



---

**The Prague Bulletin of Mathematical Linguistics**  
**NUMBER 110 APRIL 2018 71-84**

---

## **Search for the Relation of Form and Function Using the ForFun Database**

Marie Mikulová, Eduard Bejček, Eva Hajičová, Jarmila Panevová

Charles University, Faculty of Mathematics and Physics, Institute of Formal and Applied Linguistics,  
Prague, Czechia

---

### **Abstract**

The aim of the contribution is to introduce a database of linguistic forms and their functions built with the use of the multi-layer annotated corpora of Czech, the Prague Dependency Treebanks. The purpose of the Prague Database of Forms and Functions (ForFun) is to help the linguists to study the form-function relation, which we assume to be one of the principal tasks of both theoretical linguistics and natural language processing. We demonstrate possibilities of the exploitation of the ForFun database.

This article is largely based on a paper presented at the 16th International Workshop on Treebanks and Linguistic Theories in Prague (Bejček et al., 2017).

---

### **1. Introduction**

The study of the relation between (linguistic) forms and their functions or meanings (terms known from Saussure's structural linguistics (Saussure, 1916) as the relation between "signifié" and "signifiant") is one of the fundamental tasks of linguistics, with important implications for natural language understanding. As Katz (1966, p. 100) says, to understand the ability of natural languages to serve as an instrument to the communication of thoughts and ideas we must understand what it is that permits those who speak them consistently to connect the right sounds with the right meanings. This, however, is obviously not an easy task as the relation between form and function is a many-to-many relation. At present, the availability of richly annotated corpora helps the linguist to analyze the given relation in its variety, and it is a challenging task to provide linguists with useful tools for their study.

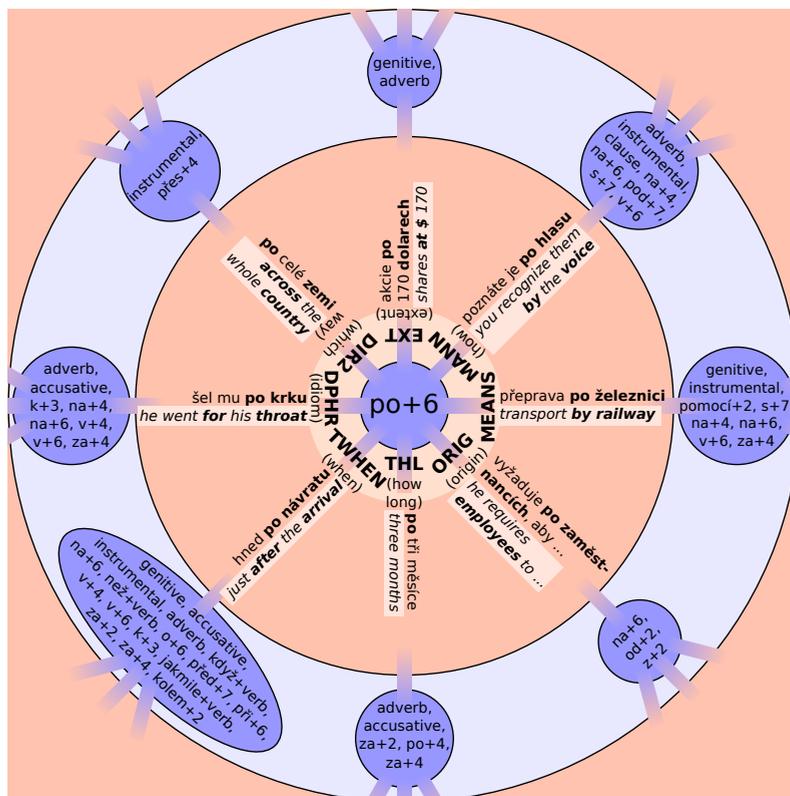


Figure 1. Many-to-many relation between forms and functions demonstrated on prepositional case *po* + Locative.

One of the most useful types of corpora for this task are treebanks based on a stratificational (multi-layer) approach, where the form-function relation may be understood as a relation between units of two layers of the system. The aim of this paper is to introduce a database of language forms and their linguistic functions built with the use of the multi-layer annotated corpora of Czech, the Prague Dependency Treebanks (PDTs), with the purpose to help the linguists to study the form-function relation. We offer a new database ForFun which gives a possibility to search in a user-friendly way all forms (almost 1 500 items) used in PDTs for particular functions and vice versa to look up all functions (66 items) expressed by the particular forms.

