

# Towards a lexical database of Dutch taboo language

Gerhard B van Huyssteen<sup>1</sup>, Carole Tiberius<sup>2</sup>

<sup>1</sup> Centre for Text Technology (CtexT), North-West University, Potchefstroom, South Africa

<sup>2</sup> Instituut voor de Nederlandse Taal, Leiden, The Netherlands

E-mail: Gerhard.VanHuyssteen@nwu.ac.za, Carole.Tiberius@ivdnt.org

## Abstract

Over the past 45 years, at least eighteen Dutch paper-based dictionaries of taboo-language (or taboo-related language) have been published (i.e., as visible works of lexicography). However, none of these are available as (linked) lexical data that could be integrated in natural language processing (NLP) tools and applications (i.e., as invisible works of lexicography). In this paper, we describe the development of a comprehensive lexical database of taboo language (LDTL) for Dutch (TaboeLex) that can be integrated in NLP tools and applications. TaboeLex will be made available as open data, i.e., as a freely available, structured, annotated lexicon that can be linked to other data in the future. The paper focusses on the first phase of the project, namely, to define and design TaboeLex.

**Keywords:** Dutch; lexical database; swearword; taboo language

**Warning:** This paper contains content that may be offensive or upsetting.

## 1. Introduction

Despite giant strides that have been made over the past thirty years in digitalising and automating lexicographic work, resources for specialised purposes and non-mainstream languages are still often neglected. As a case in point, even though at least eighteen Dutch paper-based dictionaries of taboo words (see 2.1 for a definition) have been published over the past 45 years (i.e., as visible works of lexicography), none of these are available as (linked) lexical data that could be integrated in natural language processing (NLP) tools and applications (i.e., as invisible works of lexicography).

Lexical databases of taboo language (LDTLs) are specialised digital resources that could be used as sources of linguistic and extralinguistic knowledge in many natural language processing (NLP) systems (see 2.2). Although such an LDTL could be simply a wordlist, for our purposes we consider an LDTL a digital collection of linguistic constructions that has been annotated or enriched in some way (e.g., with part-of-speech information, offensiveness ratings, meanings), and that is structured (e.g., encoded in XML). Most often, the primary use of LDTLs is to recognise words that could be potentially offensive to a specified community of language users (e.g., children). Despite their immediate practical value, and despite the fact that “much work has been done on abusive language detection in general”, much remains to be learned about “lexical knowledge for the detection of abusive language” (Wiegand et al., 2018), as

well as about the development and implementation of LDTLs for languages other than English.

In this paper we will report on the first phase of a project<sup>1</sup> to develop a Dutch LDTL (**TaboeLex**) consisting of potentially offensive constructions (words, word groups, expressions) as linked open data (i.e., a freely available, structured, annotated lexicon that could be linked to other data in future). In section 2, we will give a definition of what we mean by taboo language, and we will set the scope of TaboeLex. Section 3 then describes the design of the database. Section 4 concludes the paper, outlining future work.

## 2. Definition and scope of TaboeLex

### 2.1 Taboo language

Referring to the term *swearing*, Stapleton et al. (2022: 2) point out that “precise definitions and criteria are sometimes difficult to pin down [..., e.g.,] whether swear words can be used with literal (as opposed to figurative) meaning”. For purposes of this project, we define *taboo language* as linguistic constructions that are potentially offensive to some users in some contexts; constructions are form-meaning pairings on a morphological, lexical or syntactic level (see Goldberg (2006) for an extended view). We therefore use *taboo language* as a hypernym to include other phenomena and/or synonyms like *swearing*, *cursing/cussing*, *maledicta*, *profanity*, *blasphemy*, *obscenity*, *vulgarity*, *euphemisms and dysphemisms*, *verbal abuse*, *verbal sparring*, *(racial) slurs*, *terms of abuse*, *insults*, *offensive language*, *dirty language*, etc.

Our definitions and categories are all based on an extensive review of literature from various disciplines that aim to define taboo language, identify types of taboo language, sources of taboo language, etc. Most influential were Hirsch (1985), Hoeksema (2019), Jay (2018), Jay and Janschewitz (2008), Lewandowska-Tomaszczyk et al. (2021), Ljung (2011), Ruitenbeek et al. (2022) and Van Sterkenburg (2019), while the following books were also formative in our thinking about taboo language: Andersson and Trudgill (1990); Jay (1992, 2000); McEnery (2006); Montagu (1967); Pinker (2007). To inform us on the values of attributes, we also scrutinised the tags and definitions in GSW (2007) and Van Sterkenburg (2001), in order to create curated lists of possible values (see 3.2).

Some features of taboo words that are relevant to this project, include the following:

---

<sup>1</sup> Ethical clearance for the research project was obtained through the Language Matters Ethics Committee of the North-West University (ethics number: NWU-00632-19-A7).

- **Morphosyntactic type:** Taboo constructions include linguistic material on various morphosyntactic levels of independence and compositionality; these types are implemented in TaboeLex as an element `<headwordType>`. In addition to words, it also provides for sub-word items (like affixes), reduced forms (like initialisms), and multiword expressions (MWEs) (see 3.1 for values and examples).
- **Taboo domain:** Much work has been done to identify and delineate the source or reference domains of taboo language, such as religion, sex, scatology, animals, death, disease, etc. Within the scope of this paper, suffice to note that a taboo ontology will be declared as part of the `<denotatum>` element, which is a child element of the `<sense>` element (see 3.2).
- **Taboo type:** While the literal vs. figurative meaning requirement for taboo constructions are still being debated, we take the stance that both constructions with literal meanings, and constructions with figurative meanings could be taboo. For example, while neutral, scientific terms (i.e., orthophemisms) like *penis* and *vagina* could be considered by most people in most contexts as non-taboo, they could still be offensive to some people in some contexts, e.g., they might be dysphemistic in front of one’s grandparents at a Christmas dinner, or in a geography class for grade 5 learners.

This adds a layer of complexity to the development of LDTLs, since homonymous and polysemous constructions need to be handled appropriately. For example, *emmer* refers mostly to ‘bucket’ (container) – see for example the abridged Dutch dictionary, and the multilingual dictionaries in VDO (2021). However, in some rather obscure cases *emmer* could also refer to ‘an inferior person, specifically a prostitute’ (i.e., as an abusive term), or ‘female genitalia’ (i.e., as an obscenity), as reflected in the more comprehensive, unabridged Dikke Van Dale (DVD Online, 2022). This feature of taboo language is practically resolved by introducing the element `<tabooType>` that can be added to any sense of an entry (see 3.2).

- **Tabooness:** Tabooness ratings of constructions will differ between different social groups and are subject to change over time. It is therefore not only essential that constructions should be rated in terms of their observed tabooness in or for certain groups, but also that such ratings should be re-evaluated regularly. For example, it is the task of the British public regulator for communication services, Ofcom, to determine public attitudes towards offensive language on TV and radio, specifically when children are particularly likely to be listening (roughly speaking between 06:00 and 19:00) (Ipsos MORI, 2021a: 3). To this effect, they commission research reports roundabout every five years (Ipsos MORI, 2016, 2021b; Synovate UK, 2010; The Fuse Group, 2005) to determine which words are to be considered mild, moderate, or strong (Ipsos

MORI, 2021a: 4). Similar (but not necessarily comparable) investigations have been done for Dutch in 1998, 2001, 2007, and 2018 (Van Sterkenburg, 2001, 2008, 2019). The element `<tabooValue>` will capture this knowledge with attribute values on a scale ranging from `highlyTaboo` to `notTaboo`; see section 3.2 for other potential values.

