

# Lexical Access, a Search-Problem

Michael Zock (1), Didier Schwab (2) and Nirina Rakotonanahary (2)

(1) LIF-CNRS, TALEP, 163, Avenue de Luminy

(2) LIG-GETALP, University of Grenoble

zock@free.fr, didier.schwab@imag.fr, damanidaddy@msn.com

## Abstract

Our work is confined to *word access*, that is, we present here our ideas of how to improve electronic dictionaries in order to help language producers (speaker/writer) to find the word they are looking for. Our approach is based on *psychological findings* (representation, storage and access of information in the human mind), observed *search strategies* and typical *navigational behavior*.

If one agrees with the idea that lexical access (word finding) is basically a search problem, then one may still want to find out *where* and *how* to search. While the space, i.e. the *semantic map* in which search takes place is a *resource problem*,— any of the following could be used: dictionary, corpus, thesaurus, etc. or a mix of them,— its exploration is typically a *search problem*. Important as it may be, the building of a high quality resource is not the focus of this work, we rely on an existing one, and while we are concerned with its quality, we will be mostly concerned here with search methods, in order to determine the best.

## 1 Problem: find the needle in a haystack

One of the most vexing problems in speaking or writing is that one knows a given word, yet one fails to access it when needed. This kind of search failure, often referred to as *dysnomia* or *Tip of the Tongue-problem*, occurs not only in communication, but also in other activities of everyday life.

Being basically a search problem it is likely to occur whenever we look for something that exists in real world (objects) or our mind: dates, phone numbers, past events, peoples' names, or you just name it.

As one can see, we are concerned here with the problem of words, or rather, how to find them in the place where they are stored: the human brain, or an external resource, a dictionary. Our work being confined to lexical access, we would like to develop a *semantic map* and a *compass* to help language producers to find the word they are looking for. More precisely, we try to build an index and a navigational tool allowing people to access words no matter how incomplete their conceptual input may be. Our approach is based on psychological findings concerning the *mental lexicon* (Aitchison, 2003; Levelt et al., 1999), i.e. *storage* and *access* of information in the human mind, observed *search strategies* and typical *navigational behavior*.

## 2 Consider the following elements before attempting an engineering solution

Before conceiving a roadmap leading to an engineering solution it may be useful to consider certain points. The list here below is by no means complete, neither is the following discussion. Nevertheless we believe that the following points are worth consideration: features of the mental lexicon, how to build and use the resource, searching, ranking and weights, interface problems. For reasons of space constraints we will touch briefly only upon some of these points.

Our main goal is the enhancement of electronic dictionaries to help speakers or writers to find

quickly and intuitively the word they are looking. To achieve this target we take inspiration in the findings concerning the human brain (structure, process) when it tries access words in the mental lexicon.

## 2.1 The mental lexicon, a small-world network?

While *forms* (lemma) and *meanings* (lexical concepts, definitions) are stored side by side in paper dictionaries (holistic presentation), the human brain stores them differently. The information concerning meaning, forms and sound is distributed across various layers. Lexical fragmentation or information distribution is supported by many empirical findings,<sup>1</sup> and while this fact is arguably the reason accounting for word access problems, it is probably also the explanation of the power and the flexibility of the human mind which generally manages to find in no time the right term after having searched for it in a huge store of words.

While it is still not entirely clear what is stored, or whether anything is stored at all<sup>2</sup> coming close to the kind of information generally found in dictionaries, it does seem clear though that the structure of mental lexicon is a multidimensional network in which the user navigates. "Entries in the lexicon are not islands; the lexicon has an internal structure. Items are connected or related in various ways...There are item relations *within* and *between* entries." (Levelt, 1989). While the former relate *meanings* and *forms*: syntactic (part of speech), morphological, phonological information, the latter connect lexical entries.<sup>3</sup> In sum,

<sup>1</sup>*Speech errors* (Fromkin, 1980), studies on *aphasia* (Dell et al., 1997; Blanken et al., 2004) or *response times* i.e. *chronometric studies* (Levelt et al., 1999), *neuroimaging* (Shafto et al., 2007; Kikyo et al., 2001), *eye movements*, (Roelofs, 2004), experiments on *priming* (Schvaneveldt et al., 1976) or the *tip of the tongue problem* (TOT) (Brown and McNeill, 1996).

<sup>2</sup>An important feature of the mental lexicon lies in the fact that the entries are not *accessed* but *activated* (Marslen-Wilson, 1990; Altmann, 1997). Of course, such a detail can have far reaching consequences concerning knowledge representation and use, i.e. structure and process.

<sup>3</sup>These are typically the kind of relations we can find in WordNet (Fellbaum, 1998), which happens to be quite rich in this respect, but relatively poor with regard to intrinsic, i.e. intralexical information.

lexical networks store or encode the information people typically have with regard to words, and finding the needed information, amounts to enter the graph at some point,— in the case of writing or speaking, usually a node dedicated to meaning,— and to follow the links until one has reached the goal (target word). While computer scientists call this kind of search 'navigation', psychologists prefer the term 'activation spreading'. While not being exactly the same, functionally speaking they are equivalent though.

