

Manual for Morphological Annotation

Instructions for Annotators

Jiří Hana, Hana Hanová

In cooperation with
Jan Hajič, Barbora Hladká, Emil Jeřábek

CKL MFF UK Praha



Contents

1	Introduction	9
2	Lemma and tag structure	11
2.1	Lemma structure	11
2.1.1	Derivational Information	11
2.1.2	Semantic Information	12
2.2	Tag Structure	12
2.2.1	Positional tags	12
2.2.1.1	1 – Part of speech	13
2.2.1.2	2 – Detailed part of speech	14
2.2.1.3	3 – Gender	18
2.2.1.4	4 – Number	19
2.2.1.5	5 – Case	19
2.2.1.6	6 – Possessor’s Gender	19
2.2.1.7	7 – Possessor’s Number	19
2.2.1.8	8 – Person	20
2.2.1.9	9 – Tense	20
2.2.1.10	10 – Degree of Comparison	20
2.2.1.11	11 – Negation	20
2.2.1.12	12 – Voice	21
2.2.1.13	15 – Variant	21
2.2.2	Compact tags	21
2.2.3	Informal abbreviations	21
3	Names	23
3.1	Personal names	24
3.1.1	von, van, etc.	24
3.1.2	Chinese names	24
3.1.3	Korean names	25
3.1.4	Foreignized Czech names	25
3.2	Compound names (names consisting of other names)	26
3.3	Horses, DJ’s etc.	27
3.4	Sport clubs, etc.	27
3.5	Other	28
3.5.1	Geographical names	28
3.5.2	Initials	29
3.5.3	Institutions, companies	29
3.5.4	Sporting and other events	30
3.5.5	Televisions	31
3.5.6	News, Magazines	31
3.5.7	Song names, etc.	31

4	Abbreviations	33
4.1	Gender	33
4.2	Normal abbreviations	33
4.3	Isolated letters	34
4.4	RM-systém, samopal SA-58	35
4.5	Units of measurements	35
4.6	Authors abbreviations	36
4.7	Academic titles	36
5	Colloquial Czech	37
5.1	Cos, jaks, kdys	37
5.2	Suffix -é in plural of neuter	37
6	Foreign words and phrases	39
6.1	Part of speech	39
6.1.1	English noun clusters	39
6.1.2	Examples	39
6.2	Articles	40
6.2.1	Citation use	40
6.2.2	Word use	40
6.2.3	Examples	40
6.3	Nouns	40
6.3.1	Citation use	40
6.3.1.1	English	40
6.3.2	Word use	41
6.3.2.1	English	41
6.4	Verbs	41
6.4.1	Citation use	41
6.4.1.1	English verbs	41
6.4.2	Word use	41
6.5	Slovak language	41
7	Errors	43
7.1	Characters	43
7.2	Separators, etc.	43
8	Hard to decide	45
8.1	až	45
8.2	jak	45
8.3	málo	47
8.4	moc	47
8.5	proto	48
8.6	svůj	48
8.7	tak	48
9	Sólokapři	51
9.1	Date and time	52
9.2	Numbers, numerals and quantifiers	52
9.3	Hyphenated composites	52

10 Insertion	55
10.1 Possessive adjectives	55
10.2 Words ending with -ismus, -izmus	55
10.3 Strange and unique things	56
10.4 Other	57
11 Errors in PDT 1.0	59

PREFACE

We are pleased to publish the first version of the manual for morphological annotation of Czech sentences. We believe that such guidelines can be of use to the users of Prague Dependency Treebank 1.0 (PDT 1.0), as well as for preparation of new data.

Let us recall the most important steps we passed in order to get about two million morphologically annotated words (PDT 1.0). At the very beginning, we put together a team of eight annotators – we did introduce them to a system of morphological tags we designed to describe Czech morphological properties; we also introduced them a morphological analyzer for processing isolated words we use (as a preprocessing step), and, last but not least, we did rely on their knowledge of Czech morphology they have acquired while studying at secondary school, i.e. we did not offer them any annotation guidelines.

One can assume that this strategy is too hazardous – how to deal with discrepancies the annotators produce to ensure the consistency of annotation? First, two annotators annotated each text file. Then, by a “blind” automatic procedure (no matter what word is processed – just comparing two strings) we detected words annotated differently. Consequently, the only one annotator (as a member of just two-member team) handled these cases and, also, checked the morphological annotations against the syntactic-analytical annotations. This way we replaced the absence of annotation guidelines by sequential elimination of discrepancies across both the morphological and syntactic-analytical levels of annotation.

Along the way we were writing this annotation manual. It is not intended as a comprehensive guide to the morphological annotation of Czech sentences (in contrast to the manual for syntactic-analytical annotations). The authors concentrate “only” on those cases which caused the most ambiguities and problems while annotating PDT 1.0. The ongoing effort is directed to the treating of not-yet-solved problematic cases in accord with the conventions of automatic morphological analyzer.

The morphological annotation of PDT 1.0 was carried out in the framework of experimental verification of the definition of formal representation of the analysis of Czech sentences (the project GAČR 405/96/0198, “Formal representation of language structures”). The material obtained in this way (data) is used in many domains of research in computational linguistics, above all as basic (training) data in projects of the automatic language analysis, the MŠMT research project MSM113000006, the “Laboratory for Language Data Processing” (the MŠMT project VS961510) and the Center for Computational Linguistics (the MŠMT project LN00A063). These data have been also used as verification material for various partial projects within the complex program GAČR 405/96/K214 (“Czech Language in Computer Age”). The “Center for Computational Linguistics” project financially supported work on these morphological annotation guidelines.

We are grateful to Petr Pajas – this document “as it is” would not appear without his XML and LaTeX skills.

Typographical conventions

Vertical bar on the outer side of the page is used to highlight comments we make or suggestions we propose.

Gray is used to highlight something what should be checked.



Chapter 1

Introduction

Sometimes, the writer uses the word incorrectly – e.g. a name of a woman as a name of a man, surname as a first name, etc. it is necessary to annotate the real usage not the should-be usage.

Maybe it should be somehow marked, if we encounter it.

To get an idea what a foreign name, etc. mean it is useful to try to find using an internet portal, in an encyclopedia, on a map, etc. During annotation, we have found the following internet links useful:

Portals

<http://www.seznam.cz> – for Czech products, companies

<http://search.seznam.cz/search.cgi?mod=f&hlp=y> – for Czech companies

<http://www.google.com>

<http://www.altavista.com> (shop section for various searching products)

Encyclopedias

<http://www.britannica.com>

<http://www.encyclopedia.com>

<http://www.encarta.msn.com>

Dictionaries

<http://dictionary.oed.com/entrance.dtl> – Oxford English Dictionary

<http://slovník.seznam.cz> – various dictionaries

Maps

<http://mapy.atlas.cz> – Czechia

<http://www.mapquest.com/maps> – U.S.A and the world

Chapter 2

Lemma and tag structure

2.1 Lemma structure

Lemma in PDT 1.0 has two parts. First part, the lemma proper, has to be a unique identifier of the lexical item. Usually it is the base form (e.g. infinitive for a verb) of the word, possibly followed by a number distinguishing different lemmas with the same base forms. Second part (optional) is not part of the identifier and contains additional information about the lemma, e.g. semantic or derivational information.

Note: There is a convention that if lemmas use numbers to distinguish lexical items with the same base form, they all have to use them- i.e. instead of sets of lemmas {X, X-1, X-2} or {X, X-2, X-3}, there should be a set {X-1, X-2, X-3}

Note: The lemmas having different semantic suffixes should have different numbers.

In this manual we behave as the annotator. We try to mark such improper numbers by roman font (other part of the lemma is in italics). For example *stop* in *akce Stop million* will be marked as *stop-1*;m and not *stop-1*;m).

Table 2.1: Examples

Whole lemma	Lemma proper	Second part
<i>Chemik</i>	<i>chemik</i>	
<i>maso</i> ^(jídlo_apod.)	<i>maso</i>	^(jídlo_apod.)
<i>Bonn</i> ;G	<i>Bonn</i>	;G
<i>vazba-1</i> ^(obviněného)	<i>vazba-1</i>	^(obviněného)
<i>vazba-2</i> ^(spojení)	<i>vazba-2</i>	^(spojení)
<i>Martinův-1</i> ;Y^(*4-1)	<i>Martinův-1</i>	;Y^(*4-1)

2.1.1 Derivational Information

The morphological component used in PDT 1.0, handles only inflection, not derivations – it means lemmas are rather shallow. However, sometimes the lemma contains information about lemmas it is derived. For example lemmas of possessive adjectives contain information about the noun they are derived from (*otcův* ← *otec*). The information is encoded in the following way – how many characters you have to remove from the end, and what string you have to add to get the deeper lemma. Only the proper lemmas are both input and output of this process.

Following examples illustrate this:

kardinál $v_{-}S_{-}^{(*2)}$ – remove two letters: *kardinál*

Karl $l_{-}Y_{-}^{(*3el)}$ – remove 3 characters, add "el": *Karel*

přijetí $-2_{-}^{(např._{.}návřh)}_{-}^{(*5mout-2)}$ – remove 5 characters, add "mout-2": *přijmout-2*

Martin $l_{-}Y_{-}^{(*4-1)}$ – remove 4 characters, add "-1": *Martin-1*

Other examples:

Soros $l_{-}S_{-}^{(*2)}$

chlapc $l_{-}^{(*3ec)}$

Mách $l_{-}S_{-}^{(*2a)}$

Hlink $l_{-}S_{-}^{(*4a-1)}$

podání $_{-}^{(něco_{-}[někomu]_{-}[někam])}_{-}^{(*3at)}$

prohlášení $_{-}^{(*4sit)}$

protiprávnost $_{-}^{(*3ý)}$

2.1.2 Semantic Information

Some lemmas (esp. names) contain suffixes expressing semantic information about their use, etc.:

G – geographical name: *Praha, Ústí nad Labem*

Y – given (first) name, formerly used as default: *Petr, John*

S – surname (last name): *Dvořák, Zelený, Agassi, Bush*

E – name of a nationality: *Čech, Kolumbijec*

R – name of a product: *Tatra* (the car),

K – name of a company: *Tatra* (the company)

m – default – names of mines, stadiums, guerilla bases, etc; also used for functional words in names.

