as reported in Table 1, and an agreement with the statements reported in Table

2. The agreement was expressed in a 1-5 Likert-type scale as follows:

1 = strongly disagree, 2 = disagree, 3 = neither agree nor disagree, 4 = agree, 5 = strongly agree

|  | **REQUESTED INFORMATION** |
|---|---|
| **Information about user** | Age: 18-35, 36-50, 51-70, 70+ <br> Sex: M or F <br> Education level: Secondary school, High school, Bachelor/Master, PhD <br> Computer skills: None, Poor, Average, Good, Excellent <br> Employment information: Employed, Unemployed, Retired <br> Contact with health system: None, Sporadic, Frequent, Continuous <br> Interest towards eHealth web apps: None, Poor, Average, Good, Excellent <br> Medical knowledge: None, Poor, Average, Good, Excellent |

Table 1. Questionnaire on user information.

| **CATEGORY** | **STATEMENTS** |
|---|---|
| **Perceived usefulness** | Using SIMPLE can improve the comprehension of medical texts |
| **Perceived ease of use** | I find the SIMPLE interface clear and easy to understand |
|  | I find easy to access additional information (consumer terms and definitions) provided by SIMPLE |
|  | I think that SIMPLE can be used without prior knowledge |
|  | I think that the explanations provided by SIMPLE are useful and understandable |
| **Trust** | As I understand it, I believe that SIMPLE provides me with correct information about medical texts |
| **Performance** | I find that SIMPLE provides an immediate answer to the input |

| | |
|---|---|
| **Intention to use** | I am willing to continue using SIMPLE |
| | I would recommend SIMPLE to my friends and acquaintances |
| **System features need** | **Input:** I think the feature of text input is necessary |
| | **Output:** I think the feature of providing consumer terms is necessary |
| | **Output:** I think the feature of providing term explanations is necessary |

Table2. Questionnaire on statements.

### 2.2.2 Subjective tests with expert users

As a further step on subjective evaluation, we have decided to use 160 Italian radiology reports provided by the IRCCS-ISMETT hospital with the purpose of verifying the accordance between automatic (SIMPLE) and human evaluations. IRCCS-ISMETT is a highly-specialized hospital located in Palermo whose acronym stands for 'Istituto di Ricovero e Cura a Carattere Scientifico - Istituto Mediterraneo per i Trapianti e Terapie ad Alta Specializzazione'.

IRCCS-ISMETT also provided a panel of ten healthcare professionals from different fields, such as Hepatology, Diabetology, Neurology, Biology and Psychology. The experts were specifically chosen among healthcare professionals with different specializations, in order to ensure a heterogeneous set of expertise for the experimental phase. Since radiologic reports had been chosen, experts of this field were not included to avoid that skills of experts on the specific medical field could influence the choice of medical terms. By doing so, the experts could be considered as expert users, i.e., people with sufficient skills to ensure full comprehension of the text and related terminology although with no specific expertise on the medical field under examination.

A specific web application was built on top of SIMPLE so to create a test environment for the experts and the 160 radiology reports were distributed among them in such a way that each report was tested by three different experts. As a consequence, each expert examined 48 reports and he/she was asked to:

- read each of the 48 radiology reports assigned to him/her;

- highlight those terms that he/she considered technical (medical);

- check the translation to a consumer term provided by SIMPLE, when available;

- confirm the given translation or suggest a more suitable one and, when no translation was available, add a new one.

### 2.2.3 Objective tests based on term familiarity

We performed objective tests by taking ten medical reports and computing the *term familiarity* of the medical terms and consumer terms, i.e., the number of google results related to each term [2, 33]. In particular, we measured the average of the term familiarity of the medical terms in the reports (found by SIMPLE) and compared it with the one obtained by the consumer terms provided by SIMPLE.