- **Context dependence:** Whether a construction is taboo or not, is not only dependent on the situational and/or textual contexts (e.g., whether the derogatory meanings of *emmer* are activated or not), but also on the social context. The word *rambam* (‘undefined, imaginary illness’) appears only in taboo constructions, like *krijg de rambam* (‘get an illness’), but is not considered taboo in most social contexts. The prototypicality rating (`<tabooPrototypicality>`) will – to a large extent – account for situational, textual, and social contextual dependence of taboo constructions. Words that are taboo in all contexts (e.g., *oetlul* ‘jerk, wanker’) will get the value `alwaysTaboo`, while words that are rarely used in the taboo sense (like *emmer*), will have the value `rarelyTaboo` – see 3.2 for other potential values.
- **Intention and effects:** From a sociopragmatic point of view, taboo language is often defined as language with an expressive/emotive function (Jay, 2020: 39). Hirsch (1985) therefore made a strong case that a taxonomy of taboo language should be based first and foremost on the speech acts (Austin, 1962; Searle, 1969, 1979) in which expressions occur. Following this general approach, we therefore provide for three pragmatic-specific elements, viz. `<speechAct>` for the type of speech act, `<illocution>` for the speaker’s intention, and `<perlocution>` for the effect on the hearer (see 3.2).

## 2.2 Lexical databases of taboo language

We define LDTLs as digital, structured, enriched collections of linguistic constructions that are potentially offensive to some users in some contexts (e.g., in children’s books). When implemented in NLP systems as simple look-up lists (gazetteers) for filtering of results, they might sometimes also be called *blacklists*, *greylists*, *swearword stop lists*, or *profanity filters* (e.g., Shutterstock, 2020). Two prominent examples of LDTLs are the following:

- Hurtlex is a lexicon of 1,156 Italian “hate words” that were “linked to synset-based computational lexical resources such as MultiWordNet and BabelNet” (Bassignana et al., 2018).
- Taboo Wordnet is an online, synset-based Japanese resource that could “help detection systems regulate and curb the use of offensive words online” (Choo & Bond, 2021). It consists of 2,095 words with 912 synsets, and it is linked to the Open Multilingual Wordnet.

Besides proprietary lists that are not accessible in the open-data domain, there are also numerous data sets for various taboo-related domains available (see Nakov et al., 2021; Rosenthal et al., 2020; Wiegand et al., 2021; Wiegand et al., 2019; Wiegand et al., 2018; Zampieri et al., 2019b; Zampieri et al., 2020 for overviews of available material). The different tagging schemas of more than 60 such data sets have been compared by Lewandowska-Tomaszczyk et al. (2021), with the aim to create an ontology basis for offensive language identification, while also getting insight in how the concept *offensive* is understood across different projects. They use the term *offensive language* similar to how we use *taboo language* (see 2.1) as a superordinate term for all kinds of language phenomena (Lewandowska-Tomaszczyk et al., 2021: 7). Their proposed ontology of offensive language, together with their methodology for the detection of such language, hold the potential to play an important standardisation role with regards to the treatment of taboo language in the context of Linguistic Linked Open Data (LLOD). In the next phase of our project, their ontology will therefore be the first point of reference to which we will compare our own ontology.

Of utmost importance is that re-usability should be a compulsory design requirement of any LTDL. To make the data re-usable for multiple purposes in several different applications, the database should ideally be rich with as much information as possible – either in the database itself, or otherwise through links to other existing resources. By using subsets of data, or a selection of elements, attributes and/or values, the data could be used in a variety of practical NLP applications like some of the following:

- Offensive language identification (Zampieri et al., 2020) has been a prevailing topic in NLP for a number of years, especially with a view on hate speech, cyber-bullying and abuse detection on social media platforms (Akiwowo et al., 2020; Davidson et al., 2017; Fišer et al., 2018; Jarquín-Vásquez et al., 2020; Korotkova & Chung, 2023; Li et al., 2023; Mostafazadeh Davani et al., 2021; Nakov et al., 2021; Narang et al., 2022; Pradhan et al., 2020; Roberts et al., 2019; Rosenthal et al., 2020; Schmidt & Wiegand, 2017; Teh et al., 2018; Waseem et al., 2017; Zampieri et al., 2019a). The identification of taboo language is also an important aspect of sentiment analysis (Byrne & Corney, 2014; Cachola et al., 2018), especially since the speech acts and language associated with sentiment analysis can oftentimes be more subtle or indirect, e.g., by using humour (Ahuja, 2019; Ahuja et al., 2018; Bansal et al., 2020; Meaney et al., 2021), or irony and sarcasm (Frenda et al., 2022; Husain & Uzuner, 2021).
- More recently the evaluation of large language models for biased and toxic language (Osoba & Welser IV, 2017; Schäfer, 2023; Wiegand et al., 2019) have been pushed to the fore with the public availability of OpenAI’s GPT-4 and ChatGPT models. However, from a linguistic and user interface design perspective, our understanding of the implementation of these models in conversational artificial intelligent agents (e.g., speech assistants and chatbots), and especially the relation with taboo language, is still in its infancy.

- LDTLs have been used for many years in applications of text filtering; see Zhou (2019) for an elaborate evaluation of some of these, as well as his own improved implementation. These include, inter alia:
  - **predictive text filtering**, e.g., for search engines, keyboards on mobile phones, online text editors, etc.;
  - **suggestion filtering**, e.g., for spelling checkers and electronic dictionaries (especially dictionary apps for children) that should not suggest swearwords as corrections for ordinary typos;
  - **taboo language censoring**, i.e., redacting, modifying, replacing or removing a word in a text that matches a word in the LDTL; implemented typically as part of parental control software for text, audio, and video (see Porutiu (2023) for an overview and marketing reviews of a number of these applications);
  - **content filtering**, e.g., social media algorithms that (semi-)automatically delete posts or ban users, like Facebook’s profanity filter for Facebook Page, or spam filters used in email applications. Other examples of content filtering include e-lexicography tools for choosing good dictionary examples (Kilgarriff et al., 2008), or computer-assisted language learning systems that automatically selects suitable texts for learners (Belaid, 2016).

### 2.3 Dutch resources of taboo language

Dutch has a rather long tradition in taboo language research, going back to at least 1834 with an history-focused article by J.F. Willems titled *On some old Dutch curses, oaths and exclamations* [translated – the authors] (Willems, 1834). However, the first specialised printed dictionary focusing on language from a taboo domain only appeared in 1977 (EW, 1977). Since then, at least seventeen other printed dictionaries (or dictionary-like books) on various aspects of taboo language have been published (DBG, 1991/2021; GSW, 2007; GT, 1997; HEW, 1988; KDV, 1998; LNS, 1989; LOS, 1990; Lutz-van Elburg, 1990; Lutz-van Elburg & Jager, 1989; NSW, 1984; Van der Gucht et al., 2018; Van der Meulen et al., 2018; Van Lichtenvoorde & Van Lichtenvoorde, 1993; Van Sterkenburg, 2001; WAON, 2013; WEPCT, 2001; WPTG, 2020-2023). Of these, only three are available as digital data: GSW (2007); Van Sterkenburg (2001); WPTG (2020-2023). Since WPTG (2020-2023) is a general dictionary of slang, and therefore also contains many non-taboo constructions, we only use data from the other two dictionaries as candidate taboo constructions for TaboeLex.

