

Verb Frames Extracted from Dictionaries

Hana Skoumalová

Abstract

In this work, an electronic lexicon of Czech verbs is presented. The lexicon contains valency frames of approx. 15,000 Czech verbs, and its purpose is to enrich information contained in other electronic dictionaries. The trend of recent years is to make large-scale reusable sources which can be combined with other sources. This work shows how the lexicon cooperates with an existing morphological lexicon and how it can be used in various NLP systems.

This article contains a substantial part of my dissertation thesis *Czech syntactic lexicon*, which was defended at Charles University in 2001. The full version of the dissertation can be found on the web (Skoumalová, 2001). In this article, some errors are corrected, and the development of the dictionary since finishing the thesis is briefly described.

1 Introduction

In the era of computers, language processing has gained a form different from what was known before. Vast amounts of data are available and computers can process them in a reasonably short time, but they need adequate tools for their work. Beside grammar rules they also need lexicons which they can understand.

In this work, an electronic lexicon of Czech verbs is presented. The use of the lexicon in Natural Language Processing (NLP) makes special demands on it. It differs from “human” lexicons in that all information must be explicit or deducible by exactly formulated rules of derivation.

While sketching the format of the dictionary, interesting theoretical problems were encountered, which are discussed in this work. The lexicon contains valency frames of approx. 15,000 Czech verbs, and its purpose is to enrich information contained in other electronic dictionaries.

1.1 Terminological remarks

Various authors differ in their understanding of the term *subject*. We will consider a subject only such a member of a frame which is in Nominative and with which the main verb agrees. Our criterion is the question for a subject: *kdo, co* (who_{Nom}, what_{Nom}). This means that we will not take Genitive in such constructions as *vody ubývá* (water_{Gen} diminishes) as subject. On the other hand, a clause or an infinitive can be a subject, as we can ask the above question; in such a case the verb shows agreement with neuter singular.

In the text, we will use the terms *actants* and *inner participants* as synonymous. *Actant* is Tesnière’s term, while FGD uses *inner participant*, but their meaning is so close that they are often interchanged.

We will also use the terms *animate* and *animacy*. For purposes of this work we will divide nouns into two groups: personal and non-personal. The former can be represented by the pronoun *kdo* (who), the latter by the pronoun *co* (what). Sometimes we will refer to personal nouns as to animate ones and to non-personal as to inanimate.

1.2 Theoretical background

When describing the role of verbs in the language, all authors agree on the necessity to describe syntactic properties of verbs in the dictionary. But they differ in the understanding of what kind of information should be included. Dictionaries for practical usage (language dictionaries for human readers, or machine dictionaries for grammar checking or shallow parsing) contain usually only the surface information.

Dictionaries that serve more sophisticated purposes must contain also information on the argument structure, and the relations between the two layers of linguistic description. The two views of the

dictionary differ also in their understanding of what belongs to the verb frame. The classical lexicologists collect all *typical* complementations while the theoreticians discriminate between *arguments* and *adjuncts*. The arguments are listed in subcat lists and grammar rules check whether all of them are present in the sentence. The adjuncts, on the other hand, are not obligatory, can occur more than once in a sentence, and they are not listed in the dictionary entries.

The dictionary described in this work is meant to provide for the automatic processing of the Czech language. The algorithms for the language processing do not necessarily have to be based on a linguistic theory, but we believe that with a theory we can develop algorithms that are efficient and elegant because they are linguistically adequate. The results of these algorithms, on the other hand, can serve to a linguistic theory as a feedback which helps to improve the theoretical description.

For this work we decided to utilize the Functional Generative Description (FGD) developed by Sgall, Hajičová and Panevová (Sgall, Hajičová, and Panevová, 1986), and especially the part dealing with the verb frames (Panevová, 1974-75; Panevová, 1980).

1.2.1 An overview of FGD

In FGD, several levels of language description are distinguished. For purposes of this work, we will only work with two of them—the *tectogrammatical* level and the *morphemic* level. To be able to express certain relations we will also need the notion of *subject*.

Each level has its own units, basic and compound. The compound units are formed from the basic ones with the help of *C-relations*. The translation between two neighbouring levels is provided by *R-relations*. The basic units on tectogrammatical level are *semantemes* (lexical units), *functors* (syntactic units) and *grammatemes*. The compound units are *propositions*. The functors also serve as the C-relations with the help of which the propositions are constructed (Sgall, 1967).

There are two types of functors—*inner participants* (Tesnière's *actants*) and *free modifications*. A verb frame denotes which functors are required by a certain semanteme (verb lemma). A frame can contain up to five inner participants (*Actor*, *Patient*, *Addressee*, *Origin* and *Effect*) and any number of free modifications. Some of the inner participants can be *optional* (also called *facultative*), which means that they do not need to be present in the sentence—neither on the tectogrammatical nor morphemic level. Other participants are always *obligatory*. However, they can be realized as *general*—the structure on the tectogrammatical level then contains a general participant, which is not realized as such on the morphemic level. Whether a participant is optional or obligatory, and whether an obligatory participant can be realized as general can be tested by a question test (Panevová, 1980). Let us imagine the following dialogue:

- (1) *Petr čte. Co? Nevím.*
Petr is reading. What? I don't know.

The answer 'I don't know' is acceptable, as the the speaker does not need to know what Petr's reading is, but it must be something which is usually read (a newspaper, a book, etc). This shows that Patient in the frame of the verb *číst* (read) can be general. On the other hand, in dialogue (2), the answer 'I don't know' is nonsensical. This shows that Actor is an obligatory participant in the frame of the verb *přijít* (come).¹

- (2) *Už přišel. Kdo? *Nevím.*
(He) has already come. Who? I don't know.

In example (3) the sentence is actually ungrammatical, if the participant is omitted—this is clear evidence that the participant is obligatory.

- (3) **Petr daroval.*
Petr donated.

Free modifications normally are not members of a frame, but they can become members as *obligatory free modifications*:

¹The fact that the surface realization of Actor in this sentence is omitted is caused by another phenomenon: Czech is a so called pro-drop language and thus a personal pronoun in the position of a subject can be omitted. Morphological markers of the person and number (in the past tense also of the gender) are present also in the verb form and thus the personal pronoun is redundant (Karlík, 2000).

- (4) a. *Jan se choval jako blázen.*
Jan behaved like a fool.
- b. **Jan se choval.*
Jan behaved.

In some cases, when the modification is known from the context, it can be omitted on the surface; such free modification is called *obligatory and deletable free modification*. For testing whether a free modification is an obligatory member of a frame the question test can be used again. In the sentence in (5) the question test proves that the direction is an obligatory and deletable free modification of the verb *přijít* (come, arrive).

- (5) *Petr přišel. Kam? *Nevím.*
Petr arrived. Where? I don't know.

In other theoretical models (Daneš, Hlavsa, and others, 1987; Grepl and Karlík, 1989; Karlík, Nekula, and Rusínová, 1995), the repertory of participants is wider: instead of Actor the authors speak about Agent, Causer, Experiencer, etc. Patient is more or less a synonym of the direct object and Recipient a synonym of the indirect object. In FGD, Actor and Patient are determined by syntactic criteria rather than by semantic ones (cf. Tesnière (1959)), and other participants are determined semantically:

- (6) 1. If the verb frame contains only one participant, this participant is Actor.
2. If the frame contains two participants, one of them is Actor and the other is Patient. In most cases, Actor is the subject of the active construction, but there are some exceptions to this rule, which will be discussed later.
3. If the verb frame has more than two participants, the roles of Actor and Patient must be occupied, and the other participants occupy the roles of Addressee, Effect or Origin. The decision about which participant bears which role is based on the semantics of the participants.

The basic units on the morphemic level are *semata*, and the compound units are *morphemes* and *formemes*—units which combine prepositions with morphological cases.

The lexicon in FGD contains semantemes, their functors and grammatemes. In our informal example, parentheses denote whether a functor is obligatory, obligatory deletable or optional:

- (7) *spát* (sleep) Act
pojídat (devour) Act Pat
těšit_se (look forward) Act Pat Gram:{Refl[se]}
darovat (donate) Act Pat (Addr)

Beside it, the lexicon should also define the R-relation which translates every functor and grammateme to the morphemic level. After this addition, the lexicon will have the following format:

- (8) *spát* Act[Noun+Nom]
pojídat Act[Noun+Nom] Pat[Noun+Acc]
těšit_se Act[Noun+Nom] Pat[Noun+Acc+na] Gram:{Refl[se]}
darovat Act[Noun+Nom] Pat[Noun+Acc] (Addr[Noun+Dat])

When we try to compare FGD with other linguistic theories we can make a parallel between functors and θ -roles. In the other theories, subcat lists are viewed as primary syntactic structure attached to lexical entries and the θ -roles are mapped onto the subcat list by some kind of mapping function. In FGD an opposite assumption is made: the tectogrammatical functors form a primary syntactic structure of a verb and the surface forms are their counterparts on the morphemic level which are translated by R-relation from the functors.

Beside this, the θ -roles differ from the repertory of participants in FGD. Not only are their names different, but also their distributions to single verbs. An *Actor* in FGD can be marked as *Agent* or *Bearer* or *Experiencer* in other theories, etc.

If we use FGD as the background theory of a dictionary, we will be unable to transfer the lexicon to another theoretical framework 'as is'; we tried, however, to make it possible to extract the subcat lists. More details on this issue can be found in (Skoumalová, 2001).

For utilizing the tectogrammatical information, we would have to find a mapping function which would have to take into consideration also the semantics of single verbs, which will be the subject of further research.

2 Using existing resources

When we try to create a new electronic dictionary, it is of course possible to start from scratch, but it is more efficient to use existing resources. Printed dictionaries usually contain syntactic information, but unfortunately this information is meant for human readers, and very often it is assumed that the reader knows the rules that apply in usual cases, and only exceptions are listed. Beside this, the information is not encoded in a formal way that could be understandable to a machine.

There exists a Czech dictionary of verbs (Svozilová, Prouzová, and Jirsová, 1997) which contains the verb frames encoded in a formal way. But its size is quite limited (ca 600 verbs) and the information concerns only the surface frames. Nevertheless, this dictionary can serve as an aid to creators of an electronic dictionary.

One of the first attempts at making an electronic dictionary of verb frames was made in the project RUSLAN (Oliva, 1989). This project was focused on machine translation from Czech to Russian and the format of the lexicon was adapted for this purpose; it contained the Czech word stem and its Russian translation, Czech and Russian morphological information, the Czech surface frame and its translation to the Russian surface frame. The domain of the translated texts were programming manuals, which affected the coverage of the lexicon. Another drawback (caused by limited computational resources) was the small size of the lexicon—it contained ca 10,000 entries (including all word classes). The work invested in this project was useful for gaining experience with natural language processing rather than for creating working software.

Another small lexicon was created for the purposes of the project LaTeSlav (Avgustinova et al., 1995). This was a project for creating grammar-checkers for two Slavic languages (Czech and Bulgarian). In fact, there were two lexicons for Czech, as the project split into two branches. Both the lexicons contained a small number of entries which had very rich syntactic information, but unfortunately they were “hardwired” in the software and it would not be easy to extract them for other purposes.

The most promising source of valency frames is a dictionary created at Masaryk University by Karel Pala and his team (Pala and Ševeček, 1997). This dictionary was compiled from several printed dictionaries, and the valency frames were taken mainly from SSJČ (1989). We used this dictionary as a source of surface frames and enhanced them with information at the tectogrammatical level.

2.1 Source data

Our dictionary contains ca 15,000 verbs with surface frames. The original format called BRIEF contains lemma, starting delimiter of the list of frames (<v>) and the list itself (see example in 9a). (9b) translates this notation to a readable form.

- (9) a. `agitovat <v>hPc4,hPc3-hPc4,hPTc4r{pro},hPTc3r{proti}`
b. **agitovat** (to agitate) *koho (komu), pro koho, proti komu*

In BRIEF format, frames are separated by commas, and single members of a frame are separated by dashes. The obligatoriness is not marked, but a frame can be repeated several times, with and without the optional, deletable or generalizable members. In example (9) this is the case of the frame *koho (komu)*.

BRIEF encoding is described in Horák (1998a) and Horák (1998b). Here, we only provide a short overview of attributes and values used in the dictionary. Every member of a frame is described by a list of attributes and their values. We can understand these attributes and their values as grammatemes occurring on the tectogrammatical level.

2.1.1 The attributes used in the lexicon and their values

h — ‘Semantic’ feature. This attribute has rather heterogeneous values. Single values are only applicable for certain word classes and thus they include implicit information on the part of speech as well. The values are:

- P — Person (only applicable for nouns and pronouns); this value actually stands for ‘case questions’ *kdo* (who), *koho*, etc.
- T — Thing (only nouns and pronouns); it stands for ‘case questions’ *co* (what), *čemu*, etc. The values P and T can be grouped together.
- R — Long reflexive pronoun *sebe*, *sobě*, etc.
- Q — Quality (only adjectives).
- M — Amount (only numbers).
- L — Location (only adverbs).
- A — Direction where (only adverbs).
- F — Direction from (only adverbs).
- D — Which way (only adverbs).
- W — When (only adverbs).
- c — Morphological case. This attribute is only applicable for nominal word classes, and so it only occurs if the h attribute has one of the values P, T, R, or Q. The values are 1, 2, 3, 4, 6 and 7.
- r — Preposition. This attribute can only occur after a morphological case. The value is the preposition itself closed in curly brackets: *r{na}*, *r{o}*, *r{vzhledem k}*, etc.
- s — Clause or infinitive. The values are:
- I — Infinitive.
 - C — Clause attached by the conjunction *až* (when).
 - D — Clause attached by the conjunction *že* (that).
 - F — Clause attached by the conjunction *jestli, zda* (if, whether).
 - P — Clause attached by the conjunction *ať* (let).
 - R — Clause attached by a relative expression *co* (what), *kteřý* (which), *kdo* (who), *kolik* (how many), etc.
 - U — Clause attached by the conjunction *aby* (so that).
 - Z — Clause attached by the conjunction *jak* (how).
- e — Negation (in a clause). The values are A (affirmative) and N (negative). The affirmative value is the default and it is not marked in the lexicon.²
- i — Idiom. The value is a string closed in curly brackets. The string contains words forming the idiom and a case marker for the variable part. If there are possible variants in the fixed part, they are put in parentheses and separated by commas, or they are separated by a vertical bar. The variants in the variable part are separated by a vertical bar. Examples:
- | | |
|--|-------------------------------|
| <code>brát <v>i{pod ochranu do ochrany <koho>}</code> | (take sb under protection) |
| <code>dávat <v>i{konzert hru film}</code> | (put concert, play, movie on) |
| <code>házet <v>i{přes palubu <koho co>}</code> | (throw sb over board) |
| <code>chovat <v>i{(přátelství, zášť, nenávisť) <ke komu>}</code> | (feel friendship, hatred) |
- v — Constraint applied for a single valency frame. The constraint is an attribute with a required value, or an attribute with a prohibited value, preceded by `^`. Currently, only `v{eN}` is used, for verbs whose negated forms have different valency frames:

²This attribute is mainly used together with a clause attached by the conjunction *aby* (so that)—`sUeN`, e.g. *bát se* (fear), *varovat* (warn), etc. Though this is a typical usage, the affirmative clause cannot be excluded. After a simple query in the Czech National Corpus (Kocěk, Kopřivová, and Kučera, 2000) we found eight affirmative clauses (out of ca 230 occurrences of the verb *bát se* with the conjunction *aby*), e.g. *Po volbách se úředníci bojí, aby přežili ... změnu dnešního ministra ...* (After elections, clerks are afraid whether they will survive the change of the current minister ...).

hledět <v>hPTc4r{na}, hPc3, hTc2r{do}, hPc3-hTc2r{do}, v{eN}hTc3r{k} (not to look at st)
páchnout <v>hTc6r{po}, hTc7, v{eN}hTc2r{do} (not to set foot on st)
znát <v>hPTc4, v{eN}hTc2z{jen se zápořem}, hTc4-hPTc6r{na} (not to know--Genitive of negation)

z — Comment in curly brackets (see the example above).