As every day language experience shows, things may go wrong, we lack information, hence we get blocked. Yet when trying to complete the puzzle we do not start from scratch, we rely on existing information, which, in terms of the network metaphor means that we start from (information underlying) a word being close to the target word.<sup>4</sup>

It is interesting to note, that our lexical graphs seem to have similar characteristics as *small-world networks*. These latter are a type of graph in which most nodes, eventhough not being direct neighbors, can be reached via a small number of clicks, about 6, regardless of the starting point. This property of networks, where objects, or the nodes standing for them, are highly connected has first been described by Frigyes Karinthy (1929) a Hungarian writer, to be tested then many years later by a social psychologist (Milgram, 1961). Nodes can be anything, people, words, etc. If they represent people, than edges specify their relationship, i.e. the very fact that they know each other, that they are friends, etc. Given this high connectivity, anything seems to be at the distance of a few mouse clicks. Hence, it is easy to connect people or to find out who entertains with whom what kind of relationship. Obviously, there is a striking similarity to our lexical graphs, and the small-world feature has been tested by mathematicians, who concluded that the distance for words is even smaller than in the original Milgram experiments, namely 4 rather than 6. Indeed, (Motter et al., 2002) and colleagues could show that more than

<sup>4</sup>As TOT experiments have shown (Brown and McNeill, 1996), people always know something concerning the target word (meaning, form, relation to other words), hence finding a word in such a situation amounts to puzzle-completion.

99 percent of the word pairs of their corpus could be connected in 4 steps at the most.

## 2.2 Building the resource

There are two elements we need to get a clearer picture of: the nature of the *resource* (semantic map), and the *search method* i.e. the way to explore it. Concerning the resource, there are many possible sources (dictionary, thesaurus, corpora, or a mix of all this) and many ways of building it. Since our main goal is the building of an index based on the notion of word relations (triples composed of two terms and a link), the two prime candidates are of course *corpora* and *association lists* like the ones collected by psychologists. While the former are raw data, containing the information in a more or less hidden form, the latter (often) contain the data explicitly, but they are scarce, subject to change, and some of the links are questionable.<sup>5</sup>

**Corpora:** Concerning the resource the following points deserve consideration: *size*, *representativity* and *topic sensitivity*.

- *Size or coverage:* While size or coverage are critical variables, they should not be overemphasized though, trading quantity against quality. We need to define the meaning of quality here, and whether, when or how lack of quality can be (partially) compensated by quantity. In other words, we need to define thresholds. In the absence of clear guidelines it is probably wise to strive for a good balance between the two, which again assumes that we know what quality means.
- *Representativity:* Obviously, the system we have in mind is only as good as the data we use, i.e. the purity/accuracy and representativity of the word/feature-association lists.

<sup>5</sup>This flaw is due to the experimental protocol. Subjects are asked to give the first word coming to their mind right after a stimulus. Not having been asked to specify the link it is the experimenter who does so. Yet, many word pairs,— say, cat and dog,— allow for various links (love, tease, chase, etc.), and it is not obvious at all which is the one intended by the user. This problem could have been avoided to a large extent if the instruction had been, "build a *sentence* containing the following word". Another potential problem may be due to the distance between the source and the target word: the link may be mediated.

No single set of data (dictionary, corpus, thesaurus) will ever suffice to capture the knowledge people have. While it would be unrealistic to try to model the semantic map of everyone, it is not unreasonable to try to reach an average user, say someone who has been to school and is a computer literate. If we want to capture the world-knowledge of this kind of user (target), then we must beware that it is contained in the material we use, since our resource will be based on this data. Hence, taking as corpus only the newspapers read by an elite (say, *Le Monde*, in France), will surely not suffice to capture the information we need, as it will not relate information ordinary citizens, say sport fans, are familiar with or interested in. In sum, we need to take a wide variety of sources to extract then the needed information. While there is shortage of some document types needed, there are nevertheless quite a few sources one may consider to begin with: Wikipedia, domain taxonomies, topic signatures, (Lin and Hovy, 2000), a database like (<http://openrdf.org>), etc.

- *Topic sensitivity*

Weights are important, but they tend to change dynamically with time and the topic. Think of the word 'piano' uttered in the contexts of a 'concert' or 'household moving'. It is only in this latter case that this term evokes ideas like *size* or *weight*. The dynamic recomputation of weights as a function of topic changes requires that the system be able to recognize the topic changes, as otherwise it might mislead the user by providing of inadequate weights. For some initial work see (Ferret and Zock, 2006).

**Association lists:** Psychologists have built such lists already decades ago (Deese, 1965; Schvaneveldt, 1989). Similar lists are nowadays freely available on the web. For example, for English there is the Edinburgh Associative Thesaurus <sup>6</sup> and the compilation done by Nelson and his colleagues in Florida <sup>7</sup>. There are also some re-

<sup>6</sup><http://www.eat.rl.ac.uk/>

<sup>7</sup><http://cyber.acomp.usf.edu/FreeAssociation/>

sources for German (see <sup>8</sup> or <sup>9</sup>), for Japanese,<sup>10</sup> and probably many other languages.