### 3. Results

Twenty-five users (with no specific expertise in the medical field) used SIMPLE and filled out the questionnaire reported in Table 1 and Table 2 of Section 2.2.1. SIMPLE was mainly used for understanding diagnoses, exam results, package inserts and web medical content. Fig. 3 shows the statistics on the requested user information (Table 1). We tried to have a heterogeneous audience but, in the end, a majority of young males performed the tests. The education level was quite high as well as the computer skills. The contacts with the health systems happened occasionally and the interest towards e-health resulted quite various. Finally, the medical

knowledge was mostly low or very low confirming the status of 'non-medical experts' of the
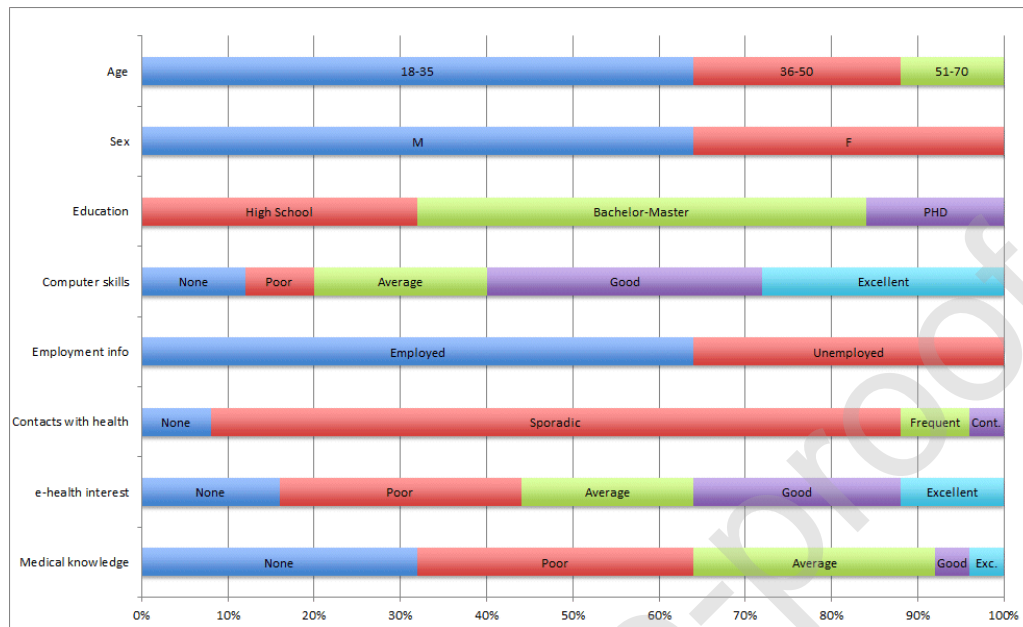
users.



Figure 3. User information statistics.

Fig. 4 shows the score shares obtained by SIMPLE for the different statements presented in

Table 2 by grouping them in three classes, i.e., 1-2, 3, and 4-5. Most users found SIMPLE easy

to use and responsive. The users reported that SIMPLE helped them understand medical texts

and found it useful in order to have access to both consumer terms and explanations. They

declared their willingness to continue using SIMPLE and recommend it to other people. The

results, though obtained with a small group of people, are encouraging and show the

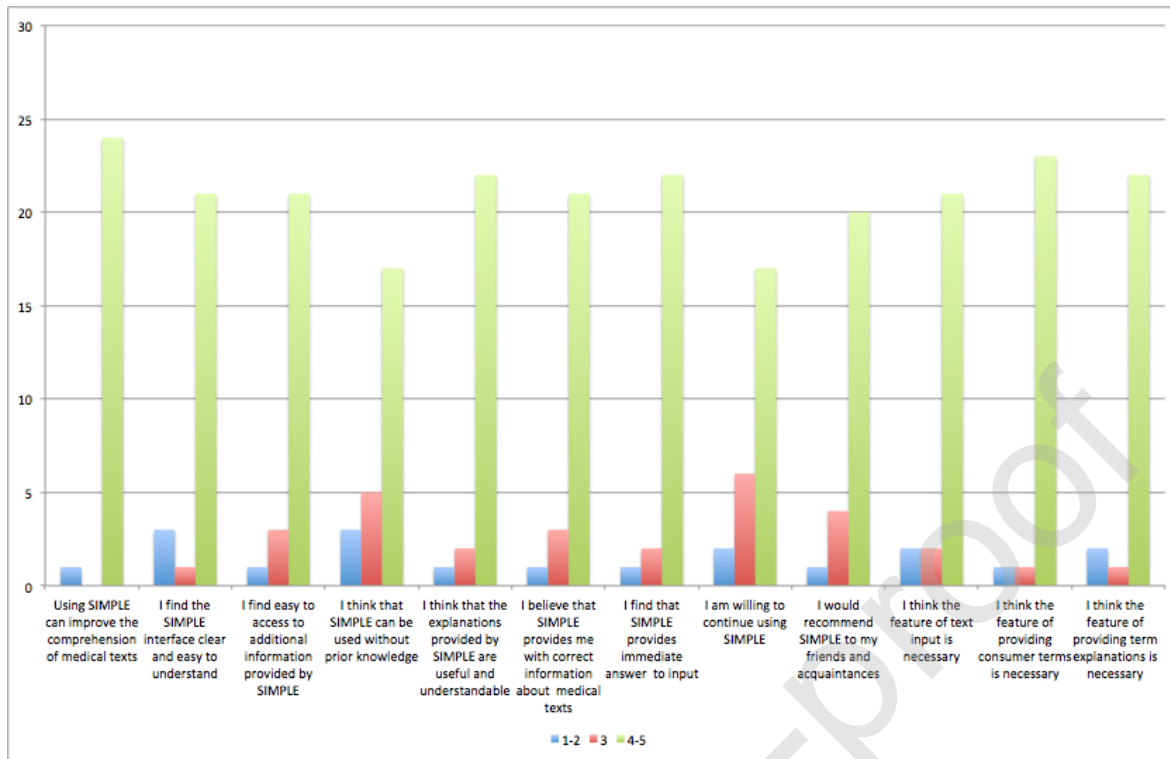effectiveness of SIMPLE and its simplicity of use.