2.2 Tag Structure

2.2.1 Positional tags

A positional tag is a string of 15 characters. Every position encodes one morphological category using one character (mostly upper case letters or numbers).

Position	Name	Description
1	POS	Part of speech
2	SubPOS	Detailed part of speech
3	Gender	Gender
4	Number	Number
5	Case	Case
6	PossGender	Possessor's gender

continued on next page

continued from previous page

Position	Name	Description
7	PossNumber	Possessor's number
8	Person	Person
9	Tense	Tense
10	Grade	Degree of comparison
11	Negation	Negation
12	Voice	Voice
13	Reserve1	Reserve
14	Reserve2	Reserve
15	Var	Variant, style

Some of the characters encode aggregation of more atomic values – for example: 'X' – means any value, 'Y' means masculine animate ('M') or inanimate ('I'). Dash ('-') means no value (e.g. tense for nouns).

Not all combinations of tag values are possible. There is about 4K tags¹.

Examples:

hraniční: AAIS4-----1A----- standard adjective, masc. inanimate, singular, accusative, positive

potok: NNIS4-----A----- noun, masc. inanimate, singular, accusative, positive

karikaturistou: NNMS7-----A----- noun, masc. animate, singular, instrumental, positive

ODS: NNFXX-----A---8 noun, feminine, any number, any case, positive, abbreviation

podle: RR--2----- preposition (non vocalized), requiring genitive

volen: V_SYS---XX-AP--- verb, passive participle, masculine, singular, any person, any tense, positive, passive

2.2.1.1 1 – Part of speech

Value	Description
A	Adjective
C	Numeral
D	Adverb
I	Interjection
J	Conjunction
N	Noun

continued on next page

¹See also: http://ufal.mff.cuni.cz/pdt/Morphology_and_Tagging/Doc/docc0pos.pdf, for quick reference: http://ufal.mff.cuni.cz/pdt/Morphology_and_Tagging/Doc/hmptagqr.html.

continued from previous page

Value	Description
P	Pronoun
V	Verb
R	Preposition
T	Particle
X	Unknown, Not Determined, Unclassifiable
Z	Punctuation (also used for the <i>Sentence Boundary</i> token)

2.2.1.2 2 – Detailed part of speech

Further subcategorizes POS. The POS value is uniquely specified by SubPOS value.

Table 2.4: SUBPOS

Value	Description	POS
#	Sentence boundary	
	Word <i>krát</i> (lit.: <i>times</i>)	C – numeral
,	Conjunction subordinate (incl. <i>aby, kdyby</i> in all forms)	J – conjunction
}	Numeral, written using Roman numerals (XIV)	C – numeral
:	Punctuation (except for the virtual sentence boundary word ###, which uses the 2.4 SUBPOS #)	
=	Number written using digits	C – numeral
?	Numeral <i>kolik</i> (lit. <i>how many/how much</i>)	C – numeral
@	Unrecognized word form	X – unknown
^	Conjunction (connecting main clauses, not subordinate)	J – conjunction
4	Relative/interrogative pronoun with adjectival declension of both types (<i>soft</i> and <i>hard</i>) (<i>jaký, který, čím, ...</i> , lit. <i>what, which, whose, ...</i>)	P – pronoun
5	The pronoun <i>he</i> in forms requested after any preposition (with prefix <i>n-</i> : <i>něj, něho, ...</i> , lit. <i>him</i> in various cases)	P – pronoun
6	Reflexive pronoun <i>se</i> in long forms (<i>sebe, sobě, sebou</i> , lit. <i>myself / yourself / herself / himself</i> in various cases; <i>se</i> is personless)	P – pronoun

continued on next page

continued from previous page

Value	Description	POS
7	Reflexive pronouns <i>se</i> (2.8 CASE = 4), <i>si</i> (2.8 CASE = 3), plus the same two forms with contracted -s: <i>ses</i> , <i>sis</i> (distinguished by 2.11 PERSON = 2; also number is singular only) This should be done somehow more consistently, virtually any word can have this contracted -s (<i>cos</i> , <i>polívkus</i> , ...)	P – pronoun
8	Possessive reflexive pronoun <i>svůj</i> (lit. <i>my/your/her/his</i> when the possessor is the subject of the sentence)	P – pronoun
9	Relative pronoun <i>jenž</i> , <i>již</i> , ... after a preposition (<i>n-</i> : <i>něhož</i> , <i>niž</i> , ..., lit. <i>who</i>)	P – pronoun
A	Adjective, general	A – adjective
B	Verb, present or future form	V – verb
C	Adjective, nominal (short, participial) form <i>rád</i> , <i>schopen</i> , ...	A – adjective
D	Pronoun, demonstrative (<i>ten</i> , <i>onen</i> , ..., lit. <i>this</i> , <i>that</i> , <i>that ... over there</i> , ...)	P – pronoun
E	Relative pronoun <i>což</i> (corresponding to English <i>which</i> in subordinate clauses referring to a part of the preceding text)	P – pronoun
F	Preposition, part of; never appears isolated, always in a phrase (<i>nehledě</i> (<i>na</i>), <i>vzhledem</i> (<i>k</i>), ..., lit. <i>regardless</i> , <i>because of</i>)	R – preposition
G	Adjective derived from present transgressive form of a verb	A – adjective
H	Personal pronoun, clitical (short) form (<i>mě</i> , <i>mi</i> , <i>ti</i> , <i>mu</i> , ...); these forms are used in the second position in a clause (lit. <i>me</i> , <i>you</i> , <i>her</i> , <i>him</i>), even though some of them (<i>mě</i>) might be regularly used anywhere as well	P – pronoun
I	Interjections	I – interjection
J	Relative pronoun <i>jenž</i> , <i>již</i> , ... not after a preposition (lit. <i>who</i> , <i>whom</i>)	P – pronoun
K	Relative/interrogative pronoun <i>kdo</i> (lit. <i>who</i>), incl. forms with affixes -ž and -s (affixes are distinguished by the category 2.16 VAR (for -ž) and 2.11 PERSON (for -s))	P – pronoun
L	Pronoun, indefinite <i>všechmen</i> , <i>sám</i> (lit. <i>all</i> , <i>alone</i>)	P – pronoun
M	Adjective derived from verbal past transgressive form	A – adjective
N	Noun (general)	N – noun
O	Pronoun <i>svůj</i> , <i>nesvůj</i> , <i>tentam</i> alone (lit. <i>own self</i> , <i>not-in-mood</i> , <i>gone</i>)	P – pronoun

continued on next page

continued from previous page

Value	Description	POS
P	Personal pronoun <i>já, ty, on</i> (lit. <i>I, you, he</i>) (incl. forms with the enclitic <i>-s</i> , e.g. <i>tys</i> , lit. <i>you're</i>); gender position is used for third person to distinguish <i>on/ona/ono</i> (lit. <i>he/she/it</i>), and number for all three persons	P – pronoun
Q	Pronoun relative/interrogative <i>co, copak, cožpak</i> (lit. <i>what, isn't-it-true-that</i>)	P – pronoun
R	Preposition (general, without vocalization)	R – preposition
S	Pronoun possessive <i>můj, tvůj, jeho</i> (lit. <i>my, your, his</i>); gender position used for third person to distinguish <i>jeho, její, jeho</i> (lit. <i>his, her, its</i>), and number for all three pronouns	P – pronoun
T	Particle	T – particle
U	Adjective possessive (with the masculine ending <i>-ův</i> as well as feminine <i>-in</i>)	A – adjective
V	Preposition (with vocalization <i>-e</i> or <i>-u</i>): (<i>ve, pode, ku, ...</i> , lit. <i>in, under, to</i>)	R – preposition
W	Pronoun negative (<i>nic, nikdo, nijaký, žádný, ...</i> , lit. <i>nothing, nobody, not-worth-mentioning, no/none</i>)	P – pronoun
X	(temporary) Word form recognized, but tag is missing in dictionary due to delays in (asynchronous) dictionary creation	
Y	Pronoun relative/interrogative <i>co</i> as an enclitic (after a preposition) (<i>oč, nač, zač</i> , lit. <i>about what, on/onto what, after/for what</i>)	P – pronoun
Z	Pronoun indefinite (<i>nějaký, některý, číkoli, cosi, ...</i> , lit. <i>some, some, anybody's, something</i>)	P – pronoun
a	Numeral, indefinite (<i>mnoho, málo, tolik, několik, kdovíkolik, ...</i> , lit. <i>much/many, little/few, that much/many, some (number of), who-knows-how-much/many</i>)	C – numeral
b	Adverb (without a possibility to form negation and degrees of comparison, e.g. <i>pozadu, naplocho, ...</i> , lit. <i>behind, flatly</i>); i.e. both the 2.14 NEGATION as well as the 2.13 GRADE attributes in the same tag are marked by – (Not applicable)	D – adverb
c	Conditional (of the verb <i>být</i> (lit. <i>to be</i>) only) (<i>by, bych, bys, bychom, byste</i> , lit. <i>would</i>)	V – verb
d	Numeral, generic with adjectival declension (<i>dvojí, desaterý, ...</i> , lit. <i>two-kinds/..., ten-...</i>)	C – numeral
e	Verb, transgressive present (endings <i>-e/-ě, -íc, -íce</i>)	V – verb