The research question we follow by constructing the database can be illustrated e.g. by the example of the Czech preposition *po* + Locative case of a noun (translated to English as *along, on, about, at, ... + noun*) in Figure 1. The dark colour indicates the

forms, the light colour the functions, identified in the PDTs by the functors attached to the nodes representing the given item (see below Section 2).<sup>1</sup> The prepositional case *po* + Locative (see the inner circle) may express the following eight functions (see the middle circle): *TWHEN* (when), *THL* (how long), *ORIG* (origin), *MEANS*, *MANN* (manner), *EXT* (extent), *DIR2* (direction which way), *DPHR* (idiomatic meaning). Each of these functions, in turn, may be expressed by a number of forms (see the outer circle) one of which is *po* + Locative. Thus for example, the function labelled *THL* (how long) may be expressed by an adverb, or Accusative of a noun (prepositionless case), or prepositional cases *za* + Genitive, *za* + Accusative, *po* + Accusative, and, of course, by the already mentioned *po* + Locative. In Figure 1, only a few functions of *po* + Locative are displayed; for a full list of 32 functions see their list in Table 3.

## 2. Multi-layer Architecture of Prague Dependency Treebanks

PDTs (on which our ForFun database is based) are complex linguistically motivated treebanks based on the dependency syntactic theory of the Functional Generative Description (see Sgall et al. 1986). The original annotation scheme has the following multi-layer architecture:<sup>2</sup>

- **morphological layer:** all tokens of the sentence get a lemma and a (disambiguated) morphological tag,
- **surface syntax layer** (analytical): a dependency tree capturing surface syntactic relations such as subject, object, adverbial; a (structural) tag reflecting these relations is attached to the nodes as one component of their (complex) labels,
- **deep syntax layer** (tectogrammatical) capturing the semantico-syntactic relations: on this layer, the dependency structure of a sentence is a tree consisting of nodes only for autonomous meaningful units (function words such as prepositions, subordinating conjunctions, auxiliary verbs etc. are not represented as separate nodes in the structure, their contribution to the meaning of the sentence is captured within the complex labels of the autonomous units). The types of dependency relations are captured by means of the so-called functors.

Functors (66 in total) are classified according to different criteria. The basic subdivision is based on the the valency criterion, which divides functors into the argument functors and adjunct functors. There are five arguments: Actor/Bearer (ACT), Patient (PAT), Addressee (ADDR), Origin (ORIG) and Effect (EFF). The repertory of adjuncts is

<sup>1</sup>Throughout the paper, we use the term *functor* for the label of the type of the dependency relation between the governor and its dependent; in the dependency tree structure representing the sentence on the deep (underlying, tectogrammatical; see Section 2) layer this label is a part of the complex label attached to the dependent node. The term *prepositional case* is used for a combination of a preposition and a noun or a nominal group in a morphological case. In the figures and tables, morphological cases are indicated by numbers, i.e. 2 for Genitive, 3 for Dative, 4 for Accusative, 6 for Locative, 7 for Instrumental. When the noun or nominal group is not accompanied by a preposition, we use the term *prepositionless case*.

<sup>2</sup>The PDTs annotation scenario is described in detail in Mikulová et al. (2006) and Hajič et al. (2017).

much larger than that of arguments. Their set might be divided into several subclasses, such as temporal (TWHEN for “when?”, TSIN for “since when?”, TTILL for “till when?”, TPAR for “during what time?”, THL for “how long?”, THO for “how often?”, TFHL for “for how long?”, TFRWH for “from when?”, and TOWH for “to when?”), local (LOC for “where?”, DIR1 for “where from?”, DIR2 for “which way?”, DIR3 for “where to?”), causal (CAUS for “cause”, AIM for “aim”, INTT for “intention”, COND for “condition”, CNCS for “concession”), functors for manner (MANN for general “manner”, MEANS for “means or instrument”), and other functors for other adjuncts (such as ACOMP for “accompaniment”, EXT for “extent”, INTF for “intensifier”, BEN for “benefactor”, etc.). For a full list of all dependency relations and their labels see Mikulová et al. (2006).

The nodes on a lower layer are explicitly referenced from the corresponding closest (immediately higher) layer. These links allow for tracing every unit of annotation all the way down to the original raw text. For the ForFun database, we use the annotations of the nodes on the deep syntactic layer and their counterparts on the morphological layer, which has made it possible to retrieve the relations between functions (expressed on the deep layer by functors) and forms and vice versa.