While association lists are generally built manually, one can also try to do so automatically or with the help of people (see section 5 in (Zock and Bilac, 2004)). JeuxdeMot (JdM), a collectively built resource focusing on French being an example in case.<sup>11</sup>

### 2.3 Searching

The goal of searching is more complex than one might think. Of course, ultimately one should find the object one is looking for,<sup>12</sup> but the very process should also be carried out quickly and naturally. In addition we want to allow for recovery in case of having taken the wrong turn, and we want to avoid looping, that is, walking in circles, without ever getting closer to the goal. Last, but not least we want to make sure that stored information can also be accessed.

That this is less obvious than it might seem at first sight has been shown by (Zock and Schwab, 2008). Taking two resources (WN and Wikipedia) that contain both a given target word, we wanted to see whether we could access it or not. The target word was ‘vintage’. In order to find it we provided two access keys, i.e. trigger words: ‘wine’ and ‘harvest’. Combining the two produced a list of 6 items in the case of WN and 45 in the case of Wikipedia, yet, while the latter displayed the target word, it was absent from the list produced by WN. This example illustrates the fact that our claim concerning storage and access is well founded. Having stored something does by no means guarantee its access.

In the next sections we will present a small experiment concerning search.

### 3 System architecture

To allow for word access, we need at least two components: an index, i.e. a resource, representing or encoding the way words are connected

(database or semantic network encoding associative relations between words) and an efficient search algorithm to find the needed information, in our case, words.

In other words, since search requires a map or a resource in which to search and a good algorithm to perform the search, we are keen in finding out how different resources (for example, Wikipedia, WordNet or JeuxdeMots) and various search algorithms might affect efficiency of word access. While there is a link between (the quality of) the resource and the searching, we will separate the two, focusing here mainly on the search algorithms and possible ways to evaluate them.

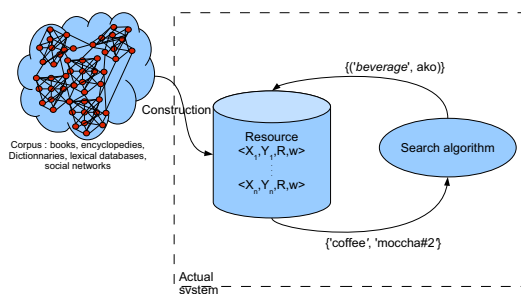


Figure 1: Overall architecture of the system: the resource (association matrix) and a set of search algorithms

### 4 Corpora and resources

Resources are built from corpora which can be of various kinds: news, books, social media, encyclopedias, lexical databases,... They can be general or specific, representing a particular domain. ‘Text genre’ may of course have an impact on what we can expect to retrieve. Obviously, one will not take a database of stock exchange news if one is looking for words referring to tennis- or fishing equipment.

To build our resource, we relied on WordNet (WN).<sup>13</sup> The resource can be seen in various ways: as a semantic network, an association

<sup>8</sup><http://www.schultheimwalde.de/resource.html>

<sup>9</sup><http://www.coli.uni-saarland.de/projects/nag/>

<sup>10</sup><http://www.valdes.titech.ac.jp/terry/jwad.html>

<sup>11</sup><http://www.lirmm.fr/jeuxdemots/rezo.php>

<sup>12</sup>This poses special requirements concerning the organization, indexing and ranking of the data, i.e. words. We will not get into these issues here.

<sup>13</sup>Note, that one may consider WN (Fellbaum, 1998) not only as a dictionary, but also as a corpus. Actually we used precisely this kind of corpus for building our resource.

matrix, or as a list or database of 4-tuples composed of terms, links and weights. These elements can be represented in the following way,  $\langle X, Y, R_Z, w \rangle$ , where  $X$  and  $Y$  are terms or arguments of a given link, whose name expresses the type of relationship holding between them [ $R_Z$  (synonyme, antonyme, hyperonyme, collocation,...)]. Links and terms have a certain weight  $w$  which can be crucial for navigation and information display (interface). Words may be grouped into clusters,<sup>14</sup> and the clusters as well as their elements may be presented in the descending order of the weight: frequent terms being shown on top of the list.

Weights can be calculated in various ways: mutual information, co-occurrences, etc. For example, for corpora, they can be seen as the number of times two items co-occur in a given window (sentence, phrase, paragraph,  $n$  words before or after, ...). Unfortunately, this kind of information is not always available in current resources. For example, WN relates terms (or senses), but does not assign them any particular weight. Yet, this information is very important and might be added via some learning method.