continued on next page

continued from previous page

Value	Description	POS
f	Verb, infinitive	V – verb
g	Adverb, forming negation (2.14 NEGATION set to A/N) and degrees of comparison 2.13 GRADE set to 1/2/3 (comparative/superlative), e.g. <i>velký, za\-jí\ -ma\ -vý, ..., lit. big, interesting</i>	
h	Numeral, generic; only <i>jedny</i> and <i>nejedny</i> (lit. <i>one-kind/sort-of, not-only-one-kind/sort-of</i>)	C – numeral
i	Verb, imperative form	V – verb
j	Numeral, generic greater than or equal to 4 used as a syntactic noun (<i>čtvero, desatero, ..., lit. four-kinds/sorts-of, ten-...</i>)	C – numeral
k	Numeral, generic greater than or equal to 4 used as a syntactic adjective, short form (<i>čtvery, ..., lit. four-kinds/sorts-of</i>)	C – numeral
l	Numeral, cardinal <i>jeden, dva, tři, čtyři, půl, ...</i> (lit. <i>one, two, three, four</i>); also <i>sto</i> and <i>tisíc</i> (lit. <i>hundred, thousand</i>) if noun declension is not used	C – numeral
m	Verb, past transgressive; also archaic present transgressive of perfective verbs (ex.: <i>udělav</i> , lit. <i>(he-)having-done</i> ; arch. also <i>udělaje</i> (2.16 VAR = 4), lit. <i>(he-)having-done</i>)	V – verb
n	Numeral, cardinal greater than or equal to 5	C – numeral
o	Numeral, multiplicative indefinite (<i>-krát</i> , lit. <i>(times): mnohokrát, tolikrát, ..., lit. many times, that many times</i>)	C – numeral
p	Verb, past participle, active (including forms with the enclitic <i>-s</i> , lit. <i>'re (are)</i>)	V – verb
q	Verb, past participle, active, with the enclitic <i>-ě</i> , lit. (perhaps) <i>-could-you-imagine-that?</i> or <i>but-because-</i> (both archaic)	V – verb
r	Numeral, ordinal (adjective declension without degrees of comparison)	C – numeral
s	Verb, past participle, passive (including forms with the enclitic <i>-s</i> , lit. <i>'re (are)</i>)	V – verb
t	Verb, present or future tense, with the enclitic <i>-ě</i> , lit. (perhaps) <i>-could-you-imagine-that?</i> or <i>but-because-</i> (both archaic)	V – verb
u	Numeral, interrogative <i>kolikrát</i> , lit. <i>how many times?</i>	C – numeral
v	Numeral, multiplicative, definite (<i>-krát</i> , lit. <i>times: pětkrát, ..., lit. five times</i>)	C – numeral

continued on next page

continued from previous page

Value	Description	POS
w	Numeral, indefinite, adjectival declension (<i>nejeden, tolikátý, ..., lit. not-only-one, so-many-times-repeated</i>)	C – numeral
y	Numeral, fraction ending at <i>-ina</i> ; used as a noun (<i>pětina, lit. one-fifth</i>)	C – numeral
z	Numeral, interrogative <i>kolikátý, lit. what (at-what-position-place-in-a-sequence)</i>	C – numeral

Table 2.5: Obsolete values

Value	Description
!	Abbreviation used as an adverb
.	Abbreviation used as an adjective
~	Abbreviation used as a verb
;	Abbreviation used as a noun
3	Abbreviation used as a numeral
x	Abbreviation, part of speech unknown/indeterminable

2.2.1.3 3 – Gender

Value	Description
F	Feminine
H	{F, N} – Feminine or Neuter
I	Masculine inanimate
M	Masculine animate
N	Neuter
Q	Feminine (with singular only) or Neuter (with plural only); used only with participles and nominal forms of adjectives
T	Masculine inanimate or Feminine (plural only); used only with participles and nominal forms of adjectives
X	Any
Y	{M, I} – Masculine (either animate or inanimate)
Z	{M, I, N} – Not feminine (i.e., Masculine animate/inanimate or Neuter); only for (some) pronoun forms and certain numerals

2.2.1.4 4 – Number

Value	Description
D	Dual , e.g. <i>nohama</i>
P	Plural, e.g. <i>nohami</i>
S	Singular, e.g. <i>noha</i>
W	Singular for feminine gender, plural with neuter; can only appear in participle or nominal adjective form with gender value Q
X	Any

2.2.1.5 5 – Case

Table 2.8: CASE

Value	Description
1	Nominative, e.g. <i>žena</i>
2	Genitive, e.g. <i>ženy</i>
3	Dative, e.g. <i>ženě</i>
4	Accusative, e.g. <i>ženu</i>
5	Vocative, e.g. <i>ženo</i>
6	Locative, e.g. <i>ženě</i>
7	Instrumental, e.g. <i>ženou</i>
X	Any

2.2.1.6 6 – Possessor's Gender

Value	Description
F	Feminine, e.g. <i>matčín, její</i>
M	Masculine animate (adjectives only), e.g. <i>otců</i>
X	Any
Z	{M, I, N} – Not feminine, e.g. <i>jeho</i>

2.2.1.7 7 – Possessor's Number

Value	Description
P	Plural, e.g. <i>náš</i>
S	Singular, e.g. <i>můj</i>

2.2.1.8 8 – Person

Table 2.11: PERSON

Value	Description
1	1st person, e.g. <i>píšu, píšeme</i>
2	2nd person, e.g. <i>píšeš, píšete</i>
3	3rd person, e.g. <i>píše, píšou</i>
X	Any person

2.2.1.9 9 – Tense

Value	Description
F	Future
H	{R, P} – Past or Present
P	Present
R	Past
X	Any

2.2.1.10 10 – Degree of Comparison

Table 2.13: GRADE

Value	Description
1	Positive, e.g. <i>velký</i>
2	Comparative, e.g. <i>větší</i>
3	Superlative, e.g. <i>největší</i>

2.2.1.11 11 – Negation

Table 2.14: NEGATION

Value	Description
A	Affirmative (not negated), e.g. <i>možný</i>
N	Negated, e.g. <i>nemožný</i>

2.2.1.12 12 – Voice

Value	Description
A	Active, e.g. <i>píšící</i>
P	Passive, e.g. <i>psaný</i>

2.2.1.13 15 – Variant

Table 2.16: VAR

Value	Description
-	Basic variant, standard contemporary style; also used for standard forms allowed for use in writing by the Czech Standard Orthography Rules despite being marked there as colloquial
1	Variant, second most used (<i>less frequent</i>), still standard
2	Variant, rarely used, bookish, or archaic
3	Very archaic, also archaic + colloquial
4	Very archaic or bookish, but standard at the time
5	Colloquial, but (almost) tolerated even in public
6	Colloquial (standard in spoken Czech)
7	Colloquial (standard in spoken Czech), less frequent variant
8	Abbreviations
9	Special uses, e.g. personal pronouns after prepositions etc.

2.2.2 Compact tags

For most (but not all cases) just omit the dashes from positional tags. For more information, see ²

2.2.3 Informal abbreviations

In certain cases (including some places in this manual), the following tag abbreviations are used. Most of them are self-evident (dashes and rarely used fields dropped), as you can see in the following list:

Ngnc – noun; NFS1 = NNFS1-----A-----
 Aagnc – adjective; AAXXX = AAXXX-----1A-----
 Db – adverb; Db = Db-----
 Dg – adverb; Dg = Dg-----1A-----
 Dgd – adverb; Dga2 = Dg-----2A-----
 J^ – conjunction; J^ = J^-----

²http://ufal.mff.cuni.cz/pdt/Morphology_and_Tagging/Doc/compact_tags.pdf

CHAPTER 2. LEMMA AND TAG STRUCTURE

J , - conjunction; J , = J , -----

Rc, RRc - preposition, RR7 = RR--7-----

RVc - vocalized preposition, RV7 = RV--7-----

TT - particle; TT = TT-----

Ng-8, NNgXX-8 - noun abbreviation; NFXX-8 = NNFXX-----A----8

AX-8, AAXXX-8 - adjective abbreviation; AAXXX-8 = AAXXX-----1A----8

Db-8 - adverb abbreviation; Db-8 = Db-----8

Rc-8, RRc-8 - preposition abbreviation; RR7-8 = RR--7-----8

Chapter 3

Names

Proper names (either directly or the lemmas they consist of) have suffixes marking the category of that name:

G – geographical name: *Praha, Ústí nad Labem*

Y – given (first) name, formerly used as default: *Petr, John*

S – surname (last name): *Dvořák, Zelený, Agassi, Bush*

E – name of a nationality: *Čech, Kolumbijec*

R – name of a product: *Tatra* (the car),

K – name of a company: *Tatra* (the company)

m – default – names of mines, stadiums, guerilla bases, etc; also used for functional words in names.

The lemma should start with upper case if the word is always in upper-case in names (*Tatra* is always in uppercase, but *banka* not).

Keeping this categorization in the same level as lemmas is quite unsustainable and very unsuitable.

1. In theory every word can occur in any category. For example, *new* in *New York* (G), *New Jersey Devils* (sport club – K), *New Jersey Devils cards* (product – R), etc. Because it would explode the lexicon, usually common words (besides *new* and alike, all functional words) have only two lemmas: one for common words and one for all names (using default category m). But such approach is highly unsystematic and works only for a small corpus.
2. But even then, the system is not used consistently – some functional words even do not have the above mentioned two versions (normal and m). For example, *nad* in *Ústí nad Labem*, should be *nad-2-;m*, but there is no such a lemma. Similarly *a* in *a.s.* should be annotated differently when part of a name of some company and when not, but it is not.
3. Moreover lemmas having different categories are formally not connected, e.g. if you see *Martin-3-;K*, you do not know if it is derived from *Martin-1-;Y* (*Martin* + *Martin, s.r.o*) or *Martin-2-;G* (*DS Martin, a.s.*).

Proposed solution:

- The G, K, etc. categories should be independent of the morphology, and should be assigned to phrases on a different level. This would also require some enhancement to the annotation tool DA.
- Only the words that always (i.e. >90% or so) belong to some category would be marked with that category (e.g. *new* has no special suffix, *England* has G).
- Personal name would be annotated as separate lemmas *Petr Pánvička – Pánvička-;S* not (*pánvička*)_S.

- Names containing name of a person, where the original link is not perceived (usually geographical names that do not contain possessive construction) have separate entries. N.B. that current guidelines require all the following words to be annotated with lemmas containing G (incl. *ostrov*, and *úžina*!!)

Columbus (town in Ohio) – *Columbus-2-;G* not (*Columbus-1-;S*)_G
Martin (town in Slovakia) – *Martin-2-;G* not (*Martin-1-;S*)_G
Beringova úžina – (*Beringův-;S-^(*)2*) *úžina*_G not *Beringův-2-;G* *úžina-2-;G*
Ostrov Sergeje Kirova – (*ostrov Sergej-;Y* *Kirov-1-;S*)_G
Kirov (town in Russia) – *Kirov-2-;G* not (*Kirov-1-;S*)_G

3.1 Personal names

Some names are sometimes declined, sometimes not (*Bill* – *o Bill Clintonovi*, *o Billu Clintonovi*, *o Billovi*). The tag for nondeclined form is NgXXA.