Figure 4. Score shares for the different statements about SIMPLE.

For what concerns the subjective tests with expert users, presented in Section 2.2.2, Table 3 shows the results of the experiments in terms of number of unique words of all reports (either technical or non-technical), unique words indicated by the IRCCS-ISMETT experts as medical terms, unique words found by SIMPLE in its metavocabulary and unique words for which SIMPLE has a consumer term and/or an explanation. Furthermore, we counted the stop words (commonly used words) and recomputed the different items eliminating the stop words (last two columns of the table).

| | Total unique words | % (with respect to total unique words) | Stop words | Total unique words - stop words | % (with respect to total unique words - stop words) |
|---|---|---|---|---|---|
| | | | | | |

| | | | | | |
|---|---|---|---|---|---|
| All words | 8,065 | - | 151 | 7,915 | - |
| Words indicated by ISMETT experts | 1,104 | 13.69 | 21 | 1,083 | 13.68 |
| Words found by SIMPLE in metavocabulary | 4,590 | 56.91 | 128 | 4,462 | 56.37 |
| Words found by SIMPLE with consumer term and/or explanation | 3,883 | 48.14 | 117 | 3,766 | 47.58 |

Table 3. Experimental results in terms of total words and words indicated by experts and

SIMPLE.

By considering the total unique words without stop words, the experts indicated 1083 words as

technical ones out of a total of 7,915 unique words, whereas SIMPLE found 4,462 words in its

metavocabulary and highlighted 3,766 words, i.e., the ones for which it had a consumer term

and/or an explanation. The difference in the number of medical terms indicated by the experts

and SIMPLE can be partially explained by the fact that the experts would not consider many

terms being technical (given their expertise) whereas SIMPLE, by automatically using medical

vocabularies, would take whatever terms are in those vocabularies.

Table 4 shows the matching between the unique words indicated by the experts and the ones

found by SIMPLE. Again, we considered both the unique words found by SIMPLE in its

metavocabulary and the ones for which there is a consumer term and/or an explanation.

Moreover, we also counted the stop words and recomputed the different items eliminating the stop words (last two columns of the table).

| | Unique matching words | % (with respect to words indicated by experts) | Stop words | Unique matching words - stop words | % (with respect to words indicated by experts - stop words) |
|---|---|---|---|---|---|
| Experts & SIMPLE | 900 | 81.52 | 21 | 879 | 81.16 |
| Experts & SIMPLE with consumer term and/or definition | 822 | 74.55 | 21 | 801 | 74.05 |

Table 4. Matching words between experts and SIMPLE.