### 3. List of available Prague Dependency Treebanks

For Czech, the following four treebanks are available, each of them contains data of a different source. The Prague Dependency Treebank version 3.5 (PDT 3.5),<sup>3</sup> the newest edition of the core Prague Dependency Treebank, consists of articles from Czech daily newspapers. A slightly modified scenario was used for the annotation of the Prague Czech-English Dependency Treebank 2.0 (PCEDT 2.0),<sup>4</sup> the Prague Dependency Treebank of Spoken Czech 2.0 (PDTSC 2.0),<sup>5</sup> and the PDT-Faust corpus. In contrast to the original PDT project, in these treebanks, the morphological and surface syntactic annotations were done automatically, and the manually annotated deep syntactic layer does not contain some special annotations. However, the annotation of functors, which is important for our research of the form-function relation, has been done manually in all treebanks.

In the parallel PCEDT 2.0 (Hajič et al., 2012), the English part consists of the Wall Street Journal sections of the Penn Treebank (Marcus et al., 1993), and the Czech part, which is used in the ForFun database, was manually translated from the English original. PDTSC 2.0 (Mikulová et al., 2017b) contains dialogs from the Malach project<sup>6</sup> (slightly moderated testimonies of Holocaust survivors) and from the Companions

---

<sup>3</sup><https://ufal.mff.cuni.cz/pdt3.5>

<sup>4</sup><https://ufal.mff.cuni.cz/pcedt2.0/>

<sup>5</sup><https://ufal.mff.cuni.cz/pdtsc2.0>

<sup>6</sup><https://ufal.mff.cuni.cz/cvbm/vha-info.html>

project<sup>7</sup> (two participants chat over a collection of photographs). PDT-Faust is a small treebank containing short segments (very often with vulgar content) typed in by various users on the reverso.net webpage for translation.

It is obvious (see Table 1) that the Prague Dependency Treebank family provides rich language data for our purpose, i.e. for the study of the relation of forms and their functions since every content word there is assigned one of those 66 functors. Altogether, the treebanks contain around 180 000 sentences with their morphological, syntactic and semantic annotation.

	PDT 3.0	PCEDT 2.0	PDTSC 2.0	Faust	Total
Tokens	833 195	1 162 072	742 257	33 772	2 771 296
Sentences	49 431	49 208	73 835	3 000	175 474

Table 1. Volume of data in Prague Dependency Treebanks

## 4. Prague Database of Forms and Functions

ForFun 1.0, the Prague Database of Forms and Functions (Mikulová and Bejček, 2018), is a rich database of syntactic functions and their formal realizations with a large amount of examples coming from both written and spoken Czech texts. Since the database is extracted from the PDTs (see Section 3), it takes over the list of syntactic functions as well as the terminology (they are called *functors*).

ForFun is provided as a digital open source accessible to all scholars via the LINDAT/CLARIN repository.<sup>8</sup>

### 4.1. Design

We have already mentioned that in general the relation between forms and functions is a many-to-many relation. As such, it has to be explored from both sides: a given form has several functions and any of these functions may again be realized by several forms (the given one among them). When such relations have to be explored, ForFun is a perfect choice, since it is designed exactly for this kind of traversing through data.

Although the annotated example sentences are the same, they can be retrieved by asking either for their forms or for their functions. The ForFun database provides two entry points (cf. Figures 2 and 3):

<sup>7</sup>[http://cordis.europa.eu/project/rcn/96289\\_en.html](http://cordis.europa.eu/project/rcn/96289_en.html)

<sup>8</sup><http://hdl.handle.net/11234/1-2542>

**do+2**

DIR3 (9415x)

PoS	corpus	examples	occurs
v (7414x)	FAUST		73
	PCEDT	<ul style="list-style-type: none"> <li>• Zpět v centru stihli šéfové v hotelu pár schůzek, aby se opět nalodili <b>do autobusu</b>. (do autobusu–autobus) ×</li> <li>• Rapanelli nedávno řekl, že vláda prezidenta Carlose Menema, který nastoupil <b>do úřadu</b> 8. července, cítí, že × významné snížení jistiny a úroku je jediný způsob, jak může být problém s dluhem vyřešen. (do úřadu–úřad)</li> <li>• Dostihová dráha míří od Chile přes Rakousko až <b>do Portugalska</b>. (do Portugalska–Portugalsko) ×</li> <li>• ...</li> </ul>	1703
	PDT	<ul style="list-style-type: none"> <li>• Dotace se promítají <b>do cen</b> energií, prodávaných ostatním spotřebitelům. (do cen–cena) ×</li> <li>• Drahá energie pak konečně donutí odběratele investovat <b>do úspor</b> paliv. (do úspor–úspora) ×</li> <li>• ...</li> </ul>	2034
	PDTSC		3604
	n (1618x)	FAUST	
	PCEDT		592

