

DEVELOPING A SPELL CHECKER AS AN EXPERT SYSTEM

Sandor Dembitz*, Petar Knezevic, Mladen Sokele**

*University of Zagreb, Faculty of Electrical Engineering and Computing, Unska 3, Zagreb, Croatia

**Croatian Telecom, Folnegoviceva 5, Zagreb, Croatia

Abstract: We discuss introductory some basic concepts of knowledge engineering, expressing our feeling that there is something missing. According to our opinion, the missing link, able to clarify all the basic concepts, should be the concept of energy, not in use in knowledge engineering practice and theory at all. Since energy is an extremely abstract concept, which led to misunderstandings even in physics, we humbly introduce it by describing an existing expert system, called Hascheck. Hascheck is a spell checker for Croatian language, implemented as learning (semi)automaton, accessible via E-mail, and being in public use for more then 5 years. We also present here a mathematical model of Hascheck's learning. The model led to so-called cognoelectrical analogy, which allowed some energy considerations about learning, and knowledge.

1 Introduction

Knowledge management is currently attracting a great deal of interest in scientific and business communities. However, is knowledge engineering a scientifically well founded discipline? There are some reasons to doubt this.

"Expert systems has to be an applied discipline, not a theoretical one. In application, real-word needs such as knowledge-base maintenance and learning need to be accommodated. We need to see expert systems as primarily a discipline to do with supporting knowledge applications; primarily an applied discipline, primarily concerned with the logic of the knowledge and only secondarily a technological discipline. This needs a modification of viewpoint on the part of many expert systems practitioners." (Taylor, 1999; p. 8)

The quoted thoughts of R.M.Taylor are bringing knowledge engineers in the position of ancient Roman civil engineers, when they were constructing bridges. It is not such a bad position. Not knowing the Newton axioms, they have constructed many bridges, some of which are still in function, from Spain to Near East, from England to Africa. But, for modern engineers this position is very weird; they are taught to think axiomatically. It is not so easy, when knowledge is concerned.

On the other hand, there are some attempts trying to formalise fundamentals of knowledge engineering. How it works will be illustrated by quoting only one

of these attempts, (Uschold, 1998). Introductory M.Uschold says: "Very importantly, *this is a descriptive exercise, not a normative one*" (p. 5). A few pages later he says: "In particular, readers are expected to know how the term KNOWLEDGE is being used, and what INFERENCE is. They should be familiar with fundamental notations of LANGUAGE in general and KNOWLEDGE REPRESENTATION LANGUAGES in particular, as well as what is meant by a KNOWLEDGE BASE" (p. 8).

Do we all mean the same when using terms like *knowledge, reasoning or language*? Surely not! The authors of these lines, when using *language*, primarily think of their own Croatian language.

Is there any knowledge and reasoning out of language? Of course there is. Artists producing paintings or music have such knowledge, and such reasoning. But, is there any sharable, common knowledge and reasoning out of natural languages? We doubt it.

Regarding Uschold's and others papers about theoretical foundations of knowledge engineering we got a strong feeling that there is something missing, something that could connect, physically, all basic weakly defined concepts. Being engineers, we believe that we have found the missing link in the concept of energy.

"Knowledge is power", people say. Everyone will agree that learning is difficult. Thus, in everyday life we describe knowledge and knowledge acquisition, or learning, by energy terms. Nowadays many researchers and engineers are trying to build up learning automata, or knowledge based systems. Do they consider their attempts trough energy analysis? Usually not, because of lack of scientific concepts which could support such considerations.

Even in physics energy is an extremely abstract concept which led to many misunderstandings: only *perpetuum mobile* is to be mentioned as an example. Therefore our paper has to be regarded as a humble - but practical - attempt, pointing where to research further in order to enable energy considerations about knowledge and learning.

At the very beginning we had a simple - for many years now commonly regarded as unattractive - knowledge engineering task: to collect word types and build up an acceptable spell checker for Croatian language. Our approach was unconventional: we

decided to implement our checker as a telematic service embedded in E-mail, so the users had to send their texts to a given address and to wait for automatic reply in a form of checker's report. Texts sent for checking have demonstrated to be a valuable source for acquiring new knowledge for the checker. This provoked our interest for learning in technical systems.

The simplicity of our original task had two faces: positive and negative. We dealt with common knowledge, correctly and incorrectly written Croatian words, so everyone having some Croatian language competence could judge us, whether our product functions well or not. The positive side of our job was that we did not have to cope much with structures, relations and similar, what embarrasses most knowledge engineering tasks. This circumstance allowed us to go fairly deep in modelling learning. The results will be presented here.