One of the most prominent or most used look-up lists of Dutch taboo words (so to see), is the Dutch version of the *List-of-Dirty-Naughty-Obscene-and-Otherwise-Bad-Words*

(RolfBly, 2020).<sup>2</sup> This list was derived from the Dutch section of *The Alternative Dictionaries* (TAD, 2004), although it is not clear when this was done, and by whom it was done. RolfBly (2020) consists of 190 constructions: 165 one-word constructions, and 25 MWEs. While this list will be used in a next phase of the project as one of the baselines for evaluation, several potential problems with the list could already be identified:

- The list is not free of linguistic errors. These include:
  - four spelling errors (i.e., *\*johny* > *johnny*; *\*pijpbekkieg* > *pijpbekkie*; *\*tongzoeng* > *tongzoen*; *\*triootjeg* > *triootje*);
  - six errors related to obsolete orthographic forms due to spelling reforms in Dutch (i.e., *\*boerelul* > *boerenlul*; *\*bokkelul* > *bokkenlul*; *\*krentekakker* > *krentenkakker*; *\*kuttelikkertje* > *kutlikkertje*; *\*paardekop* > *paardenkop*; *\*paardelul* > *paardenlul*);
  - one compound that should be written as one word (i.e., *\*trottoirprostituée* > *trottoirprostituée*);
  - an ephemeral word that only exists in TAD (2004) and its derivatives (i.e., *hoempert*, apparently meaning ‘hard excrement’).
- The list contains only lemmas, e.g., *op+sodemieter · en* (up+tumble · INF ‘to fuck off’), and no other word forms, e.g., *op+ge · sodemieter · d* (PTCP). This is particularly problematic for purposes of look-up lists in applications using predictive text filtering, and suggestion filtering (see 2.2). In such applications, the input text cannot be lemmatised first, since filtering needs to happen in real-time and on the fly.
- The MWEs are only presented as lemmas, e.g., *op z’n sodemieter gev · en* (on his carcass give · INF ‘to beat the hell out of him’). There is therefore no indication of:
  - orthographic variants, e.g., related to the example above, *zijn/zn/zun* instead of *z’n*, the latter of which does not appear in the 5.9-billion-word nlTenTen20 corpus (Sketch Engine, 2020);
  - morphosyntactic variants, e.g., again related to the above example, *op zijn* (3SG.M) *sodemieter* accounts for only roughly half the cases in the nlTenTen20 corpus; *zijn* is followed by *hun* (3PL), *mijn* (1SG), *ons* (1PL), *de* (DET), and *her* (3SG.F);

---

<sup>2</sup> An older version (2014) of the list is available at <https://github.com/chucknorris-io/swear-words/blob/master/nl>, while the list is also reproduced elsewhere on the web.

- lexical variants, e.g., *krijg · en* (‘to get’) occurs more frequently than *gev · en* (‘to give’) on the righthand side of *sodemieter* in the nITenTen20 corpus (Sketch Engine, 2020); or
  - syntactic variants, e.g., *geeft hem op zijn sodemieter* instead of *hem op zijn sodemieter geeft*.
- In addition, the MWEs are not always presented uniformly. Compare for instance the lemma *op z'n sodemieter geven* that is presented as a prepositional phrase [<sub>PP</sub> *op<sub>PREP</sub> z'n<sub>PN</sub> sodemieter<sub>N</sub>*], followed by the verb [*geven<sub>V</sub>*]. However, the lemma *reet trappen, voor zijn* has the same [PP V] structure as the former example (i.e., [<sub>PP</sub> *voor<sub>PREP</sub> zijn<sub>PN</sub> reet<sub>N</sub>*] [*trappen<sub>V</sub>*]), but is presented here as [*reet<sub>N</sub> trappen<sub>V</sub>* , *voor<sub>PREP</sub> zijn<sub>PN</sub>*]. Also, in most cases in the list, only bare verbs are added as lemmas, e.g., *bedonderen* or *belazeren* (both meaning ‘to swindle, take someone for a ride’). However, in the case of [*besodemieteren<sub>V</sub>*] (also meaning ‘to swindle, take someone for a ride’) a copula verb phrase [*besodemieterd<sub>PTCP</sub> zijn<sub>COP</sub>*] (‘to have been swindled, taken for a ride’) is provided additionally as a separate lemma.
  - Numerous polysemous constructions that are most frequently used in a non-taboo way, are included in the list. Compare for instance *achter het raam zitten*, which is an ordinary phrase for ‘to sit in a window (looking at what’s happening outside)’. However, it is also rarely used with the meaning ‘to work as a prostitute’ (TAD, 2004), or ‘to present oneself in a prostitute-like manner’ (DVD Online, 2022). Also compare *welzijn · s+mafia* (welfare · LK+mafia ‘ineffective and meddling social workers corps’) in the list, which is always used unmarked in the Dutch mainstream media.
  - Many of the examples are general slang that is not taboo at all. Compare for instance *buffelen* (‘to hit; to work hard; to wolf down food’), *huisdealer* (‘drugs dealer associated with a certain establishment’), or *kanen* (‘to eat’; associated with slang in The Hague).
  - Many others are euphemisms, like *de hond uitlaten* (‘to let the dog out’), but which can also be used as a euphemism for ‘to urinate’. Another example is *de koffer induiken* (‘to jump in one’s bed’), which is mostly used euphemistically with the meaning ‘to have sex’.
  - Numerous expected candidates, i.e., highly frequent, highly taboo constructions, are not included in the list. These include words like *debiel* (‘mentally deficient’), *trut* (‘twat, cunt’), *kanker+wijf* (cancer+woman ‘stupid bitch’), and many racial slurs. The list also excludes many English taboo words that are used frequently in Dutch, like *bitch*, *fuck*, and *bullshit*.



A much better and unproblematic list is the *GRoninger OFFensive Lexicon* (GrofLex) (Van der Veen, 2020), a Dutch lexicon of abusive lemmas based on version 1.2 of the Dutch section of HurtLex (Basile, 2020) (see below for more details on Hurtlex). It consists of 847 one-word constructions only (no MWEs). The list has been annotated with part-of speech information, as well as the offensive category (what we call *denotatum* – see 3 below) of each lemma (e.g., ethnic slurs, physical disabilities and diversity, words related to religion, male genitalia, etc.). While the list still contains polysemous constructions (like *kuiken* ‘chicken’; *kalf* ‘calf’; *druif* ‘grape’), and orthophemisms (like *pretentieus* ‘pretentious’, *fascistisch* ‘fascist’, *snob* id.), it could be used fruitfully in a next phase of the project as another baseline for evaluation.

### 3. Design of the TaboeLex lexical database

Our goal is to design an LDTL for Dutch, of which the data can be integrated into various NLP applications and tools, but which can potentially also be useful for human users, or for linguistic research. The general principles and structure of TaboeLex is in line with most existing standards and encoding formats such as Ontolex-Lemon (Cimiano et al., 2016), DMLex (Měchura et al., 2023), LMF,<sup>3</sup> and TEI Lex-0 (Tasovac et al., 2018). General aspects are briefly discussed in section 3.1, followed by those aspects that relates specifically to a LDTL in section 3.2. Figure 1 presents an illustrative example, with LDTL-specific information marked in red. The complete XML schema and documentation, plus eventually all the TaboeLex data, will be made available under a CC BY-SA 4.0 license.

