

# Side Effect Alerts Generation from EHR in Polish

Wojciech Jaworski<sup>1</sup>[0000-0002-7838-7781], Małgorzata  
Marciniak<sup>2</sup>[0000-0002-0953-758X], and Agnieszka  
Mykowiecka<sup>2</sup>[0000-0002-8939-3255]

<sup>1</sup> LekSeek Polska, Puławska 465, Warsaw, Poland

<sup>2</sup> Institute of Computer Science Polish Academy of Sciences,  
Jana Kazimierza 5, Warsaw, Poland  
{agn,mm}@ipipan.waw.pl

**Abstract.** The paper addresses the problem of extending an existing and widely used program for Polish public healthcare with a function for detecting possible occurrences of drug side effects. The task is performed in two steps. First, we extract information that binds names of drugs with side effects and their frequency. In the next step, we look for similar phrases in the list of side effect phrases. For all words in phrases, we use Polish Wordnet to find similar ones, and check if phrases with replaced words exist in the list. For long side effect phrases, which never occur in patient records, we look for simpler internal side effect phrases to generate alarms. Finally, we evaluate to what extent this action increases the efficiency of side effect alarms.

**Keywords:** drug side effects · electronic health records · Polish

## 1 Introduction

Electronic health record (EHR) systems are in common use all over the world in clinics, both for administrative services and to support the work of physicians. One of the basic required functionalities of such systems is to store information about patients' medical care: the reasons for patients' visits, the results of check-ups, tests, diagnosis, and prescribed medications. Analysis of this data along with other resources can increase the efficiency of physicians and the accuracy of undertaken decisions. Clinical Decision Support Systems (CDSSs) are quite common for English data [26], and available for other languages e.g. Swedish[8], German [14], and Korean [3]. Polish EHR systems are focused on the organizational and administrative aspects of the clinic's functioning and on collecting information about patients, but do not analyze the data. A summary of the use of EHR systems in Poland in 2016 and their perspectives is given in [4].

The slow development of clinical decision support systems in Poland is due to the poor resources for the processing of medical data and there not yet being a national standard for the storage of medical information. Medical terminology resources in Polish are limited to International Classification for Nursing Practice (INCP, <https://www.icn.ch/>); a small part of The Unified Medical Language System [15] (UMLS), i.e: Medical Subject Headings (Mesh —

a controlled biomedical vocabulary designed for medical literature indexing and searching) [27], and the International Statistical Classification of Diseases and Related Health Problems (ICD-10). Additionally, a list of drugs and supplements approved for use in Poland is publicly available. The Systematized Nomenclature of Medicine Clinical Terms [25] (SNOMED CT) is not translated into Polish. We are not aware of any general medical ontology of Polish medical data.

The paper addresses one of the EHR systems available on the Polish market, i.e. drWidget, which has been developed for 7 years and is implemented in over 16K outpatient clinics. The system collects data concerning patients' visits. Some of the information is given in a structured form such as: visit identification, doctor identification, and basic information about patient (sex, age, id). The records describing the patient's visit (interview and examination) are given in free texts. They usually contain a large number of spelling errors which makes them difficult to process. Nowadays, all prescriptions of medications are given in electronic form. Several records containing the structural description of a drug (and its dosage) can be added to a visit. This eliminates the problem of multiple variants of drug names being used in free texts. The system also provides doctors with Summaries of medical Product Characteristics (SmPCs) in text form.

In the paper, we describe a new functionality of the system, i.e., generating alerts when a symptom described in a free text about a patient visit might be a side effect of a drug being taken by the patient. Identification of situations when a drug causes undesirable secondary effects in addition to the desired therapeutic effect can help both a doctor and a patient by making it easier to make a proper diagnosis and sparing a patient unnecessary ailments. Systems which could facilitate a doctor's diagnoses can thus potentially be very advantageous for patients, but are frequently not very well received by the physicians themselves if they are not fully reliable. The general solution to this task is difficult, as it requires many aspects to be considered to recognize not only which drugs a patient takes now but what new symptoms he/she developed and for what reason. In the current version of the system, we limit generating alerts to the case when a medication prescribed in the previous patient visit has side effects which are similar to symptoms detected in the free text about the current visit.

Ultimately, the alerts will be generated in a system used by hundreds of doctors; therefore we want to minimize the number of false positives. Too many unjustified alarms would quickly lead to users ignore them [1].

## 2 Related Work

Medical texts, both scientific and clinical, constitute a vast amount of data which can be mined for different kinds of information. In [9], text mining was described as an emerging tool for leveraging underutilised data sources that can improve pharmacovigilance, including the objective of Adverse Drug Event (ADE) detection and assessment.