Section 2 describes the checker. Section 3 brings a mathematical model of learning process based on data collected during 5 years of checker's life. Finally, Section 4 introduces a cognoelectrical analogy that allows some energy considerations about knowledge and learning.

2 Hascheck - A Learning (Semi)automaton

"Spelling is one of the best examples I've seen of the need for prototyping: build something small, try it, see how useful it is in practice, then modify and extend" (Bentley, 1985; p. 460). This idea led to Croatian Academic Spelling Checker, called Hascheck (an acronym derived from its full Croatian name: Hrvatski akademski spelling checker). Hascheck is the first Croatian spell checker; some others Croatian checkers were developed later (Sokele, 1997). Hascheck functions as a telematic service embedded in E-mail (the address: hacheck@fer.hr). The service has been in operation since March 21, 1994.

Hascheck's initial dictionary counted less than 100,000 common word types. Here we must explain the term word and how we use it. A lemmatised word is what one finds on the left side of a conventional dictionary. By applying its morphology, a natural language produces word forms for each lemmatised word. Not all produced word forms, regarded as alphabetic strings, are necessarily distinct. In English, for example, the noun *work* and the verb *work* are two distinct word forms, because they are two distinct lemmatised word, but one word type. Further, it is necessary to distinguish common words (common nouns, verbs, adjectives etc.) from names (proper nouns, acronyms, foreign words in original orthography etc.), because of their different behaviour in texts. Finally, words in text are tokens. But, not all tokens are words. Some tokens are nonwords

(misspellings or typos), and the spell checker's job is to point them out.

With a small initial dictionary we encountered several problems. Croatian is a Slavic language, as well as Russian, Polish, Bulgarian and others. They all belong to the group of highly inflected languages with a great variety of different word forms - and word types - for each lemmatised word. A Russian estimation sets the lowest limit of dictionary size for an acceptable spell checker at 1,000,000 word types; the upper limit is 100,000,000 word types (Dolgopov, 1986). Having a dictionary 10 times smaller than the estimated lowest limit, we had to find, in order to make the Hascheck report useful for the users, a way how to divide unrecognised tokens from a text on legal word types and nonwords, respectively. An idea how to treat this problem was found in a paper from the very beginning of practical spell checking (Morris&Cherry, 1975).

The AT&T Bell Labs' spell checker *typo* - not in use anywhere for many years now - used the concept of string peculiarity. This concept - despite of its first technical implementation - meant that strings like *approachment* and *apprchment* have to be treated differently. For any person having some knowledge of English the first one seems much more English-like than the second one does. Is it possible to build up a program able to make a similar distinction?

An acceptable solution was found by applying binary n-grams ($n=3,4,5,6$) derived from common word dictionary. It is a classical AI technique used mostly for error corrections (Riseman&Hanson, 1974; Ullmann, 1977). Finally a very human-like fuzzy evaluation of strings not contained in the dictionary appeared, classifying them in several classes of peculiarity. The solution proved to be language-dependant. We experimented both in Croatian and English (Dembitz, 1993), but we never managed to put more than 50% of unknown English word types in the class of the lowest peculiarity; the best result in Croatian was around 80%. By all these attempts the rate of typos and misspellings in the same class varied between 10-15% for both languages.

Using our classifying algorithm we were able to bring the most part of nonword types at the beginning of Hascheck report, opposite to the legal, but in the dictionary not contained Croatian word types, which were pushed to the end of the report. Such selectivity, taking into account a small dictionary volume, i.e. a poor coverage of incoming texts at the beginning of service life, was of great importance for accepting Hascheck as a useful service.

It is the variety of word endings that make a language highly inflected. Text sent for checking offered a valuable source for acquiring word types unknown to Hascheck. Inspired by (Weischedel et al., 1993) we developed a tagging algorithm for spell checking

purposes. The tagging algorithm is applied on collections of Hascheck reports recorded for learning purposes. The tagging for common words is applied on less peculiar strings. Strings beginning with capitals, or having all capital letters, are assumed to be potential names and are treated separately. The tagging for names is applied on all classes.

Our word-guessing method succeeds in tagging correctly more than 70% of unknown Croatian word types. This data refer to common words (common nouns, verbs etc.), while the efficiency of the algorithm on name types (proper names, acronyms etc.) is around 50%. In both cases the noise - percentage of typos and misspellings in a set of tagged strings - is under 10%.