---

<sup>3</sup> <https://www.iso.org/standard/68516.html>

```

<lexicographicResource title="TaboeLex" language="ndl">
  <entry id="debiel-word-n">
    <headword>debiel</headword>
    <headwordType>word</headwordType>
    <partOfSpeech tag="noun" />
    <variantForm>dubiel</variantForm>
    <patternForm />
    <linkExternal gigantMolex="12324" />
    <sense>
      <denotatum>entity [person] [mental ability/health]</denotatum>
      <definition language="eng">mentally deficient person</definition>
      <example>
        <text>Mensen laat je toch niet zo opnaaien door die achterlijke
          debiel.</text>
        <source>nlTenTen20-23694165</source>
      </example>
      <tabooType value="dysphemism">epithet</tabooType>
      <tabooValue value="highlyTaboo"></tabooValue>
      <tabooPrototypicality value="alwaysTaboo"></tabooPrototypicality>
      <speechAct>
        <member value="insult">
          <member value="name-calling">
            <member value="abuse">
              </member>
            </member>
          </member>
        </member>
      </speechAct>
      <illocution>
        <member value="anger">
          <member value="disrespect">
            <member value="contempt">
              </member>
            </member>
          </member>
        </member>
      </illocution>
      <perlocution>
        <member value="offensive">
          <member value="derogatory">
            <member value="insulting">
              </member>
            </member>
          </member>
        </member>
      </perlocution>
      <relation type="synonym">
        <member idref="debiel-word-n" />
        <member idref="idiot-word-n" />
      </relation>
    </sense>
  </entry>

```

Figure 1: Sample entry for *debiel* ('retard; retarded')

### 3.1 General design

Following our definition of constructions as form-meaning pairings, each taboo construction in the database is defined by aspects related to form, and aspects related to meaning. Regarding form, we use common elements like `<headword>`, `<headwordType>`, `<partOfSpeech>` (of the headword), and `<variantForm>` (e.g., for variants like *f\*ck*, *f@ck*, *fark*, etc. for the English loanword *fuck*). The element `<headwordType>` could be extended in future to provide more detailed subcategories, but currently has the following primary values (with Dutch examples):

- subword: for affixes (e.g., *·erik* in *bang·erik* (scared·NMLZ ‘coward’)), and affixoids (e.g., *kanker÷* ‘cancer’ used as an intensifier in *kanker÷homo* ‘bad gay man’)<sup>4</sup>;
- reductionForm: for initialisms like *WTF*;
- word: for the uninflected form of words, e.g., *neuk·en* (fuck·INF ‘to fuck’); and
- MWE: for multiword expressions like:
  - word groups, e.g., *kwark blaffen* (‘to ejaculate (male)’), where neither *kwark* (‘curd’), nor *blaffen* (‘to bark’) is taboo, but their combination in a word group is;
  - construction idiom, e.g., *krijg X* (‘get X’), used as an imprecation, where X can be various illnesses; and
  - fixed expression, e.g., *Ik kan kakken en pissen en u gemakkelijk missen* (‘I can shit and piss without missing you at all’).

The rationale behind the element `<patternForm>` is to include some kind of pattern representation for each headword: on the one hand to allow for the automatic identification of the headword in corpus data (cf. Gantar & Krek, 2022; Odijk, to appear), and on the other hand to deal with the flexibility and variation that many MWEs exhibit. For single words (see Figure 1), the pattern representation is the same as `<headword>`. For verbal MWEs, the pattern representation is a finite sentence, similar to the way in which patterns are being described in the Corpus Pattern Analysis approach of Hanks (2013). However, rather than using semantic types in the argument slots, we use dummies such as *iemand* ‘someone’, and *iets* ‘something’. See also the recently compiled DUCAME<sup>5</sup> (*DUtch CAnonicalised Multiword Expressions*) resource, and the pattern descriptions in the project *Woordcombinaties*<sup>6</sup>.

The last aspect related to the form of an entry, involves the representation of all related word forms of a lemma, e.g., the verb *neuk·en* (‘to fuck’) has the grammatical forms *neuk* (1SG), *neuk·t* (2/3SG), *neuk·te* (SG.PST), *neuk·ten* (PL.PST), and *ge·neuk·t* (PTCP). Moreover, a comprehensive LDTL should ideally not only include grammatical forms, but also compounds (like *vuist+neuk·en* (fist+fuck·INF ‘to fist fuck’)), and derivations (like *neuk·er* ‘fucker’). This morphological information will be resolved in TaboeLex by means of links (`<linkExternal>`) to another lexical database, viz.

---

<sup>4</sup> We use the following notations: middle dot ( *·* ) for affix boundaries; divide symbol ( *÷* ) for affixoid boundaries; plus symbol ( *+* ) for compound boundaries.

<sup>5</sup> <https://surfdrive.surf.nl/files/index.php/s/2Maw8O0QTPH0oBP>

<sup>6</sup> <https://woordcombinaties.ivdnt.org>

GiGaNT-Molex<sup>7</sup>, the modern part of the computational lexicon of the Dutch language, compiled by the Dutch Language Institute. Because it is linked to GiGaNT-Molex, the full inflectional paradigms and word-formation families of the headwords need not be stored in TaboeLex itself. Instead, this information can be retrieved dynamically from GiGaNT-Molex, if required. This also pertains to MWEs, which are included in GiGaNT-Molex as a whole, and with individual components linked to the appropriate lemmas. This element could also be used in future to link TaboeLex data to other resources, such as thesauri, translation dictionaries, etc.

All information related to the meaning side of a construction are accommodated under the <sense> element. While most of its children elements are taboo-specific (see 3.2), three common elements are included, viz. <definition> (in English); <example>, including the <text> and reference to the <source>; and <relation> to represent lexical relations like synonyms and antonyms.

### 3.2 LDTL-specific design feature

Various elements, attributes, and/or values that are specific to LDTLs have been added to the design. These are all part of the <sense> element since their values can vary depending on which sense of the word is involved; see the information in red in Figure 1. The taboo-specific elements are the following:<sup>8</sup>

- <denotatum>: The denotata on a superordinate level are: event; relation; state; entity; locale; process. Subtypes provide for constructions related to specific domains; for example, the exonymic epithet *kaas+kop* (cheese+head ‘Dutch person’) will have the value entity [person] [inhabitant, citizen], while a euphemistic verb like *drukk · en* (press · INF ‘to defecate’) will be process [body] [substance] [excretion].
- <tabooType>: We distinguish four main taboo types on lexicopragmatic grounds, viz.:
  - orthophemism (e.g., *penis*);
  - euphemism (e.g., *klok-en-hamer-spel* clock-and-hammer-game ‘penis’);
  - dysphemism (e.g., *paal* pole ‘penis’); and

---

<sup>7</sup> <https://ivdnt.org/corpora-lexica/gigant/>

<sup>8</sup> Since it is impossible in terms of space restrictions to list all possible values for all elements or attributes here, these will be made available as part of the XML schema and documentation; suffice to present here some illustrative examples.

- witticism, i.e., constructions that were created originally with the purpose to be humorous (e.g., *sperma+spuiter* sperm+gusher ‘penis’).

Additionally, we also provide for constructions that can be both euphemistic or dysphemistic, like *aap* (monkey ‘penis’); these are eu-/dysphemism. Following Hoeksema (2019), we also have a category *rudeImperative*, for expressions like *sterf aan bloedpoep* (‘die of bloody diarrhoea’). Including subcategories (not discussed here) and a category *other* (for miscellaneous cases), `<tabooType>` has a total of 16 values.