TTILL (1910x)

PoS	corpus	examples	occurs
n (209x)	FAUST		4
	PCEDT		108
	PDT		74
	PDTSC	<ul style="list-style-type: none"> <li>• To byla škola jenom <b>do páté třídy</b>. (do třídy–třída) ×</li> <li>• Mluvím o době <b>do devatenácti let</b>, kdy jsem dospívala a byla pořád ještě v Turnově. (do let–rok) ×</li> <li>• ...</li> </ul>	23
	adj (60x)	FAUST	
	PCEDT		42

Figure 2. A screenshot of the ForFun web interface: From Form to Function.

- The user can choose one of almost 1 500 formal realizations of sentence units (i.e. prepositionless and prepositional cases, subordinated and coordinate conjunctions, adverbs, infinitive and finite verb forms, etc.) and obtains all functions it can represent.
- The user can choose one of 66 syntactic functions (i.e. LOC, TTILL, CAUS etc.) and obtains all forms used to express it.

The view can be always switched from a list of forms to a list of functions of one of them and vice versa.

For each form-function relation there are plenty of examples in the form of a sentence with the highlighted expression representing the relation. All these examples are sorted by various criteria:

**DIR3**  
v (23386x)

form	corpus	examples	occurs
#adv (4357x)	FAUST	<ul style="list-style-type: none"> <li>• ti dva ředitelé vzhledli <b>nahoru</b> na střechu budovy opery (<b>nahoru</b>)</li> <li>• ...</li> </ul>	41
	PCEDT	<ul style="list-style-type: none"> <li>• R. Hormats říká, že "nikdo nechce, aby se Američané sbalili a odjeli <b>domů</b>". (<b>domů</b>)</li> </ul>	577
do#2 (7414x)	FAUST		73
	PCEDT	<ul style="list-style-type: none"> <li>• Zpět v centru stihli šéfové v hotelu pár schůzek, aby se opět naladili <b>do autobusu</b>. (<b>do autobusu</b>–autobus)</li> <li>• Rapanelli nedávno řekl, že vláda prezidenta Carlose Menena, který nastoupil <b>do úřadu</b> 8. července, cítí, že × významné snížení jistiny a úroku je jediný způsob, jak může být problém s dluhem vyřešen. (<b>do úřadu</b>–úřad)</li> <li>• Dostihová dráha míří od Chile přes Rakousko až <b>do Portugalska</b>. (<b>do Portugalska</b>–Portugalsko)</li> <li>• ...</li> </ul>	1703
	PDT	<ul style="list-style-type: none"> <li>• Dotace se promítají <b>do cen</b> energií, prodávaných ostatním spotřebitelům. (<b>do cen</b>–cena)</li> <li>• Drahá energie pak konečně donutí odběratele investovat <b>do úspor</b> paliv. (<b>do úspor</b>–úspora)</li> <li>• ...</li> </ul>	2034
	PDISC		3604
#vfin (55x)	PCEDT	<ul style="list-style-type: none"> <li>• "Nemůže udělat nic pro to, aby se dostala zpět <b>tam</b>, kde <b>byla</b>," říká její právník James Bierbower. (<b>tam byla</b>–být)</li> <li>• Stejně jako právníci v nepřátelském prostředí akvizice jde i dítě <b>tam</b>, kde <b>jsou</b> peníze. (<b>tam jsou</b>–být)</li> <li>• ...</li> </ul>	12
	PDT	<ul style="list-style-type: none"> <li>• "Já si myslím, že Martina má jít <b>tam</b>, kam <b>patří</b>, všechno chce svůj čas," říká maminka. (<b>tam patří</b>–patřit)</li> <li>• Chci vrátit právo <b>tam</b>, kde <b>bylo</b> před padesáti lety," říká poslanec Svoboda. (<b>tam bylo</b>–být)</li> <li>• 0 Až nyní jsem si uvědomil, že v tenise jsem se dostal <b>tam</b>, kam jsem chtěl. (<b>0</b>–dostat_se)</li> </ul>	10
	PDISC	<ul style="list-style-type: none"> <li>• Ať <b>jsem</b> přišla, kam <b>přišla</b>, nikdo mě nemohl zaskočit. (<b>jsem přišla</b>–přijít)</li> <li>• Podíváme se, kde nás to <b>zajímá</b>. (<b>zajímá</b>–zajímat)</li> </ul>	33
adj (441x)			
do#2 (1933x)	FAUST	<ul style="list-style-type: none"> <li>• Sledovací systém je zabudovaný <b>do pásu</b> za účelem vedení pásu schodů, který neustále táhne schody od spodního nástupiště zpět nahoru v nekonečné smyčce. (<b>do pásu</b>–pás)</li> </ul>	2
	PCEDT	<ul style="list-style-type: none"> <li>• Tvrdí, že mnoho vozidel zařazených <b>do tříd</b> komerčních lehkých nákladních vozů převezve ve skutečnosti více osob než nákladu, a tudíž by měla mít stejné bezpečnostní prvky jako auta osobní. (<b>do tříd</b>–třída)</li> <li>• Společnost Armstrong očekává uzavření prodeje jednotky barev koncem listopadu a prodej jednotky na koberec v prosinci, s příjmy zahrnutými <b>do výsledků</b> čtvrtého nebo prvního čtvrtletí. (<b>do výsledků</b>–výsledek)</li> <li>• Záliba televize v dramatických konfliktech podporuje nadměrné používání sloganů vyvolávaných <b>do megafonů</b>, × militantní gestikulace, obviňujících plakátů a dalších taktik působících na city. (<b>do megafonů</b>–megafon)</li> <li>• ...</li> </ul>	89
	PDT	<ul style="list-style-type: none"> <li>• Milevsko: jméno tesaře <b>do žuly</b> (<b>do žuly</b>–žula)</li> <li>• ...</li> </ul>	79
	PDISC	<ul style="list-style-type: none"> <li>• Ještě se vrátím k tomu, že táta byl <b>době</b> války pravníkem nasazený <b>do Německa</b>. (<b>do Německa</b>–Německo)</li> </ul>	23