The frames do not contain subjects as the printed dictionaries usually do not list them. For an automatic processing of language, however, this information is necessary. We can make a simple assumption that the subject will be a noun in Nominative but there are exceptions to this rule. We will discuss this in more detail in Section 4.

3 Content of the final lexicon

In this section, a detailed description of phenomena recorded in the lexicon is given, as well as a thorough description of the encoding of all the linguistic information. First we will give a formal description of the format of a frame and then we will explain the meaning of single fields. After that we will describe in depth what kinds of reflexive and reciprocal verbs we distinguish in the Czech language and how we encode them in the lexicon. Then we will deal with diatheses covered by the lexicon and finally we will discuss the so called *equi* and *raising* verbs.

3.1 Format of a lexical entry

A lexical entry contains a lemma and its frame.³ The term *frame* usually denotes all types of complementations of a verb in one meaning. The existence of another frame then signals a new meaning. There are, however, variants of surface realizations of functors—in such a case we do not introduce a new meaning but we merge the variants into one frame. In our lexicon, the frame contains all the variants merged together, and in addition it also includes information on possible diatheses. As it is not always possible to accommodate all the combinations of surface realizations and diatheses into one frame, we may be forced to split one meaning into several lexical entries. The identification of one lexical meaning is then provided by indices (different from the indices from the morphological lexicon) attached by ~. Examples of lexical entries are shown in (10).

(10) adresovat	Act [Noun+Nom] Pat [Noun+Acc] \	
	Addr [Noun+Dat Noun+Acc+na Noun+Acc+pro] \	
	PeriphPass ReflPass	(address)
stát-2~1	Act [Noun+Nom] Gram: {Refl[se]} NoPass	(happen)
stát-2~2	Act [Noun+Nom] Pat [Noun+Ins] Gram: {Refl[se]} NoPass	
		(become)
stát-3~1	Act [Noun+Nom] ReflPass	(stand)
stát-3~2	Act [Noun+Nom] Pat [Noun+Acc+o] ReflPass	(long for)
stát-3~3	Act [Noun+Nom] Pat [Noun+Ins+za] Gram: {Refl[si]} NoPass	
		(be convicted)
stát-4	Act [Noun+Nom] Pat [Num+Acc] NoPass	(cost)
učit~1	Act [Noun+Nom] Gram: {Refl[se]} NoPass	(learn)
učit~2	Act [Noun+Nom] Pat [Noun+Acc] Addr [Noun+Acc] NoPass	
učit~2	Act [Noun+Nom] Pat [Noun+Dat] Addr [Noun+Acc] \	
	PeriphPass ReflPass	(teach)

The verb *adresovat* (address) has only one lemma in the morphological lexicon and only one meaning. The verb *stát* has three different lemmas in the morphological lexicon—one for the reflexive verb *stát*

³As we expect our lexicon to be used together with the morphological lexicon created by J. Hajič (Hajič, 1994) the lemmas must be identical with the lemmas of the morphological lexicon. This means that lemmas must contain the same indices as the morphological lexicon (e.g. *stát-2* (happen), *stát-3* (stand), *stát-4* (cost), etc.). Furthermore, lemmas of reflexive verbs do not contain the reflexive particle (e.g. *stát se* (happen) will have the lemma *stát-2*).

se and two for the non-reflexive verb *stát* (the reflexive verb *stát si* is morphologically covered by the non-reflexive verb *stát*). The reflexive verb is split into two entries with two different meanings in our lexicon (*stát-2~1* and *stát-2~2*), the meanings of the non-reflexive verbs are partly differentiated by the indices from the morphological lexicon, so we have to decide which of the “morphological” meanings will be split. The verbs *učít se* and *učít* have only one entry in the morphological lexicon, but we have to introduce two meanings for them. The second meaning (*učít*—teach) must itself be split into two frames, as the frame variant with two Accusatives does not allow for the formation of a passive, while the variant with Accusative and Dative allows for the formation of both periphrastic and reflexive passive.⁴

The frame is separated from the lemma by a tabulator. A frame has the following format:

`<voice><reflexivity><subject>?[<functor><grammatemes>]*<diathesis>+`

A frame starts with a voice marker, which is obligatory. Then follows a marker for reflexivity, which is also obligatory. The subject marker may be missing, as there exist verbs without a subject. After the subject marker, a list of functors and their corresponding grammatemes follows. This list can be empty, as we suppose that there are verbs with an empty frame (the obvious candidates, meteorological verbs, however, do not belong to this category, as they need an obligatory modification of the location; e.g. *pršet (kde)*—rain). The frame ends with markers of possible diatheses.

In the following sections, single parts of a frame will be described in detail.

3.1.1 Voice

The voice marker shows whether the frame concerns the active voice or the passive voice of the verb. The passive frames are listed only rarely, as normally they are “derived” from the active frames. The marker occupies one position and currently the following characters are used:

R — active frame

P — irregular passive frame

All frames in example (10) will have the marker R. The missing passives of the verb *učít (matematika je učena, matematika se učí*—mathematics is taught) will be encoded in a frame starting with P.

3.1.2 Reflexivity

The reflexivity abbreviation marks the type of reflexive particle; reflexive pronouns are treated as a value of the grammateme *semantic features* (see below). The possible values are:

-- — no reflexive particle

SE — reflexive tantum with particle *se* or reflexive passive

DE — derived reflexive with particle *se*

se — reflexive with optional particle *se*

SI — reflexive tantum with particle *si*

DI — derived reflexive with particle *si*

si — reflexive with optional particle *si*

The term *reflexive with optional particle* denotes verbs that can occur with or without the reflexive particle in the same meaning, and both these possibilities are grammatical.

3.1.3 Subject

The subject marker points to the member of the frame which is the subject (if the construction has a subject, otherwise this marker is missing). For an active frame, it points to the subject of an active sentence. When a passive frame is derived from the active one, this pointer changes so that it points to the subject of a passive sentence. In a passive frame, this pointer must point to the subject of a passive sentence.

⁴In fact, the variant with Patient realized by Accusative also allows passives, but only if the Addressee is general. We will show later how we encode passive which needs special treatment.

- s[i1] — Actor is the subject
- s[a1] — the subject is raised from Actor's frame

3.1.4 Functor

Functor is a one-character abbreviation of the functor on the tectogrammatical level. All the values are listed in Appendix A. Here we list only abbreviations of inner participants.

- 1 — Actor
- 2 — Patient
- 3 — Addressee
- 4 — Origin
- 5 — Effect
- 0 — no participant; used in frames of raising verbs

3.1.5 Grammatemes

The list of grammatemes determines the morphemes on the morphemic level. There can be several possible surface realizations which are separated by a vertical bar (|). The notation of grammatemes is taken from the source dictionary, but the repertory is enhanced by some features not previously taken into consideration. The grammatemes are given below:

- h — 'semantic' features; their description is given above in Section 2 and in Appendix A. We added the value S for a short reflexive pronoun and we allow grouping of all four nominal values together (hPTSR). More details are discussed below in Section 3.2. Another value which we added is the value Z for pronouns which can stand for a clause, in a sentence. We also added the value G for general participants and E for deleted (empty, erased) deletable modifications in certain secondary frames. Another value which was added is C for the direct speech.
- c — morphological case; possible values are 1, 2, 3, 4, 6, 7
- r — preposition; prepositions are enclosed by curly brackets ({na}, {o}, etc.)

The following grammatemes were added:

- n — number; the values are S and P for singular and plural, respectively. This grammateme was added to the original BRIEF attributes because of the proper treatment of reciprocal verbs (see Section 3.2).
- x — reciprocal coreference; the value points to a functor which is coindexed with the functor containing this grammateme. It was added because of reciprocal verbs.
- a — subject raised to object position; the value points to the embedded clause from which the subject was raised
- q — subject- or object-control
- p — "patient" control
- t — "addressee" control
- d — diatheses of embedded infinitive; the values are identical with values of the "main" frame
- l — required lexeme
- m — modality marker

Their meaning will be explained in the further text.

The whole list of grammatemes is closed in brackets whose shape determines whether the participant (functor) is obligatory, general, obligatory and deletable, or optional:

- [] — obligatory

() — obligatory inner participant which can be realized as general participant, or obligatory and deletable free modification

< > — optional

In FGD, only participants and obligatory free modifications are considered to belong to a verb frame. In practical applications, however, it may be useful to include also free modifications which occur frequently with a given verb. M. Straňáková (Straňáková-Lopatková, 2001) introduced the term *quasi-valency* for such free modifications and we will mark them as optional free modifications. The term *quasi-valency* will be used in one more meaning: it will denote a free modification which only allows some of the surface realizations typical for that free modification. (More exactly, the term *quasi-valency* is to be used in the latter sense, while in the former sense the term *typical* is more appropriate.)

3.1.6 Diatheses

Many of the diatheses, especially passive constructions, are derived regularly, as will be shown in Section 3.3. This is why we do not list all of them in the lexicon but we rather mark single frames with a sign showing which types of diatheses can be derived from the active frame. We adopted special marks for single types of diatheses and we concatenate them to strings.

% — periphrastic passive can be derived

- (11) a. *Nájemníci_{Act} žádají správcovou_{Addr} [o přístup na dvůr]_{Pat}.*
Tenants_{Nom} ask caretaker_{Acc} for access to yard.
- b. *Správcová_{Addr} je (nájemníky_{Act}) žádána [o přístup na dvůr]_{Pat}.*
Caretaker_{Nom} is (tenants_{Ins}) asked for access to yard.

\$ — reflexive passive is possible

- (12) a. *[O tom]_{Pat} právě mluvíme.*
About it_{Loc} just now speak_{1Pl}.
- b. *[O tom]_{Pat} se právě mluví.*
About it_{Loc} SE just now speaks.
'It is being spoken about just now.'

@ — no passive is possible (most reflexives tantum)

- (13) a. *Strašidel_{Pat} se nebojíme.*
Ghosts_{Gen} SE fear_{Neg1Pl}.
'We don't fear ghosts.'
- b. * *Strašidel není báno.*
Ghosts_{Gen} is_{Neg} feared.
- c. * *Strašidla nejsou bána.*
Ghosts_{Nom} are_{Neg} feared.
- d. * *Strašidel se nebojí.*
Ghosts_{Gen} SE fears_{NegSgNeut}.

The sentence (13d) is of course grammatical if we understand it as an active sentence with dropped personal pronoun.

— constructions with *mít* (they are discussed in Section 3.3)

- (14) a. *Maminka slíbila Pěťovi hračku.*
Mummy_{Nom} promised Pěťa_{Dat} toy_{Acc}.
- b. *Pěťa má slíbenou hračku.*
Pěťa_{Nom} has promised_{Prtcp1FemAcc} toy_{FemAcc}.

~ — constructions with *dostat*

- (15) a. *Maminka slíbila Pěťovi hračku.*
 Mummy_{Nom} promised Pěťa_{Dat} toy_{Acc}.
 b. *Pěťa dostal slíbenou hračku.*
 Pěťa_{Nom} got promised_{Prtcpl Fem Acc} toy_{Fem Acc}.
 c. *Učitelka vynadá neposlušným dětem.*
 Teacher_{Nom} scolds disobedient children_{Dat}.
 d. *Neposlušné děti dostanou vynadáno.*
 Disobedient children_{Nom} get_{Fut} scolded.

* — another type of construction with *mít*. Linguists consider this construction to be rather a special verb tense (Hausenblas, 1963) or they include it in a system of aspects (Panevová, 1971). We will discuss this in Section 3.3.

- (16) a. *Kuchařka uvařila oběd.*
 Cook_{Fem Nom} cooked lunch_{Acc}.
 b. *Kuchařka má oběd uvařen.*
 Cook_{Fem Nom} has lunch_{Masc Acc} cooked_{Prtcpl Masc Acc}.
 c. *Kuchařka má uvařeno.*
 Cook_{Fem Nom} has cooked_{Prtcpl Neut}.

The whole frame then looks as in (17):

- (17) a. *akumulovat* R--s [i1] 1(hPTc1)2[hTc4] % \$ (accumulate st)
 b. *kazit~2* RDEs [i1] 1[hTc1] @ (decay)
 c. *přihlásit~1* R--s [i1] 1(hPc1)2[hPTSrC4] A [hTc2r{do} | hTc4r{na}] % \$
 (enroll sb/st where)
 d. *vyhrát~3* R--s [i1] 1(hPc1)2[hTc4] 4<hPc6r{na}> % \$
 (win st of sb)
 e. *tázat* P--s [i3] 1(hPc7)2(sF | sR | hPTc4r{na}) 3[hPc1] (ask)

The frame in (17a) is a frame of a transitive verb. The frame has two participants, Actor (1) and Patient (2). Patient is obligatory ([]), while Actor can be general (()). The Actor is realized as a noun (a person or a thing) in Nominative (hPTc1), Patient is realized as a noun (a thing) in Accusative (hTc4). The subject of an active sentence is Actor (s[i1]). Both periphrastic (%) and reflexive (\$) passives are possible.

The frame in (17b) is a frame of a derived reflexive (DE). The frame contains only obligatory Actor which is realized as a noun (a thing) in Nominative (1[hTc1]). There is no possibility of passive voice (@).

The frame in (17c) is a frame of transitive verb with quasi-valency. The Patient can be also realized by a reflexive pronoun (both short and long form—hPTSr). The quasi-valency is a free modification with the meaning *where*, but not all realizations of this meaning can be applied. For example preposition *pod* (under) plus a noun in Accusative are unacceptable.

In (17d) we can see a frame with obligatory Patient and generalizable Origin.

The frame in (17e) is an example of an irregular passive frame. The generalizable Actor is realized as a noun in Instrumental, Patient as Accusative with the preposition *na* and Addressee as Nominative.

3.2 Reflexivity

In this section, we will look closer on reflexive verbs. In Czech, there is a reflexive pronoun *se* which has several different forms for different cases which can be stressed (long) or unstressed (short). There also exist two reflexive particles which are homonymous with the unstressed reflexive pronoun in Dative (*si*) and Accusative (*se*). In linguistic theory, we distinguish several types of reflexive verbs, but in the lexicon some distinctions will be omitted. We base our work on the taxonomy by K. Králíková (1981) in Table 1, but we will adapt it slightly.

In the lexicon, Dative of possession is not listed as it does not belong to a verb frame (it is treated as a free modification Beneficiary). The reflexive passive belongs among diatheses and will be treated by the respective rules. The “independent category” will be treated as a diathesis as well.

		<i>se</i>	<i>si</i>
<i>se (si)</i> is a complementation of the verb	true reflexive	<i>mýt se</i>	<i>koupit si jízdenku</i>
	reciprocal	<i>milovat se</i>	<i>psávat si</i>
	dative of possession	Ø	<i>držet si klobouk</i>
<i>se (si)</i> changes the meaning of the non-reflexive verb	passive	<i>obilí se mlátí</i>	Ø
	derived lexical meaning	<i>větev se zlomila</i> <i>vrátit se, učít se</i>	<i>zlomit si ruku</i> <i>sednout si</i>
	independent category	<i>ta kniha se dobře čte</i> <i>chce se mi spát</i>	Ø
<i>se (si)</i> is a particle	reflexive tantum	<i>smát se</i>	<i>stěžovat si</i>

Table 1: Taxonomy of reflexive verbs

3.2.1 True reflexive with *se*

True reflexive with *se* is a verbs with reflexive pronoun in Accusative. The pronoun occupies a place of a participant and expresses the coreference of this participant with subject. In most cases it is possible to use the stressed form of the pronoun as well, though the meaning is not fully synonymous.