Considering the unique matching words without stop words, SIMPLE found 879 words in its metavocaburary and highlighted 801 words out of the 1083 words indicated by the experts, showing an accordance value between automatic and human evaluations respectively of 81.16% and 74.05%. The latter percentage could increase (80% and more) by considering that we are in the process of adding further consumer terms and explanations to the system. Moreover, the obtained result is 'improved' by considering a certain degree of variability among experts' evaluation since they did not always agree in evaluating what terms were medical and what terms were not (considering that each report was seen by three different experts).

The second step of the experiment was to evaluate the matching between experts and SIMPLE for what concerns the consumer terms. To this end, we considered the total number of terms for

which SIMPLE provided a consumer term, equal to 1,157. The experts confirmed the SIMPLE translation for 688 terms of those 1,157 terms and, for the remaining 469 terms, they provided different translations. As a consequence, we obtained a matching percentage of 688/1,157 = 59.5% that can be seen as a good result considering the high variability of choice of consumer terms by the experts.

For what concerns the objective tests based on term familiarity and described in Section 2.2.3, Figure 5 presents a 100% stacked bar chart that summarizes the results related to n=10 Italian medical reports. The numbers on the left of each horizontal bar are the terms that SIMPLE found on each of the 10 reports (e.g. 20 for the first medical report). Each blue bar reports the average term familiarity (numbers of google results in millions) of the medical terms on the medical report while the red bar reports the average term familiarity (numbers of google results in millions) of the corresponding consumer terms. The x-axis shows, for each stacked bar, the relative percentage of the two data series (medical terms and consumer terms) where the total of each bar always equals to 100%.

We can notice that the consumer terms are significantly more familiar than the related medical terms. In particular, the consumer terms are, on average, eighteen times more familiar than the relative medical terms showing, once more, the effectiveness of SIMPLE in simplifying the medical terms.
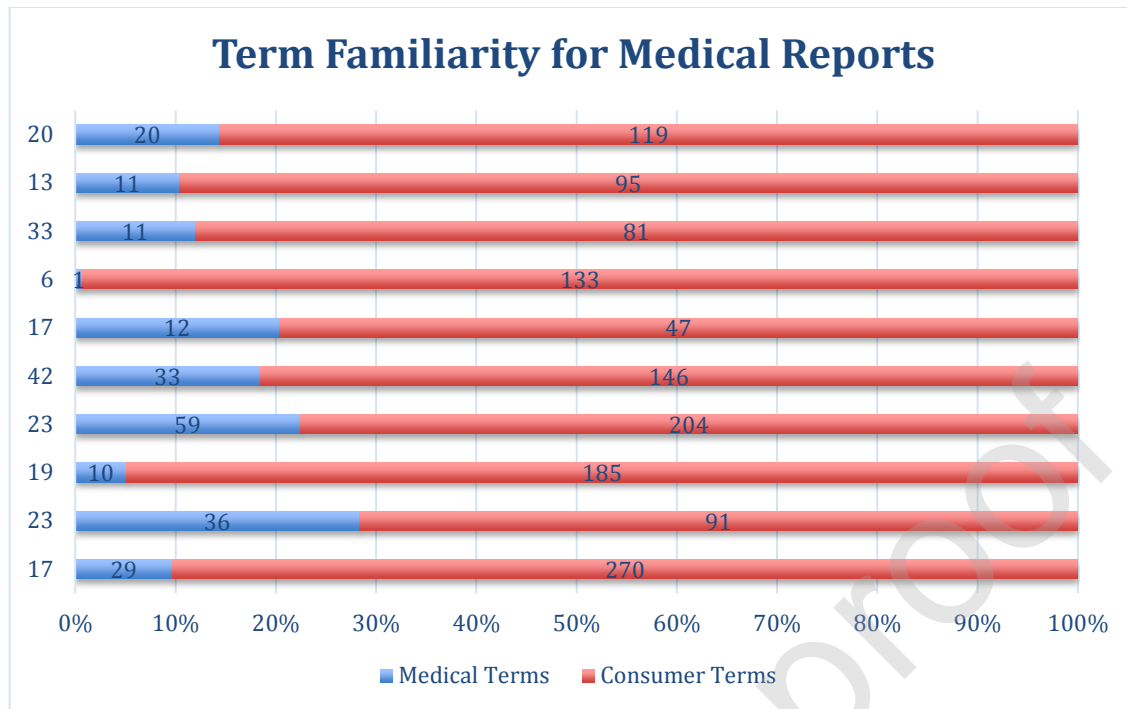
Figure 5. Term-familiarity averages related to medical and consumer terms for ten Italian

reports.