Figure 3. A screenshot of the ForFun web interface: From Function to Form.

- the word class of the parent node,
- the particular forms for the function or particular functions for the form, and
- the source of data (written, spoken, translated texts and texts from internet).

The number of examples available in the database is displayed for each pair form + functor, or functor + word class, each combination functor + form + word class and each specified 4-combination (form + functor + word class + source). Either first ten

examples or all of them are displayed on demand. On top of that, examples can be also first filtered by their source, which allows the user to hide e.g. all forms used only in spoken language or use only sentences from written corpora.

An illustration of how the result of user's search for the functions of the prepositional case *do* + Genitive looks like is given in Figure 2. In the upper part of the screenshot of the ForFun web interface, there are 9 415 occurrences in all PDTs of the form *do* + Genitive representing the functor DIR3. The occurrences of *do* + Genitive are divided according to their heads (be it a v(erb) or a n(oun), see the first column); their distribution within particular treebank is given in the second column followed by real examples from the corresponding treebank. A few of them are displayed on demand whereas many (see the last column) stay hidden. In the lower part of Figure 2, the same form *do* + Genitive in the function TTILL is exemplified in the same style. Note that Figure 2 presents only a part of the full response obtained from the ForFun database for the given query. The other functions of *do* + Genitive (PAT, EXT, EFF and others) are also not included in this shortened sample. (The list of all functions expressed by *do* + Genitive is in Table 3.)

For the opposite direction "from function to form" see the screenshot in Figure 3, where (among others) the same sentences for *do* + Genitive as the functor DIR3 can be found searching for all representations of the functor DIR3. Other forms include a finite verb (#vfin) or an adverb (#adv).

## 4.2. Volume

The ForFun database contains 2.2 million examples altogether for all forms (and the same number from the function point of view), split approx. 3:1 between written and spoken text (see Table 2). Each example is one sentence long.<sup>9</sup> They can be examined from the function side (66 functors) or the form side (1 469 forms). All examples are split into 13.5 thousand of 4-combinations (form + functor + word class + source), each with 163 examples in average.

While the average number is high, median is only two examples. The reason is that there is a long tail of 4-combinations used very rarely. These occurrences with very low frequencies in the data are one of the main benefits of the large volume of database, but they have to be used carefully. Every result has to be always understood solely as an input for a subsequent research, as ForFun may contain errors (caused by annotators as well as speakers/writers) considering its volume.