The final result of our research and development was a spell checker functioning as a learning (semi)automaton. Unrecognised tokens from a text are classified by the classifier as more or less peculiar. Tagger processes the less peculiar. According to capitalisation, the tagger separately treats potential common words and potential names. It produces collections of new word types that should be learned. After a minor human supervision and correction (therefrom prefix *semi-* when describing Hascheck) new word types are added to Hascheck name and common word dictionary, respectively. Only accepted common word types are used for updating n-gram base, if necessary.

Hascheck is implemented to be open to all word processors. The service is fairly quick in responding: a book of 100,000 tokens waits about 5 minutes for the Hascheck reply. During 5 years of public life Hascheck received more than 2,000 texts for processing, or a text corpus amounting to more than 12,000,000 tokens. As it was mentioned earlier, Hascheck started functioning with a dictionary counting less than 100,000 common word types. In 5 years its knowledge increased to more than 300,000 Croatian common word types, and more than 80,000 name types.

3 Learning Process

Hascheck learns words to improve its covering of Croatian texts. In order to keep track of learning, an elementary statistical record follows each processing. On the basis of these records two variables, text coverage (TC) and learning index (LI), are calculated:

$$TC = \left(1 - \frac{\text{NumberOfUnrecognizedTokens}}{\text{TextVolume}}\right) \cdot 100 \quad (1)$$

$$LI = \left(\frac{\text{NumberOfWordTypesLearned}}{\text{TextVolume}}\right) \cdot 100 \quad (2)$$

To describe the Hascheck learning process in time it was necessary to find a connection of TC and LI,

respectively, to the volume of previously processed corpus (PPC). PPC, the total amount of text processed by Hascheck before the checking of a particular text file, is the only acceptable measure of Hascheck's maturity. PPC is the natural substitution for time in analytical modelling of Hascheck's learning.

We took 20 LI values obtained by proofreading the Croatian Lexicon (Dembitz&Sokele, 1998). The Croatian Lexicon was split into 20 text files, which came in for checking separately, with long time spans, during the period between November 1994 and May 1997. The Lexicon proofreading was an extremely well controlled process. Therefore, we were sure that statistical data, obtained from this process, are very reliable.

Using PPC as the time variable, we looked for functions that best fit the data. The number of free parameters in all tested functions was limited to 3 or less. The best fitting was obtained by using:

$$LI = a + (100 - a) \cdot e^{-\frac{PPC - \Delta t}{\tau}} \quad (3)$$

The construction of the function (3) needs an explanation. Learning index LI, as defined in (2), cannot exceed the value of 100. The extreme value of LI can be reached only at the very beginning of learning, when the "time", $PPC - \Delta t$, equals 0. The "time" shift Δt is also easy to explain. Hascheck needed some research, development and programming, and processing of a certain amount of text in order to acquire basic Croatian word knowledge, all these expressed through Δt parameter, before it could become a public service capable of receiving the first user's text, and giving a satisfactory response to it.

The result of fitting is presented in Fig. 1, together with optimal values of free function parameters a , Δt and τ , and with obtained correlation coefficient r . It is worth mentioning that even researchers in physics, when experimenting with some fundamental physical law, would be satisfied with the correlation of 0.975.

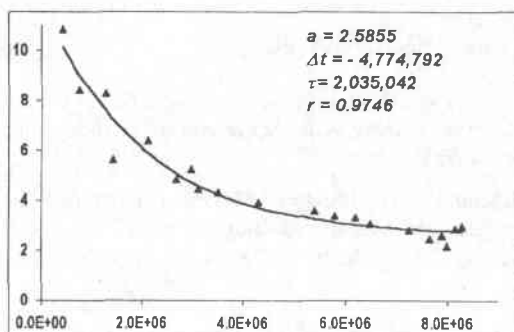


Fig. 1. Fitting of learning index LI
(X-axis=PPC [tokens]; Y-axis=LI)

Independently of learning indexes we took 20 TC values, representing average text coverage calculated on the basis of a 3-month period of processing. These 20 points have a total span of 5 years, or 12,000,000 of tokens, when counted in PPC. By fitting of these data we fixed the parameters Δt and τ on values obtained in previous fitting. The function we chose was a natural choice:

$$TC = b \cdot \left(1 - e^{-\frac{PPC - \Delta t}{\tau}} \right) \quad (4)$$

The result of this fitting is presented in Fig. 2. The fitting was again very well, since the correlation coefficient (parameter r in Fig. 2) was near 0.92.

