

Identifying *Condition-Action* Sentences Using a Heuristic-based Information Extraction Method

Reinhardt Wenzina and Katharina Kaiser

Institute of Software Technology and Interactive Systems
Vienna University of Technology
Favoritenstrasse 9-11, 1040 Vienna, Austria

Abstract. Translating clinical practice guidelines into a computer-interpretable format is a challenging and laborious task. In this project we focus on supporting the early steps of the modeling process by automatically identifying conditional activities in guideline documents in order to model them automatically in further consequence. Therefore, we developed a rule-based, heuristic method that combines domain-independent information extraction rules and semantic pattern rules. The classification also uses a weighting coefficient to verify the relevance of the sentence in the context of other information aspects, such as effects, intentions, etc. Our evaluation results show that even with a small set of training data, we achieved a recall of 75 % and a precision of 88 %. This outcome shows that this method supports the modeling task and eases the translation of CPGs into a semi-formal model.

Keywords: computer-interpretable clinical guidelines; medical guideline formalization; Many-Headed Bridge (MHB); UMLS; Information Extraction (IE)

1 Introduction

Clinical Practice Guidelines (CPGs) are defined as “systematically developed statements to assist practitioners and patient decisions about appropriate healthcare for specific circumstances” [6]. They include recommendations describing appropriate care for the management of patients with a specific clinical condition. An important part of CPG contents refers to the procedures to perform, often formulated together with specific conditions that have to hold in order to execute an activity.

CPGs are published as textual guidelines, but in order to deploy them in some kind of computerized tool (e.g., a reminder system or a more complex decision-support system) they have to be represented in specialized languages (see [15, 7] for a comparison and overview). Although different authoring/editing tools are often associated with these languages, authoring is a labor-intensive task that requires comprehensive knowledge in medical as well as computer science.

There have been several approaches to ease the modeling process, amongst others by introducing intermediate representations that provide more semi-structured and less formal formats. One of them is MHB, the Many-Headed Bridge

[21], that tries to bridge the gap between the guideline text and its corresponding formalized model. It falls in the category of document-centric approaches [22] and is devised to produce a non-executable XML document with the relevant CPG fragments, starting from the original text. The knowledge of a CPG is thereby represented in a series of chunks that correspond to a certain bit of information in the CPG (e.g., a sentence, part of a sentence, more than one sentence). The information in a chunk is structured in various dimensions, e.g., control flow, data flow. To additionally support the modeling task for non-IT experts, MHB was further developed and split into MHB-F (free-text version) and MHB-S (semantically enriched version) [19]. MHB-F now provides a very simplified structure to make the modeling even for non-IT experts feasible and to leave modeling details to knowledge engineers in a later step.

In order to ease the laborious modeling, parts of the task should be automated by applying information extraction methods. In this work we will focus on the identification of *condition-action* sentences that form a prominent aspect of the process flow in CPGs. The discovery of such combinations is not a trivial one. On the one hand, *condition-action* sentences are rarely of the form ‘*if condition then action*’, but require more sophisticated identification methods. On the other hand, conditions may refer to effects, intentions, or events and not activities, and these combinations have to be sorted out by our method. Table 1 shows a few example sentences in regard to their MHB-F aspects.

Table 1. Examples of sentences and their categorization in MHB-F

sentence	MHB-F aspect
An episiotomy should be performed if there is a clinical need such as instrumental birth or suspected fetal compromise.	decision based activity
Women with pain but no cervical changes should be re-examined after two hours.	decision based activity
Women should be informed that in the second stage they should be guided by their own urge to push.	clinical activity
The partogram should be used once labour is established.	background information
Administration of inhaled steroids at or above 400 mcg a day of BDP or equivalent may be associated with systemic side-effects.	effect
Legend: activity, condition, effect, explanation	

In order to identify *condition-action* sentences we propose a rule-based method using a combination of linguistic and semantic information. Furthermore, we introduce a weighting coefficient called *relevance rate* (rr) that shows whether a sentence is relevant for modeling.

The following section gives a short overview on the usage of information extraction methods as well as knowledge-based approaches for guideline modeling. Our method is explained in Section 3 and subsequently evaluated and discussed in Sections 4 and 5.