Some authors doubt about the group of *true reflexive* verbs. It was proposed already by B. Havránek (1928), that *se* in such constructions as *mýt se* (wash self) is not a pronoun (representing a member of a verb frame), but rather a reflexive particle. The group of true reflexive verbs would contain only a couple of constructions like *vidět se v zrcadle* (see oneself in a mirror), *udělat se samostatným* (make oneself independent), etc. This view is supported nowadays by K. Oliva (2000) who shows the behaviour of the particle *se* in opposition with the long form of the pronoun *sebe* and with the short personal pronouns:

- (18) a. $_i$ *Umyl se_i celý_i.*
 — Washed_{3Sg} SE whole_{Nom.}
- b. $_i$ *Umyl sebe_i celého_i.*
 — Washed_{3Sg} self_{Acc} whole_{Acc.}
- c. $_i$ *Umyl ho_j celého_j.*
 — Washed_{3Sg} him_{Acc} whole_{Acc.}

Oliva claims that the verb frames with stressed and unstressed forms of the pronoun *se* are actually two different frames. The verb with unstressed form of the pronoun behaves like reflexive tantum and the pronoun is in fact a particle.

For us, the important criterion is whether the form *se* (or *si*) can be replaced by the long form *sebe* (*sobě*), and whether the constructions with the short reflexive pronouns are similar to constructions with other (short) pronouns. If we adopted the view that *se* is a particle with no representation on the tectogrammatical level we would get two different descriptions of sentences which we consider nearly synonymous. Therefore, we do not go as far as Oliva and still consider the short form to be a pronoun (not a particle), but we are aware of the fact that the short and long forms of the pronoun are not always replaceable and thus, in the lexicon, both possibilities must be explicitly mentioned. We enhanced the repertory of ‘semantic’ features and added the feature S for the short form of the reflexive pronoun. The frames for the verb *umýt* then will have the following form:

- (19) *umýt* R--s [i1]1 (hPc1)2 [hPTrc4]3 (hPSrc3)%\$⁵

⁵The notation in (19) also allows realizations

- (20) a. * *umýt si se*
 wash self_{Dat} self_{Acc}
- b. ? *umýt si sebe*
 wash self_{Dat} self_{Acc}

	<i>koulovat (se)</i> (snowball)	<i>hašteřit se</i> (wrangle)	<i>soutěžit</i> (compete)
	A kouluje B B kouluje A		
reciprocal	AB se koulují	AB se hašteří	AB soutěží
	A a B se koulují	A a B se hašteří	A a B soutěží
	A s B se koulují	A s B se hašteří	A s B soutěží
	A se kouluje s B	A se hašteří s B	A soutěží s B

Table 2: Three types of reciprocal verbs

3.2.2 True reflexive with *si*

True reflexive with *si* is a verb with a reflexive pronoun in Dative. The pronoun occupies the place of a participant and expresses the coreference of this participant with the subject. The reflexive pronoun in Dative also has a short and a long form (*si* and *sobě*), which can be used in the same constructions. The treatment of these verbs is similar to the treatment of true reflexives with *se*.

- (21) a. *Každý_{Act,i} si_{Addr,i} koupí jízdenku_{Pat}.*
 Everyone SI buys ticket.

b. `koupit R--s[i1]1(hPTc1)2[hTc4]3(hPSRc3)%$`

3.2.3 Reciprocal verbs with *se*

Common definition of reciprocal verbs says that a reciprocal verb is a verb with *se*, where the reflexive pronoun has the meaning ‘each other’. Similarly to the situation with true reflexives, Actor is identical with other participant (usually Patient) and the reflexive pronoun expresses this. The difference is that there must be at least two bodies participating in the action and their roles are cross-linked. In fact, there are two actions occurring at the same time, in one of them the participant *i* is Actor and participant *j* is Patient and in the other action the roles are exchanged.

When we examine so called reciprocal verbs closer we discover that there are three types of them. The first type (represented by the verb *koulovat (se)*) was described in the previous paragraph. The second type is reflexive tantum with reciprocal meaning (inherently reciprocal verb). The reciprocal meaning is manifested by obligatory participant with the surface form *s kým* (with whom). The third type is a “plain” verb with reciprocal meaning. The three types are shown in Table 2.

All these types were described by J. Panevová (1999), with a proposal how to encode the information in a lexicon. Her work, however does not suggest structures for sentences with reciprocal verbs. We try to make a proposal of the structures and we will compare them to structures proposed in (Hajičová, Panevová, and Sgall, 2000). Our proposal is shown in (22):

- (22) a.
$$\begin{array}{c} \textit{koulovat} \\ \hline \textit{chlapci}_{Act,Pat,RECP} \quad \textit{se}_{Act,Pat,RECP} \end{array} \quad \textit{Chlapci se koulují.}$$
 Boys_{Nom} SE snowball.

c. ? *umýt se sobě*
 wash self_{Acc} self_{Dat}

d. ? *umýt sobě sebe*
 wash self_{Dat} self_{Acc}

which can be handled by a general rule of grammar saying that two (short) reflexive pronouns cannot occur as realizations inside one verb frame.

- b.
$$\begin{array}{c} \text{koulovat} \\ \text{COORD}_{Act,Pat,RECP} \quad \text{se}_{Act,Pat,RECP} \\ \text{Petr} \quad \text{Pavel} \end{array}$$
- Petr a Pavel se koulují.*
Petr_{Nom} and Pavel_{Nom} SE koulují.
Petr s Pavlem se koulují.
Petr_{Nom} with Pavel_{Ins} SE snowball.
- c.
$$\begin{array}{c} \text{koulovat_se} \\ \text{Petr}_{Act,RECP} \quad \text{Pavel}_{Pat,RECP} \end{array}$$
- Petr se kouluje s Pavlem.*
Petr_{Nom} SE snowballs with Pavel_{Ins}.
- d.
$$\begin{array}{c} \text{koulovat_se} \\ \text{Petr}_{Act,RECP} \quad \text{GNRL}_{Pat,RECP} \end{array}$$
- ? *Petr se kouluje _GNRL.*
Petr_{Nom} SE snowballs.

It may be surprising that *se* is treated as a pronoun in (22a) and (22b), and as a particle in (22c) and (22d). This is a result of applying the criteria which we used already for the true reflexives:

- (23) a. *Chlapci koulují sebe (navzájem).*
 Boys snowball self (each other).
 ‘Boys snowball each other.’
- b. *Petr a Pavel koulují sebe (navzájem).*
 Petr and Pavel snowball self (each other).
- c. *Petr s Pavlem koulují sebe (navzájem).*
 Petr with Pavel snowball self (each other).
- d. **Petr kouluje sebe s Pavlem.*
 Petr snowballs self with Pavel.

In (23a)–(23c), the short form of the pronoun *se* is replaced by the long form *sebe* (which indicates that it is really a pronoun), while in (23d) this replacement is not possible (which indicates that *se* is a particle).

We decided to introduce a new grammateme *x* whose value is the coreferential functor. In example (24b), there is Actor in plural,⁶ and Patient is realized by the reflexive pronoun *se* (both short and long form). It is marked as reciprocally coreferential with Actor. In example (24c), Patient has morphological realization by Instrumental+*s*, and it is also reciprocally coreferential with Actor.

- (24) a. R--s[i1]1[hPc1]2[hPc4]%% (koulovat)
- b. R--s[i1]1[hPc1nP]2[hSRc4x1]@ (koulovat se)
- c. RDE1[hPc1]2(hPRc7{s}x1)@ (koulovat se)

The frame in (24b) corresponds to the sentences (22a), (22b) and (23a)–(23c). The frame in (24c) corresponds to the sentences (22c) and (22d).

The frames of of inherently reciprocal verbs will have only two forms corresponding to (24b) and (24c). They are not listed here, but readers can find them in (Skoumalová, 2001).

From the above description it follows that there is no need to introduce a new mark for “reciprocal” *se* as it is possible to use other markers.

3.2.4 Reciprocal verbs with *si*

Reciprocal verbs with *si* are similar to reciprocal verbs with *se* and they are treated similarly. The difference is that the functors assigned to the participants of the action are Actor and Addressee. The types of reciprocal verbs with *si* are shown in Table 3. In (25), we show frames of reciprocal verbs with *si*.

⁶The plural here means *semantic* plural, not grammatical. It can be realized as a noun in plural, or as a coordination or as a noun with meaning of a group, e.g. *třída* (class).

	<i>povídat (si)</i> (chat, imperf.)	<i>popovídat si</i> (chat, perf.)
	A povídá B o ... B povídá A o ...	
reciprocal	AB si povídají o ...	AB si popovídají o ...
	A a B si povídají o ...	A a B si popovídají o ...
	A s B si povídají o ...	A s B si popovídají o ...
	A si povídá s B o ...	A si popovídá s B o ...

Table 3: Reciprocal verbs with *si*

- (25) a. R--1(hPc1)2[hTc6r{o}]3(hPc3)\$ (povídat)
 b. R--1[hPc1nP]2[hTc6r{o}]3[hSRx1]@ (povídat si)
 c. RDI1[hPc1]2[hTc6r{o}]3(hPc7{s}x1)@ (povídat si)
 d. RSI1[hPc1nP]2[hTc6r{o}]3[x1]@ (popovídat si)
 e. RSI1[hPc1]2[hTc6r{o}]3(hPc7{s}x1)@ (popovídat si)

3.2.5 Reflexive tantum with *se*

Reflexive tantum with *se* is a verb which has an obligatory reflexive particle *se*. This particle has no representation on the tectogrammatical level.

- (26) a. *Helena_{Act} se směje všemu_{Pat.}*
Helena_{Nom} SE laughs everything_{Dat.}
 ‘Helena laughs at everything.’
 b. *František_{Act} se nebojí ničeho_{Pat}*
František_{Nom} SE fears nothing_{Gen.}
 ‘František is not afraid of anything.’

Frames of verbs from above examples will look as follows:⁷

- (27) a. smát RSE1[hPc1]2<hPTRc3>@
 b. bát RSE1[hPc1]2(hPTRc2|sD|sF|sU)@

3.2.6 Derived reflexive verbs with *se*

This category contains verbs which behave like reflexive tantum but they have origin in true reflexive verbs. Their lexical meaning, however, changed so that they cannot be understood as true reflexives any more. For example the verb *rozčítit se* (get angry) could be understood as true reflexive, as it is possible to say *rozčítit sám sebe* (make angry oneself), but the meaning is different (as the translation also shows). Beside it, the verb *rozčítit se* has only Actor in its frame, while *rozčítit koho/sebe* has Actor, Patient and Addressee. The verb *rozčítit* then will have two meanings with two frames, as shown in example (28).

- (28) a. rozčítit~1 R--1[hPTc1]2(hTc7)3[hPTRc4]@
 b. rozčítit~2 RDE1[hPc1]@

⁷In the frame of the verb *bát se* the realization by infinitive is missing. This is because the infinitive needs special treatment—raising or control must be marked. This will be discussed in Section 3.4 and thus we did not want to obscure this example.

3.2.7 Reflexive tantum with *si*

Reflexive tantum with *si* is a verb which has an obligatory reflexive particle *si*.

- (29) a. *Nájemníci_{Act} si stěžují [na správčovou]_{Pat}*
Tenants_{Nom} SI complain about caretaker.
b. *stěžovat* RSI[hPc1]2[hPTRc4r{na}]@

3.2.8 Derived reflexive verbs with *si*

This category is similar to derived reflexive verbs with *se*.

- (30) a. *Děti_{Act} si hrají [na indiány]_{Pat}*
Children_{Nom} SI play at indians_{Acc}.
b. *hrát* RSI[hPc1]2<hPTRc4r{na}>@

3.2.9 Reflexive with optional *se*

This is a verb with reflexive particle *se* which is not obligatory. It is usually true for such verbs that the reflexive particle is optional for some meanings, and obligatory or impossible for others.

- (31) a. *Na co_{Pat} (se) koukáš?*
On what_{Acc} (SE) look_{2Sg}?
'What are you watching?'
b. *koukat~1* Rse1[hPc1]2(hPTRc4r{na})\$
c. *Kouká ti_{Ben} podolek_{Act}.*
Looks you_{Dat} shirt-tail_{Nom}.
'Your shirt-tail is showing.'
d. * *Kouká se ti podolek.*
Looks SE you_{Dat} shirt-tail_{Nom}.
e. *koukat~2* R--1[hTc1]@

3.2.10 Reciprocal verb with optional *se*

Some of the the reflexive verbs with optional *se* can also be inherently reciprocal.

- (32) a. *Vy už (se) spolu nekamarádíte?*
You_{2PlNom} already (SE) together hobnob_{Neg}?
b. *Já (se) s Jirkou kamarádím!*
I_{Nom} (SE) with Jirka_{Ins} hobnob.
c. *kamarádit* Rse1[hPc1nP]2[x1]@
d. *kamarádit* Rse1(hPc1)2[hPc7{s}x1]@

3.2.11 Reflexive with optional *si*

Reflexive with optional *si* is a verb with reflexive particles *si* which is not obligatory.

- (33) a. *Aleš_{Act} (si) myslí, [že Jiřina nepřijde]_{Pat}.*
Aleš_{Nom} (SI) thinks that Jiřina comes _{FutNeg}.
b. *Aleš_{Act} si to_{Pat} nemyslí.*
Aleš_{Nom} SI it_{Acc} thinks_{Neg}.
c. *Co_{Pat} (si) myslí Aleš_{Act}?*
What_{Acc} (SI) thinks Aleš_{Nom}?

In example (33) we can see that the verb *myslet si* does not require the particle obligatorily if it is complemented by a clause. It requires the particle, however, if the complementation is realized by a pronoun.

On the other hand, the particle *si* cannot occur if we use the verb in its intransitive meaning or in the meaning ‘have in mind’.

- (34) a. *Myslím, tedy jsem.*
 Think_{1Sg}, then am.
 ‘Cogito, ergo sum.’
- b. * *Myslím si, tedy jsem.*
 Think_{1Sg} SI, then am.
- c. *CoPat tím myslíš?*
 What_{Acc} it_{Ins} think_{2Sg}?
 ‘What do you mean by it?’
- d. * *CoPat si tím myslíš?*
 What_{Acc} SI it_{Ins} think_{2Sg}?

The verb *myslet (si)* then will need several frames which will express the behaviour of the particle *si*.

- (35) a. *myslet~1* Rsi1[hPc1]2[sD]@
 b. *myslet~1* RSI1[hPc1]2[hZc4]@
 c. *myslet~2* R--1(hPc1)\$
 d. *myslet~3* R--1(hPc1)2[hZc4|sD]I(hTc7)&
 e. *myslet~4* R--1(hPc1)2[hPTc4r{na}]&

3.2.12 Reflexive passive

Reflexive passive is a construction with the particle *se*. It is one of the possible passive constructions in Czech. This construction is usually derived from the basic active frame and therefore the passive frames are not listed in the lexicon separately.

- (36) *Brána se zavírá v devět hodin.*
 Gate SE closes at nine o'clock.

This construction will be discussed in detail in Section 3.3.2.

3.2.13 Mediopassive

Mediopassive constructions are a kind of reflexive passive and they will be described later in Section 3.3.

- (37) *Z této látky se šije dobře.*
 From this fabric_{Gen} SE sews well.
 ‘This is good fabric for sewing.’

In our lexicon these constructions will be treated as reflexive passives. The discussion about this type of construction follows in Section 3.3.2.

3.2.14 Ambiguity of reflexive verbs

Very often, reflexive verbs have several meanings; they appear as a true reflexive, reciprocal verb, derived reflexive verb, reflexive tantum or reflexive passive. We can find tests which help the lexicographers to discover all possible meanings of a certain reflexive verb. The lexicon should contain all the variants, even though they may cause ambiguity in syntactic analysis or other application. The disambiguation very often depends on semantics of the participants, and so we cannot formulate syntactic constraints which would solve it.