## 5 Search Algorithms

### 5.1 Definitions

Informally, a search algorithm is a method allowing to retrieve a list of terms (candidate target words presented in a given order) from a list of pairs containing the cue- or trigger words and their relations. For example figure 1 shows that the pair  $\langle \text{'beverage'}, \text{'ako'} \rangle$  allows the retrieval of  $\{\text{'coffee'}$  and  $\text{'moccha\#2'}\}$ , two potential target words. More formally, let's define

$$\begin{aligned} f(\{(t_1, R_1), (t_2, R_2), \dots, (t_m, R_m)\}) \\ = \{(T_1, w_1), (T_2, w_2), \dots, (T_m, w_m)\} \end{aligned} \quad (1)$$

where  $t_1, t_2, \dots, t_m$  are keys, cue- or trigger-words,  $R_1, R_2, \dots, R_m$  the type of relation,  $T_1, T_2, \dots, T_m$  candidate target-words and  $w_1, w_2, \dots, w_m$ , the associated weights. The curly brackets  $\{\}$  represent the fact that we have an ordered set.

<sup>14</sup>All words having the same link will be stored and presented together. For example, 'cat' and 'dog' are likely to fall in the category 'animal', while 'hammer' and 'screwdriver' will fall in the category 'tools'.

Indeed, the 'trigger word-relation' pairs are ordered, that is, they parallel the order in which these terms were given as input at query time. The candidate target words are ordered in terms of confidence, ranking which may vary with respect to a given search algorithm. In this paper, confidence is based on the weights in the resource. Of course, one could imagine other ways to define or compute it.

Note that we can have  $R_1 = R_2$ , provided that we do not have at the same time  $t_1 = t_2$ . For instance, while it is possible to have  $\{\langle \text{'island'}, \text{'instance'} \rangle, \langle \text{'island'}, \text{'ako'} \rangle\}$ , one cannot have at the same time  $\{\langle \text{'island'}, \text{'instance'} \rangle, \langle \text{'island'}, \text{'instance'} \rangle\}$

We present in the next sections various ways to use the resource and various search algorithms. In these experiments, we tried to use *direct* and *indirect links* (mediated associations) contained in the tuples and to establish linearly the weight as a function of the position of the word in the list of the trigger words.

### 5.2 General algorithm

In the general algorithm, we consider that our candidate *target words* are at the intersection of the sets corresponding to the pairs of *trigger words* and their *relations*.

$$\begin{aligned} f(\{(t_1, R_a), (t_2, R_b)\}) = \\ f(\{(t_1, R_a)\} \cap f(\{(t_2, R_b)\})) \end{aligned} \quad (2)$$

We will now show, step by step, how  $f(\{(t, R)\})$  is affected by various uses (direct vs. indirect use) and orderings. This will yield 4 kinds of search algorithms.

### 5.3 The use of the tuples

To illustrate our algorithms, let us consider the following resource:

$\langle \text{'mouse'}, \text{'rodent'}, \text{'ako'}, 3 \rangle; \langle \text{'rodent'}, \text{'animal'}, \text{'ako'}, 4 \rangle;$   
 $\langle \text{'rat'}, \text{'rodent'}, \text{'ako'}, 1 \rangle; \langle \text{'rat'}, \text{'animal'}, \text{'ako'}, 2 \rangle;$

#### 5.3.1 Direct use

In this case, we rely only on the direct links  $\langle t, T, R, w \rangle$  of the resource *Res*:

$$f(\{(t, R)\}) = \{(T, W) | \langle t, T, R, w \rangle \in Res\} \quad (3)$$

that is, in the case of direct use, the search algorithm fed with the trigger word  $t$  and the relationship  $R$  found all target words  $T$  contained in the tuple  $\langle t, T, R, w \rangle$  of the resource  $Res$ . The computation of the weight  $W$  is defined in 5.4. For example,  $\langle mouse \rangle$ , would yield  $\langle rodent \rangle$ , while  $\langle rat \rangle$ , would trigger  $\langle animal \rangle$  and  $\langle rodent \rangle$ .

### 5.3.2 Indirect use

In order to boost recall this algorithm takes also indirect links into account.

$$f(\{(t, R)\}) = \{(T, W) | \langle t, T, R, w \rangle \in Res\} \cup \{(T, W) | \langle t, X, R, w_1 \rangle \in Res \text{ and } \langle X, T, R, w_2 \rangle \in Res\} \quad (4)$$

Hence, we consider neighbor words of the neighbors of the trigger words.<sup>15</sup> Again, for  $\langle mouse \rangle$ , we get  $\langle rodent \rangle$  and  $\langle animal \rangle$ , while for  $\langle rat \rangle$ , we continue to get  $\langle animal \rangle$  and  $\langle rodent \rangle$ .

## 5.4 Weighting

### 5.4.1 Basic Weighting

$$\text{For } f(\{(t_1, R_1)\}, \{(t_2, R_2)\}, \dots, \{(t_n, R_n)\}), \\ W = \sum_{t_i, T, R_i, w_j} w_j \quad (5)$$

In our example and for *direct use*, if our trigger word list is  $\{\langle mouse \rangle, \langle rat \rangle\}$  then the weight ( $W$ ) will be 4 (3 + 1) for  $\langle rodent \rangle$  and 6 (4 + 2) for  $\langle animal \rangle$ . In the case of *indirect use*, the weight will be 4 (3 + 1) for  $\langle rodent \rangle$  and 10 (3 + 4 + 1 + 2) for  $\langle animal \rangle$ .