- `<tabooValue>`: To indicate to what degree the construction is generally considered to be taboo, a Likert-like scale of values are available: *highlyTaboo*; *moderatelyTaboo*; *slightlyTaboo*; *notTaboo* (e.g., for orthophemisms). Since assignment of these values will be based on empirical research by Van Sterkenburg (2019), an additional value, *unspecified*, is required for constructions for which such empirical data is not available.
- `<tabooPrototypicality>`: The prototypicality of a taboo construction is expressed here as a value of its prominence in multiple sources as an exclusively taboo construction (more prototypical), or not (less prototypical). These values are also expressed on a Likert-like scale: *alwaysTaboo*; *oftenTaboo*; *sometimesTaboo*; *rarelyTaboo*. In addition to an *unspecified* value like above, a sixth value is required for constructions that are euphemistic.
- `<speechAct>`: We distinguish 32 values that can be used to complete the leading sentence: “This lemma is mostly/often used in/as an act of ...”. These values range from very specific (e.g., *blasphemy* or *self-malediction*), to general (e.g., *expressivenessNegative*), and include also values for “positive” speech acts (e.g., *expressionPhysicalSensationPositive*). A sense can be assigned multiple values.
- `<illocution>`: A total of 60 illocutionary intentions have been identified in the literature. They complete the leading sentence: “This lemma is often used to express ...”, with values like *dislike*, *disgust*, *admiration*, *aestheticAppreciation*, *horror*, etc. Again, a sense can be assigned multiple values.
- `<perlocution>`: To complete the leading sentence: “This lemma is often used to be / perceived as being ...”, we distinguish 16 values like *offensive*, *politicallyCorrect*, *racist*, *jocular*, etc. As with the previous two elements, multiple values can be assigned to a sense.

## 4. First evaluation, conclusions, and future work

To populate TaboeLex as far as possible automatically in the next phase of the project, we will compile a combined list of candidate taboo constructions based on existing Dutch taboo-language dictionaries, which will then be curated based on corpus material. In addition, we will use the labels that are used for taboo constructions in these dictionaries as a second set of seed terms in a bootstrapping fashion to extract increasingly more data from existing resources, specifically dictionaries developed and maintained by the Dutch Language Institute. Thus, we divide the labour between two separate sets of seed terms: a list with macrostructure seed terms, and one with microstructure seed terms.

The list of macrostructure seed terms (or candidate <headword> list) consists of two primary sources, and several secondary sources. The two primary sources are two well-known, published dictionaries that are available as open, unparsed, digital data,<sup>9</sup> viz. GSW (2007) and Van Sterkenburg (2001). We consider them primary, since they are assumed to be authoritative on whether a given construction is taboo or not. Our secondary sources are considered less authoritative, since they are lists that are generally available (and editable) from the internet. These include a list of lemmas tagged as *pejoratief* (pejorative) and *scheldwoord* (swearword) from Wiktionary (Wiktionary (NL), 2023), and a list of Dutch sexual slang and euphemisms from Wikipedia (Wikipedia (NL), 2023). RolfBly (2020) and Van der Veen (2020) will not be included in the candidate list, so that we can use them as part of our quality assessment.

An initial seed list shows that most taboo constructions only occur in one of the primary or secondary sources – see Table 1. The small overlap between the two printed dictionaries (GSW, 2007; Van Sterkenburg, 2001) can be ascribed to their different coverage of semantic domains: while GSW (2007) includes all kinds of taboo words from a variety of domains, Van Sterkenburg (2001) is more focused on taboo constructions related to oaths, curses, and (self-)maledictions. Similarly, the relatively small overlap between the other lists may also be due to a difference in focus, scope or aim of the respective lists, with the greatest overlap (37,4%) between GSW (2007) and (Wiktionary (NL), 2023).

---

<sup>9</sup> We distinguished between *parsed digital data* (e.g., a lexical/lexicographic database); *unparsed digital data* (e.g., a word document with systematic formatting); and *digital documents/files* (e.g., scanned PDF documents). All these types can be *open* (i.e., available for research and development), or *proprietary* (i.e., not available).

	n	Van Sterkenburg (2001)	Wiktionary (NL) (2023)	Wikipedia (NL) (2023)
<b>GSW (2007)</b>	2,619	82 (4,2%)	382 (37,4%)	68 (6,7%)
<b>Van Sterkenburg (2001)</b>	1,973		31 (3,0%)	74 (7,3%)
<b>Wiktionary (NL) (2023)</b>	1,022			44 (4,3%)
<b>Wikipedia (NL) (2023)</b>	1,015			
<b>Total</b>	6,629			
<b>Unique</b>	5,295			

Table 1: Lemma lists, with number and percentage of shared lemmas between lists

The resulting candidate list will be further populated by extending it with headwords from other lexical resources that are labelled with one of the microstructure seed terms, i.e., constructions that occur either as tags in existing dictionaries (not only taboo dictionaries), or in the definitions of such dictionaries. In English, these would include stylistic tags like *vulgar* or *obscene*, and orthophemisms like *male genitalia* or *faeces*. Our initial list of microstructure seed terms is based on the tags and definitions used in GSW (2007) and Van Sterkenburg (2001). Initial results show that some tags do indeed result in new candidates, but that manual inspection of the results is needed. For example, a label such as *straattaal* (lit. street language 'bad language') produced one result in the ANW<sup>10</sup>, i.e., *straatbijbel* (lit. street Bible, a version of the Bible meant for young people), which is indeed written in a type of informal, street language, but which is clearly not a potential taboo construction.

To check the validity of the taboo constructions in our candidate list, we will check the constructions on the list against corpus data (and this information will be included in the database). A small pilot study shows that simply checking for occurrences in a corpus is not enough. The frequency counts require manual inspection of the data as some candidate constructions do occur in the corpus, but not as taboo constructions. For instance, all occurrences of *God in de hoge hemel* ('God in the highest heaven') and *God vergeef me* ('God forgive me') in the nlTenTen20 corpus can be considered as non-offensive. Furthermore, (normalised) frequencies can differ substantially between different types of corpora. As taboo constructions may be more likely to occur in certain types of texts than others, this is not unexpected but needs to be considered when interpreting frequency data. Moreover, the fact that a taboo construction does not occur in the corpus data does not automatically imply that it should be removed from the list.

<sup>10</sup> <https://anw.ivdnt.org/search>

Once TaboeLex is populated with the curated list of taboo constructions, the lexicographic editing process will start. The very first step will be to validate the ontology of our annotation schema against (a) other similar ontologies, notably the one of (Lewandowska-Tomaszczyk et al., 2021); and (b) real-world data. Editing will therefore be done in a modular way, gradually refining not only the annotation schema, but also the amount of information for each taboo construction in the database.