3.3 Diatheses

Another lexical information useful for language processing is the information about diatheses. The most important marked diatheses are passive constructions. In Czech there exist two syntactic constructions with passive meaning: the periphrastic passive formed by an auxiliary verb *být* (be) and passive participle, and reflexive passive formed by indicative and the reflexive particle *se*. As both these passives are derived regularly from the active constructions, we will only list the information of what *type* of passive is acceptable for a certain verb and its frame, and we will not list all the passive constructions in our lexicon. Of course, there are exceptions—passive constructions which are derived by exceptional rules—such passives must be listed explicitly (but there will be only single cases of such passive constructions).

Beside periphrastic and reflexive passive, there exist also other types of diatheses which we consider regular. For example, constructions with support verbs *dostat* (get) and *mít* (have) are very frequent. The possibility of marking these types of diatheses in the lexicon will also be discussed.

In our lexicon, we only consider such derived constructions in which the surface syntactic structure is changed. Such constructions as

- (38) a. *Bolest probudila Pavla.*
Pain_{Nom} woke Pavel_{Acc}.
b. *Marie probudila Pavla.*
Marie_{Nom} woke Pavel_{Acc}.

differ in the semantics of subjects. In (38a), the subject has the semantic role of Causer, while in (38b), the subject is Agent (Daneš, Hlavsa, and others, 1987; Štícha, 1984; Grepl and Karlík, 1998). In the FGD approach, however, both subjects are Actors. Both the constructions are identical on the surface level and they only differ in the lexical setting of the subject.

In this article, we do not provide the full discussion with other authors, we refer the reader to Skoumalová (2001).

3.3.1 Periphrastic passive

The verb is in the form of the periphrastic passive, the predicate agrees with subject in person, gender and number:

- (39) a. *Petr čte knihu.*
Petr_{Nom} reads book_{Acc}.
b. *Knih je čtena.*
Book_{FemNom} is read_{PrtcplFemSg}.

This construction is usually formed from transitive verbs (i.e. verbs with object in Accusative), but there are exceptions. Not all transitive verbs can be passivized (e.g. *mít* ‘to have’, *dostat* ‘to get’, etc.), and on the other hand, some verbs without an Accusative object can form passive:

- (40) a. *Úřad vyhověl jeho žádosti.*
Office_{Nom} granted his application_{Dat}.
‘The office granted his application.’
b. *Jeho žádosti bylo (úřadem) vyhověno.*
His application_{FemDat} was (by office_{Ins}) granted_{PrtcplSgNeut}.
‘His application was granted (by the office).’

The subject slot of the passive construction is either filled by the original Accusative object (typically Patient), or it is empty (if the active construction did not contain any Accusative). In the case when the subject is empty or it is a clause (finite or non-finite) the verb shows agreement with neuter singular.

The original subject (Actor) changes its case to Instrumental; Actor in these sentences can be general, and thus it can be omitted on the surface level.

- (41) a. *Knih byla napsána slavným autorem.*
Book_{Nom} was written famous author_{Ins}.
‘The book was written by a famous author.’

- b. *Bazén byl vypuštěn.*
Swimming pool was emptied.

There is another possible surface form of Actor: the prepositional phrase *od* (from) + Genitive, but this form cannot be used with all verbs—here, the semantics of the verb and its participants plays a role:

- (42) a. *Pepík je bit od otce.*
Pepík is beaten from father.
- b. **Kniha byla napsána od slavného autora.*
Book was written from famous author.

The conditions in which this construction can be used will be examined in the future work. Here, we assume that Actor can only change to Instrumental.

Before we start describing the algorithm, we have to make one more important remark: when we speak about a change of the structure we always work with an *instance* of a verb frame. The verb frame is an abstract set of all possible realizations, and we can only make a diathesis of a chosen member of this set.

The algorithm for deriving the frame of the periphrastic passive is described here:

The verb form changes to periphrastic passive.

If there is a nominal object in Accusative in the frame, it becomes subject (in Nominative). The subject marker changes so that it pointed to the new subject.

If the object in Accusative is a clause or the infinitive, it becomes the subject, with a special kind of agreement (3rd person, singular, neuter).

If there is no object in Accusative the passive has empty subject, with the same kind of agreement as the infinitive or clause subject. The subject marker is deleted.

If our frame instance contains only the subject on the surface, this type of passivization is prohibited.

The original subject becomes a generalizable member which is realized by Instrumental.

All other members of the frame stay in their positions.

There are some exceptions to the above rules. The first group of exceptional verbs are ditransitive verbs (verbs with two Accusatives in the frame). We have found only two such verbs in Czech:

stát koho co - to cost sb sth

This verb does not have the passive.

učit koho co_{acc}/čemu_{dat} - to teach sb sth

If we choose the frame with Accusative and Dative, no problems occur. But in the frame with two Accusatives, one of them must be omitted (both can be generalized) before we create the passive construction:

- (43) a. *Děti jsou učeny (matematice).*
Children are taught (to mathematics_{Dat}).
- b. **Děti jsou učeny matematiku.*
Children are taught mathematics_{Acc}.
- c. *Matematika je učena.*
Mathematics is taught.
- d. **Matematika je učena děti.*
Mathematics is taught children_{Acc}.

The periphrastic passive is marked by % in the lexicon, and the entries of the verb *učit* will look as follows:

- (44) a. *učit*~2 R--s[i1]1(hPc1)2(hTc3)3(hPc4)%\$

- b. učít~2 R--s[i1]1[hPc1]2[hTc4]3(hPc4)@
 c. učít~2 P--s[i2]1(hPc7)2[hTc1]@
 d. učít~2 PSEs[i2]1(hG)2[hTc1]@

Another exceptional group of verbs are reflexives tantum which can have passive forms. The member of the frame which undergoes the change into subject is not a member in Accusative but in Genitive:

- (45) a. *Soudce se tázal svědka, zda něco viděl.*
 Judge_{Nom} SE asked witness_{Gen} if he saw anything.
 b. *Svěděk byl (soudcem) tázán, zda něco viděl.*
 Witness_{Nom} was (judge_{Ins}) asked if he saw anything.

This group of verbs is not very numerous. It contains verbs *tázat se* (and its prefixed variants), *obávat se*, and perhaps some more. It is a question whether we should introduce new rules for this type of passive or rather store these passive frames as exceptions:

- (46) a. tázat RSEs[i1]1[hPc1]2[sF|sR|hPTZc4{na}]3(hPTc2)@
 b. tázat P--s[i3]1(hPc7)2(sF|sR|hPTZc4{na})3[hPTc1]@

The periphrastic passive is felt as rather formal, bookish or obsolete in modern Czech, especially the passive with expressed Actor. Unlike its English counterpart, Czech passive is rarely used for changing the topic-focus articulation—for this purpose the change of the word order is employed. The passive construction is mainly used, if the speaker wants to avoid saying who/what Actor is, or if Actor is general. In both these cases, however, the reflexive passive is used more often.

3.3.2 Reflexive passive

In this construction, the verb changes its form to reflexive passive form, the participant in Accusative (if present) becomes the subject, and Actor becomes general.

- (47) a. *Bábovka se peče.*
 Cake SE bakes.
 ‘The cake is being baked.’
 b. *Do města se jde tudy.*
 To town SE goes this way.
 ‘This is the way to the town.’

The example in (47a) is the real reflexive passive, derived from a transitive verb, while the sentence in (47b) is an impersonal active construction, derived from an intransitive verb. We mark both these constructions as reflexive passive as the algorithms for deriving them are very similar.

The reflexive passive is sometimes indistinguishable from the intrinsic or true reflexive. The sentence

- (48) *Děti se učí dobře.*
 Children SE teach well.
 ‘Children are easy to teach.’ or ‘The children learn well.’

has two readings, as the verb *učit* ‘to teach’ in reflexive passive has the same form as the reflexive verb *učit se* ‘to learn’. This ambiguity is inherent in the language and we will not try to solve this problem in the lexicon.

The algorithm for deriving the reflexive passive frame is nearly identical with the algorithm for the periphrastic passive:

The verb changes its form to a reflexive passive form.

If there is a nominal object in Accusative in the frame, it becomes subject (in Nominative). The subject marker is changed so that it points to the new subject.

If the object in Accusative is a clause or the infinitive, it becomes the subject, with a special kind of agreement (3rd person, singular, neuter).

If there is no object in Accusative the passive has an empty subject, with the same kind of agreement as the infinitive or clause subject. The subject marker is deleted.

The original subject is generalized (and thus omitted on the morphemic level).

All other members of the frame stay in their positions.

The rules for handling the ditransitive verbs *stát* ‘to cost’ and *učit* ‘to teach’ are the same as at the periphrastic passive: *stát* cannot be passivized and with the verb *učit*, the frame to be passivized can contain only one Accusative (see 44).

- (49) a. *Děti se učí (matematice).*
 Children_{Nom} SE teach (to mathematics_{Dat}).
- b. **Děti se učí matematiku.*
 Children SE teach mathematics_{Acc}.
- c. *Matematika / Matematika se učí od první třídy.*
 Mathematics_{Nom/Dat} SE teaches from first grade.
- d. **Matematika / Matematika se učí děti.*
 Mathematics_{Nom/Dat} SE teaches children_{Acc}.

The reflexive passive of *učit*, however, is homonymous with the reflexive verb *učit se* ‘to learn’, and thus it is difficult for a Czech speaker to understand the examples in (49a) and (49b) in the passive meaning. As an active sentence with the verb *učit se*, (49b) is correct.

Reflexive passive is marked by \$ in the lexicon and an example of a lexical entry was given in (44a).

For the proper treatment of the verb *učit* we also have to add an irregular frame for the reflexive passive:

- (50) *učit~2* PSEs[i2]1(hG)2[hTc1]@

The reflexive passive is used especially in cases when Actor is general and the periphrastic passive cannot be used:

- (51) a. *Tady se hodně čte.*
 Here SE much reads.
 ‘Here, people read a lot.’
- b. **Tady je hodně čteno.*
 Here is much read_{Prtcpl}.
- c. *Matematice se učí od první třídy.*
 Mathematics_{Dat} SE teaches from first grade.
- d. ?*Matematice je učeno od první třídy.*
 Mathematics_{Dat} is taught from first grade.

3.3.3 Mediopassive

This construction is very similar to the previous one—some linguistic books actually do not distinguish between them. In mediopassive, Actor is present (though it can be general) and an adverb like *dobře* (well), *špatně* (badly), *snadno* (easily), etc. (i.e. free modification of Manner), must be present in the construction. This type of passive was described by M. Dokulil (1941) as a special case of description of the way something is done. P. Karlík (1995) considers this construction a special case of the subject diathesis of the type agent–patient where the agentive role is put to the background and the agent is getting a role of Experiencer.

Examples:

- (52) a. *Matematika se mi učí snadno.*
 Mathematics_{Nom} SE me_{Dat} learns easily.
 ‘It’s easy for me to learn/teach mathematics.’

- b. *Z této látky se šije dobře.*
 From this fabric SE sews well.
 ‘It’s easy (for anyone) to make clothes from this fabric.’

The algorithm for deriving the mediopassive frame is nearly identical with the algorithm for the periphrastic passive:

The verb form is changed in a reflexive passive form.

If there is a nominal object in Accusative in the frame, it becomes subject (in Nominative). The subject marker is changed so that it points to the new subject.

If the object in Accusative is a clause or the infinitive, it becomes the subject, with a special kind of agreement (3rd person, singular, neuter).

If there is no object in Accusative the passive has an empty subject, with the same kind of agreement as the infinitive or clause subject. The subject marker is deleted.

The original Actor (subject) changes its surface realization to Dative.

All other members of the frame stay in their positions.

We do not introduce a separate mark for the possibility of deriving mediopassive as we believe that there is a general rule: any frame of an imperfective verb which can be transformed to reflexive passive can also be transformed to mediopassive. The information on reflexive passives is contained in our lexicon, and the information on aspect is contained in the morphological lexicon. If it turned out that the above rule does not hold we can introduce a new mark.

There is, however, a verb that needs special treatment: the verb *chtít* can have a reflexive form *chtít se* where Actor has the form of Dative. We will call this construction mediopassive, but it requires a separate entry in the lexicon. As this verb requires an infinitive in its frame we will show the encoding of the frame in Section 3.4.2.

3.3.4 Constructions with *mít* and *dostat*

In this type of construction, a Dative member of the frame (typically Addressee) becomes the subject of a construction with the support verb *mít* or *dostat* and the main verb occurs in the predicate as a passive participle in Accusative. If the main verb has an Accusative object (typically Patient), the participle agrees with it in gender and number. If the Accusative object is missing, the participle has the form of singular neuter. Actor (the original subject) becomes an optional member of the frame in the form of *od* + Genitive:

- (53) a. *Obec přidělila žadatelům byty.*
 Municipality_{Nom} granted applicants_{Dat} flats_{Acc}.
- b. *Žadatelé mají/dostali (od obce) přiděleny byty.*
 Applicants_{Nom} have/got (from municipality) granted_{PrtcplAcc} flats_{Acc}.
 ‘Applicants were granted flats (by municipality).’
- c. *Otec vynadá Pepíkovi.*
 Father_{Nom} will scold Pepík_{Dat}.
- d. *Pepík dostane vynadáno (od otce).*
 Pepík_{Nom} will get scolded (from father_{Gen}).
 ‘Pepík will be scolded (by the father).’
- e. *Vnučka babičce uvařila.*
 Granddaughter_{Nom} grannie_{Dat} cooked.
 ‘Granddaughter has cooked for grannie.’
- f. *Babička má uvařeno.*
 Grannie_{Nom} has cooked_{PrtcplNeutSg}.
 ‘(The meal) has been cooked for grannie.’

Some verbs allow both of the two support verbs, while others allow only one of them (*mít/dostat přiděleno, dostat/*mít vynadáno, *dostat/mít uvařeno*). This is why we introduced two marks—one for each of the support verbs.

The algorithm for deriving the verb frame of this construction follows:

An object in Dative (Addressee, Patient, or Beneficiary) becomes subject (in Nominative). The subject marker is changed accordingly.

Actor becomes an optional member of the frame of the form *od* + Genitive.

An object in Accusative (if exists) can become general.

All other members of the frame stay in their positions.

Frames of the verbs which allow this diathesis are in the following example (the diathesis with *mít* is marked by # and the diathesis with *dostat* is marked by *):

- (54) a. přidělit R--s[i1]1(hPc1)2[hPTc4]3[hPc3]%%\$#*
 b. vynadat R--s[i1]1(hPc1)2[hPc3]%%\$*
 c. uvařit R--s[i1]1(hPc1)2[hPTc4]3<hPc3>%%\$#~

3.3.5 Resultative construction with *mít*

There is one more construction with the support verb *mít*. This is not really a diathesis, as Actor remains as subject and the change on the surface only affects the verb form. It is rather a kind of resultative tense, which corresponds to English perfective constructions. K. Hausenblas (1963) ranks this construction to verb tense, while J. Panevová (1971) considers it a kind of aspect. We decided to include this construction among other diatheses because they are derived regularly and we have no other means how to create these constructions.

- (55) a. *Upeču bábovku.*
 Bake_{1SgFut} cake_{FemAcc}.
 b. *Bábovku už mám upečenu/upečenou.*
 Cake_{FemAccSg} already have_{1Sg} baked_{PrtcplFemAccSg / AdjFemAccSg}.
 c. *Už mám upečeno.*
 Already have_{1Sg} baked_{PrtcplNeutAccSg}.

In this derivation, the frame remains the same as in the base form. The only operation in forming this construction is changing the predicate.

All the above constructions can only be derived from perfective verbs, as they express a result.

This diathesis is marked by ~, and an example of a verb frame allowing this diathesis is in (54c).