3.1.1 von, van, etc.

For names (e.g *Ludwig van Beethoven*) the *van*, etc. phrase is perceived as a surname – annotate it that way. For other it is still perceived as geographical name (e.g *Kryštof Harant z Polžic a Bezdružic*). Of course the borderline is fuzzy.

Examples:

Ludwig van Beethoven – *Ludwig-;Y* *van-2-,t-^(v_hol.-jménech)* *Beethoven-;S*
Vincent van Gogh – *Vincent-;Y* *van-2-,t-^(v_hol.-jménech)* *Gogh-;S*
Kryštof Harant z Polžic a Bezdružic – *Kryštof-;Y* *Harant-;S* *z-1* *Polžice-;G* *a-1* *Bezdružice-;G*
Brigida z Háje – *Brigida-;Y* *z-1* *Háj-;G*

3.1.2 Chinese names

Usage

The surname precedes the given name. In most cases, the whole name is used (not just the family name). The thing is complicated by the fact, that many Chinese living abroad often change the order of their name or use their given name as a surname, etc. The discussion below can help you to determine, which part of a name is the given name and which part is the surname. If you are in doubt annotate them all as given names (Y).

That was the original recommendation, but probably annotating them as S would be better, because they are often used that way (You can say *Clinton for Bill Clinton*, but you cannot say *Po for Po Li*).

Surnames

There are relatively few surnames in China (200 most common surnames account for >96% of all surnames). Most of them consist of one syllable (*Wang*, *Li*, *Chen*, etc.) Only few surnames consist of two syllables (*Ou-yang*, *Mo-qi*, *Si-ma*, *Pu-yang*). Married women do not get their husband's surname.

Given names

Mostly two syllables, often connected with a dash (however sometimes separated by a space). Some can be widely used, some can be unique. Often it is impossible to say whether it is a name of a male or a female. The second syllable is usually used in informal addressing. The first syllable can be shared by all siblings. In traditional China a person had several given names during his/her life.

Most common Chinese surnames (in Pinyin):

Cai, Ceng-Zeng, Chen, Chen-Shen, Deng, Gao, Guo, He, Hu, Huang, Li, Liang, Lin, Lü, Ma, She, Sun, Tang, Wang, Wu, Xie, Xu, Yang, Ye, Zhang, Zhao, Zheng, Zhu

Links

<http://www.wlu.edu/hhill/names.html> – Chinese names explained

<http://www.geocities.com/Tokyo/3919/atoz.html> – Alphabetical Index of Chinese Surnames (incl. Pinyin, Anglicized and other versions)

3.1.3 Korean names

Korean names behave similarly as Chinese names. Surname precedes given name. Given name of most Koreans consists of two parts, in Latin alphabet often connected with a dash. Most common Korean surnames are (45% of the population): *Kim, Lee* (often spelled as *Rhee, Yi* or *Li*), *Park*.

Examples:

Yang Sung-jin – S: *Yang, Y: Sung, jin*

Yang Sungjin – S: *Yang, Y: Sungjin*

Kim Il-Sung (former dictator of North Korea) – S: *Kim, Y: Ir, Sen*

Kim Ir-Sen (= *Kim Il-Sung*) – S: *Kim, Y: Il, Sung*

He Wung – S: *He, Y: Wung*

3.1.4 Foreignized Czech names

Sometimes you can encounter names that are Czech in their origin, but are somehow altered to fit other languages (diacritics is omitted, female and male surnames are the same – e.g. *Judy Sedivy*).

Use the following guidelines to decide the lemma and tag for such a name:

- a name that does not distinguish female and male variant, should have just one lemma and three different tags (gender M, F, X¹)

Peter Janda – *Janda_;*S + NNMXX-----A----- or NNMS1-----A-----

Jane Janda – *Janda_;*S + NNFXX-----A-----

Jane a Peter Janda – *Janda_;*S + NNXXX-----A-----

- a name that has the same spelling as in Czech, should use the Czech lemma *Jane Janda* – *Janda_;*S + NNFXX-----A-----
- a name with altered spelling has its own lemma (with ,t suffix) *Judy Sedivy* – *Sedivy_;*S-,t + NNFXX-----A-----

¹If {M,F} gender is introduced, the tag NN{M, F}XX-----A----- should be used.

3.2 Compound names (names consisting of other names)

All lemmas of autosemantic words in compound names must have the category determined by the whole name (e.g. K, R). The lemmas of functional words contain default type category (m).

The problem is that a name of one type can occur as part of a name of a different type:

New England – G

New England Association of Chemistry Teachers – K

New England Association of Chemistry Teachers Journal – R

England is G noun in the first, K adjective in the second and R adj. in the third name.

If the lemma of the category you need does not exist and you have to insert a new one, do not care about numbering of lemmas, somebody else will do it (it would be impossible to ensure that the numbers were unique across all annotators). That means, if there is another lemma having just a different category (e.g. there is *England_;G* available, but you need *England_;R*), just change the category label.

Using the above-proposed separation² of morphology and name categorization, the *New England* example would be annotated quite easily (only *England* is marked by a category (G) by the morphological analyzer, the rest is done by some other kind of tool):

(*new England_;G*)_G

((*new England_;G*)_G association of chemistry teacher)_K

(((*new England_;G*)_G association of chemistry teacher)_K journal)_R

If the annotator did not recognize the components of the name (e.g. it is in Burmese), (s)he would annotate just the highest level.

The categorization is sometimes quite tricky – you do not know, whether to consider a phrase a name or a name plus normal word:

Nobelova nadace – *Nobelův_;K nadace_;K*³

Nobelův stůl (e.g. in a museum) – *Nobelův_;S stůl*

Nobelova cena – hard to say (m vs. normal), decided: *Nobelův_;S cena*.

Examples:

Brownův pohyb – *Brownův_;S*

Cena J. Debrau – *Debrau_;S cena*

Mérieuxův ústav – *Mérieuxův_;K ústav* (Should be *ústav_;K* but is not)

Divadlo J. Grossmana – *divadlo_;K J-4_;B_;K Grossman_;K*

příloha Kolumbus (in Lidové noviny) – *Kolumbus_;m*

v Dobrovského ulici nejezdí ... – *Dobrovský_;G*

v Dobrovského nejezdí ... – *Dobrovský_;G*

poliklinika Dobrovského (unofficial, it is located in D. Street) – *Dobrovský_;G*

Using the separation of morphology and name categorization, this is quite easy:

Nobelova nadace – (*Nobelův_;S nadace*)_K

Nobelův stůl (e.g. in a museum) – *Nobelův_;S stůl*

Nobelova cena – easy to say: (*Nobelův_;S cena*)_m.

²See the beginning of 3.

³The lemmas have different numbers (e.g. *Nobelův-1_;S*, *Nobelův-2_;K*).

Examples:

Brownův pohyb – *Brownův-;S pohyb*

Cena J. Debrau – (*Debrau-;S cena*)_m

Mérieuxův ústav – (*Mérieuxův-;S ústav*)_k

Divadlo J. Grossmana – (*divadlo J-0-:B-;Y Grossman-;S*)_k

příloha Kolumbus (in *Lidové noviny*) – (*příloha Columbus-;S*)_m

Dobrovského ulice – (*Dobrovský-;S ulice*)_G

v Dobrovského – (*Dobrovský-;S*)_G

poliklinika Dobrovského (unofficial, it is located in D. street) – (*poliklinika (Dobrovský-;S)*)_k

3.3 Horses, DJ's etc.

Horses have all kind of names (e.g. *Vinná réva*, *Deprivace*, *He Shall Reign*, *La Paloma Monitor*, *Frýdlant*, *Gold End*, *Lučina*, *Green Peace*, *Areál*, *First*, *Bounty*), and quite often you do not know if it is female or male (sometimes even female like names belong to a male horse). One clue is, that in an Oak (a horse contest type), all horses are young mares – females.

In PDT 1.0 the names of horses were mostly not annotated correctly – simply any available name was selected (Otherwise, a new lemma with category Y would have to be inserted in each case: e.g. *Deprivace* would be *Deprivace-;Y*, annotated as *deprivace*, *He Shall Reign* annotated as normal English phrase: *he-,t, shall-,t reign-,t*).

In our opinion, if the Y category were independent of the lemma, the horse name should be annotated correctly.

Similar problem is with the names of musical groups and DJ's. For famous groups and DJ's enter separate lemmas, for others use normal available lemmas.

3.4 Sport clubs, etc.

Name of the town in the club name: if only the town is noted, it is annotated as a geographic name (G), if the whole name of the club is noted, it is annotated as an institution (K). It is analogous to countries. (*Česko vs. Německo* are annotated as G)

Examples:

Cheb vs. Plzeň – *Cheb-;G Plzeň-;G*

SKP Union Cheb vs. Plzeň – *SKP-;B-;K Union-;K⁴ Cheb-;K Plzeň-;G*

Of course, it can be a problem to know it with foreign clubs. If you do not know, annotate it as an institution (K).

Examples:

Chelsea – part of London, UK

Chelsea – *Chelsea-;G*

Chelsea FC – *Chelsea-;K FC-1-:B-;K-;w-^(...)*

Ferencvaros – part of Budapest, Hungary

Ferencvaros – *Ferencvaros-;G*

⁴In PDT 1.0, the lemma is *Union*, but it should *Union-;K*

Ferencvaros TC – *Ferencvaros*;K TC-6-;B-;K
Sparta – *Sparta-2*;K
Sparta Praha – *Sparta-2*;K Praha-;K
Viktorie Žižkov – *Viktoria-2*;K-^(jméno_sport.klubu) Žižkov-;K
Udinese – *Udinese*;K-t + NNNXX-----A-----

It is the adjective of Udine (town in NE Italy), the official name of the football club is *Udinese Calcio* (*calcio* = *football*). However in Czech, the name is perceived as a noun and as the name of that club, therefore it is probably better to use it in that way:

To determine, whether something is a name of a town or a club, you can try to find that name on a map (eg. <http://www.expedia.com/pub/agent.dll?qscr=mmfn>) and also find the club (eg. <http://www.soccerage.com>).