---

<sup>9</sup>One sentence typically contains many different functions and serves for many examples (once for each of its parts).

examples from written text	1 608 061
examples from spoken text	593 400
examples altogether	2 201 461
number of functions	66
number of forms	1 469
number of 4-combinations	13 514
avg. examples for a function	33 355
avg. examples for a form	1 500
avg. examples for a 4-combination	163
max. number of examples for a function	490 121
max. number of examples for a form	370 586
max. number of examples for a 4-combination	97 469

Table 2. Volume of the ForFun database

## 5. Possibilities of the Exploitation of the ForFun Database

To display the richness of the material we work with, we present several examples connected with the studies of the form-function relation what the user can find out in the ForFun database.

### 5.1. Multi-functionality of Forms

A rather straightforward use of the ForFun database is to retrieve which functions can be expressed by the particular form and which forms can express the particular function. Table 3 contains seven prepositional cases with the highest number of functions they express: *na* + Accusative, *v* + Locative, *k* + Dative, *za* + Accusative, *do* + Genitive, and *po* + Locative (cf. Figure 1).

### 5.2. Functions with the Most Limited List of Forms

Table 4, by contrast with Table 3, displays those functions that are expressed by the smallest number of forms (not only prepositional cases, but also other possible forms). We can observe that the HER (heritage), CONTRD contradiction, and TFRWH (from-when) functions are expressed exclusively by a single form. E.g. functor HER (heritage) is expressed exclusively by the form *po* + Locative, but HER belongs to many functions (32 in total) which are expressed by *po* + Locative (cf. their list in Table 3).

prep.	number	list of functors
<i>na+4</i>	42	ACT ADDR AIM APP ATT BEN CAUS COMPL COND CPHR CPR CRIT DIFF DIR1 DIR3 DPHR EFF EXT ID INTF INTT LOC MANN MAT MEANS MOD ORIG PAT PREC REG RESL RESTR RHEM RSTR SUBS TFHL TFRWH THL TOWH TPAR TTILL TWHEN
<i>v+6</i>	36	ACMP ACT AIM APP ATT BEN CAUS COMPL COND CPR CRIT DE- NOM DIR2 DIR3 DPHR EFF EXT ID LOC MANN MAT MEANS MOD PAT PREC REG RESL RESTR RHEM RSTR SUBS TFHL THL THO TPAR TWHEN
<i>k+3</i>	34	ACMP ACT ADDR AIM APP ATT BEN CAUS COMPL CPHR CRIT DIR1 DIR2 DIR3 DPHR EFF EXT ID INTT LOC MANN PAR PAT PREC REG RESL RESTR RHEM RSTR TOWH TPAR TSIN TTILL TWHEN
<i>za+4</i>	33	ACMP ACT AIM APP BEN CAUS CNCS COMPL COND CPHR DIR1 DIR3 DPHR EFF EXT HER ID INTT LOC MANN MEANS ORIG PAT PREC REG RSTR SUBS TFHL TFRWH THL THO TPAR TWHEN
<i>na+6</i>	33	ACT ADDR AIM APP ATT BEN CAUS COND CPR CRIT DIR2 DIR3 DPHR EFF EXT ID INTT LOC MANN MEANS ORIG PAR PAT PREC REG RESL RESTR RSTR TFHL THO TOWH TPAR TWHEN
<i>do+2</i>	33	ADDR AIM APP ATT BEN COMPL COND CPHR DIR1 DIR3 DPHR EFF EXT INTT LOC MANN MEANS MOD OPER PAR PARTL PAT REG RESL RSTR TFHL THL THO TOWH TPAR TSIN TTILL TWHEN
<i>po+6</i>	32	ACT AIM APP CAUS COND CPR CRIT DIR2 DIR3 DPHR EXT HER ID INTT LOC MANN MAT MEANS ORIG PAR PAT REG RSTR SUBS TFHL THL THO TOWH TPAR TSIN TTILL TWHEN

Table 3. The prepositional cases with the highest number of functions.

### 5.3. Absolute Frequency of Forms and Functions (in both written and spoken texts)

An observation of frequency has an important place in the description of language because it quantifies linguistic choices made by speakers and writers. For each form and function, ForFun provides information about raw frequency in all PDTs as well as in each corpus separately. The users can search quickly and in a user-friendly way which formal means are the most frequent in Czech sentences and which ones are rarely used. See Table 5 for five most frequent prepositional cases in comparison with the class of adverbs and the clause with the conjunction *že* 'that'.

