

Later, many researchers extend and improve the SOUNDEX algorithm. For example [39] combine the letter - grouping idea of SOUNDEX with the minimum edit distance (MED) measure and describe algorithm called EDITEX which assigns smaller weight for replacing letters belonging to the same group.

In our modified minimum edit distance ratio algorithm (MMEDR) [30] we also assign smaller weights to the transformations that replace phonetically similar letters.

Guy [11] describes an algorithm for identification of cognates in bilingual wordlists based on the recurrent sound correspondences. It estimates the probabilities of phoneme correspondences by using a variant of chi-square statistic on a contingency table, which holds how often given two phonemes co-occur in words of the same meaning. The algorithm used only simple one-to-one phoneme correspondences.

The ALINE algorithm ([17]) is capable to identify phonetic similarity between a pair of phonetically transcribed words. It assigns a similarity score to pairs of transcribed words by decomposing phonemes into elementary phonetic features, such as place of articulation, manner of articulation, voice, etc. Features are assigned a weight based on their relative importance. Feature values are encoded as numbers between 0 and 1. The similarity score is then computed by a dynamic-programming algorithm that finds the optimal sequence of operations insert/delete, substitute, and expand/compress.

Following these ideas of Kondrak [19] also used sound correspondence to identify cognates between languages. His algorithm was initially designed for extracting non-compositional compounds from bitexts but it is also able to find complex sound correspondences in bilingual wordlists, not just simple one-to-one phoneme correspondences.

Phonetic similarity can be measured on the basis of the phonetic transcription of the words: first the words are transcribed as a sequence of sounds represented by characters and then the orthographic similarity between these sequences is measured. Transcription allows measuring phonetic similarity between languages using different alphabets. For example in our modified minimum edit distance ratio algorithm (MMEDR) [30] we perform phonetic transcription to replace Russian letters with their Bulgarian equivalents.

Kondrak and Dorr [16] combine several phonetic and orthographic approaches in their study of the identification of confusable drug names and report high accuracy. They conclude that a simple average of several orthographic similarity measures outperforms all individual measures on the task of the identification of confusable drug names.

Manual  $T-r$  transformation Rules . Rather than applying directly some string similarity measure like MEDR or LCSR some studies first apply a set of transformation rules that reflect some typical cross-lingual transformation patterns observed for given pair of languages . This is absolutely necessary when the languages do not use exactly the same alphabet which requires some letters from the first language to be replaced with letters from the second . This idea can be further developed to replace not just single letters but also syllables , endings and prefixes .

For example in [ 30 ] we apply a set of manually constructed transformation rules for replacing Bulgarian with Russian endings , replace double consonants with single and replace Russian-specific letters with their Bulgarian equivalents . After that we use a modification of MEDR algorithm that assigns weights for the replace operations reflecting some regular phonetic changes between Bulgarian and Russian .

Manually constructed transformation rules between English and German words ( like replacing the letters *k* and *z* by *c* and changing the ending - *t ä t* by - *ty* ) are exploited also by [ 15 ] for expanding a list of cognates .

Learning  $T-r$  transformation Rules . The idea of learning automatically cross-lingual transformation rules that reflect the regular phonetic changes between a pair of languages has been exploited by number of researchers . Such techniques follow naturally the idea of using manually constructed transformation rules .

Tiedemann [ 36 ] used various measures to learn the regular spelling transformations between English and Swedish from a set of known cognate pairs . His best performing string similarity measure algorithm NMmap uses LCSR algorithm to identify the non-matching parts of two strings and statistically assigns weights corresponding to the probability for transforming between them .

The algorithm proposed by Mulloni and Pekar [ 27 ] extracts automatically from a list of known cognates a set of rules that capture regularities in the orthographic transformations between given two languages . These transformations are substitutions of a sequence of letters from the first language with a sequence of letters in the second language identified through the minimum edit distance algorithm . Special characters are added at the word boundaries to allow capturing of rules that transform the start , the middle and the end of the words . For each rule chi-square statistics is calculated and most regular rules are truncated and used while the others are ignored . Finally the transformation rules are applied as a preprocessing step and after that the normalized minimum edit distance is calculated as a similarity measure .

Mitkov et al. [26] use very similar methodology. They collect and score the transformation rules the same way like Mulloni and Pekar [27] but do not account word boundaries as special case. Once the rules are collected and scored by chi-square statistics they apply the rules on candidate pair of words and use LCSR to calculate their similarity.

All of the above techniques use positive examples of cognate pairs to learn regular transformation rules. Unlike them Bergsma and Kondrak [2] use positive and negative examples of cognate pairs to learn positive or negative weights on substring pairings in order to better identify related substring transformations. Starting from minimum edit distance they obtain an alignment of the letters in the given strings and extract corresponding substrings consistent with the alignment. Finally a support vector machine (SVM) is trained by using sets of positive and negative cognate examples and the SVM is used to discriminatively classify given two words as cognates or not.

2.2. Statistical Approach for False Friend Identification. There is no much research concerning extracting false friends directly from text corpora. Most methods (like [26] and [32]) first extract cognates and false friends candidates using some measure of orthographical or phonetic similarity and later try to distinguish between true cognates and false friends.