Using the above-proposed ⁵ separation of morphology and name categorization, this looks much more consistent:

Cheb vs. Plzeň – (*Cheb*-;G)_K (*Plzeň*-;G)_K
SKP Union Cheb vs. Plzeň – (*SKP*-;B-;K *Union*-;K *Cheb*-;G)_K (*Plzeň*-;G)_K
Ferencvaros – (*Ferencvaros*-;G)_K
Ferencvaros TC – (*Ferencvaros*-;G TC-6-;B-;K)_K
Chelsea – (*Chelsea*-;G)_K
Chelsea FC – (*Chelsea*-;G FC-1-;B-;K-;w-^(...))_K
Viktorie Žižkov – (*Viktoria-2*-;K-^(jméno_sport.klubu) Žižkov-;G)_K
Sparta – *Sparta-2*;K
Sparta Praha – (*Sparta-2*;K Praha-;G)_K
Udinese – (*Udinese*-;G-t)_K
Udinese Calcio – (*Udinese*-;G-t *calcio*-;t)_K

The name of the sport club often contains some abbreviation. Some are common and present in the analyzer's lexicon (e.g. *FC*, *AC*) some are quite unusual (e.g. *EV*, *ERC*, *EC*, *ERC*, *EG*, *VS*, *AS*). If they are not present in the lexicon, entering them, suffixing the lemma by -B-;K-;w and using NNNXX-----A----8 as tag,

3.5 Other

Insisting on inclusion of name categories (K, R, etc.), implies explosion of number of lemmas. We follow each examples section by analogous examples using the above- proposed separation of morphology and name categorization (see 3.2).

3.5.1 Geographical names

Streets

We suppose that the word *ulice*, etc. is always present, even if elided on the surface.

Examples:

Dlouhá – *dlouhý*-;G+ AAFS1-----1A-----
Dlouhá ulice – *dlouhý*-;G+ AAFS1-----1A----- *ulice* + NNFS1-----A-----
Palackého, Dobrovského, etc. – *Palacký*-;G, *Dobrovský*-;G+ NNMS2-----A-----

⁵See the beginning of the 3.

Examples:

Dlouhá – (dlouhý)_G + AAFS1-----1A-----
Dlouhá ulice – (dlouhý ulice)_G or (dlouhý)_G ulice + AAFS1-----1A----- NNFS1-----A-----

Palackého, Dobrovského, etc. – (Palacký-;S)_G, (Dobrovský-;S)_G + NNMS2-----A-----

Towns

Words in one-word names consisting that were originally adjectives are annotated as nouns.

Examples:

Hluboká – Hluboká-;G + NFS1
Dobrá Voda – dobrá-;G + AFS1 Voda-;G^(součást_názvu_Odolena_Voda) + NFS1
Ohrada u Hluboké – Ohrada-;G + NFS1 u-;m + RR2 Hluboká-;G + NFS2

Examples:

Hluboká – Hluboká-;G⁶ + NFS1
Dobrá Voda – (dobrá voda)_G + AFS1 NFS1
Ohrada u Hluboké – (ohrada u Hluboká-;G)_G + NFS1 RR2 NFS2

3.5.2 Initials

A separate character for aggregate gender {M, F} would be good (for initials following a letter in newspaper, an initial before a foreign last name, foreign names, etc.).

3.5.3 Institutions, companies

This category contains for example companies, foundations, shops, clubs, sport clubs, restaurants, etc. All autosemantic words in names of restaurants have lemmas with K. The exceptions are functional words that are annotated as default type (m)

Examples:

Porcela Plus: Plus-3 + TT

Restaurants**Examples:**

Bar Viola – bar-2-;K, Viola-2-;K
U Medvídků – u-2-;m, medvídek-2-;K
La cambusa – Le-1-;m-,t^(franc.člen-jako-souč.jmen_a_názvů)⁷, cambusa-;K-,t
Restaurant HaPi – restaurant-2-;K HaPi-;K

⁶Frequent names of towns and names when POS changes, have separate entries. Therefore not (hluboká)_G

⁷In the current morphological lexicon, the m is missing.

Čínská restaurace Jin Jiang – čínský-2;K, restaurace-2;K, jin-2;K, jiang-2;K,t
 restaurace Jin Jiang – restaurace-1, jin-2;K, jiang-2;K,t
 Francouzská restaurace v Obecním domě – francouzský-2;K, restaurace-2;K, v-2;m obecní-2;K
 dům-2;K
 Hospůdka U vylitýho mrože – hospůdka-2;K u-2;m vylitý-2;K mrož-2;K

Examples:

Bar Viola – (bar, Viola;Y)_K or (bar, viola)_K (select anyone, if you do not know the orig.)
 U Medvídků – (u-1, medvídek)_K
 La cambusa – (le-1,t^(franc.člen), cambusa,t)_K
 Restaurant HaPi – (restaurant, HaPi;K)_K
 Čínská restaurace Jin Jiang – (čínský, restaurace, jin, jiang,t)_K
 restaurace Jin Jiang – (restaurace, jin, jiang,t)_K
 Francouzská restaurace v Obecním domě – (francouzský restaurace, v-1 (obecní dům)_K)_K
 Hospůdka U vylitýho mrože – (hospůdka, u-1, vylitý, mrož)_K

3.5.4 Sporting and other events

All events should receive special lemmas with m. However, if it is registered as a company and used in that meaning, then it should be K. If not certain use m.

Examples: ⁸

Paris Indoor – Paris-2;m,t Indoor;m,t + NNNXX-----A-----
 US Open – US-3:B;t + AAXXX----1A---8 Open-1;t;m AAXXX----1A-----⁹
 akce Stop milion – stop-1;m milion'1000000;m_m

Pohár mistrů – pohár;m mistr;m
 Mistrovství světa – mistrovství;m svět;m

Examples:

Paris Indoor – (Paris-2;G,t Indoor,t;m)_m
 US Open – (US-2:B,t^(americký) Open-1;t;m)_m
 akce Stop milion – (stop-1 milion'1000000)_m
 Pohár mistrů – (pohár mistr)_m
 Mistrovství světa – (mistrovství svět)_m

⁸Many of these entries are not in the lexicon, therefore the actual numbers can be different once it is there. See note in 2.1, e.g. *mistrovství*: *mistrovství-1*, *mistrovství-2;m*, *mistrovství-3;R*, etc.

⁹We think, it is perceived as noun, probably inanimate, in Czech.

3.5.5 Televisions

Generally televisions are annotated as institutions (K). Only, if a company runs several channels, then the channels are annotated as products (R); but it is currently used only with Czech(oslovak) public television (ČT1, ČT2 and F1).

Examples:

ČT – ČT_:B_:K

ČT1 – ČT1_:B_:R

Nova – Nova_:K

NBC – NBC-4_:B_:K

CNN – CNN-1_:B_:K_:y_:b_:t

3.5.6 News, Magazines

All autosemantic word in names of news or magazines have lemmas with R. Currently, some of the newspapers are in the lexicon as institutions (e.g. *Sme*), this is not correct. Foreign names are often used as in plural, even if in the original there are in singular.

Examples:

Sme – *Sme*_:R_^(*noviny*) + NNXX

Zeitung – *Zeitung*-1_:R_,t_^(*souč. názvu něm. novin*) + NFPX or NISX

3.5.7 Song names, etc.

Names of songs, TV programs etc. are annotated as normal words. The only reason is practical – it would cause explosion of the lexicon. If the categories and morphology are separated (see beginning of 3), these items can be annotated as R or m.

Chapter 4

Abbreviations

For discussion about inserting abbreviation not present in the morphological lexicon, see 10

4.1 Gender

Abbreviations can be used with different genders (e.g. *ODS* – feminine (*strana*) or neuter). Any abbreviation can have neuter gender. If the gender cannot be disambiguated by the context, use the gender used elsewhere in article. If the author mixes genders or there are no disambiguating contexts, use the gender inherent gender of the abbreviation. In Czech, is usually easy to determine – it is the gender of the head of unabbreviated equivalent (e.g. *ODS* – *strana* → f). With foreign abbreviations it is much more problematic, different people use different genders (e.g. because of different translation). If you are not certain which of the gender is most widely used, use the default neutrum.

Examples:

UK – F (*univerzita*)

FBI (*Federal Bureau of Investigation*) – I (*úřad*), N (default or *byro*), F (probably á la *CIA*), *MP* (pl., referring to the members of the *FBI*)

CIA (*Central Intelligence Agency*) – F (*agentura*)

4.2 Normal abbreviations

Normal abbreviations have sometimes as a lemma the abbreviation (and sometimes the original unabbreviated word. Usually the former method is used for abbreviation that are more common than the unabbreviated word (and for abbreviation of multi word expressions). But it is not always true.

For discussion about determining the gender of an abbreviation, see 4.1

Examples:

např.: *například*_:B + Db-----8

P.S.:

*post-2*_:B_,t_^(*lat.,-po,-např.-P.S.*) + RR--X-----8

*scriptum*_:B_,t_^(*př.-P.S.*) + NNNXX-----A---8

n.L.: *nad-1*_:B¹+ RR--7-----8, *Labe*_:B_;G + NNNS7-----A---8

¹Should be *nad-2*_:B_,m, but is not.

r. 1998: *rok*_:B + NNIXX-----A----8

r.: *režie*_:B + NNFXX-----A----8

rež.: *režie*_:B+ NNFXX-----A----8

4.3 Isolated letters

Note: The following is still not official.

Isolated letters (e.g. *A-konto*) are handled as abbreviations. The only exception is if they are not in the name (*zápas skupiny B*). Many of the annotations suggested below are still not offered by the morphological analyzer. Moreover, sometimes the morphological analyzer is constrained to offer appropriate lemma and tag only if the letter is followed by a dot. Should be repaired.