3.4 Verbs with the infinitive in their frames

For this group of verbs, we have to describe not only the frame of the verb, but also the interaction between the higher verb and the lower verb (the infinitive)—which members of the frames they share, what kinds of derived frames are allowed for both the infinitive and the governor, and other constraints that hold for both the verbs.

These verbs are usually divided into two subclasses: **raising** and **equi** (or **control**) verbs. In both cases the subject (or rarely an object) of the infinitive is the subject or an object of the higher verb, but there is a difference between the two deep structures.

Raising verb: The subject of the infinitive becomes (is raised as) the subject or an object of the governor, but it does not belong to the governor's frame.

Equi verb: Certain participant of the governor is coreferential with a participant of the dependant. On the surface level, such a participant is present only once, but in the deep structure, there are two slots (one in every verb's frame) which are coreferential.

Many authors were concerned with these kinds of verbs; this topic is worked up well for English (Chomsky, 1986; Dalrymple et al., 1995; Pollard and Sag, 1994), for Czech, we will proceed from Paněvová (1996) but our conclusions will be different in some cases.

First we will show the difference between the two types of verbs in examples of tree structures. We will use one raising verb (*zdát se*—seem) and one equi verb (*snažit se*—try) for explanation.

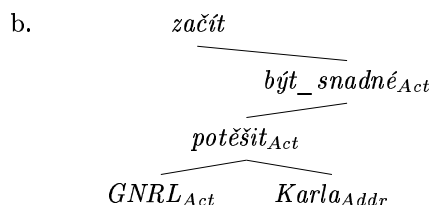
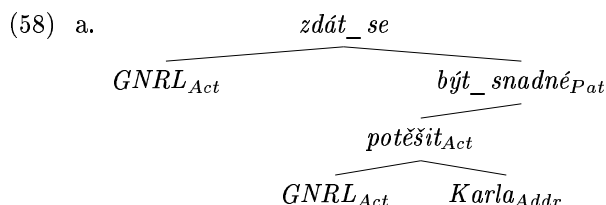


For English, certain tests were developed which should show whether a verb is raising or equi. We will try to find similar tests for Czech.

Raising verbs can have an infinitive in the subject position.

- (57) a. *Potěšit Karla se zdá být snadné.*
'To please Karel seems to be easy.'
- b. *Potěšit Karla musí být snadné.*
'To please Karel must be easy.'
- c. *Potěšit Karla začne být snadné.*
'To please Karel starts to be easy.'
- d. **Potěšit Karla zkouší být snadné.*
'To please Karel tries to be easy.'
- e. **Potěšit Karla chce být snadné.*
'To please Karel wants to be easy.'

In the above examples, the raising verb *zdát se* has two participants (general Actor in Dative and Patient), and the verb *začít* has only one participant (Actor); the subject of the upper verb is raised from the frame of the construction *být snadné* (be easy). Let us show it on graphs of sentences (57a) and (57c):



Modal verbs like *muset* (must) are treated as grammatemes in FGD and therefore they cannot have any participants. A trivial corollary of this fact is that the subject of the modal verb *must* be participant of the infinitive.

- (59) a. $\begin{array}{c} \text{být_snadné}_{\text{debitiv e}} \\ \diagup \quad \diagdown \\ \text{potěšit}_{\text{Act}} \\ \diagup \quad \diagdown \\ \text{GNRL}_{\text{Act}} \quad \text{Karel}_{\text{Addr}} \end{array}$ *Potěšit Karla musí být snadné.*
Please_{Inf} Karel must be_{Inf} easy.
- b. $\begin{array}{c} \text{být}_{\text{debitiv e}} \\ \diagup \quad \diagdown \\ \text{Jirka}_{\text{Act}} \quad \text{veselý} \end{array}$ *Jirka musí být veselý.*
Jirka must be_{Inf} merry.

Raising verbs can have a clausal subject, while equi verbs cannot:

- (60) a. *Že je chytrý, se zdá být zřejmé.*
'That he is clever seems to be obvious.'
- b. *Že je chytrý, musí být zřejmé.*
'That he is clever must be obvious.'
- c. *Že je chytrý, začne být zřejmé.*
'That he is clever starts to be obvious.'
- d. **Že je chytrý, zkouší být zřejmé.*
'That he is clever tries to be obvious.'
- e. **Že je chytrý, chce být zřejmé.*
'That he is clever wants to be obvious.'

The content of a passive sentence is the same as of the active one for raising verbs.⁸

- (61) a. *Doktor musí vyšetřit babičku.*
Doctor must examine_{Inf} grannie.
- b. *Babička musí být vyšetřena od doktora.*
Grannie must be examined from doctor.
- c. *Doktor se pokusil vyšetřit babičku.*
Doctor SE tried examine_{Inf} grannie.
- d. *Babička se pokusila nechat se vyšetřit od doktora.*
Grannie SE tried let_{Inf} SE examine_{Inf} from doctor.
- e. *Karel viděl doktora vyšetřit babičku.*
Karel saw doctor examine_{Inf} grannie.
- f. *Karel viděl babičku nechat se vyšetřit od doktora.*
Karel saw grannie let_{Inf} SE examine_{Inf} from doctor.
- g. *Karel nařídil doktorovi vyšetřit babičku.*
Karel ordered doctor examine_{Inf} grannie.
- h. *Karel nařídil babičce nechat se vyšetřit od doktora.*
Karel ordered grannie let_{Inf} SE examine_{Inf} from doctor.

Another test, which can be applied, checks the number of participants of the upper verb, and their surface realization. This number should not depend on the lexical setting of the infinitive. And also the surface realization of a certain participant should not depend on the lexical setting of another participant. If we considered for example that the verb *začít* (start) is an equi verb whose subject is coreferential with the subject of the embedded infinitive we would need several frames:

- (62) a. $\begin{array}{c} \text{začít} \\ \diagup \\ \text{pršet}_{\text{Act}} \end{array}$ *Začalo pršet.*
Started rain_{Inf}.

⁸It is impossible to use periphrastic passive in certain constructions, but we can paraphrase the passive construction by using a support verb *nechat* (let), and we can understand the construction as a kind of passive.

- b.
$$\begin{array}{c} \text{začít} \\ \diagup \quad \diagdown \\ \text{Tomáš}_{Act,i} \quad \text{pracovat}_{Pat} \\ \quad \quad \quad \diagup \\ \quad \quad \quad \text{COR}_{Act,i} \end{array}$$
 Tomáš začal pracovat.
Tomáš started work_{Inf}.
- c.
$$\begin{array}{c} \text{začít} \\ \diagup \quad \diagdown \\ \text{pršet}_{Act,i} \quad \text{být_jasné}_{Pat} \\ \quad \quad \quad \diagup \quad \diagdown \\ \quad \quad \quad \text{COR}_{Pat,i} \quad \text{všichni}_{Act} \end{array}$$
 Že prší, začalo být jasné všem.
That rains started be_{Inf} clear all_{Dat}.

The verb from (62a) would have a frame with Actor realized as an infinitive. The verb from (62b) would have a frame with Actor realized by a noun in Nominative and Patient realized by an Infinitive. The verb from (62c) would have also Actor and Patient in its frame, but Actor would be realized by a clause attached by *že*. We can see that we could continue and find even more different frames for the equi verb *začít*. On the other hand, if we suppose that the verb *začít* is a raising verb we get rid of the problem with many frames. The frame only contains Actor (the infinitive) and the subject is raised from Actor's frame. It can be even empty if the infinitive has no subject.

- (63) a.
$$\begin{array}{c} \text{začít} \\ \diagdown \\ \text{pršet}_{Act} \end{array}$$
 Začalo pršet.
Started rain_{Inf}.
- b.
$$\begin{array}{c} \text{začít} \\ \diagdown \\ \text{pracovat}_{Act} \\ \diagdown \\ \text{Tomáš}_{Act} \end{array}$$
 Tomáš začal pracovat.
Tomáš started work_{Inf}.
- c.
$$\begin{array}{c} \text{začít} \\ \diagdown \\ \text{být_jasné}_{Act} \\ \diagdown \quad \diagup \\ \text{pršet}_{Pat} \quad \text{všichni}_{Act} \end{array}$$
 Že prší, začalo být jasné všem.
That rains started be_{Inf} clear all_{Dat}.

We used a similar consideration for the so-called Slavic Accusative in sentences with verbs of perception. We believe that sentences in (64) have identical content:

- (64) a. *Petr viděl doktora vyšetřit babičku.*
Petr saw doctor_{Acc} examine_{Inf} grannie_{Acc}.
- b. *Petr viděl doktora, jak vyšetřuje babičku.*
Petr saw doctor_{Acc} how examines grannie_{Acc}.
- c. *Petr viděl, jak doktor vyšetřuje babičku.*
Petr saw how doctor_{Nom} examines grannie_{Acc}.

The verb *vidět* has only two participants, in our model, and the above sentences could be expressed by the structure in (65):

- (65)
$$\begin{array}{c} \text{vidět} \\ \diagup \quad \diagdown \\ \text{Petr}_{Acc} \quad \text{vyšetřit}_{Pat} \\ \quad \quad \quad \diagup \quad \diagdown \\ \quad \quad \quad \text{doktor}_{Act} \quad \text{babička}_{Pat} \end{array}$$

Now, we have tools for judging equi and raising verbs and we can start describing single lexical entries.

3.4.1 Raising verbs

First, we will deal with **subject raising verbs**. This group of verbs contains mainly the modal and aspectual verbs. As it was said above, modal verbs are considered grammatememes in FGD and thus they cannot have own argument structure. On the surface level, however, they impose certain constraints on the infinitives. These constraints must be encoded in the lexicon and that is why we introduce lexical entries for these verbs.

In examples in (66) we show various constructions of raising verbs; the members of the infinitival clauses are enclosed in brackets and the trace of the raised element is marked by an underscore.

- (66) a. *Petr_{Act,i} smí []_i odejít].*
Petr_{Nom} may leave_{Inf}.
- b. *Začalo [pršet].*
Started rain_{Inf}.
 'It started raining.'
- c. *Petr_{Pat,i} musí []_i být pochválen].*
Petr_{Nom} must be_{Inf} praised_{Prtpcl}.
- d. *Musí []_i se zabít] dvě mouchy_{Pat,i} [jednou ranou].*
Must SE kill_{Inf} two flies_{Nom} one hit_{Ins}.
 'Two flies must be killed by one hit.'
- e. *Bábovka_{Pat,i} [se] začala []_i péci].*
Cake_{Nom} SE started bake_{Inf}.
 'The cake started to be baked.'
- f. *Únosce_{Addr,i} musí []_i dostat slíbeno výkupné_{Pat}].*
Kidnapper_{Nom} must get_{Inf} promised_{Prtpcl} ransom_{Acc}.
 'The kidnapper must be promised the ransom.'
- g. *Kuchařka_{Act,i} [už] musí []_i mít uvařeno].*
Cook_{Nom} already must have_{Inf} cooked_{Prtpcl Neut Sg}.
 'The cook must have already cooked (everything).'
- h. *[Tady se ti_{Act}] musí [sedět nepohodlně].*
Here SE you_{Dat} must sit_{Inf} uncomfortably.
 'This must be an uncomfortable seat for you.'

We can see in the above examples that the infinitive can occur in various diatheses. The infinitive can occur in both periphrastic and reflexive passive and in the construction with the verb *dostat*; the mediopassive and the active construction with the verb *mít* are only possible with the verb *muset* (must) in the meaning of high probability. It seems that the governor can only occur in active voice, but in the spoken language, we can find evidence that it can also occur as reflexive passive; we will ignore these cases, however, because we would introduce another source of ambiguity to the lexicon.

As modal verbs have no representation on the tectogrammatical level we have to find a notation of these lexical entries that respects this theoretical constraint and gives all necessary information. In (67) we can see several examples of both modal and non-modal verbs.

- (67) a. *muset~1* R--s[a0]0[sId%\$#mD]@
- b. *muset~2* R--s[a0]0[sId%\$*~]@
- c. *začít* R--s[a1]1[sId%\$]@
- d. *zdát* RSEs[i2]1(hPc3)2[hTc1|sD]@
- e. *zdát* RSEs[a2]1(hPc3)2[sI1{být}|hQc1]@

The frame of the modal verb *muset* (67a) contains only one "argument" (0[sId%\$ mD]) whose functor is marked by 0 (zero). This notation was adopted for sentence complementations which do not belong to frame of a given verb. Attributes enclosed in brackets represent constraints imposed on the surface

forms. In (67a) these attributes have the following values: infinitive (sl) which can occur in periphrastic and reflexive passive (d%\$),⁹ and the modality feature *debitive* (mD). The subject of the construction is raised from the infinitival clause (s[a0]). The verb *muset* can occur only in active voice (@).

The frame of the verb *muset* in the meaning of high probability (67b) is very similar the frame of the modal verb. It differs in constraints imposed of diatheses of the embedded infinitive (%\$*~) and in a missing modality marker. The aspectual (phase) verb *začít* is a verb with one participant (Actor: 1[sld%\$]) which is realized by an infinitive. The infinitive can occur also in periphrastic or reflexive passive, and the verb *začít* can only occur in active voice. The subject of the verb *začít* is raised from the infinitival clause.

The verb *zdat se* has been already discussed. In (67d) we can see the frame of the verb with nominal and clausal objects (and with an “inherent” subject), in (67e) the frame with infinitive.

Object raising verbs are such verbs that have an infinitive in the frame and the subject of this infinitive becomes an object of the higher verb. This group contains the verbs of perception:

- (68) a. *Vidím ho_i __i přicházet.*
 I see him to be coming.
 ‘I see him coming.’
- b. *?Vidím ho_i __i být tázána.*
 ?I see him to be asked.
 ‘I see him being asked.’
- c. *?Cítím bábovku_i __i péct se.*
 ?I smell cake to bake SE.
 ‘I can smell that a cake is being baked.’

The passive constructions are questionable with this group of verbs; a further research on a text corpus will be necessary. In the current version of the lexicon the possibility of creating the passive voice is suppressed. The frame is encoded this way:

- (69) a. vidět R--s[i1]1(hPc1)2(hPTSc4|sD|sZ)&
 b. vidět R--s[i1]1[hPc1]2[sId\$|sZd&]0[hPTSc4a2]@

For marking the source of the raised subject we use the attribute *a*. Its value points to a functor from which the subject was raised.

3.4.2 Equi verbs

This type of verbs in Czech was described by K. Svoboda (1962) and J. Panevová (1996). Svoboda does not use the term *equi* or *control*, but he distinguishes between “subject infinitives” verb and “object infinitives”. He does not distinguish raising and equi verbs, as he only considers the surface structure and grammatical functions as subject, objects, etc.

Panevová describes carefully equi verbs from the point of view of FGD. She distinguishes four types of equi verbs:

- (70) a. Subject-control (Act-Sb):
Jan_{Act,i} se bojí [__{Act,i} zůstat doma sám].
 Jan_{Nom} fears stay_{Inf} at home alone.
- b. Object-control (Addr-Sb):
Oni_{Act} mu_{Addr,i} poručili [__{Act,i} přijít].
 They_{Nom} him_{Dat} ordered come_{Inf}.
- c. Ambiguous class (Act-Sb) or (Addr-Sb):
Rodiče_{Act,i} Petrovi_{Addr,j} slíbili [__{Act,j} svézt se na poníkovi].
 Parents to Petr promised ride_{Inf} on pony.
Rodiče_{Act,i} Petrovi_{Addr,j} slíbili [__{Act,i} přestat kouřit].
 Parents to Petr promised stop_{Inf} smoke_{Inf}.

⁹These constraints represent additional constraints to those imposed by the lexical entry of a given infinitive.

d. Object-control (Pat-Sb) (the infinitive has the function of Intent):

*Sedlák vyhnal krávy*_{Pat,i} [_{Act,i} *pást se*].
Farmer drove cows_{Acc} graze_{Inf}.