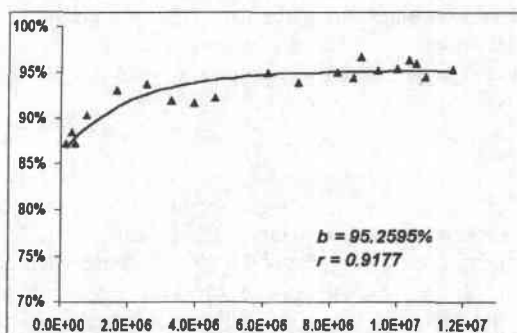


Fig. 2. Fitting of text coverage TC
(X-axis=PPC [tokens]; Y-axis=TC [%])

Functions (3) and (4) correspond to human experience. Learning words is a never-ending process, since every living language constantly produces new words; this is expressed by the small constant a in Fig. 1. "3 to 5% of word tokens are usually missing in the lexicon (*dictionary*) when tagging a real-world text" (Mikheev, 1996). The value obtained for b in Fig. 2 confirms these borders.

This correspondence with experience is giving us the faith that our mathematical model of learning is not fully out of sense.

4 Cgnoselectrical Analogy

Analogy is not the best way of scientific reasoning. However, if there is no better way of thinking, one must take this one.

Functions very similar to (3) and (4) describe the charging of a real capacitor connected to a real battery (Fig. 3) where $\tau = C \cdot R_C \cdot R_s / (R_C + R_s)$.

$$i = \frac{E}{R_C + R_s} + \frac{E}{R_s} \cdot \frac{R_C}{R_C + R_s} \cdot e^{-\frac{t}{\tau}} \quad (5)$$

$$u_C = \frac{E \cdot R_C}{R_C + R_s} \cdot \left(1 - e^{-\frac{t}{\tau}} \right) \quad (6)$$

The power function, $p(t) = i \cdot u_C$, has an extreme if $R_C > R_s$. When the extreme of power exists, it is a maximum in $t = \tau \cdot \ln[2R_C / (R_C - R_s)]$.

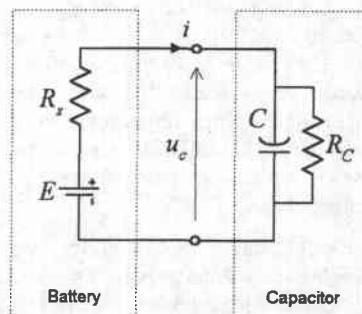


Fig. 3. Electrical equivalence of Hascheck learning process

Words are the charge for Hascheck. Flow of words, expressed by learning index LI, can be regarded as an equivalent to current i . Analogously, the text coverage can be regarded as voltage u_C . Users, with their natural language competence, are the battery; Hascheck is the capacitor. We call this the *cgnoselectrical analogy*.

Having equivalence between LI and i , and between TC and u_C , it is normal to construct and analyse a *power of learning function* (PoL):

$$PoL = \left[a + (100 - a) \cdot e^{-\frac{PPC - \Delta t}{\tau}} \right] \cdot b \cdot \left(1 - e^{-\frac{PPC - \Delta t}{\tau}} \right) \quad (7)$$

The function (7) has an extreme when the "time" has a value of (parameters are taken from Fig. 1):

$$PPC - \Delta t = \tau \cdot \ln \frac{200 - 2a}{100 - 2a} = 1,465,315 [\text{tokens}] \quad (8)$$

Now we must go back to the consideration of the value of "time" shift $\Delta t = -4,774,792$ [tokens], obtained by fitting. As we have already said, this value expresses the fact that Hascheck needed an acquisition of basic word knowledge before it could become a useful service. However, the amount of the "time" shift, nearly 5 Mtokens, cannot be verified by the text corpus used for it, only 800,000 real tokens, processed by Hascheck before it was offered to public use, can be quoted (Dembitz, 1993). The rest has to be considered as another type of time in knowledge acquisition, the time needed for research, development and programming of Hascheck.

Therefore we split the entire life of Hascheck in two periods: its virtual time, approximately 2τ long, when Hascheck was only an idea, and its real time, approximately 7τ long, when Hascheck was coping with real texts.