You have to select (or insert) the lemma according to the semantic category:

*K-0*_:B;Y – first (and most middle) names

*K-4*_:B;K – names of institutions

*K-5*_:B;G – geographical names

*K-6*_:B;R – names of products

*K-7*_:B;m – other names (sporting events, etc)

*K-9*_:B;S – last (and some middle) names

*k-8*_:B.^(*ost_zkratka*) – other abbreviations (not names)

*k-3*_.^(*označení_pomocí_písmene*) – other letters (not abbreviations, not in names)

Frequent abbreviations have their own lemmas, for example *V* – *V-1*'*volt*_:B or *k*: *ABC k.s.* – *komanditní*_:B.^(*jen_komanditní_společnost*).

Tag selection (or insertion):

- noun: gender is known: NNgXX—A—8 ($g \in \{MFIN\}$)
- noun: gender is unknown: NNXXX-----A----8
- adjective: AAXXX-----1A---8 or AAgXX—1A—8
- others: X@-----1 (variant of X@----- for one letter words)

Examples:

A: *A-mužstvo* – *a-3*_.^(*označení_pomocí_písmene*) + AAXXX-----1A-----

d: *odst. 1 písm. d* – *d-3*_.^(*označení_pomocí_písmene*) + NNNXX-----A-----

A: *16 A* – *A-1*'*ampér*_:B + NNIXX-----A----8

A: *A konto* (or *A-konto*) – *A-6*_:B;R + AAXXX-----1A-----

a: *ABC a.s.* – *akciový*_:B.^(*jen_akciová_společnost*) + AAXXX-----1A---8

s: *na s. 128* – *strana-4*_:B.^(*v_knize_rukopise...*) + NNFXX-----A----8

1. It is hard to decide, whether an isolated letter is an abbreviation or a label using a letter (e.g. *a-3*). For example, *B* in *B-konto* can be from *bežný*, but *A* in *A-konto* probably means better than *B*. Maybe not, maybe yes, who knows. What is important, we mostly annotate texts written by people that do not know. Therefore it would be reasonable to merge these two possibilities together. Maybe annotate all single letters as abbreviation, the possible exception could be labels of paragraphs and cases (*odst. 1 písm. d*) or *za a*).

2. Letters in similar configuration as nouns in noun cluster should be treated as nouns – they should be annotated as nouns. See also 6.1.1.
3. The category of a name (K, G, R, etc) and lemma selection should be orthogonal – see also 3.2.

Examples:

A: *A-mužstvo* – a-8-^(ost_zkratka_nebo_označení) + NNNXX-----A----8
 d: *odst. 1 písm. d* – d-3-^(př_odst_a,_za_a) + NNNXX-----A----8
 A: *16 A – A-1'ampér_ B* + NNIXX-----A----8
 A: *A konto* (or *A-konto*) – (a-8-^(ost_zkratka_nebo_označení)...)_R + NNXXX-----A--
 -8
 a: *ABC a.s.* – (... akciový_ B-^(jen_akciová_společnost))_K + AAXXX-----1A----8
 s: *na s. 128 – strana-4_ B-^(v_knize,_rukopise,...)* + NNFXX-----A----8

4.4 RM-systém, samopal SA-58

An abbreviation preceding a noun is an adjective, an abbreviation following a noun is a noun. We would suggest to annotate them all as nouns (see 6.1.1). Does it mean that *HIV* in *HIV virus* and *virus HIV* have different POS.

Examples:

RM-systém – *RM_ B_ K* + AAXXX-----1A----8 NNXXX-----A----8
samopal SA-58 – *SA-2_ B_ R* + NNXXX-----A----8
virus AH 3 B – *AH-1_ B* + NNXXX-----A----8
virus HIV- HIV_ B_ L_ U_^(lidský_virus_způsobující_AIDS) + NNXXX-----A----8

4.5 Units of measurements

Units called after some males person (*V – volt*, *A – ampér*, etc.), have inanimate gender. However, units using degrees ($^{\circ}\text{C}$, $^{\circ}\text{F}$) have masculine animate gender, because the word *stupeň* is always present (even if omitted in the written text). Absolute temperature uses as the unit called *Kelvin* (K) not *degree of Kelvin*. Therefore the unit has inanimate masculine gender. However, if the author uses it erroneously as degree, the tag as to be masculine animate.

Examples:

C: *Ráno byly 3 °C.* – *Celsius_ B* – NNMXX-----A----8^o
 C: *Ráno byly 3 C.* (read as *Ráno byly tři stupně Celsia*) – *Celsius_ B* – NNMXX-----A----8
 K: *Teplota 5000 K.* – *Celsius_ B* – NNMXX-----A----8
 K: *Teplota 5000 °K.* – *Celsius_ B* – NNMXX-----A----8^o

If the C character is preceded by some character trying to look as the degree symbol $^{\circ}$ (eg. -C, o C, O C), then you should mark it as an error – as lemma insert the degree² symbol $^{\circ}$ and as tag X@-----1. It should be converted into a punctuation mark.

²On Czech keyboards usually Shift+<key-on-the-left-from-1>, followed by Space. Or on any keyboard Alt+0176.

4.6 Authors abbreviations

The author's name abbreviations used in newspapers (e.g. *Ber, mas, jst, ...*) have lemma as the form + -99_:B_;S and tag NNXXX-----A---8. There is X for gender because usually we do not know it. If the {M,F} gender is introduced, it should be used here. These abbreviations are not present in the lexicon, therefore you have to insert them.

Examples:

ač: PRAHA (ČTK, ač) Problém Gabčíkova ... – *ač-99_:B_;S – NNXXX-----A---8*
gap: DUKOVANY (gap) Na základě posudku ... – *gap-99_:B_;SNNXXX-----A---8*

4.7 Academic titles

Titles distinguish genders – there has to be one lemma for men, and one lemma for women (*JUDr-1_:B_^(doktor_práv)* vs. *JUDr-2_:B_^(doktorka_práv)*); to keep it consistent the masculine has number 1, the feminine has number 2. We think, the titles should have the same form for women and men. Just the tag should be different, with possibility to have X if the gender is not known (e.g. a letter subscribed as Dr. A. B.)

Chapter 5

Colloquial Czech

If an official alternative to the colloquial form exist, then the the colloquial form has the same tag except a different variant ('5', '6', '7', ev. '3' – see 2.2.1.13).

Examples:

které: stavení, které – P4NP4-----5

Novákovíc: Novákovíc pes – *Novákův_;*S_^(*2) -AUXXXM-----6¹

takovejhlema: takovýhle – AAFFP7-----1A---6

hovadinama: hovadina – NNFP7-----A---6

naši: pro naši atletiku (officially short: naši) – *můj_^(přivlast.)* – PSFS4-P1-----6

5.1 Cos, jaks, kdys ...

We tagged these words as if they were without -s and added -9 at the end.

In our opinion it would be better to divide such an expression in two words (e.g. *cos* → *co* + *být*, analogous to *abych* → *aby* + *být*) and tag them like two normal words, just with some variant recognize it.

5.2 Suffix -é in plural of neuter

Should not be treated as misspelling, but annotated as (colloquial) variant of official -á forms (variant '5').

Examples:

které: stavení, které – P4NP4-----5

¹In PDT 1.0, this is sometimes obsoletely annotated as AUMS1M-----6 or NNXXX-----A---6

Chapter 6

Foreign words and phrases

General rule

1. For a longer phrase (or citations) in a foreign language, use morphology of that language (but distinguish genders M and I ??) (Hence citation use).
2. For a single word or shorter phrase use Czech morphology. (Hence word use) The borderline is fuzzy, of course.

6.1 Part of speech

Many foreign words used in Czech sentence can have different part of speech than in their original language. Usually the hint is how it behaves in different context, if it is declined as a noun, if agrees with its head, etc.

All foreign lemmas have *-t* suffix.

It would be good to somehow distinguish foreign words in word or citation use.

6.1.1 English noun clusters

All nouns in attributive use are annotated as adjectives.

That's quite problematic:

1. Virtually all English nouns can be used as attributes of other nouns
2. It is imported to Czech: *Staropramen Extraliga, Český Telecom Cup*, etc.

We think, it should be annotated as two nouns.

6.1.2 Examples

V kostele XY zpívala Musica Bohemica.

Bohemica annotated as a noun; in Latin it is an adjective.

Reason: When the phrase is declined, *Bohemica* is declined as a noun (*žena*): *pozvali Musicu Bohemicu*, **pozvali Musicu Bohemicou*

Annotation: *Musica*_{-t;}K + NFS1A, *Bohemica*_{-t;}K NFS1A

To je trochu ad hoc.

hoc is annotated as a noun; in Latin it is an adverb.

Annotation: *ad*_{-t}RRX, *hoc*_{-t}NXXXXA

In the following, the section headers refer to the categories of the foreign language.

6.2 Articles

English *an* should be a form of *a*.

Articles merged with a preposition (*fra du, ita della, deo im, aufs, zur*) are treated as prepositions (?Split into two words?)

Arabic short words (##)(?articles, ?prepositions) are treated as articles.

6.2.1 Citation use

Same as single words

Should distinguish gender, number, and/or case Therefore: TTgnc or AAgnr ??

6.2.2 Word use

Tag: TT-----

Lemma: Usually the same as the form

Originally, we wanted to treat articles as adjectives. Forms having different gender, number and/or case, would have the same lemma (*der* for forms *der, die, das, des, dem, den*). The problem is that Czech does not respect the original categories (*nebezpečný La Manche – la* in French F, in Czech the phrase is I; *Los Angeles – in Spanish pl, in Czech sg.*)

6.2.3 Examples

l' l-5,t^(př..l'Arc,_stažený_tvar_fr._členu) + TT

L' L-10^(př..L'Aqua,_stažený_tvar_fr._členu) + TT (should be m¹)

la la-2,t + TT

il il,t^(it._len) + TT

as as,t + TT (Arabic)

al al,t + TT (Arabic)

el el,t + TT (Arabic)

della della,t + RRX (sometimes incorrectly annotated as AAXXX)

am am,t + RRX

6.3 Nouns

6.3.1 Citation use

6.3.1.1 English

To keep it simple, number of English nouns is annotated in the same way in citation use as in word use. That means X is used instead S for nouns in singular. The difference is in X is used instead S for nouns in singular. The difference is in cases – in citation use, the case is always X, but word use it can sometimes be declined.