The users of ForFun can also find out the distribution of a particular function (various arguments or adjuncts) in the sentences. For both forms and functions, they can compare their absolute frequencies in written and spoken texts. In Table 6, the sub-

functor	meaning	list of forms	example
HER	heritage	<i>po+6</i>	<i>Podědila tu nemoc <b>po rodičích</b>. 'She inherited the disease <b>from parents</b>.'</i>
CONTRD	contradiction	<i>zatímco+verb</i>	<i>On byl jedináček, <b>zatímco</b> ona měla dvanáct dětí. 'He was an only child, <b>while</b> she had twelve children.'</i>
TFRWH	from when	<i>z+2</i>	<i>Ze kterého roku je tato fotka? '<b>From</b> which <b>year</b> is this photo?'</i>
TOWH	to when	<i>na+4; pro+4</i>	<i>Derby je vypsáno <b>na</b> 3. září. 'Derby is listed <b>on September</b> 3.'</i>
TSIN	since when	<i>od+2; z+2; adverb</i>	<i>V energetice pracuje <b>od roku</b> 1964. 'He has worked in energetics <b>since</b> 1964.'</i>
THO	how often	<i>adverb; Acc; Instr</i>	<i>Pořádáte přechod <b>každý rok</b>? 'Do you organize march <b>every year</b>?'</i>
TTILL	till when	<i>do+2; dokud+verb; adverb; než+verb</i>	<i>Smlouva nebyla <b>do dnes</b> podepsána. 'No contract has been signed <b>yet</b>.'</i>

Table 4. Functions with the most limited list of forms.

classification of the most frequent functors for adjuncts is presented in comparison of their presence in written and spoken texts. We see that spatial and temporal functors (see their list in Section 2) are by far the most frequently occurring adjunct types. Hypothetically, in a Czech text of 100 sentences, there would be 61 sentences containing an adjunct (or several different adjuncts) and out of these sentences there would be: 29 sentences with spatial functor(s), 26 with temporal functor(s), 12 with manner functor(s), 10 with causal functor(s) and 22 with other functor(s).

#### 5.4. Material for Detailed Linguistic Studies

In addition to valuable statistical data, the ForFun database provides an extremely rich material for detailed linguistic studies of individual language phenomena and for their description and classification, e.g., valency behavior, coordination/discourse relations, idioms and complex predicates, comparison of written and spoken texts, etc. The first linguistic studies based on the database analyze and subclassify the functors denoting space and time (Mikulová et al., 2017a, 2018). The studies perform a detailed description of subtle meanings of temporal and spatial adjuncts including a list of formal means with real examples coming from both written and spoken texts and as such demonstrate that ForFun can be used for fundamental linguistic research.

form	occurrences
<i>v</i> +6	51 682
<i>na</i> +4	22 444
<i>s</i> +7	19 747
<i>z</i> +2	19 502
<i>na</i> +6	17 870
adverb	93 824
<i>že</i> [ <i>that</i> ]+verb	26 831

Table 5. The most frequent prepositional cases.

sentences containing:	all texts	%	written texts	spoken texts
spatial functors	74 164	29	43 089	31 075
temporal functors	66 503	26	42 266	24 237
functors for manner	31 583	12	21 752	9 831
causal functors	26 569	10	18 022	8 547
other functors for adjuncts	50 425	20	35 967	14 458
no functor for adjuncts	99 564	39	60 060	39 504

Table 6. The frequency distribution of the selected group of functors

## 6. Conclusion

The ForFun database has been built as a rich and user-friendly resource for those researchers who (want to) use corpora in their everyday work and look for various occurrences of specific forms or patterns in relation to their syntactic functions etc. but they are not interested or just do not need to deal with various technical, formal and annotation issues. ForFun brings a rich and complex annotation in PDTs based on a sound linguistic theory closer to common researchers. It will be further developed, though it should be borne in mind that it is designed to provide only a limited number of most useful features, rather than a full interface to everything PDTs can offer. There are other complex tools for that<sup>10</sup> and ForFun does not aim to substitute them. In its simplicity and clarity, it is a user-friendly source of examples for various explorations especially in syntax.

<sup>10</sup>E.g. PML Tree Query <https://lindat.mff.cuni.cz/services/pm1tq/>, INESS Search <http://clarino.uib.no/iness>, etc.

## Acknowledgments

This article is largely based on a paper presented at the 16th International Workshop on Treebanks and Linguistic Theories in Prague (Bejček et al., 2017).

The research reported in the paper has been supported by the Czech Science Foundation under the projects GA17-12624S and GA17-07313S and by the LINDAT/CLARIN project of Ministry of Education, Youth and Sports of the Czech Republic (project LM2015071). This work has been using language resources developed, stored and distributed by the latter project (LM2015071).