The impact of ADE detection on patient treatment is analysed in [23], while in [19], the authors make an analysis of alert types in order to improve their

effectiveness in systems. In [24], the authors pay attention to the problem of ADE in older patients who take more drugs, as chronic diseases are more common and they experience more frequent ADEs. They review studies concerning usage of clinical decision support systems to reduce the prescribing of potentially inappropriate medications. The conclusions are that CDSSs are more effective in hospitals than in ambulatory care. But the authors expressed hope that more user friendly systems could improve their effectiveness.

A summary of 30 approaches published for ADE detection in the context of EHRs before 2017 is given in [6]. The problem of ADE was addressed in shared tasks. The Adverse Reaction Extraction from Drug Labels Track was organized during the Text Analysis Conference (TAC) in 2017. Based on annotated data, participants had to identify and normalize adverse reactions from drug labels. The best 10 teams taking part in the competition provided solutions based on machine learning methods: (Bi)LSTM, CRF, SVM, CNN; the best F1 score was 0.82 [22]. Another competition took place in 2018, when participants solved the problem of extraction of ADEs from clinical data. The organizers reported an F1 measure equal to 0.89 of the best systems “that process raw narrative text to discover concepts and find relations of those concepts to their medications” [10].

Most of the reports concern data analysis in English; however, the need for data processing in other languages is also noted. [21] describes the state of automatic patient data processing in Sweden in 2010, which is very close to the current situation in Poland. The authors claim that “the current structure, content and format of SmPCs make it difficult to incorporate them into CDSSs and link them to relevant patient information from the Electronic Health Records”. Our paper addresses a method of incorporating data from SmPCs to support doctor’s decisions based on EHRs in Polish.

### 3 Initial Drug Side Effects Identification Procedure

There is no official source of drug side effects in Polish and there are no corpora annotated with information on drug side effects such as the EU-ADR corpus [28] and that described in [7]. We are not able to construct the necessary resource with the help of UMLS as in [12]. The first step of the process was thus to construct a resource with possible drug side effects from SmPCs, which was the complete and up-to-date documentation of all prescription drugs authorized to be used in Poland. We extracted information about side effects associated with the frequency they occur for a given drug. It was carried out in two steps. First, we manually marked information on side effects in SmPCs. It allowed us to create a list of 13,347 various phrases. Most of them consisted of up to 4 tokens, but some were longer. We observed that in the actual visit description corpus, it was possible to match up to 5-tokens side effect phrases. Drug descriptions have free text form, but the fragments of possible side effects that are usually enumerated together with information about their frequency are of interest to us. It is expressed by very strict phrases, a list of which was prepared manually, e.g.: *często* ‘often’, *rzadko* ‘rare’, *bardzo rzadko* ‘very rare’. For a given drug, we

extracted pairs consisting of side effects and the frequency of their occurrence by simply matching the previously completed lists (of side effects and frequency phrases) with the drug description.

However, finding those side effects in a patient visit record is much harder. The main problems result from linguistic diversity:

- a side effect phrase may occur in inflected forms,
- a side effect may be expressed by various phrases,
- sometimes slightly more general or more precise information may be used in place of the term mentioned within the drug description,
- terminology used by a patient (who is usually not a doctor) differs from that used in drug descriptions,
- coordination is used quite often to describe side effects in SmPCs which is not common in patient records.

Side effect phrases which are listed within drug descriptions are in nominative form, so to be able to recognize all their forms in a text, we computed their inflected forms (noun phrases are declined by cases and numbers). This was done using a generator containing data from SGJP [31] and a guesser operating on the basis of the rules describing the inflection of the Polish language.

While creating an inflection model, we focused on the most productive Polish inflection rules. The model does not include irregular verbs and a small number of words that belong to other parts of irregular speech variety. They are not very numerous and their forms are just listed in the glossary attached to the model. The model also describes acronyms, frequently used words with a non-Polish spelling, and some dialect forms.

