

The Syntax-Semantics Interface of Czech Verbs in the Valency Lexicon¹

Václava Kettnerová, Markéta Lopatková & Eduard Bejček

Keywords: *valency, lexicon, alternations.*

Abstract

In this paper, an alternation based model of the valency lexicon of Czech verbs, *VALLEX*, is described. Two types of alternations (changes in valency frames of verbs) are distinguished on the basis of used linguistic means: (i) grammaticalized alternations and (ii) lexicalized alternations. Both grammaticalized and lexicalized alternations are either conversive, or non-conversive. While grammaticalized alternations relate different surface syntactic structures of a single lexical unit of a verb, lexicalized alternations relate separate lexical units. For the purpose of the representation of alternations, we divide the lexicon into data and rule components. In the data part, each lexical unit is characterized by a single valency frame and by applicable alternations. In the rule part, two types of rules are contained: (i) syntactic rules describing grammaticalized alternations and (ii) general rules determining changes in the linking of situational participants with valency complementations typical of lexicalized alternations.

1. Introduction

Information on valency characteristics of verbs, which are traditionally considered to be the center of a sentence, plays a key role in many rule-based NLP tasks such as machine translation, information retrieval, text summarization, question answering, etc. However, the valency behavior of verbs is so various that it cannot be described by general rules; instead, it must be captured for each lexical unit of a verb separately in the form of a lexical entry listed in the valency lexicon.

Prototypically, a single meaning of a verb corresponds to a single valency structure. However, in many cases, semantically similar uses of a verb can be syntactically structured in a surface sentence in a different way, see the pairs of examples (1) (a)-(1) (b) and (2) (a)-(2) (b).

- (1) (a) *Jana vymetla pavučiny z půdy.*
(b) *Jana vymetla půdu.*
Eng. (a) Jane swept cobwebs from the attic.
(b) Jane swept the attic.
- (2) (a) *Pavučiny byly z půdy vymeteny (Janou/od Jany).*
(b) *Půda byla vymetena (Janou/od Jany).*
Eng. (a) Cobwebs were swept from the attic (by Jane).
(b) The attic was swept (by Jane).

Then the question arises how it is possible to describe such changes in valency structure of verbs (usually referred to as alternations) in the lexicon. In last decades, theoretical linguistics has been paying considerable attention to alternations, see esp. Levin (1993). However, the results of this theoretical research have not been applied in the existed lexical resources yet except for, for example, *DeepDict*, Bick (2009), or *LexIt*, Lenci (2009), where this information was automatically extracted from corpora.

In this paper, we report on work in progress which is focused on a theoretically adequate and at the same time economical lexicographic description of Czech alternations. Moreover,

we attempt to propose such representation of alternations which can be applied in automated language processing. Our proposal is primarily formulated for the purpose of the description of valency in the valency lexicon of Czech verbs, *VALLEX 2.5*, see Lopatková et al. (2008); however, this phenomenon is to be solved in any valency lexicon. *VALLEX* provides information on the valency structure of verbs in their particular senses: on the number of valency complementations, on their type labeled by functors, and on their morphemic forms, see Žabokrtský and Lopatková (2007). This lexicon describes 2730 Czech verb lexemes containing about 6460 lexical units. As its theoretical background, the Functional Generative Description (FGD) is adopted. In FGD, valency – the range of syntactic elements either required or specifically permitted by a lexical unit – is related to a layer of linguistically structured meaning (so called tectogrammatical layer in Sgall et al. (1986), Panevová (1994)). *VALLEX* is available both in machine-tractable XML format and as human-readable structured web pages.

Verbs in *VALLEX* can be viewed and sorted according to various criteria (a number of lexical units of a verb, morphological forms of valency complementations, functors participating in the situation etc.). Moreover, more elaborate searches can be done using PML-TQ search engine (see <http://ufal.mff.cuni.cz/~pajas/pmltq/>). Queries are created in a graphical form and such a single query may aggregate several conditions imposed on the verb, see Bejček et al. (2010).

2. The structure of the lexicon

For the purpose of the representation of the alternations, we divide the valency lexicon into data and rule components. In the data component, each lexical unit of verb is represented by a single valency frame. Valency frames in the data component correspond to unmarked use (i.e., active use) of lexical units. Further, each lexical unit is ascribed by applicable alternations. The rule part of the lexicon contains rules determining changes in valency structure of verbs.