¹If the name categories and lemmas were independent, everything annotated as L-10 would become l-5 (see 3.2)

6.3.2 Word use

6.3.2.1 English

The nouns in singular in English have number annotated as X (English singular. is often used in Czech as plural). For nouns that are usually declined mark the case even if in base form, for nouns that are nondeclined mark it as X.

Examples:

flow: oba dva cash flow (oficiální i skutečný) ... – *flow_,t* – NNIXX-----A-----

statement: v cash flow statementu ... – *statement_,t* – NNIS6-----A-----

statement: Náš cash flow statement ... – *statement_,t* – NNIS1-----A-----

flow: Náš cash flow ... – *flow_,t* – NNIXX-----A-----

girl: Beatles: Girl – *girl_,t* – NNFXX-----A-----

girls: A teď zahrajeme písničku Girls. – *girl_,t* – NNFXX-----A-----

6.4 Verbs

6.4.1 Citation use

6.4.1.1 English verbs

The following tags are applied:

- Present non3sg: *go* VB-X---XP-AA---
- Present 3sg: *goes* VB-S---3P-AA---
- Imperative: *go* Vi-X---X--A-----
- Infinitive: *go* Vf-----A-----
- Past tense: *went* Vp-X---XR-AA---
- Passive participle *gone* Vs-X---XX-AP---

If it is hard to determine the base form usage, annotate it as infinitive. If it is hard to decide between past tense and passive participle, use past tense. In PDT 1.0, most of the verbs using base form were annotated using the default – infinitive.

Examples:

be: to be or not to be – *be_,t^(angl._být,_v_názvech_apod.)* – Vf-----A-----

do: Do it right now! – *do-2_,t* – Vi-X---X--A-----

6.4.2 Word use

Usually the tags and lemmas are the same as in citation use.

6.5 Slovak language

If a Slovak word has the same form as corresponding Czech one (e.g. prepositions), you should annotate it as if it were Czech. Otherwise it has to be annotated as any other foreign language.

Chapter 7

Errors

The text can contain errors. It is reasonable to correct some of them, preserving the original form. However, only low-level errors – spelling and morphology should be corrected (We do not want to correct Engels' text into Heidegger's). Never correct a colloquial form by an official one (e.g. *zelené města* *→ *zelená města*, *bez noh** *bez nohou*), even if the analyzer does not know the form ¹.

The errors have to be just marked, do not edit the file. Try to insert lemma and tag as if the form is correct, and use the DA support for marking errors – it inserts the text "(Chyba)" at the end of lemma or tag. If the lemma is correct, insert it after the tag, otherwise insert it after the lemma, if you do not know just insert it somewhere. If you want to add some comment, write it before the closing parenthesis, preceded by a dash (e.g. (*Chyba-nad c by měla být čárka, ne háček*)). This convention makes it easy to find the errors automatically.

7.1 Characters

Sometimes, foreign characters had been be screwed (e.g. *Fran?oise*), and therefore the morphological analyzer did not recognize the whole word. Mark it as a lemma error (do not edit the file), it has to be corrected and run thru the analyzer once more. There is a problem with letters not contained in Latin 2, they should be replaced by corresponding characters without diacritics. In the future, Unicode (2 bytes or UTF) should be considered.

7.2 Separators, etc.

Sometimes, the text contains *o* or *l* as bullets or separators. They should be marked for deletion (Press L (delete) in the lemma or tag list).

¹You have to insert a new lemma and/or tag – see 10 for more details.

Chapter 8

Hard to decide

8.1 až

až-1 + J[^]

2 až 3 (but not od 2 až do 3 – see až-3)

nabízí přiblížení až přijetí

až-2 + J ,

tak .. až: Nabízí se tak okatě, až je to hanba.

.. začnou pochybovat, až nakonec uvěří, že ..

Bylo mi 24, a byl jsem plný touhy se pomstít. Až jsem se ocitl před člověkem, který dostal zabrat víc než já.

až-3 + Db

If omitted, the sentence stays grammatical. It is often possible replaceable by *teprve*.

Dostanete až 250 mil zdarma.

kam až: Kam až půjdeš?

Až on me přesvědčil, že tomu tak bude.

Modifies functional word (should be probably TT)

až + conj: Je geolog a až pak filozof

až + prep: z Brna až do Prahy (Cf. až-1)

8.2 jak

jak-1;L[^](živočích) + NNMnc-----A-----

Obvious.

jak-2 + J ,

1. Meaning že (cannot be replaced by *jakpak*)

Jak řekl M. Zeman, bude třeba ..

Jak ukazuje vývoj posledních let, je to ..

Jak známo, ...

Skutečnost, jak už to bývá, byla trochu jiná. However, rarely it can be Db – depending on the interpretation

Viděl, jak upadla.

Meaning *Viděl, že upadla.* – J ,

Meaning *Viděl, jakým způsobem upadla.* – Db
Kamera zabírá poslance, jak otvírají krabici

2. Time, meaning *když, až, jakmile*
Přijdu, (hned) jak budu hotov^{ssč}.
Hned jak budu moct, zavolám.
3. In comparison, meaning *než, jako:*
Byl větší jak on^{ssč}
rychlý jak vítr^{ssč}
4. Condition (coll.), having the meaning *jestliže, když*
Jak budeš zlobit, nepůjdeš nikam^{ssč}

Japonskému turistovi upadla lžička, jak chtěl zmáčknout spoušť foťáku.
Poslední šanci, jak se probojovat do .., bude ..
Stát to měl spravovat zoláštním ministerstvem (jak je tomu např. v Rakousku)

jak-2 + J[^]

In the phrase *jak ... tak ...*, having the meaning of *i...i*. However cf. **jak-3** 2.
Byli tam jak odborníci, tak amatéři.

jak-3 + Db

Pronominal adverb

1. Interrogative – manner or extend (expr. *jak pak*).
Jak se jmenuješ?
Jak je to možné?
 Sometimes expressing large extend (often in exclamations).
Jak ten čas letí^{ssč}
Jak (pak) by ne^{ssč}. Japa by ne.
Líbí se ti to? – A jak!
2. Relative – marks subordinative adverbial clause (mostly manner expressing comparison, often with *tak* – however cf. **jak-2** + J[^])
Jak řekli, tak udělali^{ssč}
tak dlouho, jak je možné (tak .., jak ..)
Jak si kdo ustele, tak si lehne
3. Relative (coll.) – meaning *co, který*
ten člověk, jak jsem ti o něm říkal^{ssč}
4. Indefinite
buď jak buď (the verb is repeated)
jak kdo, jak kde, jak kdy,

??*Jak se kůže sama obnovuje, postupně vylučuje ..*

?? *Jak jsem chodil o berlích, tak jsem si zničil i druhé koleno.*

8.3 málo

Similar to *moc*.

málo-1.^(málo+_2._p.,_málo_peněz) + Ca--c-----

It has to be modified (in the shallow syntax) by a noun in genitive. Has only two forms:

málo and *mála* (only in genitive).

Máme málo zájemců.

bez mála peněz

před málo lety^{ssč}

Je jen o málo důslednější. – but Je málo důsledný. is málo-3 (Dg)

Udělal to jako jeden z mála odborníků, ..

Udělal to jako jeden z mála. – ?? not modified by anything

Udělal to jako jeden z mála, co přišli.

málo-2.^(př._to_málo_co_měl) + NNNnc-----A-----

vystačit s málem^{ssč}

vařit z mála^{ssč}

Děkuji. – Za málo. ^{ssč}

málo-3.^(málo+_příd._jm.,_př._byl_málo_důsledný) + Dg-----dA-----

Málo mluví, hodně dělá.^{ssč}

Je málo důsledný.

Ve srovnání s loňskou sezónou je to velmi málo. – you can say méně.

Zdržím se jen málo^{ssč}.

8.4 moc

Similar to *málo*.

moc-1.^(nad_někým;_politická,_vojenská;_plná,...)

Obvious.

převzít moc

moc proletariátu

udělám, co je v mé moci

mermo mocí

moc-2.^(mnoho_něčeho_[se_subst._v_gen.]) + Ca--X-----

Cannot be replaced by *velmi*. Can mean *příliš*, but is more colloquial. It has to be modified (in the shallow syntax) by a noun in genitive.

Má moc peněz.

Všeho moc škodí.

moc-3.^(velmi,_ve_spojení_s_adj.,_př._moc_hezká) + Db

Can be replaced by *velmi* (except ellipses). Modifies an adjective, adverb or verb.

Je moc hezká.

Vím to moc dobře.

Moc se snažil.

Ve srovnání s loňskem je to moc. – ellipse.

8.5 proto

proto-1^(*proto; a_proto, _ale_proto, ...*) + \mathcal{J}^{\wedge}

Coordinative conjunction expressing consequence (implication). Structure: reason \rightarrow consequence. Replaceable by *tedy*. Usually *a proto* or *a ... proto*

Nesplnil úkol, (a) proto nedostal odměnu.

Každé proč má své proto.

Německo se začalo dusit, a rozhodlo se proto omezit ...

proto-2^(*dal_mu_co_proto, _tak_proto!*) + Db

Pronominal adverb. Refers to the subordinative clause Structure: what \rightarrow reason

proto, že: Udělal to proto, že musel.

Udělal to proto, aby/že mu pomohl.

co proto: dát někomu co proto; dostat co proto

no proto: Říkal, že tam přece jen půjde – No proto! (Sometimes classified as a modal particle)

8.6 svůj

svůj-1^(*přivlast.*) + P8gnc-----v

Obvious.

svůj-2^(*být_svýj*) + AOgn-----v

Problem with tags, analyzer probably needs update.

Vzít za své.

Víme své. Víme svoje.

8.7 tak

In general:

- replaceable by *a proto* $\Rightarrow \mathcal{J}^{\wedge}$
- replaceable by *tím způsobem, stejně, zrovna* \Rightarrow Db

tak-2 + \mathcal{J}^{\wedge}

Coordinative conjunction. If one of the clause is subordinative then *tak* has the meaning of an adverb: (Cf. *Bál se, tak si pískal.* – \mathcal{J}^{\wedge} vs. *Kdyby se bál, tak si pískal* – Db)

1. A consequence — meaning *(a) proto, tedy*

Bál se, (a) tak si pískal.^{ssč}

Neudělali..., příspěvek tak budou muset vrátit.