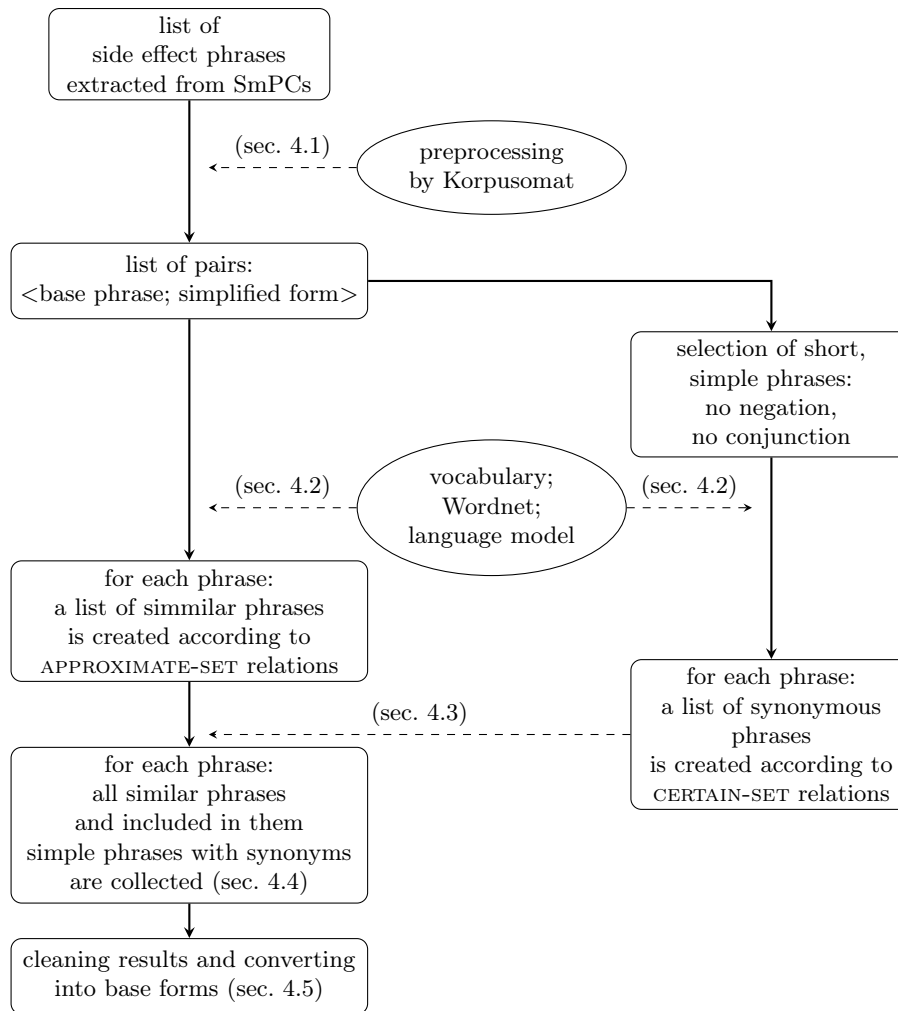
The list of symptoms is compared with the content of the interview and examination introduced by the doctor in the office program. If any symptom from the list occurs, the patient’s history is checked for whether he or she was taking the drug causing the symptom. If this is the case, a warning is displayed about the possible occurrence of side effects. Comparison with the list of previously prescribed medications is necessary due to the fact that some undesirable symptoms, e.g. ‘cough’, are also common symptoms of diseases.

## 4 Looking for Semantically Similar Side Effect Phrases

In this section, we describe a method for identifying phrases that express a similar meaning and a more general one on the side effect list. For example, in the *Oritop* leaflet, the side effect phrase *ataki lęku* ‘anxiety attacks’ is mentioned, while for *Epitoram* the phrase *napad lęku* ‘fit of anxiety’ is used. The meaning of both phrases is the same, and both can be used interchangeably in patient records. Recognition of the latter phrase should generate information about a potential undesirable effect of the first drug. Long, complex phrases never occur in patient visits, so information about the occurrence of potential side effect should be generated if a shorter, more general phrase implied by the longer one occurs. For example, the phrase *reakcje nadwrażliwości na światło słoneczne*

*i promieniowanie ultrafioletowe* ‘hypersensitive reactions to sunlight and ultraviolet radiation’ implies the following phrases: *nadwrażliwość na światło słoneczne* ‘hypersensitivity to sunlight’ and *nadwrażliwość na światło* ‘hypersensitivity to light’.

Figure 1 is a diagram of the subsequent steps of the data processing and the flow of information during the search for semantically similar side effect phrases described in the section below.



**Fig. 1.** Sequence of steps that make up the list of semantically similar phrases.

#### 4.1 Preprocessing

The side effect phrases were analysed using the Korpusomat [13] service. The phrases were tagged by the Concraft tagger [29] which uses the Morfeusz analyser [30]. The tagger gives lemmas for all words which are present in the Morfeusz dictionary while for the out-of-vocabulary words, it guesses morphosyntactic descriptions. We obtained the list of pairs consisting of a phrase and the corresponding sequence of lemmas which we hereinafter refer to as a simplified base form. As Polish is a highly inflected language, the problem of morphological variants recognition is not a simple string comparison. For matching variants, we use the method described in [16] and operate on simplified base forms. For example, for the phrase *upośledzenie czynności nerek* ‘impairment of renal function’, the simplified base form is *upośledzenie czynność nerka* ‘impairment function renal’. As we can see, only the first token in both phrase forms matches. Simplified base forms treat as equivalent phrases whose meanings are slightly different such as *świąd oka* ‘itchy eye’ and *świąd oczu* ‘itchy eyes’. They have the same simplified base form *świąd oko*. As a result, both phrases are represented as one entry in the list of side effect phrases that now has 12,514 various entries (some of the phrases have been unified).