2.1. *The data component*

The *data component* consists of word entries corresponding to verb lexemes. Lexeme is an abstract twofold data structure which associates lexical form(s) and lexical unit(s). Lexical forms are all possible manifestations of a lexeme in an utterance (e.g. perfective, imperfective and iterative verb lemmas, all their morphological verb forms, reflexive and irreflexive forms). All lexical forms of a lexeme are represented by its lemma(s).

Concerning lexical units, two parts of the verbal meaning are crucial for their delimiting. (i) *A situational meaning* reflects a situation portrayed by a verb; it is characterized by a set of *situational participants* related by particular relations. Such part of the verbal meaning is not syntactically structured, see esp. Mel'čuk (2004). (ii) The part of the verbal meaning in which the situational participants are syntactically structured is referred here to as a *structural meaning*; its components correspond to *valency complementations* Panevová (1994). Each lexical unit of a verb is characterized by both situational and structural meaning in a unique way: any change in the situational or structural meaning leads to the change of lexical unit.

In the lexicon, each lexical unit is characterized by a gloss (i.e., a verb or a paraphrase roughly synonymous with the given sense) and by example(s) (i.e., sentence fragment(s) containing the given verb used in the given sense). The core information on valency characteristics of a verb is encoded in a form of valency frames. Each lexical unit is described

by exactly one valency frame reflecting unmarked (active) use of the verb. Valency frame is modeled as a sequence of valency slots, each slot standing for a single valency complementation. The slots consist of a functor (coarse-grained semantic role), a list of morphemic form(s) and information on obligatoriness. Each relevant lexical unit is ascribed by optional attributes providing information on idiomaticity, control, and semantic class membership. Moreover, optional attributes listing alternations applicable for the particular lexical units are proposed, see Section 3.1.3 and 3.2.3.



Figure 1. VALLEX, lexeme *navracet^{impf}, navrátit^{pf}, navracívat^{iter}* 'to return/to restore'.

2.2. The rule component

The *rule component* of the lexicon consists of two sets of rules: (i) a set of formal syntactic rules determining changes in the mapping of valency complementations onto surface syntactic positions and (ii) general rules specifying changes in the linking of situational participants and valency complementations. Whereas the first type of rules (formally describing grammaticalized alternations, see below) makes it possible to obtain all possible surface syntactic manifestations of lexical units of verbs (i.e., number of complementations, their types and possible morphological forms), the second type (representing lexicalized alternations, see below) indicates the semantic relationships between different lexical units.

3. Basic types of alternations

Here we further develop the results of the theoretical research of Czech alternations and modify the proposal of their representation presented in Kettnerová and Lopatková (2010). The changes in valency structure of verbs are associated with specific relations between different surface syntactic structures of the same verb lexeme related to the same situational meaning. According to linguistic means by which these relations are expressed, we

distinguish (i) *grammaticalized alternations* expressed by grammatical means, as in (3)(a)-(3)(b) or (4)(a)-(4)(b) and (ii) *lexicalized alternations* expressed by lexical-semantic means, that is by a change of a lexical unit of a verb, e.g. (5)(a)-(5)(b). We observe that the alternations of both types are either (a) conversive, or (b) non-conversive. (a) The *conversive alternations* have the character of permutation of situational participants where the prominent surface syntactic position of subject or direct object is involved. (b) The *non-conversive alternations* are characterized by changes in the linking of a single situational participant and syntactic positions which cannot be classified as permutations. Whereas the conversive alternations play a central role in the perspectivization of a situation denoted by a verb, the non-conversive alternations represent rather peripheral means.