Byly zakázané, a tak přitahovaly

Zmizí bariéry, a tak bude možné využívat ..

Zpozdila se, a tak musela běžet.

Jsou profíci, tak ať se podle toho zařídí/

Počítá se s tím, že některé se sloučí, i tak bude třeba ..

2. A conjunction — in *jak – tak*

tak-3 + Db

1. Referring to something known, to other sentence, etc.

tak – jak: Bylo to tak, jak jsem myslel.^{ssč}

jak – tak: Jak řekli, tak udělali.

Přesně tak.

tak zvaný

Ať je to tak nebo tak ...^{ssč}

jen tak: Udělal to jen tak.

tak tak: Stihl to (jen) tak tak.

to: Stalo se tak při ..

Tak se tehdy žilo^{ssč}

Sub-Clause, *tak* Main-clause:

Když – tak: Když jsem počítal já, tak mi vyšlo velké číslo.

Pokud – tak: Pokud to není diskriminace, tak nevidím důvod ..

Dokud se člověk raduje, tak je život pěkný.

Kdyby – tak: Kdyby/Pokud by se bál, tak by si pískal.

(Cf. *Bál se, tak si pískal.* – ǝˆ)

2. Expressing amount (usually large) of a property, etc.

Kam tak rychle?^{ssč}

tak jako: Je tak velký jako já.

Zmizel z povědomí tak jako jeho pomník;

Nabízí se tak okatě, až je to hanba.

To je ale tak daleko .

tak vysoká; tak oslaben, že ...

Buďte tak laskav.^{ssč}

ani tak o ..., jako o ...: Nejde ani tak o mzdu, jako o ...

přibližně: Dostane se na burzu asi tak třetí den od ..

hned tak: Hned tak nepřijde. (koneckoců)

odmítá to, stejně tak jako ...

.. a zrovna tak hyzdit;

tak jako tak

Chapter 9

Sólokapři

Hradec Králové

Králová, G, S (Dvůr Králové) + NNFS2-----A----- It is hradec that belonged to králová

strana

na jedné straně ..., na druhé straně ...: druhý-1 (jiný), strana-1 (v prostoru) nerespektované ze strany Israele: strana-3 (u soudu, ..

stát

stane se ministrem: stát-2 (něco se přihodilo)

s=to

být sto něco udělat lemma = sto-3 (být sto), značka TT-----

tudíž

always J ##Why is there the other possibility

vážít

vážít cestu – vážít na váze? nebo ctít někoho

vedení

Everything except elektrické vedení type, is considered as form of vedení-1

Examples:

*pod vedením kamarádky – vedení-1 (*7ést-1)*

*vedení podniku – vedení-1 (*7ést-1)*

*čínské vedení – vedení-1 (*7ést-1)*

elektrické vedení – vedení

9.1 Date and time

v +

a day – accusative (4): *v sobotu, v neděli*

a month – locative (6): *v lednu, v září*

an hour – accusative (4): *ve 4 hodiny, v 6 hodin*

ve dne – locative (6) – NNIS6-----A---9 -special kind of locative that occurs only in this context (*v noci* is also in locative):

month in a date – genitive (2): *25. září, 2. října*

9.2 Numbers, numerals and quantifiers

An adjective modifying a quantified expression agrees in case with the noun not the numeral.

Examples:

za (gen) *těch* (gen) *mizerných* (acc) *deset* (gen) *korun*

Deset (nom) *nejlepších* (gen) *sportovců* (gen) *ukázalo*

1x

Lemma: as form. Tag Cv-----

Example

1x – 1x + Cv-----

4x5

Should be split into three parts. E.g. $4x5 \rightarrow 4, x, 5$

9.3 Hyphenated composites

If the hyphenated word ends with -o, and by a replacement of that -o by an adjective ending we obtain an adjective (normal or possessive), the lemma for the word is that adjective (e.g. *česko-německý* – *česko* → *český*, *Karlo-Ferdinandova* – *Karlo* → *Karlův*). Some word cannot be viewed as derived from adjectives, but rather from nouns (e.g. *rap-jazzová* – *rap* → *rap* vs. *rapovo-jazzová* – *rapovo* → *rapový*). However, the lemma of that noun cannot be used as a lemma for the hyphenated form, and a new lemma (having different number) has to be introduced.

That is extremely inconvenient (*padlý na hlavu*) – virtually any noun can be used in such a context, therefore for every noun, there have to be two lemmas – one for normal usage and one for hyphenated usage. We strongly suggest to allow any noun to have a hyphenated form in several variants – at least the bare base form and form ending in -o (variant '1').

Examples:

srbsko-černohorská – *srbský* – A2-----A-----

Univerzita Karlo-Ferdinandova – *Karlův*; *K_(*3el)* -A2-----A-----¹

¹Better: *Karlův*; *Y_(*3el)* – A2-----A----- . See 3.2

Univerzita Karel-Ferdinandova – Karel-2.;K – A2-----A-----

rap-jazzová: rap-2 – A2-----A-----

| *Better: rap – A2-----A-----*

rapo-jazzová: rap-2 – A2-----A-----

| *Better: rap – A2-----A----1*

rapovo-jazzová: rapový – A2-----A-----

Chapter 10

Insertion

If the possibilities offered by morphological analyzer are not suitable, you have to insert new lemma and/or tag. If you insert a new lemma, you have to ensure, that the lemma (lemma proper) you insert is not already used. That usually means adding unique numbers to distinguish lexical items having the same base form.

10.1 Possessive adjectives

Lemmas of possessive adjectives show how they get the noun they are derived from (see also 2.1.1). For example:

*kardinál*ŭv_^(*2) – remove two letters: *kardinál*
*Karl*ŭv_.;Y_^(*3el) – remove 3 characters, add "el": *Karel*
*Martin*ŭv-1.;Y_^(*4-1) – remove 4 characters, add "-1": *Martin-1*

Examples:

*premiér*ŭv_^(*2)
*Soros*ŭv_.;S_^(*2)
*chlapc*ŭv_^(*3ec)
*Švehl*ŭv_.;S_^(*2a)
*Mách*ŭv_.;S_^(*2a)
*Hlink*ŭv-1.;S_^(*4a-1)
*Bender*ŭv-1.;S_^(*4-1)

10.2 Words ending with -ismus, -izmus

The base form should use *-izmus* ending, the form using *-ismus* is treated as variant '1'. Currently still some entries do not follow this convention.

Examples: ¹

mechanizmus: *mechanizmus* – NNIS1-----A-----
mechanismus: *mechanizmus* – NNIS1-----A---1
exhibicionismus: *exhibicionismus* – NNIS1-----A-----

¹The examples show the desired state, in the current version of morphological analyzer they are regarded as separate lexical items (they have different lemmas)

exhibicionismus: *exhibicionismus* – NNIS1-----A---1

nacionalizmus: *nacionalizmu* – NNIS1-----A-----

nacionalismus: *nacionalismus* – NNIS1-----A---1

10.3 Strange and unique things

Transcription of pronunciation

Lemma: as the form, tag: NNXXX

Examples:

vyslovujeme "zpjev" – zpjev + NNXXX

Isolated morphemes

Lemma: as the form, tag: NNXXX

Examples:

...ve slovech končících na -ství píšeme...: ství + NNXXX

Geometry

We can meet an article about a geometric theme sometimes. It means, that there occur a lot of triangles ABC, abscissas (lines) PQ, RS, AB and so on in that article. It is necessary to create a new lemma ending 98 for every mentioned figure.

Chess codes

Lemma: The code + -1_;w. Tag NNNXX-----A---8 (neuter because *pole* is neuter)

Example

Jh8 – Jh8-1_;w + NNNXX-----A---8

Crippled forms

Lemma: the same as the form + _;t

Tag: normal if possible, otherwise NNXXX / AAXXX according to the POS

Examples:

Waklaf Hafel – Waklaf_;t + NNMS1, *Hafel_;t* + NNMS1

Gaptschikowo – Gaptschikowo_;t + NNNS1

v Gaptschikowo – Gaptschikowo_;t + NNNXX

10.4 Other

This section contains especially examples of previously inserted lemmas/tag. Some of them are already in the dictionary, however they mainly serve as an inspiration, when inserting similar things.

ad hoc

ad-x₋t + RRX, hoc-x₋t + N

pele-mele

For example as a heading in a newspaper *pele + TT, mele + TT*

zprostředkovací vs. zprostředkovat

The morphology is rather shallow. It means, for example, that lemma for *zprostředkovací* is *zprostředkovací* (precisely *zprostředkovací* $\hat{>(*2t)}$), and *zprostředkovat* as it was in the past.

Chapter 11

Errors in PDT 1.0

- *HaDivadlo_;*Y → *HaDivadlo_;*K
- Theaters – most of the theaters do not have K category (search for *Divadlo* case sensitive): *Divadlo v Celetné, Divadlo Husa*, etc.
- S/NWS/1993/mf930701:105-p4s1 – *los* should be *los-2*
- *které* should be P4NP4-----5 not corrected as an error. (look for *w <spell>které*)
- <s id="S/NWS/1992/lnd92251:095-p1s1">: *US* should be with m (*US-3_-B_;*m_;t) not K.
- *Novákovic, Perotovic, ..* should be AUMS1M-----6, but it is either AUgncM-----6 or NNXXX-----A---6 (e.g. S/NWS/1992/lnd92254:051-p5s8)
- S/NWS/1992/lnd92258:077-p91s1

	cash	flow	statement
Is	<i>cash_;</i> t AAXXX-----1A-----	<i>flow</i> NNFXX-----A-----	<i>statement</i> NNIS1-----A-----
Should be	<i>cash-2_;</i> t A2-----A-----	<i>flow-2_;</i> t AAXXX-----1A-----	<i>statement_;</i> t
Should be	<i>cash_;</i> t NNIXX-----A-----	<i>flow_;</i> t NNIXX-----A-----	