2 Background and Related Work

Guideline developers edit CPGs in a free-text format. In order to transform the medical knowledge described in a guideline into execution models a translation process is required. Moser and Miksch [14] detected prototypical patterns in free-text guidelines to bridge this gap. Serban et al. [18] proposed an ontology-driven extraction of linguistic patterns to pre-process a CPG in order to retrieve control knowledge. The evaluation showed that the modeling as well as the authoring process of guidelines was supported. Language engineering methods were used in SeReMed [5] to detect diagnoses or procedures in medical documents. The method was successfully applied to X-ray reports. These documents however, show a standardized structure and therefore are easier to handle by knowledge engineering methods than CPGs. Taboada et al. [23] identified relationships between diagnoses and therapy entities in free-text documents by matching the core information units of a sentence with a collection of predefined relationships but the quality of this matching was not rated. To implement a rule based approach to recognize medical entities the MeTAE (Medical Texts Annotation and Exploration) platform was used by Abacha and Zweigenbaum [1]. Additionally, semantic relations between each pair of these entities were identified by means of MetaMap [2]. Consequently, relations between a problem (e.g., disease) and a corresponding treatment were found. The method was applied to selected articles and abstracts of PubMed but not to CPGs. In [8] a heuristic-based approach using information extraction methods independent of the final guideline representation language was defined. This method was implemented and applied to several guidelines containing a high amount of semi-structured text. The adoption of this methodology on *'living-guidelines'* [9] showed promising results. A set of semantic patterns representing activities based on semantic relations was generated by Kaiser et al. [10] to identify medical activities in CPGs. Its effectiveness was proved by a study which showed that a large part of control flow related aspects could be identified. The relation between the activity and a corresponding condition, however, was not part of the method but is an important requirement for the future automatic translation of a guideline.

3 Methods

Condition-based medical activities are expressed in clinical practice guidelines in various ways and mostly found in single sentences. These sentences affect the clinical pathway and therefore are relevant to the computer-interpretable model of the guideline. In order to classify such a sentence as relevant, we based our approach on the following hypothesis:

1. A sentence owns a certain domain independent linguistic structure, and
2. contains recurrent domain dependent semantic key patterns.

We propose a rule-based, heuristic method using linguistic and semantic patterns to classify sentences in CPGs as relevant for describing conditional activities in order to move towards an automatic translation of such sentences into MHB in a following future step (an example is shown in Figure 1). Therefore, we analyzed a CPG document to develop a general linguistic pattern set and a semantic pattern set based on UMLS Semantic Types. These pattern sets then form the basis for the subsequent classification by calculating the *relevance rate* (rr).

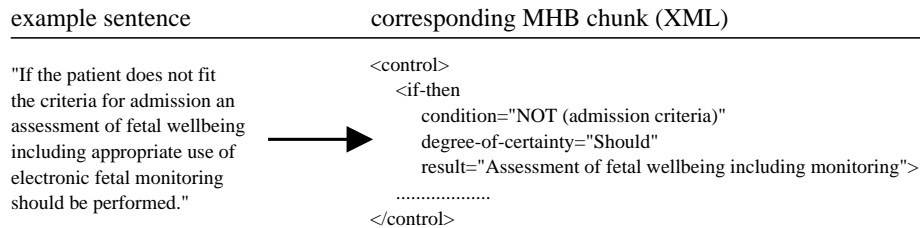


Fig. 1. A condition-action sentence represented as MHB-chunk

3.1 Knowledge Sources and Tools

The Unified Medical Language System (UMLS) [11] combines selected health and biomedical vocabularies to facilitate the standardized exchange of medical data between computer systems. It offers three different components:

1. The Metathesaurus which is an aggregation of medical terms and codes of different vocabularies (e.g., MeSH, SNOMED CT, etc.).
2. The Semantic Network which reduces the complexity of the Metathesaurus by assigning semantic types to the concepts of the Metathesaurus in order to group and define relationships among them.
3. The SPECIALIST Lexicon and Lexical Tools which provide natural language processing tools.