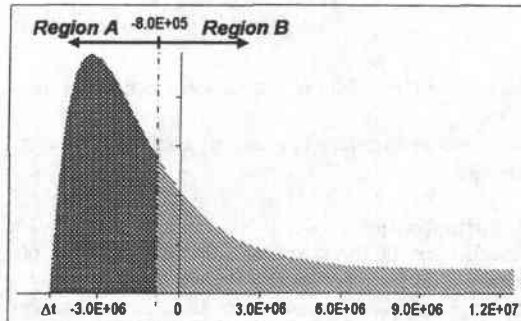


Fig. 4. The Hascheck power of learning function:
Region A = energy spent in virtual time of Hascheck's life;
Region B = energy spent in real time of Hascheck's life;
(X-axis=PPC [tokens]; Y-axis=PoL)

The Hascheck power of learning function PoL is presented in Fig. 4. It is obvious from there that the energy most demanding period was the virtual time of Hascheck's life. So it was, because all doubts, when one is faced with many problems to be solved, consume much energy. During the real time of Hascheck's life, when we had these problems mainly solved, the demand for energy dropped significantly. Here the analytic model proves our own experience.

Generally speaking, these energy considerations do not contradict common sense. Automata are built up to save human energy; users sent their text to Hascheck to save their own time needed for proofreading. In case of learning automata - or semiautomata, like Hascheck - it should be compensated by the authors' energy spent in the course of making such an automaton operable. To put it simply, learning automata function due to the power = knowledge initially supplied by their authors. Weather they will function well or not, it depends upon quality and quantity of knowledge initially supplied, and in each concrete case an estimation of these is still pure art.

5 Conclusion

The results we have presented here are consequences of a well-chosen and very unconventional approach. Being engineers, we always regarded natural language as an economic and energy balanced system. This point of view is not new. It was very modern in technical writing 50 years ago (Shannon, 1951; Zipf, 1949), but later pushed down by influence of structural and similar algebraically based approaches to language (Chomsky, 1957). Inspired by Martinet's philosophy of language (Martinet, 1960),

we decided to explore hidden natural language energy, or geometry of language, in order to solve, with minimal effort, the problem of developing a good spell checker for our highly inflected language. We believe that our achievement, presented here, might produce some impact on the renewal of energy considerations about language, and knowledge, since these two concepts are inseparable.

References

1. Bentley, J. (1985). A Spelling Checker, *Commun. of the ACM* 28(5):456-462.
2. Chomsky, N. (1957). *Syntactic Structures*, Mouton&Co., The Hague.
3. Dembitz, S. (1993). *Automatic Misspelling Detection and Telecommunication Services* (in Croatian), PhD Thesis, Faculty of Electrical Engineering, University of Zagreb.
4. Dembitz, S. & Sokele, M. (1998). Computational Proofreading of the Croatian Lexicon, *Proc. of the 9th Mediterranean Electrotechnical Conference*, pp. 1370-1374, Tel-Aviv, Israel, May 18-20, 1998.
5. Dolgopov, A. S. (1986). Automatic Spelling Correction, *Cybernetics* 22(3):332-339.
6. Martinet, A. (1960). *Eléments de linguistique générale*, Armand Colin, Paris.
7. Mikheev, A. (1996). Unsupervised Learning of Word-Category Guessing Rules, *Proceedings of the 34th Annual Meeting of the Association for Computational Linguistics*, University of California, Santa Cruz, pp. 62-70, 1996.
8. Riseman, E. M. & Hanson, A. R. (1974). A Contextual Postprocessing System for Error Correction Using Binary n-Grams, *IEEE Trans. Comput.*, C-23(5):480-493.
9. Shannon, C. E. (1951). Prediction and Entropy of Printed English, *Bell System Technical Journal*, 30(1):50-64.
10. Sokele, M. (1997). Croatian Spelling Checkers (in Croatian), *WIN.INI* 6(3):38-49.
11. Taylor, R.M. (1999). Expert Systems in the Knowledge Management Era, *Applications and Innovations in Expert Systems VI* (R.Milne et al., Eds), pp. 3-8, Springer-Verlag.
12. Ullmann, J. R. (1977). A Binary n-Gram Technique for Automatic Correction of Substitution, Deletion, Insertion and Reversal Errors in Words, *Computer Journal*, 20(2):141-147.
13. Uschold, M. (1998). Knowledge Level Modelling: Concepts and Terminology, *The Knowledge Engineering Review*, 13(1):5-29.
14. Weischedel, R., Meteer, M., Schwartz, R., Ramshaw, L. & Palmucci, J. (1993). Coping with Ambiguity and Unknown Words through Probabilistic Models, *Computational Linguistics*, 19(2):359-383.
15. Zipf, G. K. (1949). *Human Behavior and the Principle of Least Effort*, Cambridge: Addison-Wesley.