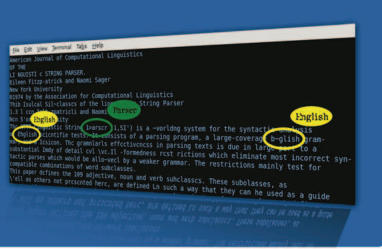# spelling correction

Hailian Jiang, Irene Kolb, Nadya Georgieva

Software Project WS 2010/11
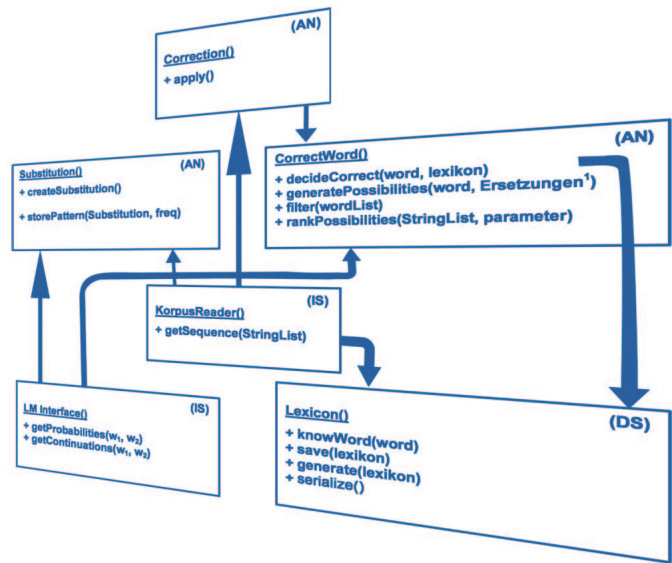
## I. Motivation and Goal:

Spelling mistakes can arise, for example, from incorrect encoding of text which ruins special character (a legacy problem). Another source is OCR (Optical Character Recognition)

Our project's goal is to improve OCR-Output
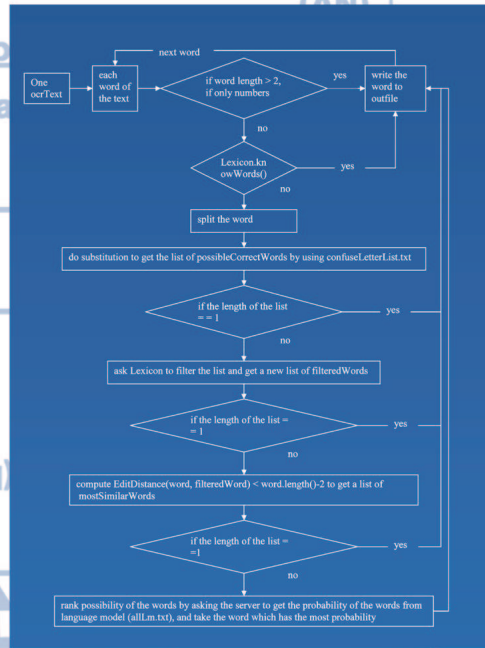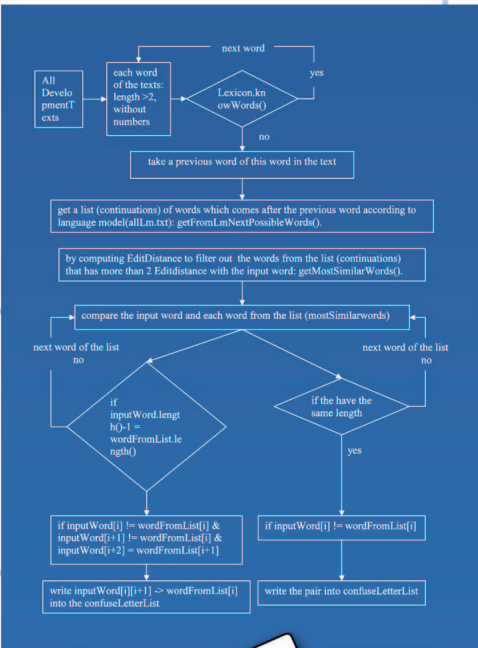
## II. Resources:

- SRILM (Stanford Research Institute Language Modeling Toolkit)
- JSBD-1.6 (Sentence Splitter Tool)
- ACL ARC (Anthology Reference Corpus)



Testdata | 1979 | 1980 | Trainingdata | 2005
Developmentdata

### Class diagram

**Correction()** (AN)
+ apply()

**Substitution()** (AN)
+ createSubstitution()
+ storePattern(Substitution, freq)

**CorrectWord()** (AN)
+ decideCorrect(word, lexikon)
+ generatePossibilities(word, Ersetzungen[1])
+ filter(wordList)
+ rankPossibilities(StringList, parameter)