## 5. References

- Ahuja, V. (2019). *Computational Analysis of Humour*. Master of Science. Hyderabad: International Institute of Information Technology.
- Ahuja, V., Mamidi, R., & Singh, N. (2018). From Humour to Hatred: A Computational Analysis of Off-Colour Humour. In M. Zhang, V. Ng, D. Zhao, S. Li, & H. Zan (eds.) *Natural Language Processing and Chinese Computing*. Cham: Springer International Publishing, pp. 144–153. [https://doi.org/10.1007/978-3-319-99501-4\\_12](https://doi.org/10.1007/978-3-319-99501-4_12).
- Akiwowo, S., Vidgen, B., Prabhakaran, V., & Waseem, Z. (eds.) (2020). *Proceedings of the Fourth Workshop on Online Abuse and Harms*. Association for Computational Linguistics. <https://aclanthology.org/2020.alw-1.0>.
- Andersson, L. G., & Trudgill, P. (1990). *Bad Language*. London: Penguin.
- Austin, J. L. (1962). *How to do things with words*. Oxford.
- Bansal, S., Garimella, V., Suhane, A., Patro, J., & Mukherjee, A. (2020). Code-switching patterns can be an effective route to improve performance of downstream NLP applications: A case study of humour, sarcasm and hate speech detection. *arXiv pre-print server*(arXiv:2005.02295v1 [cs.CL]). <https://doi.org/10.48550/arXiv.2005.02295>. (24 May 2023)
- Basile, V. (2020). *Hurtlex NL lexicon version 1.2* [GitHub repository]. <https://doi.org/https://github.com/valeribasile/hurtlex/tree/master/lexica/NL/1.2>. (24 May 2023)
- Bassignana, E., Basile, V., & Patti, V. (2018). Hurtlex: A Multilingual Lexicon of Words to Hurt. In *CLiC-it*. <https://doi.org/10.4000/BOOKS.AACCADEMIA.3085>.
- Belaid, A. M. (2016). *The Localisation of the PARSNIP Model and Authentic Materials*. Chisinau: Scholars' Press.
- Byrne, E., & Corney, D. (2014). Sweet FA: Sentiment, Swearing and Soccer. SoMuS ICMR 2014 Workshop, Glasgow, Scotland, 01 April.
- Cachola, I., Holgate, E., Preoțiu-Pietro, D., & Li, J. J. (2018). Expressively vulgar: The socio-dynamics of vulgarity and its effects on sentiment analysis in social media. Proceedings of the 27th International Conference on Computational Linguistics, Santa Fe, New Mexico, August 20-26.
- Choo, Y. H. M., & Bond, F. (2021). Taboo Wordnet. 11th GlobalWordnet Conference, Potchefstroom, South Africa, 18-21 January.
- Cimiano, P., McCrae, J. P., & Buitelaar, P. (2016). *Lexicon model for ontologies:*



*Community report.*

- Davidson, T., Warmesley, D., Macy, M., & Weber, I. (2017). Automated Hate Speech Detection and the Problem of Offensive Language. International AAAI Conference on Web and Social Media (ICWSM), Montréal, Québec, Canada, May 15–18.
- DBG: *Duizend bommen en granaten: Scheldwoordenboek van Haddock*. (1991/2021). Edited by A. Algoud. Brussel: Casterman.
- DVD Online: *Dikke Van Dale Online (Van Dale Groot woordenboek van de Nederlandse taal)*. (2022). Edited by. Utrecht: Van Dale Uitgevers.
- EW: *Erotisch woordenboek*. (1977). Edited by H. Heestermans, P. Van Sterkenburg, & J. v. d. V. Van der Kleij. Baarn: Erven Thomas Rap.
- Fišer, D., Huang, R., Prabhakaran, V., Voigt, R., Waseem, Z., & Wernimont, J. (eds.) (2018). *Proceedings of the Second Workshop on Abusive Language Online (ALW2)*. Brussels, Belgium: Association for Computational Linguistics. <https://aclanthology.org/W18-5100>.
- Frenda, S., Cignarella, A. T., Basile, V., Bosco, C., Patti, V., & Rosso, P. (2022). The unbearable hurtfulness of sarcasm. *Expert Systems with Applications*, 193, 116398. <https://doi.org/https://doi.org/10.1016/j.eswa.2021.116398>. (24 May 2023)
- Gantar, P., & Krek, S. (2022). Creating the lexicon of multi-word expressions for Slovene. Methodology and structure. XX EURALEX International Congress, Mannheim.
- Goldberg, A. E. (2006). *Constructions at work: The nature of generalization in language*. Oxford: Oxford University Press.
- GSW: *Groot scheldwoordenboek: van apenkant tot zweefteef*. (2007). Edited by M. De Coster. Antwerpen: Standaard.
- GT: *Gespierde taal. Verbaal geweld voor in het buitenland: Beknopt scheldwoordenboek Nederlands-Engels, -Duits, -Frans en -Spaans*. (1997). Edited by J. Frijters. Zutphen: Alpha.
- Hanks, P. (2013). *Lexical analysis: Norms and exploitations*. Cambridge: The MIT Press.
- HEW: *Homo-erotisch woordenboek*. (1988). Edited by A. Joustra. Amsterdam: Thomas Rap. [https://dbnl.org/tekst/jous008homo01\\_01/](https://dbnl.org/tekst/jous008homo01_01/). (24 May 2023)
- Hirsch, R. (1985). Taxonomies of swearing. In L.-G. Andersson & R. Hirsch (eds.) *Perspectives on Swearing*. Göteborg: Department of Linguistics, University of Göteborg, pp. 37–59.
- Hoeksema, J. (2019). Taboo terms and their grammar. In K. Allan (ed.) *The Oxford Handbook of Taboo Words and Language*. Oxford: Oxford University Press, pp. 160–179.
- Husain, F., & Uzuner, O. (2021). Leveraging Offensive Language for Sarcasm and Sentiment Detection in Arabic. 6th Arabic Natural Language Processing Workshop, Kyiv, Ukraine, 19 April.
- Ipsos MORI. (2016). *Attitudes to potentially offensive language and gestures on TV and radio*. <http://stakeholders.ofcom.org.uk/binaries/research/tv->

- research/Offensive-language/Offensive-Language-2016-report.pdf. (24 May 2023)
- Ipsos MORI. (2021a). *Public attitudes towards offensive language on TV and Radio: Quick Reference Guide*.  
[https://www.ofcom.org.uk/\\_\\_data/assets/pdf\\_file/0020/225335/offensive-language-quick-reference-guide.pdf](https://www.ofcom.org.uk/__data/assets/pdf_file/0020/225335/offensive-language-quick-reference-guide.pdf). (24 May 2023)
- Ipsos MORI. (2021b). *Public attitudes towards offensive language on TV and Radio: Summary Report*.  
[https://www.ofcom.org.uk/\\_\\_data/assets/pdf\\_file/0021/225336/offensive-language-summary-report.pdf](https://www.ofcom.org.uk/__data/assets/pdf_file/0021/225336/offensive-language-summary-report.pdf). (24 May 2023)
- Jarquín-Vásquez, H. J., Montes-y-Gómez, M., & Villaseñor-Pineda, L. (2020). Not All Swear Words Are Used Equal: Attention over Word n-grams for Abusive Language Identification. In K. M. Figueroa Mora, J. Anzures Marín, J. Cerda, J. A. Carrasco-Ochoa, J. F. Martínez-Trinidad, & J. A. Olvera-López, *Pattern Recognition*. Cham. [https://doi.org/10.1007/978-3-030-49076-8\\_27](https://doi.org/10.1007/978-3-030-49076-8_27).
- Jay, T. B. (1992). *Cursing in America: A psycholinguistic study of dirty language in the courts, in the movies, in the schoolyards and on the streets*. Amsterdam: John Benjamins.
- Jay, T. B. (2000). *Why we curse: A neuro-psycho-social theory of speech*. Amsterdam: John Benjamins.
- Jay, T. B. (2018). The psychology of expressing and interpreting linguistic taboos. In *The Oxford Handbook of Taboo Words and Language*. pp. 76–95. <https://doi.org/10.1093/oxfordhb/9780198808190.013.5>.
- Jay, T. B. (2020). Ten issues facing taboo word scholars. In N. Nassenstein & A. Storch (eds.) *Swearing and Cursing*. Berlin: De Gruyter Mouton, pp. 37–52. <https://doi.org/10.1515/9781501511202-002>.
- Jay, T. B., & Janschewitz, K. (2008). The pragmatics of swearing. *Journal of Politeness Research. Language, Behaviour, Culture*, 4(2). <https://doi.org/10.1515/jplr.2008.013>.
- KDV: *Krijg de vinkentering! 1001 Nederlandse en Vlaamse verwensingen*. (1998). Edited by E. Sanders & R. Tempelaars. Amsterdam: Uitgeverij Contact.
- Kilgarriff, A., Husák, M., McAdam, K., Rundell, M., & Rychlý, P. (2008). GDEX: Automatically finding good dictionary examples in a corpus. 13th EURALEX International Congress, Spain, July.
- Korotkova, E., & Chung, I. K. Y. (2023). Beyond Toxic: Toxicity Detection Datasets are Not Enough for Brand Safety. *arXiv preprint*(arXiv:2303.15110v1 [cs.CL]). <https://doi.org/10.48550/arXiv.2303.15110>. (24 May 2023)
- Lewandowska-Tomaszczyk, B., Žitnik, S., Bączkowska, A., Liebeskind, C., Mitrović, J., & Oleškevičiene, G. V. (2021). LOD-connected offensive language ontology and tagset enrichment. 42nd Conference on Very Important Topics (CVIT 2016), RWTH Aachen.
- Li, Z., Cabello, L., Yong, C., & Hershovich, D. (2023). Cross-Cultural Transfer Learning for Chinese Offensive Language Detection. *arXiv pre-print server*. <https://doi.org/arXiv:2303.17927v1> [cs.CL]. (24 May 2023)

