

United in Diversity: Dutch Historical Dictionaries Online

Katrien Depuydt

Jesse de Does

Instituut voor Nederlandse Lexicologie

*The Integrated Language Database of Dutch (ILD) is a project of the Institute for Dutch Lexicology in Leiden, which integrates corpora, computational lexica and dictionaries describing the Dutch language from ca. 500 until the present. In 2007, the dictionary component was released, already containing two major historical dictionaries of Dutch, the *Woordenboek der Nederlandsche Taal* (WNT, Dictionary of the Dutch Language, 1500-1976) and the *Vroegmiddelnederlands Woordenboek* (VMNW, Dictionary of Early Middle Dutch, 1200-1300). When, by 2009, the *Middelnederlandsch Woordenboek* (MNW, Dictionary of Middle Dutch, ~1250 - 1550) and the *Oudnederlands Woordenboek* ("ONW", Dictionary of Old Dutch, a current project at INL, to be finished in 2008, ca. 500-1200) will have been added, researchers of Dutch will have access to dictionaries covering the complete history of the Dutch language. The choice of a single application, integrating the dictionaries so that a user might query one or more dictionaries simultaneously, was a logical step because of the complementary nature of the dictionaries. The challenge was not only providing the user with optimal access to the dictionary information, but also doing so without compromising the uniqueness of each individual dictionary. We sketch the principles underlying the application.*

1. Introduction

The *Integrated Language Database of Dutch* (ILD) is a project of the Institute for Dutch Lexicology in Leiden in which corpora, computational lexica and dictionaries describing the Dutch language from ca. 500 until the present will be integrated¹. In 2007, the dictionary component was released, already containing two major historical dictionaries of Dutch, the *Woordenboek der Nederlandsche Taal* (WNT, Dictionary of the Dutch Language, 1500-1976) and the *Vroegmiddelnederlands Woordenboek* (VMNW, Dictionary of Early Middle Dutch, 1200-1300). When, by 2009, the *Middelnederlandsch Woordenboek* (MNW, Dictionary of Middle Dutch, ~1250-1550) and the *Oudnederlands Woordenboek* ('ONW', Dictionary of Old Dutch, a current project at INL, to be finished in 2008, ca. 500-1200) will have been added, researchers of Dutch will have access to dictionaries covering the complete history of the Dutch language.

Choosing for one application, integrating the dictionaries so that a user can query one or more dictionaries simultaneously, was a logical step to take because the dictionaries complement each other. The challenge was not only to give the user optimal access to the dictionary information, but also to do so without compromising the uniqueness of each individual dictionary.

2. Dictionary data model

All dictionaries were already available in digital form. We started by first analysing the content, then the level of encoding and finally the applied encoding system. A thorough analysis of the content of each dictionary revealed that, in spite of obvious differences, they are very similar as to their macrostructure: headword, section with linguistic information at entry level, section with semantic analysis of the headword and section with related entries. They differed however greatly as to their level of encoding. In the original WNT data, the sense hierarchy of the article is encoded, but

¹ For information on the project, see Kruyt (2004). When the project started in 2000, the assumption was that the oldest Dutch material would not go back further than the 8th century. The editors of the ONW (Dictionary of Old Dutch) discovered later on that there is Old Dutch material dating from around 500.

individual citations only sporadically. The opposite is true for the MNW². As for the VMNW and ONW, the situation is close to ideal: virtually every information category is distinguishable, either as a table in a relational database (VMNW) or encoded in the XML of the article (ONW). Since for each dictionary, the encoding system was different, and there was no compelling reason to use any of them, we chose to standardize the data by converting it to the XML version of TEI P4 (Text Encoding Initiative³). It is not only widely used for online publishing of dictionaries (Grimm⁴, *Mittelhochdeutsche Wörterbücher im Verbund*⁵, *Anglo-Norman Dictionary*⁶), but application to our data was pretty straightforward. More important to us was the fact that it enables both fine-grained and coarse-grained encoding. We decided to convert all available encoding in each dictionary to TEI, and we established a minimal level of encoding required for all dictionaries. Thus we did not need to impose one dictionary structure and level of encoding upon the others, but were still able to have simultaneous retrieval on the dictionaries. Achieving the minimal level of encoding implied a lot of data work for some of the dictionaries. Some additional data development was done for the sake of simultaneous retrieval. We have added a Modern Dutch (equivalent) lemma to each headword, so as to deal with the different headword spellings each dictionary has according to the language period it describes. And we mapped the indication of part of speech in each dictionary to a uniform one. Finally abbreviated variant forms, compounds and derivatives were expanded. The latter two also received encoding as a headword, so that it does not matter how a particular word is treated in a dictionary: as headword or related entry.