The open source framework for text engineering GATE (Version 6) [4] allows the combination of different text engineering components to develop reusable applications. The following components were used in our method:

- ANNIE: A set of information extraction (IE) components, distributed within the GATE system and relying on finite state algorithms and the JAPE language [3].
- OpenNLPChunker supports the detection of phrases within a parsed text.
- MetaMap Annotator: A tagger that maps biomedical texts to the UMLS Metathesaurus and discovers Metathesaurus concepts and their semantic types [2].

3.2 Manual Development of the General Linguistic Pattern Set

For the development of the general linguistic pattern set we used a chapter from an Asthma guideline developed by SIGN [17], where *chapter 4: pharmacological management* had been modeled in the semi-structured modeling language MHB-F [21] by a guideline modeling expert. We analyzed this chapter with regard to the control flow aspects and started generating an initial linguistic pattern set based on trigger words (12 occurrences for 'if' and 4 occurrences for 'should').

Table 2. Selected general linguistic patterns

#	rule type	pattern	weight(w)
1	IF	* [Ii]f *	0.5
2	IF	If {condition} {consequence}.	1.0
3	IF	If {condition}, {consequence}.	1.0
4	IF	{consequence} if {condition}.	1.0
5	IF	If {condition} then {consequence}.	1.0
6	SHOULD	* should have *	0.5
7	SHOULD	* should be *	0.5
...

In order to identify the semantic clauses of a sentence (these are the parts describing the *condition* and the *consequence*), the patterns had to be grouped into 6 different patterns for 'if' and 4 different patterns for 'should' (some straightforward patterns are listed in Table 2). Condition and consequence are distinguished according to the sentence's syntax, punctuation and its sequence of phrases. Two complex examples including conditions spread over multiple phrases are shown in Figure 2.

We assigned a weighting factor to every pattern of the set - the value 0.5 to show that only a trigger word was found and 1.0 to express that also the semantic clauses were identified. These constants can be adapted for new rule types in the future.

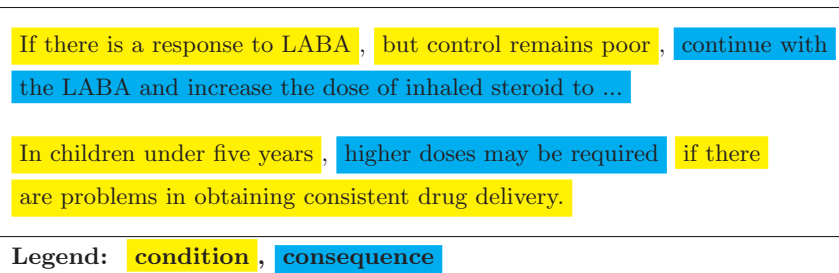


Fig. 2. Sentences with multi-part conditions

3.3 Generation of the Domain-Dependent Semantic Pattern Set

Amongst the general syntactic patterns we also used semantic patterns based on the UMLS Semantic Network. Therefore, we used the MetaMap plugin within the GATE framework to automatically identify medical concepts in our text and assign them to their corresponding UMLS concepts and semantic types (represented by four-letter abbreviations - e.g., 'popg' stands for 'Population Group'). By this way it was possible to find out the sequence of semantic types in the clauses of the sentences (see Figure 3).

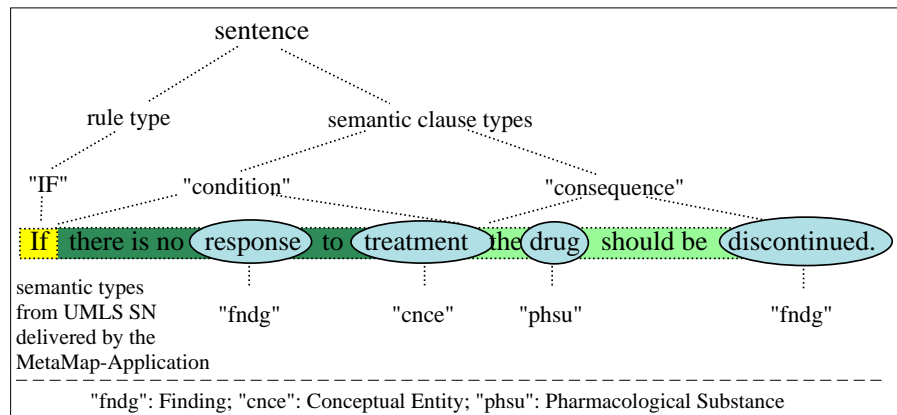


Fig. 3. Semantic abstraction of a sentence

Finally, the complete semantic abstraction of a sentence including rule type, semantic clause type and sequence of semantic types was added to the semantic pattern pool (see Table 3). A total of 32 entries was automatically generated.