- Ljung, M. (2011). *Swearing: A Cross-Cultural Linguistic Study*. Palgrave Macmillan.
- LNS: *Luilebol! Het Nederlands scheldwoordenboek*. (1989). Edited by H. Heestermans. Amsterdam: Thomas Rap.
- LOS: *Lik op stuk: Nieuw Nederlands woordenboek van agressief taalgebruik*. (1990). Edited by D. De Bleecker, P. Thomas, & J. Van Haver. Tiel: Lannoo.
- Lutz-van Elburg, I. (1990). *More Dutch you won't learn in class: not for hypocrites*. Rotterdam: Wilkerdon.
- Lutz-van Elburg, I., & Jager, P. C. W. (1989). *Dutch you won't learn in class (not for hypocrites)*. Lelystad: Zander Media Service.
- McEnery, T. (2006). *Swearing in English: Bad Language, Purity and Power from 1586 to the Present*. London: Routledge.
- Meaney, J. A., Wilson, S., Chiruzzo, L., Lopez, A., & Magdy, W. (2021). SemEval-2021 Task 7: HaHackathon, Detecting and Rating Humor and Offense. 15th International Workshop on Semantic Evaluation (SemEval-2021), Online, August.
- Měchura, M., Filip, D., & Krek, S. (2023). *Data Model for Lexicography Version 1.0* <https://docs.oasis-open.org/lexidma/dmlex/v1.0/wd01/dmlex-v1.0-wd01.html>. (24 May 2023)
- Montagu, A. (1967). *The anatomy of swearing*. New York: The Macmillan Company.
- Mostafazadeh Davani, A., Kiela, D., Lambert, M., Vidgen, B., Prabhakaran, V., & Waseem, Z. (eds.) (2021). *Proceedings of the Fifth Workshop on Online Abuse and Harms (WOAH 2021)*. Association for Computational Linguistics. <https://aclanthology.org/2021.woah-1.0>.
- Nakov, P., Nayak, V., Dent, K., Bhatawdekar, A., Sheikh, Hardalov, M., Dinkov, Y., Zlatkova, D., Bouchard, G., & Augenstein, I. (2021). Detecting Abusive Language on Online Platforms: A Critical Analysis. *arXiv pre-print server*. <https://doi.org/10.26434/chemrxiv-2021-00153>. (24 May 2023)
- Narang, K., Mostafazadeh Davani, A., Mathias, L., Vidgen, B., & Talat, Z. (eds.) (2022). *Proceedings of the Sixth Workshop on Online Abuse and Harms (WOAH)*. Seattle, Washington (Hybrid): Association for Computational Linguistics. <https://aclanthology.org/2022.woah-1.0>.
- NSW: *Nationaal scheldwoordenboek: Schelden van de Schelde tot Terschelling*. (1984). Edited by K. Laps. Amsterdam: Ploegsma.
- Odiijk, J. (to appear). *MWE-Finder: Querying for multiword expressions in large Dutch text corpora*. Berlin: Language Science Press.
- Osoba, O. A., & Welser IV, W. (2017). *An intelligence in our image: The risks of bias and errors in artificial intelligence*. Santa Monica: Rand Corporation. [https://www.rand.org/pubs/research\\_reports/RR1744.html](https://www.rand.org/pubs/research_reports/RR1744.html). (24 May 2023)
- Pinker, S. (2007). *The Stuff of Thought : Language as a Window Into Human Nature*. New York: Viking.
- Porutiu, T. (2023, 06 January). Profanity Filters: Everything You Need to Know + Our Top 5 Picks. <https://VPNOverview.com>. (24 May 2023)
- Pradhan, R., Chaturvedi, A., Tripathi, A., & Sharma, D. K. (2020). A Review on Offensive Language Detection. In M. L. Kolhe, S. Tiwari, M. C. Trivedi, & K. K.

- Mishra, *Advances in Data and Information Sciences* Singapore.
- Roberts, S. T., Tetreault, J., Prabhakaran, V., & Waseem, Z. (eds.) (2019). *Proceedings of the Third Workshop on Abusive Language Online*. Florence, Italy: Association for Computational Linguistics. <https://aclanthology.org/W19-3500>.
- RolfBly. (2020). *Dutch LDNOOBW (List-of-Dirty-Naughty-Obscene-and-Otherwise-Bad-Words)* [GitHub repository]. <https://doi.org/https://github.com/LDNOOBW/List-of-Dirty-Naughty-Obscene-and-Otherwise-Bad-Words/blob/5faf2ba42d7b1c0977169ec3611df25a3c08eb13/nl>. (24 May 2023)
- Rosenthal, S., Atanasova, P., Karadzhov, G., Zampieri, M., & Nakov, P. (2020). SOLID: A Large-Scale Semi-Supervised Dataset for Offensive Language Identification. *14th International Workshop on Semantic Evaluation (SemEval-2020)*. <https://doi.org/https://doi.org/10.48550/arXiv.2004.14454>.
- Ruitenbeek, W., Zwart, V., Van Der Noord, R., Gnezdilov, Z., & Caselli, T. (2022). “Zo Grof !”: A Comprehensive Corpus for Offensive and Abusive Language in Dutch. *Proceedings of the Sixth Workshop on Online Abuse and Harms (WOAH)* Seattle, Washington (Hybrid), July.
- Schäfer, J. (2023). Bias Mitigation for Capturing Potentially Illegal Hate Speech. *Datenbank-Spektrum*. <https://doi.org/10.1007/s13222-023-00439-0>.
- Schmidt, A., & Wiegand, M. (2017). A Survey on Hate Speech Detection using Natural Language Processing. Fifth International Workshop on Natural Language Processing for Social Media, Valencia, Spain, 3 April.
- Searle, J. R. (1969). *Speech acts*. Cambridge: Cambridge University Press.
- Searle, J. R. (1979). *Expression and meaning: Studies in the theory of speech acts*. Cambridge: Cambridge University Press.
- Shutterstock. (2020). *List of Dirty, Naughty, Obscene, and Otherwise Bad Words* [GitHub repository]. <https://doi.org/https://github.com/LDNOOBW/List-of-Dirty-Naughty-Obscene-and-Otherwise-Bad-Words/blob/master/README.md>. (24 May 2023)
- Sketch Engine. (2020). *Dutch Web corpus 2020 (nlTenTen20)*. <https://www.sketchengine.eu/>. (24 May 2023)
- Stapleton, K., Beers Fägersten, K., Stephens, R., & Loveday, C. (2022). The power of swearing: What we know and what we don't. *Lingua*, 277. <https://doi.org/10.1016/j.lingua.2022.103406>.
- Synovate UK. (2010). *Audience attitudes towards offensive language on television and radio*. [https://www.ofcom.org.uk/\\_\\_data/assets/pdf\\_file/0017/27260/offensive-lang.pdf](https://www.ofcom.org.uk/__data/assets/pdf_file/0017/27260/offensive-lang.pdf). (24 May 2023)
- TAD: *The Alternative Dictionaries*. (2004). Edited by H.-C. Holm. <http://www.notam02.no/~hcholm/altlang/>. (24 May 2023)
- Tasovac, T., Romary, L., Banski, P., Bowers, J., De Does, J., Depuydt, K., Erjavec, T., Geyken, A., Herold, A., Hildenbrandt, V., Khemakhem, M., Lehečka, B., Petrović, S., Salgado, A., & Witt, A. (2018). *TEI Lex-0: A baseline encoding for lexicographic data. Version 0.9.1* [GitHub repository].