## Bibliography

- Bejček, Eduard, Eva Hajičová, Marie Mikulová, and Jarmila Panevová. The Relation of Form and Function in Linguistic Theory and in a Multilayer Treebank. In *Proceedings of the 16th International Workshop on Treebanks and Linguistic Theories*, pages 56–63, 2017. URL <http://aclweb.org/anthology/W17-7609>.
- Hajič, Jan, Eva Hajičová, Jarmila Panevová, Petr Sgall, Ondřej Bojar, Silvie Cinková, Eva Fučíková, Marie Mikulová, Petr Pajas, Jan Popelka, Jiří Semecký, Jana Šindlerová, Jan Štěpánek, Josef Toman, Zdeňka Uřešová, and Zdeněk Žabokrtský. Announcing Prague Czech-English Dependency Treebank 2.0. In *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC 2012)*, pages 3153–3160. European Language Resources Association, Istanbul, Turkey, 2012. ISBN 978-2-9517408-7-7. URL <https://aclanthology.info/pdf/L/L12/L12-1280.pdf>.
- Hajič, Jan, Eva Hajičová, Marie Mikulová, and Jiří Mírovský. *Handbook on Linguistic Annotation*, chapter Prague Dependency Treebank, pages 555–594. Springer Verlag, Dordrecht, Netherlands, 2017.
- Katz, Jerrold J. *The philosophy of language*. Studies in languages. Harper & Row, New York, 1966.
- Marcus, Mitchell P., Beatrice Santorini, and Mary Ann Marcinkiewicz. Building a Large Annotated Corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330, 1993. URL <https://aclanthology.info/pdf/J/J93/J93-2004.pdf>.
- Mikulová, Marie and Eduard Bejček. ForFun 1.0: Prague Database of Syntactic Forms and Functions – An Invaluable Resource for Linguistic Research. In *Proceedings of the 11th International Conference on Language Resources and Evaluation (LREC 2018)*, Paris, France, 2018. European Language Resources Association.
- Mikulová, Marie, Alevtina Bémová, Jan Hajič, Eva Hajičová, Jiří Havelka, Veronika Kolářová, Lucie Kučová, Markéta Lopatková, Petr Pajas, Jarmila Panevová, Magda Razímová, Petr Sgall, Jan Štěpánek, Zdeňka Uřešová, Kateřina Veselá, and Zdeněk Žabokrtský. Annotation on the tectogrammatical level in the Prague Dependency Treebank. Annotation manual. Technical Report 30, ÚFAL MFF UK, Prague, Czech Rep., 2006.
- Mikulová, Marie, Eduard Bejček, Veronika Kolářová, and Jarmila Panevová. Subcategorization of Adverbial Meanings Based On Corpus Data. *Journal of Linguistics / Jazykovedný časopis*, 68(2):268–277, 2017a. ISSN 0021-5597.

- Mikulová, Marie, Jiří Mírovský, Anna Nedoluzhko, Petr Pajas, Jan Štěpánek, and Jan Hajič. PDTSC 2.0 – Spoken Corpus with Rich Multi-layer Structural Annotation. In *Text, Speech, and Dialogue 20th International Conference, TSD 2017*, Lecture Notes in Computer Science, pages 129–137, Cham / Heidelberg / New York / Dordrecht / London, 2017b. Charles University, Springer International Publishing. ISBN 978-3-319-64206-2.
- Mikulová, Marie, Eduard Bejček, and Jarmila Panevová. What Can We Find Out about Time and Space in the ForFun Database? In *Proceedings of the Second Workshop on Corpus-Based Research in the Humanities CRH-2*, Gerastree proceedings, pages 133–142. Austrian Academy of Science, Dept. of Geoinformation, Wien, Austria, 2018. ISBN 978-3-901716-43-0.
- Saussure, Ferdinand de. *Cours de linguistique générale*. C. Bally and A. Sechehaye, with the collaboration of A. Riedlinger, eds. Lausanne and Paris: Payot, 1916.
- Sgall, Petr, Eva Hajičová, and Jarmila Panevová. *The Meaning of the Sentence and Its Semantic and Pragmatic Aspects*. Academia/Reidel Publishing Company, Prague/Dordrecht, 1986.

**Address for correspondence:**

Marie Mikulová

mikulova@ufal.mff.cuni.cz

Charles University, Faculty of Mathematics and Physics,

Institute of Formal and Applied Linguistics

Malostranské náměstí 25, 118 00 Prague 1, Czech Republic