We will add two more types, which are quite rare but interesting. The embedded infinitive should be understood as a kind of passive, though it is in active voice:

(71) a. (Act-Addr) control:

*Anežka*_{Act,i} *chce* [_{Act} *podat knihu* _{Addr,i}].

Anežka_{Nom} wants pass_{Inf} book_{Acc}.

‘Anežka wants someone to pass her the book.’

*Anežka*_{Act,i} *chce* [_{Act} *přečíst pohádku* _{Addr,i}].

Anežka wants read_{Inf} tale_{Acc}.

‘Anežka wants someone to read her a tale.’

*Anežka*_{Act,i} *chce* [_{Act} *poučit o hudbě* _{Addr,i}].

Anežka wants instruct_{Inf} in music.

‘Anežka wants someone to instruct her in music.’

b. (Act-Pat) control:

*Plot*_{Act,i} *chce* [_{Act} *natřít* _{Pat,i}].

Fence wants paint_{Inf}.

‘The fence needs painting.’

*Pepík*_{Act,i} *potřebuje* [_{Act} *nařezat* _{Pat,i}].

Pepík needs spank_{Inf}.

‘Pepík needs spanking.’

For proper description of all the above constructions in the lexicon we also have to examine the possible diatheses of both the governor and the controlled infinitive. Let us start with (Act-Sb) control:

(72) a. *Petr*_{Act,i} *chce* _{Pat,i} *být pochválen*.

Petr wants be_{Inf} praised.

b. *Anežka*_{Act,i} *chce* _{Addr,i} *být poučena o hudbě*.

Anežka wants be_{Inf} instructed in music.

c. *Bábovka*_{Act,i} *se nechce* _{Pat,i} *péct*.

Cake SE does not want bake_{Inf}.

‘The cake refuses to get baked.’

d. *Pepík*_{Act,i} *nechce* _{Pat,i} *dostat nařezáno*.

Pepík does not want get_{Inf} spanked.

‘Pepík does not want to be spanked.’

e. *Petr*_{Act,i} *chce* _{Pat,i} *dostat/*mít slíbenou hračku*.

Petr wants get_{Inf}/?have_{Inf} promised toy.

‘Petr wants to be promised a toy.’

f. *Matka*_{Act,i} *už chce* _{Act,i} *mít uvařeno*.

Mother already wants have_{Inf} cooked.

‘Mother wants to have all cooking done already.’

We can see that the infinitive can be in passive, as well as in a construction with *mít* or *dostat*. The passivization of the governor, on the other hand, does not seem to be possible. The reason may be that the subject of the embedded infinitive is controlled by Actor which would become general in a passive construction. An exception is a mediopassive of the verb *chtít*.

- (73) *Nechce se mi_{Act,i} –_{Act,i} spát.*
 Wants_{Neg3SgNeut} SE me_{Dat} sleep_{Inf}.
 ‘I don’t want to sleep’.
- Nechce se mi_{Act,i} –_{Pat,i} být bit.*
 Wants_{Neg3SgNeut} SE me_{Dat} be_{Inf} beaten.
 ‘I don’t want to be beaten’.
- Bábovce_{Act,i} se nechce –_{Pat,i} péct (se).*
 Cake_{Dat} SE wants_{Neg3SgNeut} bake_{Inf} (SE).
 ‘The cake refuses to get baked’.

The verb *chtít* even allows reflexive passive with general Actor:

- (74) *Když se –_{Act,i} nechce –_{Act,i} pracovat, tak se nemusí –_{Act,i} jíst.*
 When SE wants_{Neg3SgNeut} work_{Inf} then SE needs_{Neg} eat_{Inf}.
 ‘If one doesn’t want to work then he doesn’t need to eat.’

Frames of two equi verbs, *bát se* (fear) and *chtít* (want) follow:

- (75) a. *bát* RSEs[i1]1[hPc1]2[hPTRc2|hPTRc4r{o}|sD|sU|sIq1d%]@
 b. *chtít~1* R--s[i1]1[hPc1]2[hTc4|sIq1d%\$#~]@
 c. *chtít~2* PSEs[i2]1(hPTc3)2[hZc4|sIq1d%\$]@

Next, we will examine the the possibility of passivization of verbs with (Pat-Sb) control.

- (76) a. *Velitelé_{Act,i} vojákům_{Addr,j} zakázali –_j chodit na pivo.*
 Commanders soldiers_{Dat} prohibited go_{Inf} for beer.
- b. *Vojákům_{Addr,j} bylo (veliteli_{Act,i}) zakázáno –_j chodit na pivo.*
 Soldiers_{Dat} was (commanders_{Ins}) prohibited go_{Inf} for beer.
- c. *Vojákům_{Addr,j} se zakázalo –_j chodit na pivo.*
 Soldiers_{Dat} SE prohibited go_{Inf} for beer.
- d. *Vojáci_{Addr,j} mají/*dostali (od velitelů_{Act,i}) zakázáno –_j chodit na pivo.*
 Soldiers_{Nom} have/*got (from commanders) prohibited go_{Inf} for beer.
- e. *Šéf_{Act,i} zabránil podřízenému_{Addr,j} –_j být povýšen.*
 Boss prevented employee_{Dat} be_{Inf} promoted.
- f. *Podřízenému_{Addr,j} bylo (šéfem_{Act,i}) zabráněno –_j být povýšen.*
 Employee_{Dat} was (boss_{Ins}) prevented be_{Inf} promoted.
- g. *Podřízenému_{Addr,j} se zabránilo –_j být povýšen.*
 Employee_{Dat} SE prevented be_{Inf} promoted.

Frames for the verbs *poručit* (order), *zakázat* (forbid) and *zabránit* (prevent) follow:

- (77) a. *poručit* R--s[i1]1(hPc1)2[sU|sIq3d@]3(hPc3)%\$#
 b. *zakázat* R--s[i1]1(hPc1)2[sU|sIq3d@]3(hPc3)%\$#
 c. *zabránit~1* R--s[i1]1(hPc1)2[sU|sIq3d%\$

The next category to be examined are the ambiguous verbs like *slíbit* (promise) or *odepřít* (refuse). First, we will examine possible diatheses of the governor.

- (78) a. *?Rodiče_{Act,i} Petrovi_{Addr,j} slíbili –_j svézt se na poníkovi.*
 Parents_{Nom} Petr_{Dat} promised ride_{Inf} on pony.
- b. *Petrovi_{Addr,j} bylo (rodiči_{Act,i}) slíbeno –_j svézt se na poníkovi.*
 Petr_{Dat} was (parents_{Ins}) promised ride_{Inf} on pony.

- c. *Petrovi_{Addr,j} se slíbilo* $_{-j}$ *svézt se na poníkovi.*
Petr_{Dat} SE promised ride_{Inf} on pony.
- d. *Petr_{Addr,j} má/dostal (od rodičů_{Act,i}) slíbeno* $_{-j}$ *svézt se na poníkovi.*
Petr has/got (from parents) promised ride_{Inf} on pony.
- e. *Rodiče_{Act,i} Petrovi_{Addr,j} slíbili* $_{-i}$ *přestat kouřit.*
Parents_{Nom} Petr_{Dat} promised stop_{Inf} smoke_{Inf}.
- f. **Petrovi_{Addr,j} bylo (rodiči_{Act,i}) slíbeno* $_{-i}$ *přestat kouřit.*
Petr_{Dat} was (parents_{Ins}) promised stop_{Inf} smoke_{Inf}.
- g. *(*Petrovi_{Addr,j} se slíbilo* $_{-i}$ *přestat kouřit.*
Petr_{Dat} SE promised stop_{Inf} smoke_{Inf}.
- h. **Petr_{Addr,j} má/dostal (od rodičů_{Act,i}) slíbeno* $_{-i}$ *přestat kouřit.*
Petr_{Nom} has/got (from parents) promised stop_{Inf} smoke_{Inf}.

The construction (78a) is rejected by some speakers, but it can be converted into passive constructions (78b)–(78d), which are admitted by all speakers. The sentence (78e) is perfectly correct, but the passivization of the controller is impossible. Only the sentence in (78g) can be accepted if we suppose Actor of the embedded infinitive to be general.

Let us now try to passivize the infinitive:

- (79) a. *Rodiče_{Act,i} Petrovi_{Addr,j} slíbili* $_{-j}$ *být pochválen.*
Parents_{Nom} Petr_{Dat} promised be_{Inf} praised.
- b. *Petrovi_{Addr,j} bylo (rodiči_{Act,i}) slíbeno* $_{-j}$ *být pochválen.*
Petr_{Dat} was (parents_{Ins}) promised be_{Inf} praised.
- c. *Petrovi_{Addr,j} se slíbilo* $_{-j}$ *být pochválen.*
Petr_{Dat} SE promised be_{Inf} praised.
- d. *Petr_{Addr,j} má/*dostal (od rodičů_{Act,i}) slíbeno* $_{-j}$ *být pochválen.*
Petr has/*got (from parents) promised be_{Inf} praised.
- e. *Rodiče_{Act,i} Petrovi_{Addr,j} slíbili* $_{-i}$ *být v práci povýšeni.*
Parents_{Nom} Petr_{Dat} promised be_{Inf} at work promoted.

Now, we can encode the frames of the verb *slíbit* (promise):

- (80) a. *slíbit*~1 R--s[i1]1(hPc1)2[hZc4|sD|sIq3d%\$##*
- b. *slíbit*~2 R--s[i1]1(hPc1)2[hTc4|sD|sIq1d%]3(hPc3)@

The constructions with (Act-Pat) control and (Act-Addr) control do not allow any kind of diathesis, so their frames will be quite simple. We only had to introduce two more attributes in the description: p for (Act-Pat) control and t for (Act-Addr) control.

- (81) a. *chtít*~3 R--s[i1]1[hPTc1]2[sIp1d@]@
- b. *chtít*~4 R--s[i1]1[hPTc1]2[sIt1d@]@

4 Algorithm for assigning the frames

In this section the automatic processing of the source data will be described. The format of the source data was described in Section 2. The desired content of the lexicon was described in Section 3. The steps which have to be done to achieve this are

1. identifying single frames
2. merging all variants of a single frame
3. marking the obligatoriness of frame members
4. assigning the functors to members
5. marking the possible diatheses

In the next sections these single steps will be described in detail.

4.1 Identifying and merging frames, marking the obligatoriness

In the source lexicon, every lemma is listed only once, even if it has several valency frames. A single valency frame, on the other hand, can have several variants (e.g. *učit koho co_{Acc}*, *učit koho čemu_{Dat}*—teach sb st). The variants of one frame are mixed with other frames and thus the first task is to separate the different frames and merge the variants. Let us show it with an example. The verb *bránit* has the following format in the source lexicon:

(82) *bránit* <v>hTc3, sI, hPc3-sUeN, hPc3-hTc6r{v}, (protect, prevent)
 hPTc4, hPTc4-hPTc3r{proti}, hPTc4-hPTc7r{před}

	A hTc3	B sI	C hPc3	D sUeN	E hTc6r{v}	F hPTc4	G hPTc3r{proti}	H hPTc7r{před}
1	+							
2		+						
3			+	+	+			
4			+		+			
5						+	+	+
6						+	+	
7						+		+

Table 4: Identifying single frames

We arrange the members of all its frames into a table (see Table 4): the rows are single “frames” from the original dictionary and the columns are single members of the frames. If there are more than one + in a column, then two or more frames share that member. We suppose that frames with a non-empty intersection are variants of one frame. We illustrate the procedure of identifying variants of single frames in Table 4: the intersecting frames are marked by the gray background. They form non-intersecting rectangles. Every grey rectangle represents one frame. Members of a single frame which never occur in one line together can be declared with high probability as variants of one member (in Table 4 we can see that items D and E are variants of one member and items G and H are variants of another member). Now, we can collapse columns with the variants, which is shown in Table 5: the frames 3 and 4 were merged into 3’ and the frames 5 and 6 into 6’.

	A	B	C	D E	F	G H
1	+					
2		+				
3’			+	+		
5					+	+
6’					+	+

Table 5: Merging frame variants

There is a small problem with single-member frames (frames 1 and 2 in our example). They can be understood as separate frames, as in the case of *mířit kam* (head somewhere), *mířit na koho* (aim at sb), or as variants of one frame, as in the case of *bádat nad čím*, *bádat o čem* (research into st). We had to make a decision whether we wanted to merge all such frames, or whether we wanted to keep them separate. We decided to “merge as much as possible” because of an easier assignment of the functors,

which will be explained in the next section. In our table, we then also merge the frames 1 and 2 into a frame with one member A|B.¹⁰

In the above table we can also see how we identify obligatory members of a frame. In lines 5 and 6', the member F is always present, while the other member G|H may be missing. Unfortunately, we are not able to say whether G|H is a general inner participant, or optional participant, or obligatory and deletable free modification, or even non-obligatory free modification, but at least the information about obligatory members of the frame should be correct. We use the square brackets for obligatory members of a frame (as was described in Section 3), and for now, we will use the parentheses for all other cases. The entry from example (82) now can be recorded as follows:

- (83) a. bránit [hTc3|sI] (bránit čemu/něco udělat) (prevent st/doing st)
 b. bránit [hPc3] [sUeN] (bránit komu, aby něco neudělal) (prevent sb from doing st)
 c. bránit [hPTc4] (hPTc3r{proti}|hPTc7r{před}) (bránit koho/co {proti komu/čemu/před kým/čím}) (protect sb/st {against sb/st|from sb/st})

As we said above, the source dictionary does not contain the so-called “left valency”, i.e. subjects. This information is usually missing in printed dictionaries, as readers are able to fill the missing information, but in an electronic dictionary which is meant for language processing, such information must be included. We will describe the process of adding the subjects in the next section.

4.2 Assigning functors

It was shown by many authors that there is no straightforward correspondence between the deep frame and its surface realization. One can, however, try to find some regularities or tendencies, and then formulate rules for assigning the functors to the surface frames. The mappings between the tectogrammatical and morphemic levels (in active voice) is shown in Figure 1.

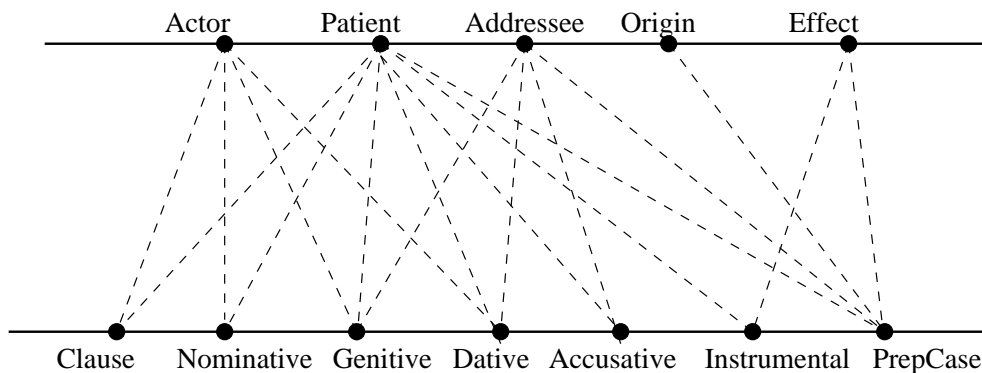


Figure 1: Mapping between TL and ML in active voice

We can see that this picture does not help much—nearly everything is possible. It is necessary to add, however, that this picture also covers all marginal frames like *líbit* RSEs[i2]1(hPRc3)2[hPTc1]@ (like, appeal) and *ubývat* R--1[hTc2]@ (dwindle).

Among all correspondences, there are some which are considered as typical. In the direction from the tectogrammatical level to the morphemic one these are Actor → Nominative, Patient → Accusative, Addressee → Dative, Effect → Instrumental, Origin → Genitive+Prep{z} (from) or Origin → Genitive+Prep{od} (from). In the opposite direction the correspondences are not so clear because of free

¹⁰A careful reader notices that the second frame should also contain Dative (hPc3) and it should in fact be merged with the third frame into one frame: *bránit* [hPc3] [sI|sUeN]. We showed here a real example from the source lexicon, where some information was missing. The correction of this type of mistake is left for the post-editor.