- (3) (a) *Recepční*_{Agent-ACT-Subj} *hostu*_{Recipient-ADDR-InObj} *přidělil*_{active} *pokoj*_{Patient-PAT-Obj} č. 11.
 (b) *Host*_{Recipient-ADDR-Subj} *dostal přidělen*_{recip} *pokoj*_{Patient-PAT-Obj} č. 11 (od recepčního)_{Agent-ACT-Adv}
- Eng. (a) The receptionist_{Agent-ACT-Subj} has allocated_{active} the guest_{Recipient-ADDR-InObj} room_{Patient-PAT-Obj} n. 11.
 (b) The guest_{Recipient-ADDR-Subj} has been allocated_{recip} room_{Patient-PAT-Obj} n. 11 (by the receptionist)_{Agent-ACT-Adv}
- (4) (a) *Jan*_{Speaker-ACT-Subj} *všechny opravy*_{Information-PAT-Obj} *domu konzultoval*_{active} *se svým otcem*_{Recipient-ADDR-InObj}
 (b) (*Jan a jeho otec*)_{Speaker/Recipient-ACT/ADDR-Subj} (*spolu*) *konzultovali*_{active} *všechny opravy*_{Information-PAT-Obj} *domu*.
- Eng. (a) John_{Speaker-ACT-Subj} consulted_{active} all house repairs_{Information-PAT-Obj} with his father_{Recipient-ADDR-InObj}
 (b) (John and his father)_{Speaker/Recipient-ACT/ADDR-Subj} consulted_{active} all house repairs_{Information-PAT-Obj} (together).
- (5) (a) *V sále*_{Location-LOC-Adv} *zní sborový zpěv*_{Bearer-ACT-Subj}
 (b) *Sál*_{Location-ACT-Subj} *zní sborovým zpěvem*_{Bearer-PAT-Adv}
- Eng. (a) The choral singing_{Bearer-ACT-Subj} sounds in the hall_{Location-LOC-Adv}
 (b) The hall_{Location-ACT-Subj} sounds with the choral singing_{Bearer-PAT-Adv}

3.1. Grammaticalized alternations

This type of alternations stem from use of specific grammatical means. The uses of a verb are characterized by the same situational and structural meaning; that is the same sets of situational participants are linked with the same sets of valency complementations in the same way. What is different for such uses of a verb is surface syntactic expression of some situational participants. Thus grammaticalized alternations are typical of the relations between different surface syntactic structures of a single lexical unit of a verb. In Czech these alternations are either of conversive, or of non-conversive character: the conversive grammaticalized alternations are connected with *diatheses*, the non-conversive alternations of the same type are associated with *reciprocity*.

3.1.1. *Diatheses*. In Czech, diatheses represent the relations between surface syntactic structures of a verb which differ in the grammatical category of voice, that is, they are associated with specific morphological meanings of the verb. Five specific meanings of Czech verbs are determined: passive, deagentive, resultative, dispositional and recipient-passive meanings, see Panevová (man.). Active voice constitutes the unmarked opposition to all the marked meanings. The use of a specific meaning of a verb results in changes in its valency

structure. These changes are conversive: they result in the permutation of valency complementations (respective situational participants): 'ACTor' is prototypically shifted from the prominent subject position into a less prominent surface position. As a result, the situation portrayed by a verb is perspectivized either from the point of view of the situational participant corresponding to 'ACTor' (usually 'Agent' or 'Causator') (the unmarked syntactic structures with active voice), or from the point of another participant (the marked structures with a certain specific morphological meaning of a verb). See the use of active and recipient-passive meaning of the verb *přidělit* 'to allocate' accompanied by the permutation of the valency complementations 'ACTor' and 'ADDRessee' (corresponding to the situational participants 'Agent' and 'Recipient', respectively) in (1)(a)-(1)(b) given above, see also schema in Figure 1.

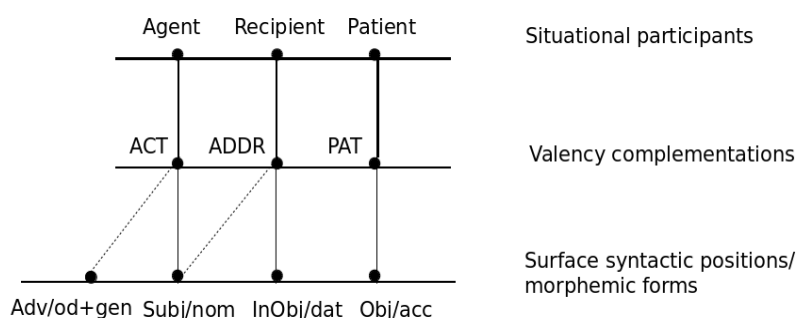


Figure 2. Recipient-passive diathesis of the verb *přidělit* 'to allocate'.