Table 3. Structure of the semantic pattern pool (selected samples)

rule type	semantic clause type	sequence of semantic types
IF	condition	[fndg][cnce]
IF	consequence	[phsu][fndg]
IF	consequence	[idcn][qlco][resa][ftcn][ftcn]
SHOULD	condition	[aggp][podg][dsyn]
SHOULD	consequence	[qnco][tmco][resa][orch, phsu][idcn]
SHOULD	condition	[podg][qlco][ftcn][orgf][qlco][gngm][phsu]

3.4 Calculation of the Relevance Rate

The relevance rate rr is a measure to find out whether a sentence contains a condition-action combination. Furthermore, it shall classify a sentence as crucial for the clinical pathway in contrast to other information aspects like intentions or explanations which are modeled in MHB-F in a different way. To find this semantic difference the syntax of the sentences as well as their containing medical semantic types must be respected. In order to calculate the relevance rate for a selected sentence its semantic abstraction has to be generated and compared with every entry in the semantic pattern pool. If the rule type and the semantic clause type are matching, the similarity of the sequences of the semantic types is calculated by using the Dice coefficient [12]. The highest value is selected for further calculation.

In general, the value of the relevance rate rr is the sum of

- the weight(w) of the applied general IE rule, and
- the sum of the maximum similarity value (s_i) for each semantic clause of the sentence divided by the number of semantic clauses (n) identified by the general IE rules.

$$rr = w + \frac{\sum_{i=1}^n \max\{s_i\}}{n} \quad (1)$$

The weight of the IE rule and the arithmetic average of the similarity values - both in the range between 0 and 1 - have the same influence on the rr .

The similarity value s_i of the semantic clause of a sentence and a matching entry in the semantic pattern pool were calculated as follows:

- If the semantic clause contains only one semantic type it is compared to those entries of the semantic pattern pool that also show only one semantic type (to receive a better accuracy). In the case that both types are equal the value for s_i is set to 1.0, otherwise
- both sequences of semantic types are interpreted as a string each and two sets of 4-letter string bigrams are composed out of them. Subsequently, these two sets are used for the calculation of the Dice coefficient.

Example:

Given are the sequence of semantic types of a new semantic clause S and the sequence of semantic types of a matching entry of the semantic pattern pool P .

S : [fndg] [orgf] [qlco] [gngm] [phsu] and P : [strd] [gngm] [phsu] [ortf].

The resulting 4-letter string bigrams are:

$S = \{“fndgorgf”, “orgfqlco”, “qlcogngm”, “gngmphsu”\}$
 $P = \{“strdgngm”, “gngmphsu”, “phsuortf”\}.$

The Dice coefficient is defined as twice the shared information over the sum of cardinalities:

$$s_i = \frac{2n_t}{n_s + n_p} \quad (2)$$

where n_t corresponds to the number of bigrams found in both sets, n_s is the number of bigrams in S , and n_p the number of bigrams in P . So the result of this example is $s = \frac{2*1}{4+3} = \frac{2}{7}$.

The general interpretation of the rr is shown in Table 4.

Table 4. Interpretation of the relevance rate

value	interpretation
$rr = 0.5$	only a trigger word or word combination were found; no semantic clause could be identified
$rr = 1.0$	an appropriate general IE pattern was found and semantic clauses could be identified
$1.0 < rr \leq 2.0$	additionally, a semantic similarity between the sequences of the semantic types was detected

4 Evaluation

The guidelines *Management of active low-risk labour - Admission for Birth* [16] (chapters 1, 2.1, 2.2, and 2.3) and *CBO Treatment of Breast Cancer* (chapter 3) were applied to evaluate the method. These guidelines were intentionally selected, because they cover a completely different medical application area in contrast to the Asthma guideline and furthermore an MHB¹ [21] model already existed, used as a “golden standard”.