We have also created a token dictionary which consists of all lemmas used in the phrases. It consists of 4629 different tokens, which includes lemmas of words and fewer than 70 other tokens as punctuation marks, digits and a combination of them which were not segmented by the tagger (e.g.: *>3-krotnej* ‘>3-times’).

#### 4.2 Use of Wordnet Relations

In order to find similar phrases within the list of side effect phrases, we use methods similar to those described in [11]. We use Polish Wordnet [20, 5] to find words and phrases similar to words from the token dictionary. For each word from the dictionary we found all words and phrases which are in the following relations: synonymy, inter-paradigmatic synonymy, and various types of derivatives e.g. pairs of nouns and adjectives *skóra – skórny* ‘skin<sub>noun</sub> – skin<sub>adj</sub>’ (CERTAIN-SET). Moreover, we select the second set which includes pairs whose meaning is more distant. This set additionally contains: hipernyms, hiponyms and fuzzynyms (APPROXIMATE-SET).

As one word may refer to several synsets and there is no effective method to select which meaning may refer to a medical topic, we collect similar words and phrases for all synsets. We then select only those words and phrases that have all elements included in the token dictionary. For example, *noga* ‘leg’ refers to 7 synsets and only two of them are connected to human anatomy. One meaning is imprecise as it identifies ‘leg’ with ‘foot’ but that probably does not cause significant errors. The dictionary condition allows us to filter out meanings that are distant such as *piłka nożna* ‘soccer’, which is called colloquially *noga* ‘leg’. The second test which we perform on the similar words/phrases collected from Wordnet is based on a distributional word2vec [17] language model. We use the model with vectors of 100 in length, calculated on the lemmatized texts of patient

visits. We select for further processing, all similar words whose distance is at least 0.1. This criterion looks very mild but we want to remove only very distant pairs of meanings such as *senność* ‘somnolence’ and its colloquial synonym *śpiączka* ‘coma’. If one or both lemmas are not in the dictionary, we accept such pairs assuming that they refer to a medical notion as they are present in the token dictionary.

After making the selection described above, we create all possible phrases where elements are substituted by all the selected synonyms for all side effect phrases. If such constructed phrase is present in the list of side effect phrases, we join them as variants. This procedure allows us to join phrases such as: *zanik skóry* ‘disappearance of skin’ and *atrofia skóry* ‘skin atrophy’; *zamazane widzenie* ‘blurred vision’, *niejasne widzenie* ‘unclear vision’, and *niewyraźne widzenie* ‘dim vision’.

Based on Wordnet, we select two sets of similar phrases that have almost certainly the same meaning and somehow similar ones. Such defined similarity might not be reciprocal and the first set consists of 1,079 phrases for which 1,303 phrases are similar. The second set consists of 2,329 entries for which 3,852 phrases are similar. The above numbers concern phrases in their simplified base forms. A large number of pairs are double counted (the relations are mainly reciprocal) and the effective number of encountered similar phrases is about half of the total.

### 4.3 Internal Phrases

As patient records contain side effect phrases up to 5 tokens, it is ineffective to look for the longer phrases which make up 20% of the side effect list. However, to make these phrases useful for generating alarms we recognize all simple phrases included in them. As a simple phrase we accept a phrase up to 5 tokens which does not include coordination (*i* ‘and’, *oraz* ‘and’, *lub* ‘or’, and characters: ‘,’, ‘/’) and negation<sup>3</sup> (*nie* ‘no’, *bez* ‘without’, *wyjatek* ‘exception’). The recognition of included phrases is limited to a very simple comparison of two bag-of-words. So phrase  $\mathcal{A}$  is included by phrase  $\mathcal{B}$  if all tokens of phrase  $\mathcal{A}$  are in phrase  $\mathcal{B}$ . To perform this comparison, we use the simplified base form of phrases. If phrase  $\mathcal{B}$  contains negation, it is shortened to the place where negation occurs. It allows us not to recognize ‘aura’ as a side effect entailed from the following phrase: *migrena bez aury* ‘migraine without an aura’. In this case, from a logical point of view, the *migraine* itself does not follow either, but a patient may complain of migraine without stating that the aura does not occur.

The comparison of bag-of-words copes with coordinated information. If a patient uses a drug with the side effect *kwasica metaboliczna i ketonowa* ‘metabolic and keto acidosis’, an alarm is generated if her/his record also contains one of the following phrases included in the coordinated one: *kwasica ketonowa* ‘keto acidosis’ *kwasica metaboliczna* ‘metabolic acidosis and just *kwasica* ‘acidosis’.

<sup>3</sup> *brak* ‘lack’ is handled differently.

Polish is a free word order language, so for example, the side effect *choroba niedokrwienne serca* ‘ischemic heart disease’ is expressed by a phrase with a different word order: *niedokrwienne choroba serca*. As they consist of the same tokens, the comparison of the bag-of-words allows us to connect them as similar.