**CONCLUSIONS**

In this paper, we have presented SIMPLE: a system that, given a medical text, automatically

finds the technical terms, translates them in consumer terms and provides additional information

with the same kind of language. We have implemented the system, built a prototype and

validated it through different subjective and objective tests. Overall, the subjective tests have

confirmed the capability of SIMPLE in spotting the medical terms in a text and providing

consumer terms and explanations that facilitate the comprehension of the text. Moreover, users

declared their willingness to continue using SIMPLE and recommend it to other people.

Objective tests have also proved the effectiveness of SIMPLE in simplifying the medical terms

by showing a much higher familiarity of consumer terms compared to the ones of the related

medical terms.

It is important to underline that the main obstacle in performing lexical simplification is the availability of medical vocabularies, thesauri and dictionaries in different languages. In fact, while English versions of lexical simplification tools are broadly developed and are considered as 'de facto' standard in the medical field, usually only a subset of such tools is made available in other languages. Thus, beside showing the ability of a lexical simplification tool to make reading of medical texts easier, the challenge is to have the tool providing accurate 'translations' between different languages for which standard sources are not readily available and alternative sources are to be used. As a consequence, the testing environment with the Italian reports does not represent the 'optimum' for SIMPLE. From this point of view, the obtained results present a higher intrinsic value since SIMPLE makes use of translations of English terms (as we have seen for the Consumer Health Vocabulary – CHV) that may be inaccurate and may easily lead to mistakes and misunderstandings.

Of course, many improvements of SIMPLE can be realized. As a first priority, the medical vocabulary needs to be expanded to be able to find more technical terms present in the texts. Very important is the completion of the Italian CHV by adding more technical terms and their consumer equivalents. Finally, other health-consumer dictionaries have to be found for increasing the number of definitions that come from medical sources.

We plan to complete the prototype and integrate it within an Electronic Health Record (EHR) or Personal Health Record (PHR) by using the HL7 "Infobutton" standard. Moreover, we are in the process of creating a mobile app for different types of smartphones that, among others, allows the user to take a picture of the text (if written on paper) and directly load it into the system for further processing.

As an extension of the system, we plan to provide the user with additional information coming directly from the Web but tailored to the specific user requirements [34-36]. Moreover, we are in the process of developing a graphical framework to build health consumer-oriented advanced services and allow users to easily deal with health information [37].

**Authors' contributions**

All authors equally contributed to the article named "Design, Development and Validation of a System for Automatic Help to Medical Text Understanding".

**Statement on conflict of interest**

The authors declare that there is no confict of interests regarding the publication of the article named "Design, Development and Validation of a System for Automatic Help to Medical Text Understanding".

**Summary table**

- Understanding medical texts is a long and complex task for laypeople that are used to surf the Internet and search for medical information and/or to consult medical consumer vocabularies/dictionaries.

- Text Simplification (TS) can be applied to the medical context in order to reduce the syntactic or lexical complexity of a text while preserving their meaning.

- There are many papers in the literature related to generic TS automation, but none of them, when applied to the medical context, considers the importance of adding information in consumer's language to facilitate comprehension of the medical text.

- Other systems that deal with TS only work with a language (e.g., English) at the time.

- We have developed a system called SIMPLE, that is able to automatically: 1) identify medical terms in a medical text by using medical vocabularies; 2) translate medical terms into consumer terms through medical-consumer thesauri; 3) provide term explanations by using health-consumer dictionaries.

- SIMPLE does not create any change in the original text by replacing the word or inserting an explanation in the text. It only provides a translation and additional info (in a tooltip) on request, leaving the user fully in charge of his/her navigation path through the original text.

- SIMPLE works with different languages. English and Italian are already implemented but other languages can easily be added.

- Subjective tests (performed with expert and non-expert users) and objective tests (using term familiarity) have proved SIMPLE efficacy and accuracy.

**References**

1.      Kokkinakis D. The journal of the Swedish medical association - a corpus resource for biomedical text mining in Swedish. In: The Third Workshop on Building and Evaluating Resources for Biomedical Text Mining (BioTxtM), an LREC Workshop. 2012, Turkey.