3. Dictionary application model

Since the dictionary application is freely accessible, we will only go into some general underlying principles. As mentioned before, we did not want to integrate the dictionaries by mere extensive linking; we wanted to enable integrated searches, respecting each dictionary's own information categories. In the application, a user can select one or more dictionaries and for instance search for a headword in the selected dictionaries simultaneously. By using the Modern Dutch lemma as the search key, this can be done without knowledge of the historical spelling. Since the dictionaries share a minimal level of encoding, simultaneous searches on other information categories within the dictionaries are also possible. When a search is not applicable to one of the dictionaries, the search field is greyed out when only that particular dictionary is selected, or, in case of combined dictionary searching, no results from that particular dictionary are obtained. Another important issue for us was to approach the dictionaries in a more corpus-like fashion, meaning that we aimed at providing the user as much relevant information as possible without forcing him/her to read through a complete dictionary article. This was necessary because of the length of the dictionary articles: an article like *water* in the WNT for instance contains over 5700 citations and more than 300 senses and subsenses, totalling 144,450 words⁷. The corpus-like approach is visible in several places. When looking for a word or words in a sense, a citation or in a full article in the simple search option, the result will be displayed as concordances with the number of results per article (cf. fig. 1).

² Both WNT and MNW first appeared on CD-ROM (not available anymore).

³ <http://www.tei-c.org>

⁴ <http://germazope.uni-trier.de/Projects/DWB>: the online version of the *Deutsches Wörterbuch* of Jacob and Wilhelm Grimm. The dictionary is also available on CD-ROM.

⁵ <http://germazope.uni-trier.de/Projects/MWV>: the online version of the *Mittelhochdeutsches Wörterbuch* of Georg Friedrich Benecke, Wilhelm Müller and Friedrich Zarncke, the *Mittelhochdeutsches Handwörterbuch* by Matthias Lexer, the *Findebuch zum mittelhochdeutschen Wortschatz* of Kurt Gärtner, Christoph Gerhardt, et. al and of the *Nachträge zum Mittelhochdeutschen Handwörterbuch von M. Lexer*. There is also a version of this application on CD-ROM.

⁶ <http://www.anglo-norman.net/>: the online version of the *Anglo-Norman Dictionary* of William Rothwell, Stewart Gregory, William Rothwell, David Trotter et al.

⁷ An average novel like *Pride and Prejudice* has about 124,000 words.

If the concordances do not provide enough information, or more context is needed, the user can always view the article the concordances were found in. Within advanced search, the user can adapt the display of the search results to the kind of research he or she wants to do. Either the result is a list of articles, like in classic dictionary applications, or he can opt for a list of dictionary meanings, citations, collocations, or head sections of a dictionary article (cf. fig. 2).

The screenshot shows the INL GTB dictionary interface in Mozilla Firefox. The search results are displayed in a table with columns: Nr., Wdb, Trefwoord, Originele spelling, Woordsoort, Frequentie, and Concordantie. The results are for the word 'water' in the WNT dictionary. A pop-up window titled 'Concordanties binnen het artikel 'WATER' in het WNT (Frequentie: 21)' shows the concordances within the article 'WATER' in the WNT dictionary. The concordances are as follows:

Nr.	Wdb	Trefwoord	Originele spelling	Woordsoort	Frequentie	Concordantie
1	VMNW	water	WATER	znw.v.,o.	1	ementen, t.w. aarde, water, vuur en lucht. De verhouding tusse
2	WNT	water	WATER	znw.(o.,v.)	21	(Zich) voor water en vuur weten te bewaren e.d. Zie Dl. XXIII, elementen enz., zooals land, lucht, vuur e.d. In deze toep. is de gedachte aan tgov. andere elementen of krachten, inz. vuur en wind. Vaak in toep. waarbij water waterketel zijn eentonig liedekke boven het vuur zingt, snieders 15, 52 (ed. 1926) [1863]. Uw Vadren, Belgen! streën met vuur. Met water, honger, pest; er Als water en vuur overeenkomen, heelemaal niet overeenkomen. Water in de eene hand en vuur in de andere hand dragen e.d. Zie Dl. Wanneer men water in een pot op vuur zet, wordt het warm, het zet zich uit beurt in de aarden test het kooltje vuur dat ze heeft gekocht voor een halven — In spreekw. enz. waarin water met vuur (of branden e.d.) genoemd wordt, vaak vaak als tegenstrijdige zaken. Vgl. ook vuur (I), III). Iem. water en vuur ontzeggen. Zie Dl. XXIII, 1401. overvallen te worden, gaf schrickelijk vuur van alle kanten, witsen, Scheepsb. 459 — De een schreeuwt om vuur, de ander om water, modderman, Bijdr. Waterberoeringen, Uitbarstingen van Vuur en opwerpingen van Bergen of Eilanden (Als) water en vuur zijn. Zie Dl. XXIII, 1406 en nog de Ik heb van mijn leven wel voor heeter vuur geseten (sey Lammert) en stont tot zijn in, en biedt als ware het tegen het vuur eenen scheidsmuur van water, of eene der Grieken: aarde, lucht, water en vuur, niet gekend hebben, dan kan men het Vriendschap is noodzakelijker dan water en vuur. Zie Dl. XXIII, 585 en nog de volg.
3	WNT	waterbad	WATERBAD (I)	znw.(o.)	3	

Figure 1. KWIC-view

The screenshot shows the INL GTB dictionary interface in Mozilla Firefox. The search results are displayed in a table with columns: Nr., Wdb, Trefwoord, Originele spelling, Woordsoort, and Betekenis. The results are for the word 'vlinder' in the WNT dictionary. The results are as follows:

Nr.	Wdb	Trefwoord	Originele spelling	Woordsoort	Betekenis
50	WNT	uil	UIL (I)	znw.(m.,v.)	..2.b (Overdr.) (Amst.) Weesjongen. Zoo genoemd naar de kleuren van de kleeding der weeskinderen in Amsterdam (zwart en rood) die overeenkomen met die van de een bep. soort vlinder: het weeskind (zie BREHM-HUIZINGA 3, 494 [1910]).
51	WNT	vijfwouter	VIJFWOUTER	znw.(m.)	1 Vlinder, kapel.
52	WNT	vleek	VLEEK	znw.(v.)	4 Opm. De bij KIL. [1599] en in enkele latere wdb. vermelde bet. 'vlinder' ('Vleke, vleken, sax. sic. j. pepel, Pzolin') berust wsch. op een verkeerde lezing van CHYTRAEUS: deze heeft Vleken, waarin de beginletter als u moet worden opgevat (DE SMET in Zs. Bijvorm van fenter; zie ald., en nog de volg. aanh. Daarnaast ook vlinder (TER LAAN).
53	WNT	vlender	VLENDER	znw.(m.)	..5.c Vand. van de ziel (onder het beeld van een vlinder).
54	WNT	vlerk	VLERK	znw.(v.,m.)	..5.c Vand. van de ziel (onder het beeld van een vlinder).
55	WNT	vlichelter(e)	VLICHELT(E)	znw.(m.,v.)	Een alleen bij KILIAAN aangetroffen benaming voor vlinder, die teruggaat op vlielter(e), onder invloed gekomen van vlichelen (vgl. PAUWELS in H. Top. Dial. 9 [1935], 344).
56	WNT	vliegwouter	VLEGGWOUTER	znw.(m.)	Een door KILIAAN en zijn navolgers vermelde etymologiseerende variant van vliwouter 'vlinder', waarmaast volgens versch. idiotica ook vliege(n)bouter schijnt voor te komen (SCHUERM. [1865-1870]; JOOS [1900-1904]; TEIRL.).
57	WNT	vliemel	VLIEMEL	znw.(m.)	Eenmaal aangetroffen als een der vele namen voor den vlinder. Misschien een combinatie van de eerste

Figure 2. Result from advanced search for definitions containing *vlinder* (butterfly).

An specially designed highlighting mechanism makes sure that there is a one to one correspondence of search results and hit highlighting in the article view (cf. fig. 3).