The IE rules for the domain independent patterns and the generation of the semantic patterns were implemented with GATE - the processing resources for

¹ MHB is the former version of MHB-F

the calculation of the relevance rate are shown in Figure 4. The IE rules were

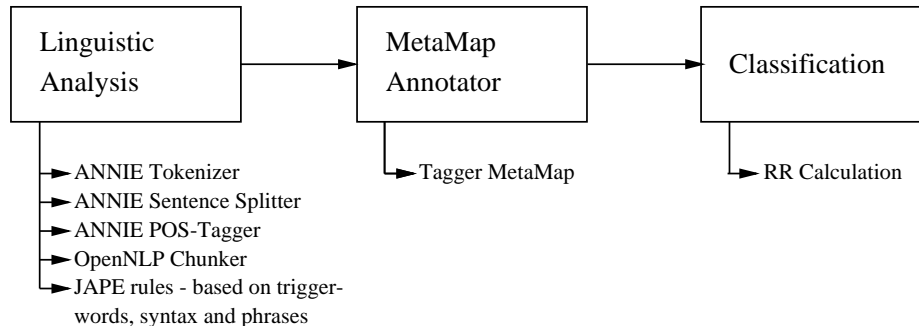


Fig. 4. Processing resources used in GATE

applied to the guidelines and the identified sentences were annotated with the semantic types retrieved from the UMLS via the MetaMap plugin. 50 sentences were found which complied to the IE rules. Twenty out of the 50 sentences had been modeled in MHB as control flow related aspects (activities based on special medical conditions) and therefore their *rr* was expected to be higher than 1. The other 30 sentences represented information about intentions and explanations in MHB with an expected *rr* lower or equal 1. The evaluation results are shown in Table 5 differentiated by the rule types.

Table 5. Evaluation results I ($rr > 1$: tp and fp; $rr \leq 1$: fn and tn)

		type: IF		type: SHOULD		total	
Gold Standard	<i>tp</i>	14	3	1	2	15	5
	<i>fn</i>	1	2	1	26	2	28

tp = true positive; fn = false negative; fp = false positive; tn = true negative

With the described method it was possible to identify 15 sentences which correctly contained control flow related aspects. Only two sentences got an incorrect *rr* higher than 1. They were not correctly rated, because they did not describe condition-based activities. Other 5 sentences got an *rr* lower than 1 although their *rr* should have been higher as they contained condition-based activities. Thus, the method had a recall of 75% and a precision of 88% (see Table 6). The negative predictive value of 79% was higher than the recall value and showed that 28 sentences had been correctly classified.

Table 6. Evaluation results II

Guideline	type: IF			type: SHOULD			total		
	REC	PRE	NPV	REC	PRE	NPV	REC	PRE	NPV
Breast Cancer	100%	67%	100%	-	-	100%	100%	67%	100%
Adm. for Birth	80%	100%	0%	33%	50%	88%	72%	93%	75%
total	82%	93%	40%	33%	50%	93%	75%	88%	85%

REC the number of correctly identified sentences over the number of the modeled sentences in the golden standard (=recall)

PRE the number of correctly identified sentences over the entire number of identified sentences (=precision)

NPV the number of correctly not identified sentences over the number of not modeled sentences in the golden standard (=negative predictive value)

Table 7. Selected sentences with an $rr \leq 1.0$

sentence	rr	reason
If, notwithstanding these procedures cervical dilatation doesn't progress, consider cesarean section after 2-3 hours of regular and painful contraction with no cervical changes.	0.5	Sentences with such a linguistic structure did not exist in the Asthma guideline → no IE rule was implemented to identify semantic clauses
If dilatation progress is not regular (<1 cm/hour in nulliparous, <1,5 cm/hours in parous) consider: - amniotomy; - oxytocin administration.	0.5	A list of resulting activities was not found in the Asthma guideline → no IE rule was implemented to identify semantic clauses
In case of abnormal FHR, monitoring should be continuous.	1.0	The sentence was wrongly categorized with a rule for 'should' because no rules for "in case of" were implemented → no semantic similarity within the semantic pool pattern was found
After umbilical cord clamping, if the second stage of labour has been physiological, the baby is given to the mother and covered.	1.0	The semantic clauses were found, but no semantic similarity occurred

Generally, the results proved that the rules of type “if” showed much better results for the precision (93%) and the recall (82%) than the ones of type “should”. Nevertheless, the negative predictive value for the latter type showed a rate of 93%.