3.1.2. *Reciprocity*. Czech non-conversive alternations are characteristic of reciprocity. In contrast to diatheses, reciprocity does not consist in the use of any specific morphological meaning of a verb (the category of voice is preserved) but it is expressed primarily by syntactic means. Reciprocalization is a syntactic operation two (or three) valency complementations (respective situational participants) – if their semantic properties allow for it – are used symmetrically. The reciprocal use of valency complementations leads to the shift of the valency complementation expressed in a less prominent surface syntactic position into the more significant syntactic position (subject or direct object) of the other symmetrically used valency complementation. As a result, whereas the prominent position is ‘multiplied’ by syntactic or by morphological means (e.g., coordination, plural), the less significant position is deleted in the reciprocal surface structure. See the reciprocity of 'ACTor' and 'ADDRessee' (corresponding to the situational participants 'Speaker' and 'Recipient') of the verb *konzultovat* 'to consult' in (2)(a)-(2)(b) given above and in Figure 2.

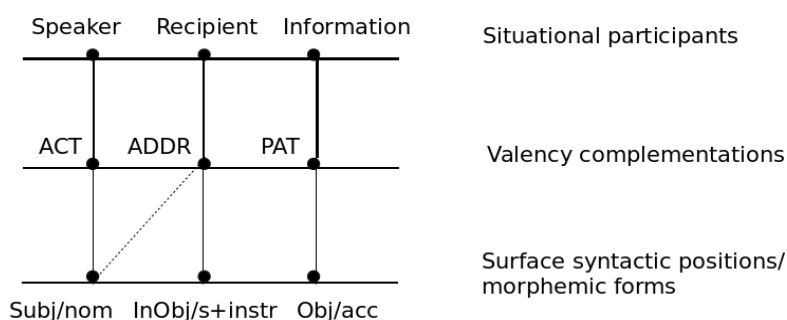


Figure 3. Reciprocity of 'ACTor' and 'ADDRessee' of the verb *konzultovat* 'to consult'.

3.1.3. *The representation of grammaticalized alternations.* Both conversive and non-conversive grammaticalized alternations are limited only to changes in morphemic forms of valency complementations which are affected by the shifts in surface syntactic positions. These changes are regular enough to be captured by formal syntactic rules. These rules are stored in the rule component of the lexicon and make it possible to derive the valency frames corresponding to the marked surface syntactic structures from the valency frames describing unmarked ones. At present, transformational rules formulated for the purposes of the description of diatheses in *PDT-VALLEX*, the lexicon of Prague Dependency Treebank, are used, see (Urešová, 2011).

As for the representations of grammaticalized alternations in the *data component*, each lexical unit is represented by a single valency frame corresponding to the unmarked use (i.e., active use). Concerning diatheses, the information on the possible application of specific morphological meanings is assigned to each relevant lexical unit in the special attribute – *diat*. Concerning reciprocity, a list of valency complementations possibly involved in reciprocal use is given in the attribute –*rcp*.

- **lemma:** *přidělovat*^{impf}, *přidělit*^{pf} 'to allocate'
- **gloss:** *dát do vlastnictví n. užívání* 'to give to ownership or usage'
- **frame:** ACT₁^{obl} ADDR₃^{obl} PAT₄^{obl}
- **example:** *učitel každému žáku přidělil učebnice* 'the teacher allocated each student a textbook'
- **diat:** recip
- **rcp:** ACT-ADDR

Table 1 gives an example of the syntactic rule for the recipient passive diathesis Recip.r describing the changes in verbal voice and the morphemic forms of the valency complementations of 'ACTor' and 'ADDRessee' of the verb *přidělit* 'to allocate' in example (1)(a)-(1)(b) above:²

Table 1. Syntactic rule for the recipient passive diathesis of the verb *přidělit* 'to allocate'.

| Type: recipient- passive | | | Commentary |
|--------------------------------|----------|--|------------|
| Action | verbform | replace (active vf → recipient passive vf) | (1) |
| | ACT | replace (nom → od+gen) | (2) |
| | ADDR | replace (dat → nom) | (3) |

Commentary:

(1) The verb form changes from active form to recipient passive form (auxiliary verb *dostat* + participle of a given lexical verb).