Fung [10] proposes methods for creating bilingual lexicons from parallel corpora and comparable corpora. His method for extracting semantically related words from sentence level aligned parallel corpus works as follows: for each word pair two binary occurrence vectors are constructed. The first vector maps the occurrences of the first word in the sentences at the left side of the parallel text. The second vector maps the occurrences of the second word in the sentences at the right side of the parallel text. Finally the correlation between these vectors is calculated and used as measure for semantic relatedness.

Brew and McKelvie [5] use sentence alignment to extract cognates and false friends directly from parallel bilingual corpora. The semantic relatedness is identified by statistical method based on collocation analysis in the aligned sentences. The orthographic similarity is measured by various string similarity algorithms. As a result the extracted candidate pairs are classified as cognates, translations, false friends, or unrelated. Their experiments are limited to verbs in English and French but their approach is capable to be applied for other languages as well.

Nakov and Pacovski [29] extract false friends directly from a parallel corpus. Their idea follows the intuition that false friends are unlikely to co-occur

Kondrak [20] extended his algorithm for measuring semantic similarity based on WordNet and used eight semantic similarity levels as binary features: gloss identity, keyword identity, gloss synonymy, keyword synonymy, gloss hyponymy, keyword hypernymy, gloss meronymy and keyword meronymy. These features are combined with a feature based on phonetic similarity and naive Bayes classifier is used to distinguish between cognates and non-cognates.

Mitkov et al. [26] proposed few methods for measuring semantic similarity between orthographically similar pairs of words used to distinguish between cognates and false friends on the basis of similarity threshold estimated on a training data set. Their first method uses comparable corpora and relies on the distributional similarity. For given pair of words a set of  $N$  most similar words are collected using skew divergence [21] as similarity function. The similarity between the words is calculated as Dice coefficient between the obtained sets. A bilingual glossary is used to check if two words can be translations of each other. Their second method extracts co-occurrence statistics for each word of interest from the respective monolingual corpus using a dependency parser. Thus verbs are used as distributional features of the nouns. Semantic vectors are created for the two sets of verbs (using skew divergence again) and similarity between them is measured by Dice coefficient and using a bilingual glossary. The first method requires a glossary of equivalent nouns while the second requires a glossary of equivalent verbs. In the same study the first method is further extended to use taxonomy data from EuroWordNet (when available). The proposed methods are shown to have different performance on different language pairs and none of them was superior to the others.

The idea of using the Web as a corpus has been exploited by many scientists working on different problems (see [14] for an overview). Some of them use Web search engines for finding how many times a word or phrase is met on the Web and extracting pointwise mutual information ([13]), whereas others directly retrieve context from the text snippets returned by the Web search engines ([31]).

The idea of retrieving information from the text snippets returned by Web search engines is used in [6]. The model they introduce is based on the idea that if two words  $X$  and  $Y$  are semantically bound, then searching for  $X$  should cause  $Y$  to appear often in the results, and vice versa: searching for  $Y$  should cause  $X$

to appear often in the results. As it is later discovered by Bollegala et al. [4], this produces incorrect zero semantic similarity for most of the processed pairs.

Bollegala et al. [4] combine retrieval of information about the number of occurrences of two words (both together and individually) from a Web search

**Lemmatization Lexicons** . We used two large monolingual morphological lexicons for lemmatization for Bulgarian and Russian .

The Bulgarian morphological lexicon [33] is created at the Linguistic Modeling Department of the Institute for Parallel Processing in the Bulgarian Academy of Sciences (BAS) and contains about 100 wordforms and 700 lemmata . Each lexicon entry consists of a wordform , a corresponding lemma , followed by morphological and grammatical information . There can be multiple entries for the same wordform , in case of multiple homographs .

The Russian morphological lexicon [33] is also created at the Linguistic Modeling Department of the Institute for Parallel Processing in the Bulgarian Academy of Sciences (BAS) . It is in the same format like the Bulgarian and contains about 1500 wordforms and 100 lemmata . Its core content is based on the grammatical dictionary of [38] .

**Bilingual Glossary** . We used a large Bulgarian - Russian electronic glossary consisting of 59582 pairs of words which are translations of each other . The glossary was adopted by scanning , parsing and processing the Bulgarian - Russian dictionary of [3] and the Russian - Bulgarian dictionary of [7] . We use the word - word translations from these dictionaries ignoring the phrase - word and phrase - phrase translations . Most of the words have multiple translations so we have a set of Russian translation words for each Bulgarian word and vice versa . This is taken into account during the comparison of the Bulgarian and Russian contextual semantic vectors as described in Section 3.4 .

**Searches in Google** . During our experiments we performed searches in Google for 557 Bulgarian and 550 Russian wordforms and collected as many as possible ( up to 10 ) page titles and text snippets from the search results . We used this text information to extract the local contexts of these words and build their contextual semantic vectors as described in Section 3.4 .

**4.2. Experiments** . This section describes the experiments performed with the statistical , semantic and combined algorithms for identification of false friends .

**Baseline** . As baseline we took the following algorithm :

- **ASC** – words pairs sorted in ascending order ( first by the Bulgarian word and second by the Russian word ) . It behaves nearly like a random function .

**Statistical Algorithms** . We performed the following experiments based on the statistical approach for identifying false friends in a parallel text :