The analyzed guidelines contained 7 sentences with control flow related information but they were not classified because their patterns did not exist in the Asthma guideline. Consequently, no corresponding general linguistic IE rule existed in the pattern pool. An extension of the general linguistic pattern set with rules for the trigger words “when”, “could” and “in case of” should be taken into consideration. Additionally, 1 condition based activity could not be found, because the semantic information was distributed over more than one sentence. In Table 7 selected examples are shown with a relevance rate lower or equal than 1.0 together with the corresponding reasons.

Even though only a small amount of training data (16 sentences) was available from the Asthma guideline, our method identified condition-based activities for control flow related aspects in a guideline document. Furthermore, it showed that the combination of domain independent information extraction rules and an automatically created semantic pattern pool leads to valuable results.

5 Conclusions and Further Work

The aim of this paper was to develop a method to identify condition-action sentences. By defining a set of linguistic patterns we split up sentences semantically - from one selected training guideline - into their clauses showing the condition and the consequence. We used the UMLS Semantic Network [13] to find out which types of medical concepts were applied in these clauses. The outcome was a semantic abstraction of every training sentence which then was stored in a semantic pattern pool. This pool facilitated the classification of new sentences regarding to their relevance to the corresponding MHB-F model expressed by the measure relevance rate (rr).

Modeling experts benefit from the method in two ways:

1. Condition-based activities in free-text guidelines, which must be modeled in MHB-F, are identified and rated.
2. These sentences are automatically split into the condition and the resulting activity.

An integration of the presented method into modeling tools will ease the work of all parties involved.

Ongoing steps will be (1) the implementation of additional information extraction rules to expand the general linguistic pattern set in order to improve the hit rate; (2) the extension of the semantic pattern pool with additional training data in order to increase the significance of the relevance rate; (3) the application of the method in the context of *‘living-guidelines’* [20]; (4) an investigation,

whether this method can support the modeling of processes also in other application areas by substituting the UMLS SN by other domain dependent thesauri; (5) the refinement of information extraction rules to make a step towards an automatic translation of condition-based activities of free-text guidelines into the modeling language MHB-F; and (6) the development of case studies to evidence the effectiveness of the method in real-world scenarios. If it is possible to tap the full potential of the presented method, the implementation of CPGs will be fostered tremendously.

Acknowledgement. This research was carried out as part of project no. TRP71-N23 funded by the Austrian Science Fund (FWF).

References

1. Abacha, A.B., Zweigenbaum, P.: Automatic extraction of semantic relations between medical entities: a rule based approach. *Journal of Biomedical Semantics* 2(Suppl 5), S4+ (2011)
2. Aronson, A.R., Lang, F.M.M.: An overview of MetaMap: historical perspective and recent advances. *Journal of the American Medical Informatics Association : JAMIA* 17(3), 229–236 (May 2010)
3. Cunningham, H., Maynard, D., Bontcheva, K., Tablan, V.: GATE: A framework and graphical development environment for robust NLP tools and applications. In: *Proceedings of the 40th Anniversary Meeting of the Association for Computational Linguistics* (2002)
4. Cunningham, H., Maynard, D., Bontcheva, K., Tablan, V., Aswani, N., Roberts, I., Gorrell, G., Funk, A., Roberts, A., Damjanovic, D., Heitz, T., Greenwood, M.A., Saggion, H., Petrak, J., Li, Y., Peters, W.: *Text Processing with GATE (Version 6)*. University of Sheffield Department of Computer Science (2011), <http://tinyurl.com/gatebook>
5. Denecke, K.: Semantic structuring of and information extraction from medical documents using the UMLS. *Methods of Information in Medicine* 47(5), 425–434 (2008)
6. Field, M.J., Lohr, K.N. (eds.): *Clinical Practice Guidelines: Directions for a New Program*. National Academies Press, Institute of Medicine, Washington DC (1990)
7. Isern, D., Moreno, A.: Computer-based execution of clinical guidelines: A review. *International Journal of Medical Informatics* 77(12), 787 – 808 (2008)
8. Kaiser, K., Akkaya, C., Miksch, S.: How can information extraction ease formalizing treatment processes in clinical practice guidelines? A method and its evaluation. *Artificial Intelligence in Medicine* 39(2), 151–163 (2007)
9. Kaiser, K., Miksch, S.: Versioning computer-interpretable guidelines: Semi-automatic modeling of 'Living Guidelines' using an information extraction method. *Artificial Intelligence in Medicine* 46(1), 55–66 (2009)
10. Kaiser, K., Seyfang, A., Miksch, S.: Identifying treatment activities for modelling computer-interpretable clinical practice guidelines. In: Riao, D., Teije, A., Miksch, S., Peleg, M. (eds.) *Knowledge Representation for Health-Care, Lecture Notes in Computer Science*, vol. 6512, pp. 114–125. Springer Berlin Heidelberg (2011)
11. Lindberg, D., Humphreys, B., McCray, A.: The Unified Medical Language System. *Methods of Information in Medicine* 32(4), 281–291 (1993)