(2) The morphemic form of 'ACTor' changes from nominative into the prepositional group *od*+genitive.

(3) The morphemic expression of 'ADDRessee' changes from dative into nominative.

3.2. *Lexicalized alternations*

Lexicalized alternations³ represent such changes in valency structure of a verb which are

associated with the change of lexical unit. They are characteristic of the uses of a verb which are characterized by the same situational meaning whereas their structural meaning is different. It implies that the same set of situational participants is mapped onto the valency complementations in a different way. As a result, the involved situational participants are differently syntactically structured on surface. The lexicalized alternations are conversive or non-conversive: whereas the conversive lexicalized alternations (referred here to as lexical-semantic conversions) – similarly as grammaticalized conversive alternations – are crucial for the perspectivization of a situation portrayed by a verb, the non-conversive lexicalized alternations are restricted to a few language specific constructions (e.g., multiple structural expression of a situational participant, or structural splitting of a situational participant).

3.2.1. *Lexical-semantic conversion.* Lexical-semantic conversion relates different surface syntactic structures based on different lexical units of the same verb lexeme. These lexical units share the same situational meaning; however, its situational participants are mapped onto a different set of valency complementations. The changes in valency frames of such lexical units can affect the number of valency complementations, their types, obligatoriness and morphemic forms. Prototypically, they lead to a permutation of some situational participants while the prominent subject or direct object position is affected. See the alternation of the verb *znít* 'to sound' in (3)(a)-(3)(b) and in Figure 4 resulting in the inverse role of the situational participants 'Bearer (of action)' and 'Location'.

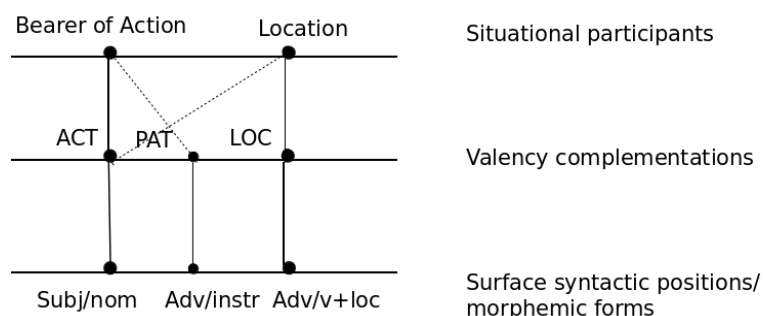


Figure 4. Lexical-semantic conversion Bearer of Action-Location of the verb *znít* 'to sound'.

3.2.2. *Czech non-conversive lexicalized alternations.* This type of lexicalized alternations can be exemplified by the multiple structural expression of a situational participant or by the structural splitting of a situational participant. In case of the *multiple structural expression of a situational participant*, the changes in valency structure of a verb arise from two possible mappings of a single situational participant onto different valency complementations. See the different linking of the situational participant 'Goal' of the verb *vyjít* 'to climb' onto 'DIR(ection)' valency complementation (4)(a) and on the 'PAT(ient)' (4)(b), respectively.

The *structural splitting of a situational participant* is typical of verbs of communication. These verbs allow one of its participants to be linked either with a single valency complementation, or with two valency complementations. See the difference in the mapping of the participant 'Information' of the verb *řici* 'to say', which is related either to 'PAT(ient)' in (5)(a) or it is split into 'PAT(ient)' and 'EFF(ect)' in (5)(b).

- (4) (a) *Horolezci*_{Agent-ACT} *vylezli na Mount Everest*_{Goal-DIR}
 (b) *Horolezci*_{Agent-ACT} *vylezli Mount Everest*_{Goal-PAT}
 Eng. (a) The mountaineers_{Agent-ACT} climbed up Mount Everest_{Goal-DIR}
 (b) The mountaineers_{Agent-ACT} climb Mount Everest_{Goal-PAT}

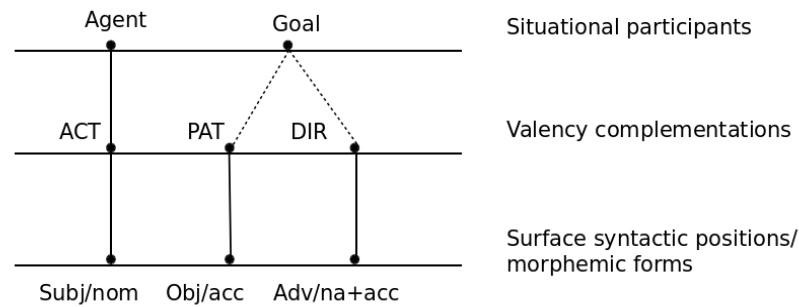


Figure 5. Multiple structural expression of 'Goal' of the verb *vylézt* 'to climb'.