### 5.4.2 Weighting based on the cue-word's position and its relation to other words

Let us suppose that the user gave in this order the following cue words  $\langle A \rangle$ ,  $\langle B \rangle$ ,  $\langle C \rangle$ . In this case we assume that  $\langle A \rangle$  is more important than  $\langle B \rangle$ , which is more important than  $\langle C \rangle$  in order to find the target word. Following this line of reasoning, we may consider the following cases:

$$\text{For } f(\{(t_1, R_1)\}, \{(t_2, R_2)\}, \dots, \{(t_n, R_n)\}), \\ W = \sum_{t_i, T, R_i, w_j} (n - i + 1) \times w_j \quad (6)$$

<sup>15</sup>Please note, we consider here only one intermediate word (two links), as this is, computationally speaking, already quite expensive.

In our example of direct use, if our trigger word list is  $\{\langle mouse \rangle, \langle rat \rangle\}$ ,  $W$  will be 7 ( $2 \times 3 + 1 \times 1$ ) for  $\langle rodent \rangle$  while  $W = 10$  ( $2 \times 4 + 1 \times 2$ ) for  $\langle animal \rangle$ . For indirect use,  $W$  is 4 ( $2 + 1 \times 1$ ) for  $\langle rodent \rangle$  and 17 ( $2 \times (3 + 4) + 1 \times (1 + 2)$ ) for  $\langle animal \rangle$ .

## 5.5 Proposed Search Algorithms

Crossing the characteristics of *weight* (direct vs. indirect) and *use* (direct vs. indirect), we get 4 possible *search methods* : direct use with *basic weighting* (A1) or *linear weighting* (A2); and *indirect use* with *basic weighting* (A3) or *linear weighting* (A4).

## 6 Evaluation

### 6.1 The problem of evaluation.

The classical *in vivo / vitro* approaches do not seem to fit here. While the former tests the system for a given application, the latter tests the system independently. Given the fact that our system has several components, we can evaluate each one of them separately. More precisely, we can evaluate the quality of the *resource* and/or the quality of the *search algorithm*. We will focus here on the search method.

### 6.2 Procedure

The basic idea is to provide each algorithm with an ordered set of *trigger words* and to see how many of them are generally needed in order to reveal the *target word*.

Another way to evaluate the quality of the search mechanism is to check at each step the *position* of the target word in the list generated by the algorithms (output).

### 6.3 Building the test corpus

Psychologists have studied the differences of monolingual and bilingual speakers experiencing the 'tip-of-the-tongue' problem (Pyers et al., 2009). Their experiments were based on 52 pictures corresponding to low-frequency names. Starting from this list we were looking for associated words. In order to build this list we used as resource the results produced by 'Jeux de Mots' (Lafourcade, 2007).<sup>16</sup>

<sup>16</sup>As mentioned by one of the reviewers, we could and probably should have used the Edinburgh Associative The-

### 6.3.1 JeuxdeMots (JdM)

JeuxdeMots, meaning in English 'word games' or 'playing with words', is an attempt to build collaboratively, i.e. via a game, a lexical resource. The goal is to find out which words people typically associate with each other and to build then the corresponding resource, that is, a lexical-semantic network. What counts as a *typical association* is established empirically. Given some input from the system –(term and link, let us say 'Americans' and 'elect as president')– the user produces the associated word –(second term, let us say 'Obama')– answering this way the question, what term *x* is related to *y* in this specific way. Once the network is built, terms should be accessible by entering the network at any point (via some input) and by following the links until one reaches the target word. This is at least the theory. Unfortunately, in practice things do not always work so well.

Actually, JdM has several flaws, especially with respect to access or search. The shortcomings are probably rooted in too heavy reliance on the notion of weight and in excessive filtering of the output, i.e. premature elimination of the candidates presented to the user, list of elements among which the user is meant to choose. Indeed, JdM presents only the highest ranked candidate. Hence, words may never make it to the critical level to be included in the set from which the user will choose the target word. Also, weights do not necessarily correlate with users' interests. This problem can be solved in interactive search, provided that the output contains a critical mass of candidates (possibly organized according to some point of view), but the problem is most likely remain if one presents only one candidate (the highest ranked term), as this latter is not necessarily the target, neither is it always a term from which one would like to continue search.

There is also a problem with the *link names*,

---

saurus (<http://www.eat.rl.ac.uk/>) as it contains authentic word associations collected from people. This point is well taken and we will consider this resource in the future as its coverage is better than our current one and it also avoids possible problems due to the translation. Though being generic, JeuxdeMots has mainly data on French, yet our tests were run on English.

i.e. (metalanguage),<sup>17</sup> though, to be fair, one must admit that identifying and naming links is a very difficult problem. Last, but not least, though more related to the quality of the resource than to the problem of search, there is a chance of user-bias. Indeed, it is not entirely clear whether people really give the first association coming to their mind, or the one fitting them best to continue the game and win more points.