12. Manning, C.D., Schütze, H.: Foundations of statistical natural language processing. MIT Press, Cambridge, MA, USA (1999)
13. McCray, A.: UMLS Semantic Network. In: Proc. of the 13th Annual Symposium on Computer Applications in Medical Care(SCAMC). pp. 503–507 (1989)
14. Moser, M., Miksch, S.: Improving clinical guideline implementation through prototypical design patterns. In: Miksch, S., Hunter, J., Keravnou, E. (eds.) Artificial Intelligence in Medicine, Lecture Notes in Computer Science, vol. 3581, pp. 126–130. Springer Berlin Heidelberg (2005)
15. Peleg, M., Tu, S., Bury, J., Ciccicarese, P., Fox, J., Greenes, R., Hall, R., Johnson, P.D., Jones, N., Kumar, A., Miksch, S., Quaglini, S., Seyfang, A., Shortliffe, E.H., Stefanelli, M.: Comparing Computer-Interpretable Guideline Models: A Case-Study Approach, . The Journal of the American Medical Informatics Association (JAMIA) 10(1), 52–68 (Jan - Feb 2003)
16. Remine: Documentation of formalized guidelines (2010), <http://www.remine-project.eu/>
17. Scottish Intercollegiate Guidelines Network (SIGN): British Guideline on the Management of Asthma. A national clinical guideline. Scottish Intercollegiate Guidelines Network (SIGN) (May 2011)
18. Serban, R., ten Teije, A., van Harmelen, F., Marcos, M., Polo-Conde, C.: Extraction and use of linguistic patterns for modelling medical guidelines. Artif. Intell. Med. 39(2), 137–149 (Feb 2007)
19. Seyfang, A., Kaiser, K.: MHB-F Specification (2011)
20. Seyfang, A., Martnez-Salvador, B., Serban, R., Wittenberg, A., Miksch, S., Marcos, M., Teije, A.T., Rosenbrand, K.: Maintaining formal models of living guidelines efficiently. In: Proc. of the 11th Conference on Artificial Intelligence in Medicine (AIME07). Springer Verlag (2007)
21. Seyfang, A., Miksch, S., Marcos, M., Wittenberg, J., Polo-Conde, C., Rosenbrand, K.: Bridging the gap between informal and formal guideline representations. In: Brewka, G., Coradeschi, S., Perini, A., Traverso, P. (eds.) European Conference on Artificial Intelligence (ECAI-2006). vol. 141, pp. 447–451. IOS Press, Riva del Garda, Italy (2006)
22. Sonnenberg, F.A., Hagerty, C.G.: Computer-interpretable clinical practice guidelines. where are we and where are we going ? Yearb Med Inform pp. 145–58+ (2006)
23. Taboada, M., Meizoso, M., Riaño, D., Alonso, A., Martínez, D.: From natural language descriptions in clinical guidelines to relationships in an ontology. In: Proceedings of the 2009 AIME international conference on Knowledge Representation for Health-Care: data, Processes and Guidelines. pp. 26–37. KR4HC'09, Springer-Verlag, Berlin, Heidelberg (2010)