- (5) (a) *Jana*_{Speaker-ACT} *řekla*, (*že její tchyně je moc hodná*)_{Information-PAT}
 (b) *Jana*_{Speaker-ACT} *řekla o své tchyni*_{Information-PAT}, (*že je moc hodná*)_{Information-EFF}
 Eng. (a) *Jane*_{Speaker-ACT} *said* (*that her mother-in-law is very kind*)_{Information-PAT}
 (b) *Jane*_{Speaker-ACT} - *said* - *about her mother in law*_{Information-PAT} - (*that - is - very kind*)_{Information-EFF}

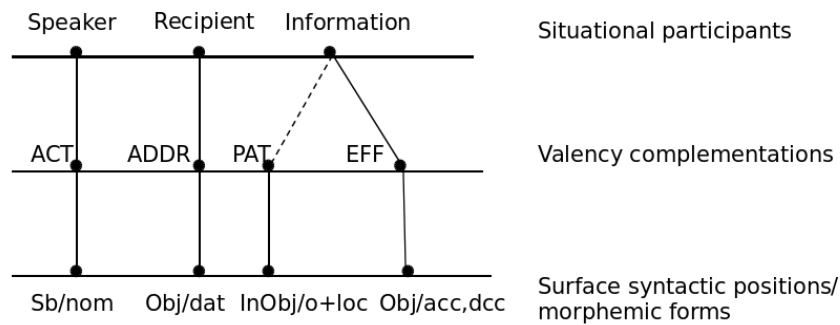


Figure 6. Structural splitting of the situational participant 'Information' of the verb *řici* 'to say'.

3.2.3. *The representation of lexicalized alternations.* In the data component, there are two lexical units related by a certain lexicalized alternation; these lexical units are represented by separate valency frames. In the rule component, general rules determining changes in the mapping of situational participants onto valency complementations are included.

For instance, in the data component of the lexicon, syntactic variants of the verb *znít* 'to sound' in the relation of lexical-semantic conversion in examples (3)(a)-(3)(b) given above are represented by the following lexical units with two different valency frames (i) and (ii) (corresponding to the use of the verb in (3)(a) and (3)(b)). These frames are interlinked by a relevant type of relation ascribe to them in the special attribute *-conv*.

- **lemma:** *znít*^{impf} 'to sound'
- **gloss:** *vydávat zvuk* 'to produce sound'
- **frame:** ACT₁^{obl} LOC^{typ}
- **example:** *v sále zněla hudba* 'music sounds in the hall'
- **conv:** bear-loc

- **lemma:** *znít*^{impf} 'to sound'
- **gloss:** *být naplněn zvukem* 'to be full of sound'
- **frame:** ACT₁^{obl} PAT₇^{obl}
- **example:** *sál zněl hudbou* 'the hall sounds with music'
- **conv:** bear-loc

The interlinking between valency complementations in the above given valency frames and the situational participants 'Bearer of Action' and 'Location' is specified by the rule given in Table 2.

Table 2. The mapping of situational participants and valency complementations of the verb *znít* 'to sound'.

| Situational participants | Valency frame (i) | Valency frame (ii) |
|--------------------------|-------------------|--------------------|
| 'Bearer of Action' | ACT | PAT |
| 'Location' | LOC | ACT |

The non-conversive lexicalized alternations (Section 3.2.2) are represented in the lexicon in the same way as conversive lexicalized alternations, that is, they are captured as separate lexical units stored in the data part of the lexicon interlinked by a relevant type of alternation. Then the rule part of the lexicon provides rules describing changes in the mapping of situational participants onto valency complementations characteristic of non-conversive lexicalized alternations.