2.      Leroy G, Endicott JE, Mouradi O, et al. Improving perceived and actual text difficulty for health information consumers using semi-automated methods. AMIA Annu Symp Proc, 2012: 522–31.

3.      Kvist M, Velupillai S. Professional Language in Swedish Radiology Reports – Characterization for Patient-Adapted Text Simplification. Scandinavian Conference on Health Informatics 2013, Copenhagen, Denmark, August 20, 2013; pp:55-59.

4.      Zeng-Treitler, Tse T. Exploring and developing consumer health vocabularies. J Am Med Inform Assoc 2006, 13(1), pp: 24-9.

5.      Kauchack D, Mouradi O, Pentoney K, et al. Text Simplification Tools: Using Machine Learning to Discover Features that Identify Difficult Text. 47th Hawaii International Conference on System Science, 06 Jan - 09 Jan, 2014, Hilton Waikoloa, Waikoloa, HI, USA.

6.      Siddadharthan A. An architecture for a text simplification system. Proceeedings of the Language Engineering Conference, 2002, 13-15 December, 2002, Hyderabad, India, pp: 64—71.

7.      Shardlow M. A Survey of Automated Text Simplification, International Journal of Advanced Computer Science and Applications, Special Issue on Natural Language Processing 2014, pp: 58-70.

8.      Damay J, Lojico G, Lu K, et al. SIMTEXT: Text Simplification of Medical Literature; 3rd National Natural Language Processing Symposium – Building Language Tools and Resources, February 17-18, 2006, Manila, pp: 34-38.

9.      Chandrasekar R, Doran C, Sinivas B. Motivations and methods for text simplification. Proceedings of COLING '96, Copenhagen, 1996, poster paper, pp: 1041-1044.

10.      Chandrasekar R, Srinivas B. Automatic induction of rules for text simplification, Knowledge-Based Systems 1997; 10: 183 – 190.

11.      Klebanov BB, Knight K, Marcu D. Text simplification for information-seeking applications. In: On the Move to Meaningful Internet Systems, Lecture Notes in Computer Science. Springer Verlag 2004, pp: 735–747.

12.      De Belder J, Deschacht K, Moens MF. Lexical simplification. In 1st International Conference on Interdisciplinary Research on Technology, Education and Communication, 25th - 27th of May 2010, Kortrijk, Belgium.

13.      Kandula S, Curtis D, Zeng-Treitler Q. A semantic and syntactic text simplification tool for health content, in AMIA Annual Symposium Proceedings. American Medical Informatics Association, 2010, pp: 366–370.

14.     Keskissarkka R. Automatic text simplification via synonym replacement, Ph.D. dissertation, Linkoping, 2012. Available at: http://liu.diva-portal.org/smash/record.jsf?pid=diva2%3A560901&dswid=-9669. Accessed: 29/09/2017.

15.     Leroy G, Endicott JE, Kauchak D, et al. User evaluation of the effects of a text simplification algorithm using term familiarity on perception, understanding, learning, and information retention. Journal of Medical Internet Research 2013; 15, (7).

16.     Jonnalagadda S, Gonzalez G. Sentence simplification aids protein- protein interaction extraction, in The 3rd International Symposium on Languages in Biology and Medicine, Hyatt Regency, Seogwiposi, Jeju Island, South Korea November 8-10, 2009, pp. 8–10.

17.     Zeng-Treitler Q, Goryachev S, Kim H, et al. Making texts in electronic health records comprehensible to consumers: a prototype translator. AMIA Annu Symp Proc. 2007; 11: 846-50.

18.     Saggion H, Gómez-Martinez E, Etayo E. Text Simplification in Simplext: Making Text More Accessible. Revista de la Sociedad Española para el Procesamiento del Lenguaje Natural 2011; 47: 341–342.

19.     Abrahamsson E, Forni T, Skeppstedt M. Medical text simplification using synonym replacement: Adapting assessment of word difficulty to a compounding language. Proceedings of the 3rd Workshop on Predicting and Improving Text Readability for Target Reader Populations (PITR) @ EACL 2014, pp: 57–65, Gothenburg, Sweden, April 26-30, 2014.

20.     Paetzold G, Specia L. Text Simplification as Tree Transduction. Proceedings of the 9th Brazilian Symposium in Information and Human Language Technology, pp: 116-125, Fortaleza, CE, Brazil, October 21-23, 2013.