Despite these shortcomings, we will use JdM as a resource as it exists not only for various languages, but is also quite rich, at least for French. Unfortunately, the English version is very poor compared to the French part.<sup>18</sup> This is why we've decided to use the French version for the test corpus.

### 6.3.2 Building the test corpus

Starting from Pyer's list, we translated each term into French and inspected then JdM in order to find the 10 most frequently connected words according to this resource. Next, we translated these terms into English, producing the list shown in the appendix.

As we will see later, in our experiment we do not have *typed relations* between the words. Actually we took from JdM what they call "associated ideas".

Nevertheless, when building the list we did have some problems. Some words do not have any, only one, or simply very few associated ideas. This is particularly true for low frequency words. This being so, we deleted them (in our case 7) from the list.

## 6.4 Description of the tool

Our tool is implemented in Java. To allow for on-line access <sup>19</sup> we use Google's Web Toolkits<sup>20</sup>. The interface is very simple, akin to Google's search engine. At the top of the page the user is invited to provide the input, i.e. the *trigger words*,

---

<sup>17</sup>The term *typical association* is underspecified to say the least.

<sup>18</sup>For example, while for English JdM has by today (july 9, 2010) only 654 relations, the French part contained 1.011.632 the very same day, and 994.889 a month ago.

<sup>19</sup><http://getalp.imag.fr/homepages/schwab>

<sup>20</sup><http://code.google.com/intl/fr/webtoolkit/>

a checkbox allows to choose *relations* and at the bottom are shown the candidate *target words*.

### 6.5 Description of the resource for the experiment

In this experiment, we use the English version of Wikipedia to build our resource. Due to corpus characteristics, only one relation is used: *neighbor* (*ngh*). We consider "words occurring in the same paragraph" as neighbours. After having deleted 'stop words' (articles, conjunction, ...) we lemmatize 'plain words' by using DELA?<sup>21</sup>

For example, a corpus containing the following two sentences "The cat eats the mouse \ The mouse eats some cheese" would yield the following resource :

$\langle \text{'cat'}, \text{'mouse'}, \text{ngh}, 1 \rangle$ ;  $\langle \text{'cat'}, \text{'eat'}, \text{ngh}, 1 \rangle$ ;  
 $\langle \text{'eat'}, \text{'cat'}, \text{ngh}, 1 \rangle$ ;  $\langle \text{'eat'}, \text{'mouse'}, \text{ngh}, 2 \rangle$ ;  
 $\langle \text{'mouse'}, \text{'cat'}, \text{ngh}, 1 \rangle$ ;  $\langle \text{'mouse'}, \text{'eat'}, \text{ngh}, 2 \rangle$ ;  
 $\langle \text{'mouse'}, \text{'cheese'}, \text{ngh}, 1 \rangle$ ;  $\langle \text{'eat'}, \text{'cheese'}, \text{ngh}, 1 \rangle$ ;  
 $\langle \text{'cheese'}, \text{'mouse'}, \text{ngh}, 1 \rangle$ ;  $\langle \text{'cheese'}, \text{'eat'}, \text{ngh}, 1 \rangle$

It should be noted that in this experiment, links are symmetrical.

$$\langle X, Y, \text{ngh}, w \rangle \rightarrow \langle Y, X, \text{ngh}, w \rangle \quad (7)$$

### 6.6 Comparison and evaluation of results

Due to time constraints, we decided to use only a small sample of words, 10 to be precise. Concerning search we have tested two parameters: the *scope* (direct vs. indirect links, i.e. associations, A1 vs. A3) and the *weight* (presence or absence, A2 vs. A4).

The results are shown in table 1, where  $\emptyset$  means that the algorithm did not find any solution, while  $\infty$  implies that the trigger word list has been fully exhausted without being able to produce the target word among the top ten candidates. Indeed, in order to be considered as a hit, the found target word has to be among the top ten.

As one can see our algorithms with indirect use (A3 and A4) never manages to find the target word. Actually, it does not fail totally. It is just that the candidate term appears very late in the list, too late to be considered. The algorithms with direct use (A1 and A2) do find the elusive word or produce a 'list' of zero candidates.