6. Conclusion

We have proposed the representation of changes in valency structure of Czech verbs in the valency lexicon *VALLEX*. We have demonstrated that whereas grammaticalized alternations can be described by syntactic rules, lexicalized alternations require rather general rules. These rules are stored in the rule component of the lexicon. In the data component, only valency frames corresponding to the unmarked use (i.e., active use) of lexical units are captured; different (morpho)syntactic uses of a single lexical unit are obtained by applying particular rules from the rule component. In case of lexicalized alternations, separate lexical units are stored in the lexicon; these lexical units are interlinked by a relevant type of alternation.

Notes

¹ This work has been using language resources developed and/or stored and/or distributed by the LINDAT-Clarin project of the Ministry of Education of the Czech Republic (project LM2010013). The research reported in this paper has been supported by the Grant of the Grant Agency of the Czech Republic No. GA P406/12/0557 and partially by the grant No. GA P406/10/0875.

² The given rule for recipient passive diathesis is simplified for better understanding. Especially, the conditions on applicability of the rule are left aside here.

³ Here we focus only on the lexicalized alternations of the same verb lexeme. The lexicalized alternations expressed by the change of verb lexeme (e.g., *koupit – prodat* 'to buy' – 'to sell') are left aside here.

References

A. Dictionaries

DeepDict. <http://gramtrans.com/deepdict/>.

LexIt. <http://sesia.humnet.unipi.it/lexit/>.

VALLEX 2.5. <http://ufal.mff.cuni.cz/vallex/2.5/>.

PDT-VALLEX. <https://ufal.mff.cuni.cz/lindat/PDTVallex.html>.

B. Other literature

Bejček, E., M. Lopatková and V. Kettnerová 2010. ‘Advanced Searching in the Valency Lexicons Using PML-TQ Search Engine.’ In P. Sojka, A. Horák, and I. Kopeček and K. Pala (eds.), *Proceedings of the 13th International Conference, TSD 2010*. Berlin / Heidelberg: Springer, 51–58.

Bick, E. 2009. ‘DeepDict – A Graphical Corpus-based Dictionary of Word Relations.’ In K. Jokinen and E. Bick (eds.) *Nordic Conference of Computational Linguistics NODALIDA 2009*. NEALT Proceedings Series, Vol. 4 (2009), Northern European Association for Language Technology (NEALT), 268–271.

Kettnerová, V. and M. Lopatková 2010. ‘The Representation of Diatheses in the Valency Lexicon of Czech Verbs.’ In H. Loftsson, E. Rögnvaldsson and S. Helgadóttir (eds.), *Proceedings of the 7th International Conference on Advances in Natural Language Processing*. Berlin / Heidelberg: Springer, 185–196.

Lenci, A. 2009. ‘Argument alternations in Italian verbs: a computational study.’ In *Atti del XLII Congresso Internazionale di Studi della Società di Linguistica Italiana*.

Levin, B. 1993. *English Verb Classes and Alternations. A Preliminary Investigation*. Chicago and London: The University of Chicago Press.

Lopatková, M., Z. Žabokrtský and V. Kettnerová 2008. *Valenční slovník českých sloves*. Praha: Karolinum.

Mel'čuk, I. A. 2004. ‘Actants in Semantics and Syntax I.’ *Linguistics* 42.1: 1–66.

Panevová, J. 1994. ‘Valency Frames and the Meaning of the Sentence.’ In P. A. Luelsdorff (ed.), *The Prague School of Structural and Functional Linguistics*. Amsterdam, Philadelphia: John Benjamin Publishing Company, 223–243.

Panevová, J. et al. Manuscript. *Syntax současné češtiny (na základě anotovaného korpusu)*. Praha: Nakladatelství Karolinum.

Sgall, P., E. Hajičová, and J. Panevová 1986. *The Meaning of the Sentence in its Semantic and Pragmatic Aspects*. Dordrecht: Reidel.

Urešová, Z. 2011. *Valence sloves v Pražském závislostním korpusu*. Prague: Institute of Formal and Applied Linguistics.

Žabokrtský, Z. and M. Lopatková 2007. ‘Valency Information in VALLEX 2.0: Logical Structure of the Lexicon.’ *The Prague Bulletin of Mathematical Linguistics* 87, 41–60.