**KorpusReader()** (IS)
+ getSequence(StringList)

**LM Interface()** (IS)
+ getProbabilities($w_1$, $w_2$)
+ getContinuations($w_1$, $w_2$)

**Lexicon()** (DS)
+ knowWord(word)
+ save(lexikon)
+ generate(lexikon)
+ serialize()

## III. The Method:

### Step 1: Create a substitution list using an unsupervised algorithm



### Step 2: Correct OCR –texts by using the created substitution list



### correction example

EHGLISH → [l> -> b, e -> c, g -> y, h -> n, m -> n, o -> e, w -> u] → [chglish, ehglisn, englisn (e -> c), english (h -> n), ehylish (g -> y)] → ENGLISH

## IV. Evaluation:

n = Number of the letters of each word
m = Number of the possible Substitutions for a letter
s = Number of the letters to be corrected
s ≤ n: Relation between s and n

Formula for the number of the possible corrected words after substitution:
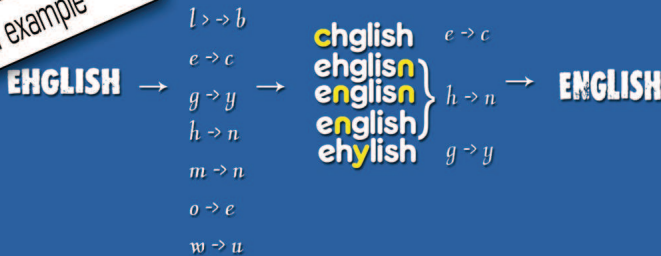
$$(m_1+1)*(m_2+1)*\ldots*(m_s + 1) - 1 = (m+1)s - 1$$

For example: for the word "Granner" : n ==7
If there are substitutions for all letters, then s=n=7
Assume that there are 9 substitutions for each letter, then the number of the possible correct words is:

$$(9+1)*(9+1)*\ldots*(9+1)*(9+1) - 1 = 10^7 - 1 = 10.000.000 - 1 = 9,999,999$$

In order to generate just one correct word, the computer has to generate 9,999,999 possible combinations between those letters. Further more, since the final result is depend on the value of the possibilities between the words after filter, it is quite possible to generate a word that has nothing to do with the word to be corrected. Therefore, the method is insufficient.

1. The substitutionList is too long: for one letter, there are 50 candidates for substitution.
   There is always outOfMemory Exception.
   Reason: There are many non-words from examples of the articles:
   For the word "lle", the following mostSimilarWords are generated:
   [the, old, are, see, lie, oke, she, one, re, cue, lcb, gee, une, due, IXI, ble, PIA, all, BIC, rule, ate, use, are].

2. The substitutionList contains all 26 letters, it means that no matter how many letters are recognized incorrectly by OCR, all letters, both correct and incorrect, of the word to be corrected will be substituted. It is too much unnecessary work.

3. If the substitutionList is reduced to be just 5 candidates for one letter, the program will run endlessly in case that the length of a word is above 8.

4. If the substitutionList is reduced to be just 2 candidates for one letter, the program will run endlessly in case that the length of a word is above 14.

## V. Optimization Suggestions:

1. Reduce the substitutionList by using a real dictionary when the mostSimilarWords are generated; the dictionary will filter out those non-words from examples of the article.
2. Define a method to filter out the part of the examples from the articles.
3. Not to consider the words that contain more than 10 letters.
4. Not to generate the substitutionList, but directly correct the word after the mostSimilarWords are generated

## VI: References:

• Steven Bird, Robert Dale, Bonnie Dorr, Bryan Gibson, Mark Joseph, Min-Yen Kan, Dongwon Lee, Brett Powley, Dragomir Radev and Yee Fan Tan. 2008. The ACL Anthology Reference Corpus: A Reference Dataset for Bibliographic Research in Computational Linguistics. In Proc. Of Language Resources and Evaluation Conference (LREC 08). Marrakesh, Morocco, May.
• Stolcke, Andreas. SRILM – An Extensible Language Modeling Toolkit. Speech Technology and Research Laboratory SRI International, Menlo Park, CA, U.S.A.
• http://www.free-ocr.com/
• http://www-speech.sri.com/projects/srilm/
• http://en.literateprograms.org/Levenshtein_distance_%28Java%29#chunk%20def:base_conditions
• http://java.sun.com/developer/technicalArticles/Programming/serialization/