The comparison of bag-of-words allows us to recognize 16,240 pairs of side effect phrases (in simplified base forms) for 7,996 entries.

#### 4.4 Unified List of Similar Phrases

For all phrases describing side effects we collect all variants in the following order. First, we generate all phrases according to the more distant Wordnet relations (APPROXIMATE-SET). Then, for each phrase and its variants, we find included phrases. In the next step, we add similar phrases for included phrases counted according to more restricted Wordnet relations (CERTAIN-SET). All repetitions are removed. The final list contains similar phrases for 9865 entries and consists of 38057 total variants.

#### 4.5 Similar Phrases Update

The list of similar phrases elaborated using the algorithm described above has two potential shortcomings. Firstly, the lists of similar terms frequently include much broader concepts. In the example below, we see ‘pain’ as equivalent to ‘muscle pain;’ which potentially will cause many false alarms. Secondly, generalization can sometime be pursued further – the two lines can be combined together:

```
bóle mięśni # ból, bolesność, bóle
‘muscles pain’ # ‘pain’, ‘ache’, ‘pains’
ból mięśni # ból, bolesność, bóle
‘muscles pain’ # ‘pain’, ‘ache’, ‘pains’
```

To overcome the above mentioned problems without manual effort, we proposed a strategy for cleaning and restructuring the list of similar terms. First, we eliminated equivalents that were shortened to just one word which is the main element of the phrase and occur independently on the symptoms list (as in the examples from the table above). Instead, we used these equivalents to exchange the first element of the phrase (removing repetitions which occur directly within one term description). The results for our example are given below:

```
bóle mięśni # ból mięśni, bolesność mięśni
‘muscles pain’ # ‘muscles pain’, ‘muscles ache’
ból mięśni # ból mięśni, bolesność mięśni
‘muscles pain’ # ‘muscles pain’, ‘muscles ache’
```

The procedure introduces some new (potentially valid) terms (in our example: ‘muscles ache’). It also introduces repetitions between terms. In this case, we use heuristics and we join terms which are lexicographically close, i.e. phrase elements are identical entirely, or at least they are identical for the first two letters, the Levenshtein distance is below 3 and words are longer than 4 letters, and



a word is not written in capital letters. These conditions allow us to cover some plural forms (such as: ‘ból mięśnia’, ‘ból mięśni’) More liberal conditions could give us improper combinations such as ‘*podbrzusze* ‘epigastrium’ and *nadbrzusze* ‘abdomen’. We also avoid identifying acronyms. The final result is:

bóle mięśni # ból mięśni, bolesność mięśni  
 ‘muscles pain’ # ‘muscles pain’, ‘muscles ache’

The problem which is not adequately solved at this moment is negation. Quite a lot of phrases in clinical notes are negated but there are no ready to use tools for Polish that are able to recognize them, similar to NegEX for English [2]. For Polish, the problem of negation in medical texts was addressed in [18]. We tested some simple methods in which we recognized several types of words introducing negation, such as *not*, *lack* and *without*, but in this particular case when texts consist mainly of noun phrases in the nominative, the simple method of recognizing only nominative forms of symptoms gave the best results. In Polish, negated phrases are in other cases, and their orthographic forms usually differ from the nominative ones. There are still some types of phrases which are incorrectly recognized, e.g. *duszność neguje* ‘shortness of breath (he/she) negates’, as they are abbreviated ways of expressing negation which are domain specific and were not identified in advance. In the future, we plan to address the problem in a more robust way.

## 5 Results and Evaluation

We applied the proposed algorithms to the set of data with 382,084 visits of 50,394 patients from different primary health care centers and specialist clinics in Poland which use the same software for data processing and storage. The documents are already segmented into various fields, but we were interested in two fields which have free text form and contain the exact text of examination and interview results written by a medical staff member (usually by physicians themselves). 11,407 of patients had only one visit registered within this data, so there were 38,987 patients left for our evaluation. The average number of visits per patient was ten, but we limited ourselves to the simple case in which we have only a description of the current and the previous visits and we are looking for any potential side effects of the drugs prescribed (newly or as a continuation of a therapy) on the last by one visit as part of the symptoms reported by a patient during the current one. Both interview and examination fields from the previous and the last visits were analyzed and searched for symptoms which can be drug side effects. Newly occurring symptoms were identified and then, for all drugs administered during the last visit, all their possible side effects were compared with this new symptom list. The results of this procedure, using three versions of the possible side effect list, are given in Table 1.