<sup>21</sup><http://infolingu.univ-mlv.fr/DonneesLinguistiques/Dictionnaires/telechargement.html>

TARGET	A1	A2	A3	A4
hive	1	1	$\infty$	$\infty$
peacock	4	4	$\infty$	$\infty$
comet	3	2	$\infty$	$\infty$
microscope	$\emptyset$	$\emptyset$	$\infty$	$\infty$
snorkel	5	5	$\infty$	$\infty$
pitcher	4	3	$\infty$	$\infty$
axe	$\emptyset$	$\emptyset$	$\infty$	$\infty$
gazebo	$\emptyset$	$\emptyset$	$\infty$	$\infty$
hoe	$\emptyset$	$\emptyset$	$\infty$	$\infty$
castle	3	3	$\infty$	$\infty$

Table 1: Comparison of the number of steps needed by each search algorithm to find the target word ie. to put it in the list of the top ten.  $\emptyset$  signals the fact that the algorithm does not find any solution, while  $\infty$  implies that eventhough the trigger word list has been totally used, it did not manage to come up with the target word among the top ten candidates.

target	1	2	3	4	5	6
A1	115	112	325	$\emptyset$	$\emptyset$	$\emptyset$
A2	118	98	273	$\emptyset$	$\emptyset$	$\emptyset$
A3	256	288	254	189	114	59
A4	234	267	262	156	115	54

Table 2: Comparison of the position of the target word at each step of the algorithm for 'microscope'

Table 2 illustrates this last point by showing the position of the *target word* with respect to one of the six *trigger words*.

While the *target word* always appears in A3-A4, A1 and A2 never produce any results beyond the 4th *trigger word*. The two experiments also show that *linear weighting* has hardly any effect on the results.

## 7 Conclusions and Future Work

We have started to characterize lexical access as a search problem. Since search requires a resource in which to search and a good algorithm to perform the search, we were interested in establishing how different *resources* –(Wikipedia (WiP), WordNet (WN), JeuxdeMots (JdM))– and various *search algorithms* might affect efficiency of word access. The focus here has been on the latter.



Next to search algorithms, we presented some methods for evaluating them. While our results are clearly preliminary and on a very small scale, we believe that the questions we have raised are of the right sort. Of course, a lot more work is needed in order to answer our questions with more authority.

## 8 Appendix

**hive** ( $t_w$ ): bee; honey; queen; cell; royal jelly; pollen; wax; group; frame; nest; ( $a_{ws}$ )

**peacock**: bird; feather; animal; spread; shout; blue; tail; disdainful; arrogant; despise;

**comet**: star; space; shooting star; astronomy; sky; galaxy; night; universe; apparition;

**microscope**: small; observe; enlarging; microscopic; observation; ocular; optical; twin; eyeglass; glasses; sea; diving; breath; ocean; mask; flipper;

**pitcher**: jug; jar; carafe; dishes; vase; ewer; container;

**axe**: cut; kill; split; fell; murder; saw; agriculture; arboriculture;

**gazebo**: pavilion; platform; viewpoint; terrace; view; architecture; house; pavillon; esplanade

**hoe** (tool): farming; tool; shovel; pick; technique; spade;

**castle**: tower; king; dungeon; fort; queen; drawbridge; prince; princess; embrasure;

**eclipse**: moon; sun; astronomy; disappearance; darkness;

**bolted joint**: door; lock; padlock; key; close; button; metal; box; house; portal; bolt;

**megaphone**: sound; loudspeaker;

**manta ray**: fish; sea; wing;

**wheelbarrow**: wheel; carry; garden; shovel; gardener; fill; push;

**dynamite**: bomb; explosion; explosive; weapon; wick; chemistry; plastic;

**compass**: direction; navigation; navy; windrose;

**chisels**: cut; scissor; prune; pair; hairdresser; paper; school; chisel;

**ostrich**: egg; bird; Australia; cassowary; emu;

**grater**: woodwork; tool; poverty; polished; dishes;

**braille**: blind; alphabet; writing;

**water well**: dig; drill; pierce;

**guillotine**: scaffold; widow; death penalty; head; reaper; decapitation; execution; torture;

**weathervane**: rooster; direction; wind; rooftop; rotation; east; south; north; west;

**churning** (butter): container; oil; milk; bottle; container; cuve; jerrycan; tank;

**carousel**: fun fair; amusement park; children; entertainment;

**canteen** (place): meal; school; eat; dessert; dish; tableland; entrance; restaurant; food; refectory; supervisor; glass; wood; tail; tooth; trapper; trunk;

**goggles**: sight; glasses; myopia; eyes; rim; twin-lens; optician; sun; vision; see; rectification; improvement; astigmatic; ophthalmologist; farsighted; longsighted; blind; binoculars; nose; optical; pair; telescope; protection;

**boomerang**: Australia; object; flying; throw; come back;

**easel**: painting; drawing; support; tripod;

**propeller**: boat; ship; plane; propulsion; curve; rotation; rolling;

**walnut**: fruit; hazelnut; almond; tree; cashew; nutcracker; oil; salad; woodwork;

**catapult**: ejection; old weapon; throw; stone; aircraft carrier; crossbow; ballista; sling;

**udder**: milk; cow; chest; nipple tit;

**gyroscope**: direction; rotation; instruments; gyrost;

**mummy**: pharaoh; Egypt; pyramid; strip; dead; embalmed; sarcophagus; fruits; funeral;

**hinge**: junction; woodwork; middle; locksmiths; assemblage;

**harmonica**: music; instrument; breath; flute; musical instruments;

**metronome**: musical, tempo, musical instruments;

**noose**: hang; boat; attach; bind; cord;

**harp**: musical instruments; zither; lute; lyre; psaltery;

**slingshot**: weapon; attack; projectile weapon; catapult;

**eiffel tower**: Paris; steel; monument; metal;

**syringe**: drug; injection; sting; nurse; sick; ill; bodycare; drug addiction;

**Words containing too little information to be usable for tests** baster, unicycle, thermos,

antlers, plunger, cleftchin, handcuffs.