21.     Paetzold G, Specia L. LEXenstein: A Framework for Lexical Simplification. Proceedings of ACL-IJCNLP 2015 System Demonstrations,July 2015; 85-90.

22.     Keskisärkkä R, Jönsson A. Automatic Text Simplification via Synonym Replacement; In Proceedings of The Fourth Swedish Language Technology Conference, Lund, (Sweden), October 25-26, 2012.

23.     Hading M, Matsumoto Y, Sakamoto M. Japanese Lexical Simplification for Non-Native Speakers. NLPTEA 2016, 92.

24.     Baorto DM,  Cimino, JJ. An "infobutton" for enabling patients to interpret on-line Pap smear reports. AMIA Annual Symposium Procs.; 2000, 47–50.

25.     Kemper D, De Fio, G, Hall L, et al. Getting patients to meaningful use: using the HL7 infobutton standard for information prescriptions. Healthwise White Paper Series. 2010.

26.     Alfano M, Lenzitti B, Lo Bosco G. An Online Multilingual Medical Vocabulary/Thesaurus/Dictionary (MED-VTD) for Facilitating Understanding of Medical Texts. Proc. of CINI Workshop on ICT for Smart Cities & Communities (I-CiTies 2015), 29-30 October 2015, Palermo.

27. Alfano M, Lenzitti B, Lo Bosco G, et al. An Automatic System for Helping Health Consumers to Understand Medical Texts. Proc. of HEALTHINF 2015b, 12-15 Gennaio 2015, Lisbona.

28. Alfano M, Lenzitti B, Lo Bosco G, et al. Facilitating text understanding for e-learning users. Proc. of International Conference on e-Learning (e-Learning'14), 11-12 Settembre 2014, Tenerife.

29. Zielstorff RD. Controlled vocabularies for consumer health. Journal of Biomedical Informatics, 2003; 36(4-5): 326–333.

30. Consumer Health Vocabulary Initiative. Available at: http://consumerhealthvocab.org/. Accessed: 29/09/2017.

31. Keselman A, Smith CA, Divita G, et al. Consumer health concepts that do not map to the UMLS: where do they fit?  Journ. Am. Med. Inform. Ass.; 2008. 15: 496-505.

32. Italian Consumer-oriented Medical Vocabulary. Available at: https://ehealth.fbk.eu/resources/italian-consumer-oriented-medical-vocabulary-icmv. Accessed: 29/09/2017.

33. Leroy, G. & Endicott, J.E., 2011. Term Familiarity to indicate Perceived and Actual Difficulty of Text in Medical Digital Libraries. , 7008(2011), pp.1–6.

34. Alfano, M., Lenzitti, B., Lo Bosco, G. 2014. A web search methodology for health consumers. In Proceedings of the 15th International Conference on Computer Systems and Technologies (CompSysTech '14), Boris Rachev and Angel Smrikarov (Eds.). ACM, New York, NY, USA, 150-157. DOI=http://dx.doi.org/10.1145/2659532.2659600.

35. Alfano, M.; Lenzitti, B.; Taibi, D. and Helfert, M. 2019. Facilitating Access to Health Web Pages with Different Language Complexity Levels.In Proceedings of the 5th International Conference on Information and Communication Technolo-gies for Ageing Well and e-Health - Volume 1: ICT4AWE, ISBN 978-989-758-368-1, pages 113-123. DOI: 10.5220/0007740301130123.

36. Alfano, M., Lenzitti, B., Taibi, D., Helfert, M. 2019. Provision of Tailored Health Information for Patient Empowerment: An Initial Study. In Proceedings of the 20th International Conference on Computer Systems and Technologies (CompSysTech '19), Tzvetomir Vassilev and Angel Smrikarov (Eds.). ACM, New York, NY, USA, 213-220. DOI: https://doi.org/10.1145/3345252.3345301.

37. Alfano, M., Lenzitti, B. Lo Bosco, G., Taibi, D. 2016. A Framework for Opening Data and Creating Advanced Services in the Health and Social Fields. In Proceedings of the 17th International Conference on Computer Systems and Technologies 2016 (CompSysTech '16), Boris Rachev and Angel Smrikarov (Eds.). ACM, New York, NY, USA, 57-64. DOI: https://doi.org/10.1145/2983468.2983473.