We can observe in Table 1 that adding similar phrases from the Wordnet database did not introduce many new concepts to II list (less than 2%), but due to newly established similarity connections, many more symptoms were identified as possible drug side effects, hence much more (3.5 times more) such alarms were

	I.		II.		III.	
	types	occ.	types	occ.	types	occ.
side effects of all drugs (lists lengths)	13,474	-	13,692	-	15,581	-
possible side effects of drugs administered during the previous visit)	7,573	-	8,620	-	8,720	-
all symptoms registered during the last visits	1,372	-	1,374	-	1,373	-
alarms: symptoms which could be drug side effects	original 143	1,309	284	4,677	215	2,027
	merged 126	-	225	-	186	-

**Table 1.** Side effects identified in the descriptions of the patient visits using different symptom lists. The first list (I) contains only symptoms extracted from textual drug descriptions. II list additionally contains terms obtained from Polish Wordnet as well as conjunct elements extracted from coordinated phrases (described in sections 3.2-3.5). The last list (III) is list II modified (as described in section 3.6) to eliminate terms which are too general, and adds more phrase equivalents. The first row of the table contains the length of the symptom lists. The second presents the number of the symptoms which are identified in the drug descriptions as possible side effects of all drugs administered to patients during the first visit from the analyzed pair. In the third row, all types of symptoms identified during the last visits are shown. The last part of the table presents the final results, i.e. the number of symptoms identified as possible side effects of drugs used by a particular patient. To make the results more comparable, the numbers of different symptom types after merging their names based on the small Levenshtein distance (using the same criteria as described above) are shown in the last line.

raised. As was expected, in a small manually checked sample, quite a few of them were judged as evidently false (e.g. ‘pain’ identified as an occurrence of a ‘back pain’) which justified our next phase of list modification. As was also expected, III list contains significantly more new elements than both I and II lists (about 15%). At the same time, the number of symptoms recognized as potential side effects in the descriptions of the drugs used by the patients is only about 1% larger than in the case of list II. And finally, this time, the final list of the possible side effect alerts is only 50% larger than in the case of list I. These results look promising, as it is more probable that this amount of additional alarms is properly supported. The exact numbers for two specific connected side effects are given in Table 2. What is interesting is that in all three cases, the numbers of types of symptoms identified in the visit descriptions are almost identical. This supports the idea that visit descriptions are written using simpler language and use more typical phrases than drug description, hence it is much easier to cover the ways physicians express symptoms. What is challenging is how to match these symptoms to side effects listed in drug descriptions which are much more formal, complicated and detailed texts.

Actual evaluation of our final solution will be possible only after deploying our method as part of the system used by physicians, which is planned. An introductory evaluation only covers evident false alarms which can be classified on

	I.		II.		III.	
	org. merged		org. merged		org. merged	
ból głowy ‘headache’	52	117	60	138	87	165
bóle głowy ‘headaches’	65	-	78	-	78	-
ból ‘pain’	42	56	778	1951	78	138
bóle ‘pains’	14	-	1173	-	60	-

**Table 2.** Examples of symptoms that are potential drug side effects recognized using different side effect lists in a description of the last visit (the names of the lists are the same as in Table 1).

the basis of general knowledge and the text itself, i.e. using a phrase in a negative scope or in the context of an improving status. Our aim was to determine the potential possibility of an undesirable symptom, while the final decision was left to the physician using the program. The evaluation was therefore performed by a person with experience in medical text annotation and not by a physician. The results on the first 20 examples are shown in the Table 3. Although the sample is small we can already observe that the coverage of symptoms is highly improved by using similar terms in list II. The further modifications make this list much more reliable. The increase in the symptom coverage after adding similar phrases is evident. It is also clear that updating the II list helps in reducing the number of false positives – symptoms wrongly recognized as possible side effects. In the case of list III, no false alarm are raised, while 2 are missing. One alarm is consequently wrongly generated by all the methods because of an error in the initial symptom list.

	correct outcome				missing alerts (FN)				erroneous recognition (FP)				F1
	nb.	%.	TP	TN	entire alerts	symptoms	entire alerts	symptoms	entire alerts	symptoms			
I.	8	.40	5	3	8	.40	3	.20	1	.05	0	.00	.45
II.	9	.45	9	0	2	.10	2	.10	6	.30	1	.05	.62
III.	17	.85	12	5	1	.05	0	.00	2	.10	0	.00	.89

**Table 3.** Manual comparison of results for 20 patients, i.e. first 20 cases for which any of the methods recognized at least one side effect. 5 from these cases are real false alarms (TN). The first part of the table presents the number and the percentage of cases in which the output of the method was correct. The number of times when the method correctly recognized the need of an alarm (TP) and the absence of such a situation (TN) is also added. Then, we give the number of alerts which were not raised at all, or were generated with incomplete lists of symptoms. The next two columns include alerts which are entirely wrong and such that have any additional (incorrect) symptom. Lastly, the F1 measure for the method is included.