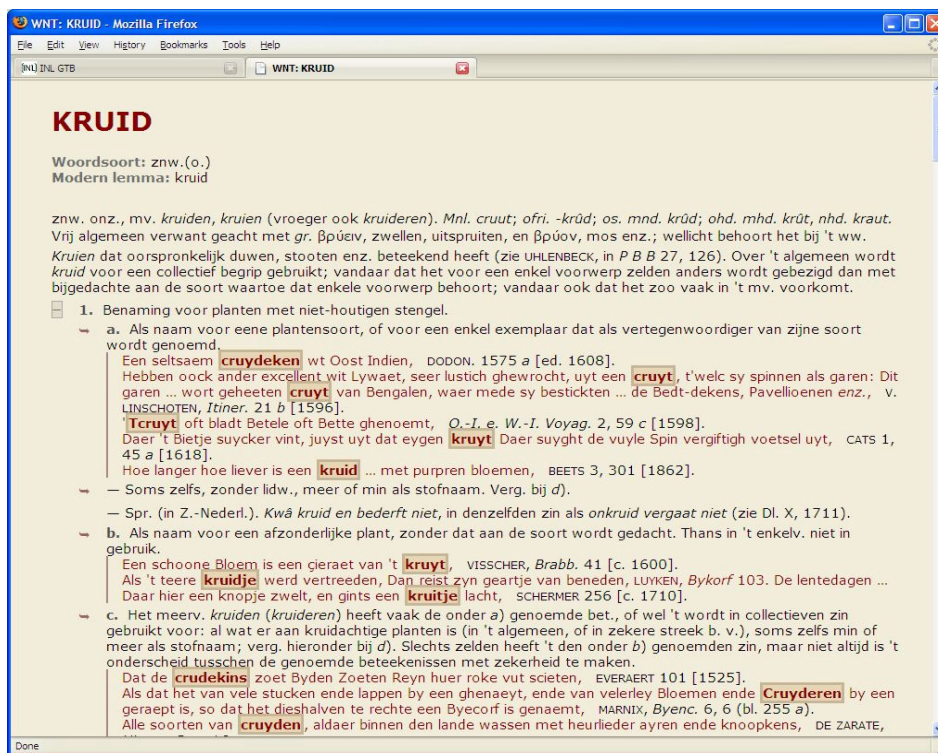


Figure 3. Highlighting of hits after regex search for variants of *kruid* (herb) in citations.

To make the dictionary application even more suitable for research purposes, all search results can be exported in HTML, XML or in CSV⁸. Searches with regular expressions are possible in every search field of the application.

4. User interface

We used the Openlaszlo platform⁹ (release 3.3.1) with flash object code to ensure cross-browser compatibility of the user-interface.

The aim of platform-independency was achieved. Though the interface was implemented before the advent of Windows Vista or even Microsoft Internet Explorer 7 and the recent beta version of IE8, no changes in the Openlaszlo part of the user interface were necessary for these platforms.

5. Search engine

The search engine was implemented independently of the user interface. A preliminary investigation showed that neither XML-based nor relational databases were adequate to the task of combining the necessary full-text search features with structured querying.

We ended up using a combination of a MySQL database¹⁰ for the storage and retrieval of the search results, and the open source Lucene search engine¹¹ (release 1.9) for full-text search functionality.

To ensure an efficient evaluation of a combined query on the entry, the sense and the quotation level, no preliminary computation of the partial results on each level is performed. Projected streams of

⁸ HTML: HyperText Markup Language; XML: eXtensible Markup Language; CSV: comma-separated values.

⁹ <http://www.openlaszlo.org>.

¹⁰ <http://www.mysql.com/>

¹¹ <http://lucene.apache.org/>

results on each level (*spans* in Lucene terminology) are combined on the level of the specified granularity of the search result. The resulting search engine can handle complex combinations of queries without loss of efficiency.

6. Access

The dictionary application is accessible without payment after a simple one time registration, providing the user with username and password. The username and password have to be entered only once for each user on a single workstation. The application URL is <http://gtb.inl.nl>.

7. Future work

The next major milestone in the development of the Integrated Language Database will be the addition of the dictionary of the Middle Dutch Dictionary and the Dictionary of Old Dutch. Since 2007, work has also started on the lexicon component of the ILD, by the development of a large integrated lexicon of the Dutch Language, a diachronic lexicon for 6th – 21st century Dutch (the so-called “GiGaNT”-lexicon). With the help of this resource (which will of course also build on existing language resources), we hope to achieve the integration of corpus and dictionary material without necessarily lemmatizing the full corpus text. This means we will have to extend the use of modern search keys from dictionary lemmata to inflected forms in running text. Given the amount of orthographic and other linguistic variation in Dutch historical documents, this is a major challenge.

References

- Kruyt, J. G. (2004). “The Integrated Language Database of 8th - 21st-Century Dutch.” In Lino, M. T.; et al. (eds.). *Proceedings of the 4th International Conference on Language Resources and Evaluation*. Paris: ELRA. 1751-1754.