- <https://doi.org/https://dariah-eric.github.io/lexicalresources/pages/TEILex0/TEILex0.html>. (24 May 2023)
- Teh, P. L., Cheng, C.-B., & Chee, W. M. (2018). *Identifying and Categorising Profane Words in Hate Speech*. Proceedings of the 2nd International Conference on Compute and Data Analysis - ICCDA 2018.
- The Fuse Group. (2005). *Language and Sexual Imagery in Broadcasting: A Contextual Investigation*.  
[https://www.ofcom.org.uk/\\_\\_\\_data/assets/pdf\\_file/0012/24015/language.pdf](https://www.ofcom.org.uk/___data/assets/pdf_file/0012/24015/language.pdf). (24 May 2023)
- Van der Gucht, F., Van der Meulen, M., Verlinde, R., & Vanbeylen, W. (2018). *Het groot Vlaams vloekboek: Slimmer schelden en vaardiger vloeken*. Tielt: Lannoo.
- Van der Meulen, M., Van der Gucht, F., Verlinde, R., & Vanbeylen, W. (2018). *Het groot Nederlands vloekboek: Slimmer schelden en vaardiger vloeken*. Tielt: Lannoo.
- Van der Veen, H. (2020). *GRoninger OFFensive Lexicon (GrofLex)* [GitHub repository].  
<https://doi.org/https://github.com/hylkevdeveen/GrofLex>. (24 May 2023)
- Van Lichtenvoorde, M., & Van Lichtenvoorde, M. (1993). *Scheldwoorden van de jaren negentig*. Helmond: Michon.
- Van Sterkenburg, P. G. J. (2001). *Vloeken. Een cultuurbepaalde reactie op woede, irritatie en frustratie* 2e ed. Den Haag: Sdu Uitgevers.
- Van Sterkenburg, P. G. J. (2008). *Krachttermen*. Schiedam: Scriptum.
- Van Sterkenburg, P. G. J. (2019). *Rot lekker zelf op: Over politiek incorrect en ander ongepast taalgebruik*. Schiedam: Scriptum.
- VDO: *Van Dale Online*. (2021). Utrecht: Van Dale Uitgevers. <https://www.vandale.nl/>. (24 May 2023)
- WAON: *Woordenboek van het Algemeen Onbeschaafd Nederlands*. (2013). Edited by H. Aalbrecht & P. Wagenaar. Houten & Antwerpen: Uitgeverij Unieboek | Het Spectrum bv.
- Waseem, Z., Chung, W. H. K., Hovy, D., & Tetreault, J. (eds.) (2017). *Proceedings of the First Workshop on Abusive Language Online*. Vancouver, BC, Canada: Association for Computational Linguistics. <https://doi.org/10.18653/v1/W17-30>.
- WEPT: *Woordenboek van eufemismen en politiek correct taalgebruik*. (2001). Edited by M. De Coster. Amsterdam: Veen/Het Taalfonds.
- Wiegand, M., Ruppenhofer, J., & Eder, E. (2021). Implicitly Abusive Language – What does it actually look like and why are we not getting there? *Conference of the North American Chapter of the Association for Computational Linguistics (NAACL): Human Language Technologies*, 576–587.  
<https://doi.org/10.18653/V1/2021.NAACL-MAIN.48>.
- Wiegand, M., Ruppenhofer, J., & Kleinbauer, T. (2019). Detection of Abusive Language: the Problem of Biased Datasets. *2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 602–608.  
<https://doi.org/10.18653/v1/N19-1060>.
- Wiegand, M., Ruppenhofer, J., Schmidt, A., & Greenberg, C. (2018). Inducing a

- Lexicon of Abusive Words – A Feature-Based Approach. Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technology, New Orleans, Louisiana.
- Wikipedia (NL). (2023). *Seksuele volkstaal en eufemismen*. Retrieved 10 April 2023 from <https://nl.wikipedia.org>.
- Wiktionary (NL). (2023). *Lemmas tagged <pejoratief> and <scheldwoord>*. Retrieved 10 April 2023 from <https://nl.wiktionary.org>.
- Willems, J. F. (1834). Over eenige oude Nederlandsche vloeken, eeden en uitroepingen. *Nederduitsche Letteroefeningen*, 218–230.
- WPTG: *Woordenboek van Populair Taalgebruik*. (2020-2023). Edited by M. De Coster. <https://www.ensie.nl/woordenboek-van-populair-taalgebruik#>. (24 May 2023)
- Zampieri, M., Malmasi, S., Nakov, P., Rosenthal, S., Farra, N., & Kumar, R. (2019a). Predicting the Type and Target of Offensive Posts in Social Media. Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT 2019), Minneapolis, Minnesota, 2-7 June.
- Zampieri, M., Malmasi, S., Nakov, P., Rosenthal, S., Farra, N., & Kumar, R. (2019b). SemEval-2019 Task 6: Identifying and Categorizing Offensive Language in Social Media (OffensEval). 13th International Workshop on Semantic Evaluation (SemEval-2019), Minneapolis, USA, 6-7 June.
- Zampieri, M., Nakov, P., Rosenthal, S., Atanasova, P., Karadzhov, G., Mubarak, H., Derczynski, L., Pitenis, Z., & Çöltekin, Ç. (2020). SemEval-2020 Task 12: Multilingual Offensive Language Identification in Social Media (OffensEval 2020). *14th International Workshop on Semantic Evaluation (SemEval-2020)*, 1425–1447. <https://doi.org/10.18653/v1/2020.semeval-1.188>.
- Zhou, V. (2019, 4 February). Building a Better Profanity Detection Library with scikit-learn. <https://VictorZhou.com>. (24 May 2023)