## 6 Conclusions

The proposed method for preparing a list of the potential side effect symptoms and their recognition in the actual visit description seem to work with an acceptable level of quality. Although the sample on which the initial evaluation was made is small, the  $F1=0.89$  seems to be quite satisfactory, e.g. compare [10]. Even if in practice it will certainly turn out to be lower, we think that it would be possible to test out the method in a real environment to observe its practical value and shortcomings. In a situation when resources for medical text processing for Polish are very limited, using a general semantic resource such as Wordnet allowed us to improve the results of symptom identification. One of the problems that is not fully solved here, but will also be addressed, is spelling corrections as patients' records contain a large number of spelling errors which affects the recognition of side effect phrases. In the current version of the system, some of the errors are already taken into account by similarity measures, but the problem needs a more general solution.

There are currently 70,000 visits a day processed by the system. Each visit potentially generates an inquiry about side effects. The tests were performed which showed that the system was able to process such a high number of questions on-line. Looking for phrases which might indicate the side effects of drugs in patients' data is executed as separate thread in the EHR system. The dictionary is organized in the TRIE structure which means a quick search is possible.

For further investigation, it would be interesting to compare lists of side effect phrases with all phrases extracted from patient records in order to find other ways of expressing symptoms and side effects in patient records. We also plan to work on eliminating alerts in cases when the symptoms are most likely related with a new illness on the basis of other new symptoms identified simultaneously.

## Acknowledgments

This work was financially supported by the National Centre for Research and Development in Poland, Grant POIR.01.01.01-00-0328/17.

## References

1. Baker, D.: Medication alert fatigue: The potential for compromised patient safety. *Hospital Pharmacy - HOSP PHARM* **44** (06 2009). <https://doi.org/10.1310/hpj4406-460>
2. Chapman, W., Bridewell, W., Hanbury, P., Cooper, G., Buchanan, B.: A simple algorithm for identifying negated findings and diseases in discharge summaries. *Journal of Biomedical Informatics* **34(5)**, 301-310 (2001). <https://doi.org/doi:10.1006/jbin.2001.1029>
3. Cho, I., Kim, J., Kim, J.H., Kim, H.Y., Kim, Y.: Design and implementation of a standards-based interoperable clinical decision support architecture in the context of the Korean EHR. *International Journal of Medical Informatics* **79(9)**, 611 – 622 (2010)

4. Czerw, A., Fronczak, A., Witczak, K., Juszczyk, G.: Implementation of electronic health records in Polish outpatient health care clinics – starting point, progress, problems, and forecasts. *Annals of Agricultural and Environmental Medicine* **23**(2), 329–334 (2016)
5. Dziob, A., Piasecki, M., Rudnicka, E.: plwordnet 4.1—a linguistically motivated, corpus-based bilingual resource. In: Fellbaum, C., Vossen, P., Rudnicka, E., Maziarz, M., Piasecki, M. (eds.) *Proceedings of the 10th Global WordNet Conference: July 23-27, 2019, Wrocław (Poland)*. pp. 353–362. Oficyna Wydawnicza Politechniki Wrocławskiej, Wrocław (2019)
6. Feng, C., Le, D., McCoy, A.: Using Electronic Health Records to Identify Adverse Drug Events in Ambulatory Care: A Systematic Review. *Applied clinical informatics* **10**, 123–128 (12/2019 2019). <https://doi.org/10.1055/s-0039-1677738>
7. Gurulingappa, H., Rajput, A.M., Roberts, A., Fluck, J., Hofmann-Apitius, M., Toldo, L.: Development of a benchmark corpus to support the automatic extraction of drug-related adverse effects from medical case reports. *Journal of Biomedical Informatics* **45**(5), 885 – 892 (2012)
8. Hammar, T., Hellström, L., Ericson, L.: The Use of a Decision Support System in Swedish Pharmacies to Identify Potential Drug-Related Problems—Effects of a National Intervention Focused on Reviewing Elderly Patients’ prescriptions. *Pharmacy: Journal of Pharmacy Education and Practice* **8** (2020)
9. Harpaz, R., Callahan, A., Tamang, S., et al.: Text mining for Adverse Drug Events: the Promise, Challenges, and State of the Art. *Drug Safety* **37**, 777–790 (2014)
10. Henry, S., Buchan, K., Filannino, M., Stubbs, A., Uzuner, O.: 2018 n2c2 shared task on adverse drug events and medication extraction in electronic health records. *Journal of the American Medical Informatics Association : JAMIA* **27**(1), 3–12 (January 2020)
11. Huang, K., Geller, J., Halper, M., Perl, Y., Xu, J.: Using WordNet synonym substitution to enhance UMLS source integration. *Artif. Intell. Medicine* **46**(2), 97–109 (2009). <https://doi.org/10.1016/j.artmed.2008.11.008>
12. Kang, N., Singh, B., Bui, Q.C., Afzal, Z., van Mulligen, E.M., Kors, J.A.: Knowledge-based extraction of adverse drug events from biomedical text. *BMC Bioinformatics* **15**, 1–8 (2014)
13. Kieraś, W., Kobyliński, Ł., Ogrodniczuk, M.: Korpusomat — a tool for creating searchable morphosyntactically tagged corpora. *Computational Methods in Science and Technology* **24**(1), 21–27 (2018)
14. Lemmen, C., Woopen, C., Stock, S.: Systems medicine 2030: A Delphi study on implementation in the German healthcare system. *Health Policy* **125**(1), 104 – 114 (2021)
15. Lindberg, D., Humphreys, B., McCray, A.: The unified medical language system. *Yearbook of medical informatics* **1**, 41–51 (1993)
16. Marciniak, M., Mykowiecka, A., Rychlik, P.: TermoPL — a flexible tool for terminology extraction. In: *Proceedings of LREC*. pp. 2278–2284. ELRA, Portorož, Slovenia (2016)
17. Mikolov, T., Sutskever, I., Chen, K., Corrado, G., Dean, J.: Distributed representations of words and phrases and their compositionality. In: *Advances in neural information processing systems*. pp. 3111–3119 (2013)
18. Mykowiecka, A., Marciniak, M., Kupść, A.: Rule-based information extraction from patients’ clinical data. *Journal of Biomedical Informatics* **42**(5), 923 – 936 (2009)
19. Page, N., Baysari, M., Westbrook, J.: A systematic review of the effectiveness of interruptive medication prescribing alerts in hospital CPOE systems to change