## References

- Aitchison, J. 2003. *Words in the Mind: an Introduction to the Mental Lexicon (3d edition)*. Blackwell, Oxford.
- Altmann, G. T. M., 1997. *The ascent of Babel: An exploration of language, mind, and understanding*, chapter Accessing the Mental Lexicon: Words, and how we (eventually) find them. Oxford University Press.
- Blanken, G., F. Kulke, T. Bormann, B. Biedermann, J. Dittmann, and C.W. Wallesch. 2004. The dissolution of word production in aphasia: Implications for normal functions. In Pechmann, T. and C. Habel, editors, *Multidisciplinary Approaches to Language Production*, pages 303–338. Mouton de Gruyter, Berlin.
- Brown, R. and D. McNeill. 1996. The tip of the tongue phenomenon. *Journal of Verbal Learning and Verbal Behaviour*, 5:325–337.
- Deese, J. 1965. *The structure of associations in language and thought*. Johns Hopkins Press.
- Dell, G.S., M.F. Schwartz, N. Martin, E.M. Saffran, and D.A. Gagnon. 1997. Lexical access in aphasic and nonaphasic speakers. *Psychological Review*, 104(4):801–838.
- Fellbaum, C., editor. 1998. *WordNet: An Electronic Lexical Database and some of its Applications*. MIT Press.
- Ferret, O. and M. Zock. 2006. Enhancing electronic dictionaries with an index based on associations. In *ACL '06: Proceedings of the 21st International Conference on Computational Linguistics and the 44th annual meeting of the ACL*, pages 281–288.
- Fromkin, V. 1980. Errors in linguistic performance: Slips of the tongue, ear, pen and hand.
- Kikyo, H., K. Ohki, and K. Sekihara. 2001. Temporal characterization of memory retrieval processes: an fMRI study of the tip of the tongue phenomenon. *European Journal of Neuroscience*, 14(5):887–92.
- Lafourcade, M. 2007. Making people play for Lexical Acquisition with the JeuxDeMots prototype. In *7th International Symposium on Natural Language Processing*, page 7, Pattaya, Chonburi, Thailand.
- Levelt, W., A. Roelofs, and A. Meyer. 1999. A theory of lexical access in speech production. *Behavioral and Brain Sciences*, 22(1):1–75.
- Levelt, W. 1989. *Speaking : From Intention to Articulation*. MIT Press, Cambridge, MA.
- Lin, C-Y and E. H. Hovy. 2000. The automated acquisition of topic signatures for text summarization. In *COLING*, pages 495–501. M. Kaufmann.
- Marslen-Wilson, W.D. 1990. Activation, competition, and frequency in lexical access. In Altmann, G.T.M., editor, *Cognitive Models of Speech Processing: Psycholinguistics and Computational Perspectives*, pages 148–172. MIT Press, Cambridge, MA.
- Milgram, S. 1961. The small world problem. *Psychology Today*, 1(1):61–67.
- Motter, A. E., A. P. S. de Moura, Y.-C. Lai, and P. Dasgupta. 2002. Topology of the conceptual network of language. *Physical Review E*, 65(6).
- Pyers, J. E., T. H. Gollan, and K. Emmorey. 2009. Bimodal bilinguals reveal the source of tip-of-the-tongue states. *Cognition*, 112(2):323 – 329.
- Roelofs, A. 2004. The seduced speaker: Modeling of cognitive control. In Belz, A., R. Evans, and P. Piwek, editors, *INLG*, volume 3123 of *Lecture Notes in Computer Science*, pages 1–10. Springer.
- Schvaneveldt, R., D. Meyer, and C. Becker. 1976. Lexical ambiguity, semantic context and visual word recognition. *Journal of Experimental Psychology/ Human Perception and Performance*, 2(2):243–256.
- Schvaneveldt, R., editor. 1989. *Pathfinder Associative Networks: studies in knowledge organization*. Ablex, Norwood, New Jersey, US.
- Shafto, M. A., D. M. Burke, E. A. Stamatakis, P. P. Tam, and L. K. Tyler. 2007. On the tip-of-the-tongue: Neural correlates of increased word-finding failures in normal aging. *J. Cognitive Neuroscience*, 19(12):2060–2070.
- Zock, M. and S. Bilac. 2004. Word lookup on the basis of associations : from an idea to a roadmap. In *Workshop on 'Enhancing and using electronic dictionaries'*, pages 29–35, Geneva. COLING.
- Zock, M. and D. Schwab. 2008. Lexical access based on underspecified input. In *COGALEX, Coling workshop*, page 8, Manchester.