	Actor	Patient	Addressee	Origin	Effect	
<i>dát</i>	(Nom)	Acc	Dat			give
<i>dostat</i>	Nom	Acc		<Gen+od>		get
<i>šít</i>	(Nom)	(Acc)	<Dat>	<Gen+z>		sew
<i>předělat</i>	(Nom)	Acc	<Dat>	<Gen+z>	<Acc+na>	remake
<i>žádat</i>	(Nom)	Acc		(Gen+od)		ask

Table 6: Prototypical frames

modifications, which have a very broad repertory of the surface realizations. Thus Accusative can represent Patient or Temporal modification, Instrumental can represent Patient (*stát se*—become), Effect (*zvolit*—elect), Means (*zaplavit*—flood), Manner (*kopat*—dig); Genitive with the preposition *od* can represent Patient (e.g. *distancovat se*—dissociate), Origin (*dostat*—get), Direction from (*odejít*—leave), Temporal modification *how long* (*spát*—sleep), Cause (*opuchnout*—swell), etc.

If we consider only frames with at least three participants¹¹ we get another picture shown in Figure 2.

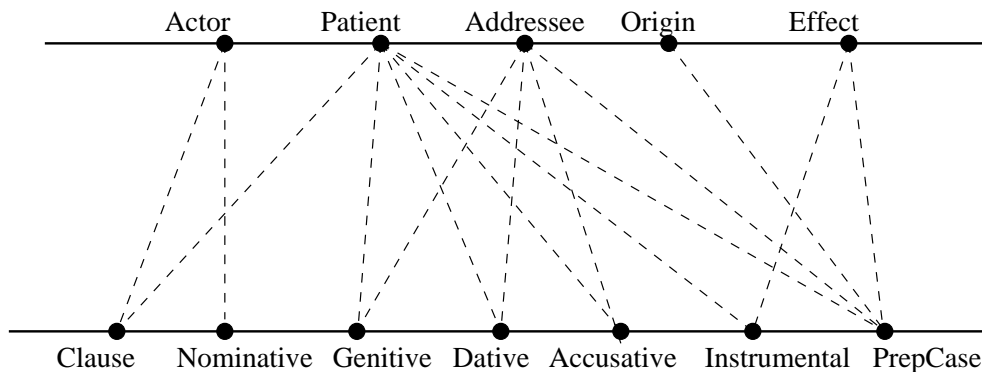


Figure 2: Mapping between TL and ML for verbs with at least three participants

Though some joins disappeared, we still cannot find a unique mapping between the tectogrammatical and morphemic level. However, we can observe that frames can be split in two groups. The first group contains verbs whose participants are only realized by typical surface forms; we call these frames *prototypical* (several examples are listed in Table 6). The other group contains verbs with *non-prototypical* frames, where at least one member is realized by a non-typical surface form (examples are in Table 7). This observation was done by J. Panevová, and an experimental algorithm for assigning the functors to surface realizations was created by Panevová and Skoumalová (1992). The algorithm checks whether a frame contains only prototypical surface forms, and if so it assigns them the corresponding functors. In Table 7, we can see that there is a possible source of problems in frames with surface forms in Accusative and Dative—their functors can be assigned the other way round than we expect. In this case we have to add one more criterion, and it is that Addressee must be “more animate” than Patient.¹² From this reason we only assume *animate* Dative as the typical realization of Addressee (hPc3 or hPTc3).

In the experiment, it was supposed that we worked only with inner participants (free modifications were filtered out), which made the task easier. In BRIEF lexicon, however, we cannot rely on getting participants only in surface frames, but on the other hand, the repertory of free modifications occurring

¹¹Frames with one or two participants are “uninteresting” as the functors are assigned after the rules listed in (6) in Section 1.2: if the frame has only one participant it is Actor, if there are two participants in the frame, they are Actor and Patient. In most cases, Actor is realized as Nominative and Patient as the “remaining” surface realization. There are some exceptional frames, as *líbit* RSEs[i]2[1[hPRc3]2[hPTc1]@ (like, appeal) or *zželet* RSE1[hPc3]2[hPTRc2]@ (take pity on sb/st) which have to be edited manually.

¹²The scale of animacy (in BRIEF notation) is hT < hPT < hP.

in the lexicon is not as wide as in the language as a whole (for example, a free modification of condition hardly occurs in a lexical entry). For this reason, we adopted a slightly different approach in the processing of BRIEF lexicon.

First, it was necessary to add the missing subjects. We did this automatically, and all frames got a subject in Nominative which was assigned the role of Actor: $s[i1]1[hPTc1]$.¹³

The second step was assigning the roles to other members of the frame. Some preparation for this was done already while merging the frames: there is a list of possible functors for every surface realization, and this list was attached to every member of the original frame.¹⁴ When we merged two members of a frame together we also made an intersection of the attached lists. An empty intersection prevented the two members from being merged. This process is shown in Table 8 on a frame of the verb *čertit se* (be angry). In BRIEF lexicon, the entry of this verb had the following form:

(84) *čertit se* $\langle v \rangle hPTc4r\{na\}, hTc4r\{pro\}, hTc7r\{nad\}, hTc3r\{kvûli\}$

Every surface realization is assigned a list of functors, as shown in the table. However, the functor ACTANT which denotes any participant is only taken in consideration if the surface realization has no variants.¹⁵ As we first try to merge all the prepositional cases into one member of the frame, we exclude ACTANT from the list. In the rest of the table, we can see that the first prepositional case ($hPTc4r\{na\}$) has an empty intersection of functors with other prepositional cases which means that it cannot be taken as their variant inside one member of a frame. The remaining surface realizations have a non empty intersection of functors containing the value CAUSE. In the resulting frame, the first prepositional case will be assigned the functors ACTANT and DIR.WHERE. Other prepositional cases will be merged into one frame member which will be assigned the functor CAUSE:¹⁶

(85) *čertit_se* $s[i1]1[hPTc1]2A[hPTc4r\{na\}] \setminus$
 $C[hTc4r\{pro\}|hTc7r\{nad\}|hTc3r\{kvûli\}]$

After the merging of participants, we get three kinds of frames: frames where every member has only one functor assigned, frames where participants are distinguished from free modifications, but some

¹³Some Czech verbs do not have a subject at all, e.g. *pršet* (rain), in some frames the subject is realized by a clause or by an infinitive, e.g. *znamenat* (mean), *zdát se* (seem), but the vast majority of Czech verbs have a nominal subject in Nominative. The exceptions will be treated by a post-editor.

¹⁴These lists of possible functors were created manually. The original lexicon was first divided into classes of frames containing a certain surface realization. These classes were analyzed and every surface realization was assigned a list of functors. Similar lists were also created for the Prague Dependency Treebank (Hajičová, Panevová, and Sgall, 2000). These PDT lists are longer because they contain all functors found in texts, not only in a lexicon. Beside it, they also contain more prepositional cases than the BRIEF lexicon.

¹⁵We do not try to assign single inner participants (Actor, Patient, etc.) in this step, we only mark whether a certain surface form can possibly represent an inner participants. Because of technical reasons we mark all potential inner participants as Patients—in a case that there is only one participant beside Actor we get Patient “for free”. In a case that there are more participants, further processing is necessary.

¹⁶For the list of abbreviations used for functors see Appendix A.4, for lists of functors attached to every surface realization see Appendix B.2.

	Actor	Patient	Addressee	Origin	Effect	
<i>zvolit</i>	(Nom)	Acc			Ins	elect
<i>hrozit</i>	(Nom)	Ins	(Dat)			threaten
<i>vystavit</i>	(Nom)	Dat	Acc			subject
<i>dědit</i>	(Nom)	(Acc)		(Loc+ <i>po</i>)		inherit
<i>hovořit</i>	(Nom)	$\langle Loc+o \rangle$	(Ins+ <i>s</i>)			speak
<i>psát</i>	(Nom)	$\langle Loc+o \rangle$	$\langle Dat \rangle$		(Acc)	write
<i>zeptat se</i>	Nom	Acc+ <i>na</i>	(Gen)			ask

Table 7: Non-prototypical frames

	hPTc4r{na}	hTc4r{pro}	hTc7r{nad}	hTc3r{kvûli}
(ACTANT)	+	+	+	
DIR.WHERE	+			
CAUSE		+	+	+
PURPOSE		+		+
WHERE			+	

Table 8: Merging frame of the verb *čertit se* (be angry)

of the free modifications are ambiguous, and frames where at least one member is ambiguous between a participant and a free modification. Approximately one third of all merged frames fall in the first category and another thousand into the second one. These frames are candidates for further processing with help of the above mentioned algorithm, and therefore they will be separated from the rest which must be left for post-editing.

Now, we will describe the process of assigning functors in the categories where participants are distinguished from free modifications. These frames fall into two subcategories: frames with at most two inner participants (i.e. Actor and Patient) and frames with at least three inner participants. The former have been handled already and we do not need to process them any further. The latter will be processed by the algorithm for assigning functors, but let us first resume the starting conditions:

We have at least three inner participants.

Actor is already assigned to the subject.

We have to decide which of the participants is Patient and what are functors of the remaining inner participants.

We will not describe the algorithm in detail, we only sketch the overall strategy. More details and a flow chart can be found in (Skoumalová, 2001).

A rule (following from the participant shifting) which must be observed after every step of the algorithm is that Patient slot must be filled. If there is only one unassigned member and the Patient slot has not been filled yet then the last member of the frame is assigned the Patient functor.

We start with searching for Origin as Origin has the narrowest set of possible surface realizations, which in addition are not “polysemous”.

Addressee assignment is ruled by the animacy of surface forms rather than the morphological cases. Animate Accusative or an animate prepositional case are realizations of Addressee rather than inanimate Dative.

The decision about Effect can be quite difficult. Beside the typical prepositional cases also Instrumental can be a surface form of Effect. We then have to take into consideration the remaining unassigned members of the frame and make a decisions about pairs of surface forms.

As was said above, approximately 7500 frames are processed by this algorithm and the program ends successfully in all cases. The remaining ca 11,000 frames must be edited manually, with help of an editor which will be created for this purpose. The editor’s work should be easier as s/he gets a (small) set of possible functors which can be assigned to every member of a frame and s/he does not have to choose from all 47 possibilities.

4.3 Marking diatheses

We made a simple assumption that

reflexive verbs cannot have any diatheses (the exception with the periphrastic passive of the verb *tázat se* was discussed above), and so they get the mark @.

intransitive verbs¹⁷ can form reflexive passive; they get the mark \$.

a verb with a member in Accusative or in an indirect case (without preposition) can form both periphrastic and reflexive passive; it gets marks %\$

a verb whose all objects are prepositional cases can form the reflexive passive; it gets the mark \$.

During the automatic processing all frames are assigned these marks and corrections will be made by the post-editor. Actors, which were added automatically to all frames, are marked as generalizable ((hPTc1)) in frames that allow for forming any passive, and they are marked as obligatory ([hPTc1]) in other frames.

5 Results and further perspectives

We have characterized a syntactic lexicon which can be used in various systems of natural language processing, especially in systems using symbolic methods (as opposed to stochastic methods). In contrast to other electronic dictionaries, we have created a large-scale lexicon, which should cover a large part of the vocabulary. The lexicon, however, still needs some editing work, but we believe that it was pre-processed in such a way that the editing work will be easy.

Recently, we started a new experiment, which should improve the quality of the lexicon: we try to extract surface frames from Prague Dependency Treebank (Hajič et al., 1999) and then we want to process them the same way as the BRIEF lexicon. This should result in several things:

1. We get some more variants of frames and we can get more accurate results of assigning functors to single members of frames.
2. We get information on diatheses which are really used in contemporary texts.
3. We get a method for extracting valency frames from text corpora.

For the last task we suppose that the corpus is preprocessed (syntactically tagged) but in fact some kind of shallow parsing or even marking the boundaries of clauses should be sufficient. The improving dictionary, on the other hand, will help to improve the parsing and tagging methods.

¹⁷The term intransitive verb here means a verb with only one participant realized as subject in Nominative.

References

- Avgustinova, Tania, Alla Bémová, Eva Hajičová, Karel Oliva, Jarmila Panevová, Vladimír Petkevič, Petr Sgall, and Hana Skoumalová. 1995. Linguistic problems of Czech. Technical report, JRP PECO 2824, Prague. Final research report.
- Chomsky, Noam. 1986. *Knowledge of Language: Its Nature, Origin, and Use*. Convergence. Praeger, Westport.
- Dalrymple, Mary, Ronald M. Kaplan, John T. Maxwell III, and Annie Zaenen, editors. 1995. *Formal Issues in Lexical-Functional Grammar*. Number 47 in Lecture Notes. CSLI, Stanford.
- Daneš, František, Zdeněk Hlavsa, et al. 1987. *Větné vzorce v češtině (Sentential paradigms in Czech)*. Number 23 in Studie a práce lingvistické. Academia, Prague.
- Dokulil, Miloš. 1941. Morfologické kategorie pasiva ve spisovných jazycích severských ve srovnání se spisovnou češtinou (Morphological categories of passive in Nordic standard languages in comparison with standard Czech). In *Hrst studií a vzpomínek: prof. dr. Ant. Beerovi jeho žáci (Handful of studies and memories: to prof. dr. A. Beer from his pupils)*. Odbočka Jednoty českých filologů v Brně, Brno, pages 77–99.
- Grepl, Miroslav and Petr Karlík. 1989. *Skladba spisovné češtiny (Syntax of Standard Czech)*. SPN, Prague, 2 edition.
- Grepl, Miroslav and Petr Karlík. 1998. *Skladba češtiny (Syntax of Czech)*. Votobia, Olomouc.
- Hajič, Jan. 1994. *Unification Morphology Grammar*. Ph.D. thesis, Charles University, Faculty of Mathematics and Physics, Prague.
- Hajič, Jan, Jarmila Panevová, Eva Buráňová, Zdeňka Urešová, Alla Bémová, Jan Štěpánek, Petr Pajas, and Jiří Kárník, 1999. *Anotace na analytické rovině: Návod pro anotátory (Annotation on analytical level: Manual for the annotators)*. Charles University, Faculty of Mathematics and Physics, Prague. URL: <http://ufal.mff.cuni.cz/projekty.html>.
- Hajičová, Eva, Jarmila Panevová, and Petr Sgall. 2000. A Manual for Tectogrammatical Tagging of the Prague Dependency Treebank. Technical Report ÚFAL/CKL TR-2000-09, Charles University, Faculty of Mathematics and Physics, Prague. URL: <http://ufal.mff.cuni.cz/pdt/pdt.html>.
- Hausenblas, Karel. 1963. Slovesná kategorie výsledného stavu v dnešní češtině (verb category of resultative in contemporary Czech). *Naše řeč*, 46:13–28.
- Havránek, Bohumil. 1928. *Genera verbi v jazycích slovanských I (Genera verbi in Slavic languages I)*. Královská česká společnost nauk, Prague.
- Horák, Aleš. 1998a. Popis formátu brief (Description of the format of the lexicon). Unpublished documentation.
- Horák, Aleš. 1998b. Verb valency and semantic classification of verbs. In Petr Sojka, Václav Matoušek, Karel Pala, and Ivan Kopeček, editors, *Proceedings of the First Workshop on Text, Speech, Dialogue — TSD'98*. Masaryk University Press, Brno, pages 61–66.
- Karlík, Petr. 2000. Hypotéza modifikované valenční teorie (Hypothesis of the modified valency theory). *Slovo a slovesnost*, LXI(3):170–189.
- Karlík, Petr, Marek Nekula, and Zdenka Rusínová, editors. 1995. *Příruční mluvnice češtiny (Handbook of Czech Grammar)*. Nakladatelství Lidové Noviny, Prague.
- Koček, Jan, Marie Kopřivová, and Karel Kučera, editors. 2000. *Český národní korpus — Úvod a příručka uživatele*. Charles University, Faculty of Philosophy, Prague. URL: <http://ucnk.ff.cuni.cz>.
- Králíková, Květa. 1981. Reflexivnost sloves z hlediska automatické analýzy češtiny (Reflexivity of verbs from the point of perspective of automatic analysis of Czech). *Slovo a slovesnost*, XLII(4):291–298.
- Oliva, Karel. 1989. *A Parser for Czech Implemented in Systems Q*. Number XVI in Explizite Beschreibung der Sprache und automatische Textbearbeitung. Charles University, Faculty of Mathematics and Physics, Prague.
- Oliva, Karel. 2000. Hovory k sobě/si/sebe/se (Discussion on sobě/si/sebe/se). In Zdeňka Hladká and Petr Karlík, editors, *Čeština—univerzália a specifika 2, (Czech—Universals and Specifics 2)*, Proceedings of the Conference held in Šlapanice u Brna, November 17–19, 1999, pages 167–171, Brno. Masaryk University.
- Pala, Karel and Pavel Ševeček. 1997. Valence českých sloves (Valency of Czech verbs). In *Sborník prací FFBU*, volume A45. Masaryk University, Brno, pages 41–54.