- prescriber behavior and improve patient safety. *International Journal of Medical Informatics* **105**, 22 – 30 (2017)
20. Piasecki, M., Szpakowicz, S., Broda, B.: *A Wordnet from the Ground Up*. Oficyna Wydawnicza Politechniki Wrocławskiej, Wrocław (2009)
  21. Rahmner, P., Eiermann, B., Korkmaz, S., Gustafsson, L., M, G., Maxwell, S., Eichle, H., Vég, A.: Physicians' reported needs of drug information at point of care in Sweden. *Br J Clin Pharmacol.* **73**(1), 115–125 (2012)
  22. Roberts, K., Demner-Fushman, D., Tonning, J.M.: Overview of the TAC 2017 adverse reaction extraction from drug labels track. In: *Proceedings of the 2017 Text Analysis Conference, TAC 2017*, Gaithersburg, Maryland, USA, November 13-14, 2017. NIST (2017)
  23. Saxena, K., Lung, B.R., Becker, J.R.: Improving patient safety by modifying provider ordering behavior using alerts (CDSS) in CPOE system. *Annual Symposium proceedings. AMIA Symposium* **2011**, 1207–1216 (2011)
  24. Scott, I.A., Pillans, P.I., Barras, M., Morris, C.: Using EMR-enabled computerized decision support systems to reduce prescribing of potentially inappropriate medications: a narrative review. *Therapeutic Advances in Drug Safety* **9**(9), 559–573 (2018)
  25. Stearns, M.Q., Price, C., Spackman, K., Wang, A.Y.: Snomed clinical terms: overview of the development process and project status. *Proceedings. AMIA Symposium* pp. 662–6 (2001)
  26. Sutton, R.T., Pincock, D., Baumgart, D.C., Sadowski, D.C., Fedorak, R.N., Kroeker, K.I.: An overview of clinical decision support systems: benefits, risks, and strategies for success. *Digital Medicine* **3**(17), 1 – 10 (2020)
  27. Ubysz, D., Fryzowska-Chrobot, I., Giermaziak, W.: Baza Tez-Mesh jako efektywne narzędzie do opracowania rzeczowego i wyszukiwania informacji z zakresu medycyny i nauk pokrewnych. *Zarządzanie Biblioteką* **11**(1), 59–73 (paź 2019)
  28. van Mulligen, E.M., Fourrier-Reglat, A., Gurwitz, D., Molokhia, M., Nieto, A., Trifiro, G., Kors, J.A., Furlong, L.I.: The EU-ADR corpus: Annotated drugs, diseases, targets, and their relationships. *Journal of Biomedical Informatics* **45**(5), 879 – 884 (2012)
  29. Waszczuk, J.: Harnessing the CRF complexity with domain-specific constraints. The case of morphosyntactic tagging of a highly inflected language. In: *Proceedings of COLING*. pp. 2789–2804 (2012)
  30. Woliński, M.: Morfeusz reloaded. In: Calzolari, N., Choukri, K., Declerck, T., Loftson, H., Maegaard, B., Mariani, J., Moreno, A., Odijk, J., Piperidis, S. (eds.) *Proceedings of the Ninth International Conference on Language Resources and Evaluation, LREC 2014*. pp. 1106–1111. ELRA, Reykjavík, Iceland (2014)
  31. Woliński, M., Saloni, Z., Wołosz, R., Gruszczyński, W., Skowrońska, D., Bronk, Z.: *Słownik gramatyczny języka polskiego*, wyd. IV (2020), <http://sgjp.pl>