- Panevová, Jarmila, 1971. *Časové a vidové kategorie predikátu (Tense and aspect categories of predicate)*, pages 23–44. In (Panevová, Benešová, and Sgall, 1971).
- Panevová, Jarmila. 1974–75. On verbal frames in functional generative description, Part I and II. *Prague Bulletin of Mathematical Linguistics*, 22:3–40,23:17–52.
- Panevová, Jarmila. 1980. *Formy a funkce ve stavbě české věty (Forms and Functions in Syntax of Czech Sentence)*. Number 13 in *Studie a práce lingvistické*. Academia, Prague.
- Panevová, Jarmila. 1996. More remarks on control. In Eva Hajičová, Oldřich Leška, Petr Sgall, and Zdena Skoumalová, editors, *Prague Linguistic Circle Papers*, volume 2. John Benjamins Publishing Company, Amsterdam/Philadelphia, pages 101–120.
- Panevová, Jarmila. 1999. Česká reciproční zájmena a slovesná valence (Czech reciprocal pronouns and verb valency). *Slovo a Slovesnost*, LX(4):269–275.
- Panevová, Jarmila, Eva Benešová, and Petr Sgall. 1971. *Čas a modalita v češtině (Tense and Modality in Czech)*. Universita Karlova, Prague.
- Panevová, Jarmila and Hana Skoumalová. 1992. Surface and deep cases. In *Proceedings of COLING '92*, pages 885–889, Nantes.
- Pollard, Carl and Ivan A. Sag. 1994. *Head-Driven Phrase Structure Grammar*. The University of Chicago Press, Chicago & London.
- Sgall, Petr. 1967. *Generativní popis jazyka a česká deklinace (Generative Description of Language and Czech Declension)*. Number 6 in *Studie a práce lingvistické*. Československá akademie věd.
- Sgall, Petr, Eva Hajičová, and Jarmila Panevová. 1986. *The Meaning of the Sentence in Its Semantic and Pragmatic Aspects*. D. Reidel Publishing Company, Dordrecht.
- Skoumalová, Hana. 2001. *Czech syntactic lexicon*. Ph.D. thesis, Charles University, Faculty of Arts, Prague. URL: <http://utkl.ff.cuni.cz/~skoumal/dissertation>.
- SSJČ. 1989. *Slovník spisovného jazyka českého (Dictionary of standard Czech)*. Academia, Prague.
- Štícha, František. 1984. *Utváření a hierarchizace struktury větného znaku (Creation and hierarchization of the structure of a sentence sign)*. Univerzita Karlova, Prague.
- Straňáková-Lopatková, Markéta. 2001. *Homonymie předložkových skupin a možnost jejího automatického zpracování (Homonymy of prepositional groups and possibility of its automatic processing)*. Ph.D. thesis, Charles University, Faculty of Mathematics and Physics, Prague.
- Svoboda, Karel. 1962. *Infinitiv v současné spisovné češtině (Infinitive in Contemporary Standard Czech)*. Rozpravy ČSAV. Academia, Prague.
- Svozilová, Naďa, Hana Prouzová, and Anna Jirsová. 1997. *Slovesa pro praxi (Verbs for practical use)*. Academia, Prague.
- Tesnière, Lucien. 1959. *Éléments de syntaxe structurale*. Klincksieck, Paris.

A Symbols used in the dictionary

A.1 Voice

R — regular frame (in active voice with possible derivations)

P — irregular passive frame

There are three more marks which are not used in the lexicon, but they are exploited in sentence patterns generated from the frames.

M — construction with support verb *mít*

D — construction with support verb *dostat*

T — resultative construction with verb *mít*

E — where (kdE)

F — diFference

G — reGard

H — Heritage

I — Intent

J — how (Jak)

A.2 Reflexivity

-- — no reflexive particle; no reciprocity

SE — reflexive tantum with particle *se* (*bát se*)

DE — derived reflexive with particle *se* (*vlévat se*)

se — reflexive with optional particle *se* (*koukat se*)

SI — reflexive tantum with particle *si* (*stěžovat si*)

DI — derived reflexive with particle *si* (*vynachválit si*)

si — reflexive with optional particle *si* (*myslet si*)

K — *reserved*

L — *reserved*

M — Means

N — Norm

O — from where (Odkud)

P — intent (Purpose, aim)

A.3 Subject

s — subject; the attribute in brackets shows the type of the subject and its value points to functor which is currently the subject

i — inherent

a — raised

Q — *reserved*

R — compaRison

S — Substitution

A.4 Functors

1 — Actor

2 — Patient

3 — Addressee

4 — Origin

5 — Effect

0 — no functor; used in frames of raising verbs

A — direction where (kAm)

B — Beneficiary

C — Cause

D — how long (jakDlouho)

T — criTerion

U — which way (kUdy)

V — accompaniment (průVod)

W — *reserved*

X — eXtent

Y — when (kdY)

Z — from when (Zekdy)

A.5 Grammatemes

h — ‘semantic’ features

- P — person
- T — thing, animal
- S — short reflexive pronoun *se* or *si*
- R — long reflexive pronoun *sebe*, *sobě*, etc.
- Z — interrogative pronoun *co* (what), demonstrative pronoun *to* (that), *všechno* (everything), etc.
- G — general participant (used in irregular passive frames and in generated sentence patterns)
- E — deleted (empty, erased) participant (used in generated sentence patterns)
- C — direct speech
- Q — quality (adjective)
- M — quantity (number, figure)
- L — location (adverb)
- A — direction where (adverb)
- F — direction from where (adverb)
- D — which way (adverb)
- W — when (adverb)

c — case

- 1 — Nominative
- 2 — Genitive
- 3 — Dative
- 4 — Accusative
- 6 — Locative
- 7 — Instrumental

r — preposition

n — number

- S — singular
- P — plural

s — clause

- I — infinitive
- C — conjunction *až*
- D — conjunction *že*
- F — conjunction *jestli*, *zda*
- P — conjunction *ať*
- R — relative expression *co*, *který*, *kdo*, ...

U — conjunction *aby*

Z — conjunction *jak*

l — required lemma

e — negation of a clause

A — affirmative (default)

N — negative

x — reciprocal coreference; the value points to a coindexed functor

a — subject raised to object position; the value points to the embedded clause from which the subject was raised

q — subject- or object-control

p — “patient” control

t — “addressee” control

d — diatheses of embedded infinitive; the values are identical with values of the “main” frame

m — modality

D — debitive (*mušet*)

H — hortative (*mít*)

V — volitive (*chtít*)

P — possibilitive (*moci*)

R — permissive (*smět*)

F — facultative (*dověst*)

A.6 Obligatority

[] — obligatory participant

() — obligatory inner participant which can be realized as general, or obligatory and deletable free modification

< > — optional participant

A.7 Diatheses

% — periphrastic passive is possible (*číst*, *stavět*)

\$ — reflexive passive is possible (*číst*, *mluvit*, *jít*)

@ — no passive (*bát se*)

— constructions with *mít* (*slíbit*)

* — constructions with *dostat* (*vyndat*)

~ — constructions with resultative *mít* (*uvařit*)

B Possible functors assigned to grammatememes

B.1 Abbreviations used in lists of possible functors

X — Unknown functor; mostly error in source data.	MEANS — Means.
PAT — Any participant. The reason why we chose this abbreviation is purely technical and it was explained in footnote 15 in Section 4.	NORM — Norm.
KAM — Direction ‘to’.	ODKUD — Direction ‘from’.
BEN — Beneficiary.	PURP — Purpose.
CAUSE — Cause.	COMPAR — Comparison.
JAKDL — Temporal modification ‘how long’.	SUBST — Substitution.
KDE — Location ‘where’.	CRIT — Criterion.
DIFF — Difference.	KUDY — Direction ‘which way’.
REGARD — Regard.	ACCOMP — Accompaniment.
HER — Heritage.	EXTENT — Extent.
INT — Intent.	KDY — Temporal modification ‘when’.
JAK — Manner.	ZEKDY — Temporal modification ‘from when’.

B.2 Lists of functors attached to every surface realization

Functors in parentheses are only taken in consideration if the surface realization has no variants. For example the prepositional case Accusative+*na* is typically a surface realization of direction, but in the frame of the verb *spoléhat na koho/co* (rely on sb/st) it is Patient.

The order of surface realization is important. A realization which is higher is listed first in brackets with variants and it are taken as a “representant” of the whole frame member.

hPc2	PAT
hPTc2	PAT
hTc2	PAT
v{eN}hTc2	PAT
hPc4	PAT
hPTc4	PAT
hTc4	PAT
sD	PAT
sF	PAT
sP	PAT
sPeN	PAT
sR	PAT
sUeN	PAT
sZ	PAT
sI	PAT INTENT KAM
sU	PAT PURP
sC	JAKDL
hA	KAM
hF	ODKUD
hL	KDE
hM	PAT
hPc1	PAT
hPTc1	PAT
hQc1	PAT
hQc7	PAT

hPc3 PAT
hPTc3 PAT
hTc3 (PAT) PURP
hPc7 PAT JAK
hPTc7 (PAT) MEANS SUBST
hRc7 PAT
hTc7 (PAT) MEANS CAUSE
hMr{na} PAT
hMr{o} DIFF
hMr{za} MEANS
hAr{do} KAM
hAr{na} KAM
hPc3r{vůči} PAT
hPc4r{o} PAT
hPc6r{o} PAT
hPc6r{po} PAT
hPc6r{při} PAT
hPc6r{v} KDE
hPc7r{mezi} (PAT) MEANS KDE
hPc7r{za} KAM
hPTc1r{jako} JAK
hPTc2r{bez} (PAT) JAK
hPTc2r{do} (PAT) KAM
hPTc2r{místo} SUBST
hPc2r{u} KDE
hPTc2r{u} (PAT) KDE
hPTc2r{vedle} KAM
hPc2r{kolem} KDE
hPc3r{proti} PAT KAM BEN
hPc7r{nad} (PAT) KDE CAUSE
hPc7r{pod} KDE
hPc7r{před} (PAT) KDE
hPTc2r{kolem} KDE KUDY
hPTc2r{od} (PAT) ODKUD
hPc2r{od} (PAT) ODKUD
hPc2r{z} PAT
hPTc2r{z} (PAT) ODKUD
hPc6r{na} PAT
hPc3r{ke} (PAT) KAM
hPc2r{do} PAT
hPc4r{mezi} (PAT) KAM
hPc4r{nad} KAM
hPc4r{na} PAT BEN
hPc4r{před} PAT
hPTc3r{ke} (PAT) KAM
hPTc3r{k} (PAT) KAM
hPc4r{za} PAT SUBST
hPc4r{pro} PAT BEN
hPc7r{s} (PAT) ACCOMP
hPTc3r{kvůli} CAUSE
hPTc3r{proti} PAT
hPTc4r{jako} (PAT) JAK
hPTc4r{mezi} KAM
hPTc4r{nad} (PAT) KAM KDE JAK
hPTc4r{na} (PAT) KAM
hPTc4r{o} PAT

hPTc4r{pod}	KAM
hPTc4r{pro}	PAT BEN
hPTc4r{před}	KAM
hPTc4r{přes}	KAM KUDY
hPTc4r{v}	PAT
hPTc4r{za}	(PAT) SUBST
hPTc6r{na}	(PAT) KDE
hPTc6r{o}	PAT
hPTc6r{po}	(PAT) HER KAM
hPTc6r{při}	PAT
hPTc6r{v}	(PAT) KDE
hPTc7r{mezi}	(PAT) KDE KUDY
hPTc7r{nad}	(PAT) KDE KUDY CAUSE
hPTc7r{pod}	KDE
hPTc7r{před}	(PAT) PURP
hPTc7r{s}	(PAT) ACCOMP
hPTc7r{za}	(PAT) KDE KAM
hRc2r{od}	ODKUD
hRc2r{ze}	PAT
hRc3r{k}	KAM
hRc4r{pod}	KAM
hRc4r{pro}	JAK
hRc4r{ze}	ODKUD
hRc7r{mezi}	PAT
hRc7r{před}	JAK
hRc7r{s}	JAK
hTc2r{bez}	JAK ACCOMP
hTc2r{během}	KDY
hTc2r{do}	(PAT) KAM
hTc2r{kolem}	KDE KUDY JAK
hTc2r{od}	(PAT) ODKUD JAKDL ZEKDY CAUSE
hTc2r{podle}	NORM CRIT
hTc2r{podél}	KDE KUDY
hTc2r{pomocí}	MEANS
hTc2r{u}	KDE
hTc2r{vedle}	KDE ACCOMP
hTc2r{z}	(PAT) ODKUD
hTc3r{kvůli}	CAUSE PURP
hTc3r{k}	PAT KAM PURP
hTc3r{proti}	PURP BEN
hTc3r{vzhledem k}	REGARD
hTc4r{jako}	(PAT) COMPAR JAK
hTc4r{mezi}	KAM
hTc4r{mimo}	KDE KAM
hTc4r{nad}	KAM
hTc4r{na}	(PAT) KAM PURP
hTc4r{o}	(PAT) KAM DIFF
hTc4r{pod}	KAM
hTc4r{po}	EXTENT JAKDL
hTc4r{pro}	(PAT) PURP CAUSE
hTc4r{před}	KAM
hTc4r{přes}	KAM KUDY MEANS JAK
hTc4r{skrže}	KUDY
hTc4r{skrz}	KUDY
hTc4r{v}	(PAT) KAM
hTc4r{za}	(PAT) KAM JAK CAUSE

hTc6r{jako v}	JAK
hTc6r{na}	(PAT) KDE JAK
hTc6r{o}	PAT KDY JAK
hTc6r{po}	KDY KAM KUDY JAK CRIT
hTc6r{při}	KDE KDY
hTc6r{v}	(PAT) ACCOMP KDE JAK
hTc7r{mezi}	(PAT) KUDY KDE
hTc7r{nad}	(PAT) KDE CAUSE
hTc7r{pod}	ACCOMP KDE KUDY CAUSE
hTc7r{před}	(PAT) KDY KDE KUDY
hTc7r{s}	(PAT) MEANS ACCOMP
hTc7r{za}	(PAT) KDE KAM
v{eN}hPTc4r{na}	PAT
v{eN}hTc2r{do}	KAM
v{eN}hTc3r{k}	PAT
hTc6	X
hPc3r{o}	X
hTc2r{v}	X
hTc7r{v}	X
hPTc4r{do}	X
hRc4r{do}	X
hRc4r{kolem}	X
hTc3r{v}	X
hTc4r{a}	X