

Univerzita Karlova v Praze
Filozofická fakulta
Ústav teoretické a počítačové lingvistiky

Jaroslava Hlaváčová

Formalizace systému české morfolgie
s ohledem na automatické zpracování
českých textů

Formalization of the Czech Morphology System
with Respect to Automatic Processing of Czech Texts

Disertační práce

Studijní program: Filologie
Studijní obor: Matematická lingvistika

Vedoucí práce: Doc. RNDr. Vladimír Petkevič, CSc.

Praha 2009

Prohlašuji, že jsem disertační práci vykonala samostatně s využitím uvedených pramenů a literatury.

Poděkování

Děkuji svému školiteli, docentu Petkeviči, za vlídné vedení a trpělivost.

Děkuji paní profesorce Panevové za konzultace ohledně lingvistických termínů a za připomínky k první verzi práce.

Děkuji své rodině za veškerou podporu. Svému muži Vaškovi za ohleduplnost. Dětem za výchovné pobídky k dokončení práce.

Práce byla podpořena grantem Informační společnosti č. 1ET100120503 poskytnutým Grantovou agenturou AV ČR a grantem č. 100008/2008 poskytnutým Grantovou agenturou UK.

Shrnutí

Přesný morfologický popis slovních tvarů je prvním předpokladem pro úspěšné automatické zpracování jazykových dat.

Systém kategorií a jejich hodnot, které se k popisu používají, jsou náplní první části práce.

Základním principem je tzv. **Zlaté pravidlo morfologie**, které říká, že každý slovní tvar by měl být v systému popsán jednoznačně. Existence variant na úrovni slovních tvarů i celých paradigmat však splnění tohoto pravidla komplikuje. Koncept variant rozšiřujeme na tzv. **mutace**, mezi které řadíme i jiné množiny slovních tvarů se stejným popisem (např. vícené tvary osobních zájmen). Mutace dělíme na **globální** pro popis na úrovni paradigmat a **flektivní** pro popis jednotlivých slovních tvarů. Toto rozdělení nám umožňuje postihnout jejich časté kombinace. Upouštíme od dělení variant (mutací) podle stylového příznaku jako neobjektivního kritéria. Při důsledném využívání hodnot kategorií **Flektivní mutace** a **Globální mutace** zůstane Zlaté pravidlo morfologie vždy splněno.

V kapitole o lemmatizaci zavádíme **vícenásobné lemma** pro popis variantních lemmat.

Podrobně se zabýváme popisem tzv. **složenin**, tedy slovních tvarů typu *zač, proň, koupilas, koliks*. Pro jejich lemmatizaci rovněž využíváme konceptu vícenásobného lemmatu. Podle slovních druhů jejich složek je dělíme na několik typů. Zabýváme se též problémem jejich vyhledávání v jazykových korpusech.

Druhá část práce popisuje systém vzorů pro popis slovních tvarů jednotlivých slovních druhů. U každého vzoru uvádíme sadu parametrů, které umožní postihnout velkou variabilitu v tvoření konkrétních paradigmat. Věnujeme se i pravidelnému odvozování příbuzných slov pomocí sufixů.

Abstract

Detailed morphological description of word forms represents one of the most important conditions of a successful automatic processing of linguistic data.

The system of categories and their values which are used for the description are the subject of the first part of the thesis.

The basic principle, so-called **Golden rule of morphology**, states that every word form has to be described by the system unambiguously. The existence of variants of word forms and whole paradigms, however, complicates the accomplishment of this rule. We introduce so called **mutations** as an extension of the variants to be able to include other sets of word forms with the same description (for instance multiple word forms of Czech personal pronouns). We divide mutations into two parts — **global** ones describing all word forms of a paradigm, and **inflectional** ones for the description on the word form level. This division enables us to express their various combinations. We do not use features of style for the mutation division, for they are subjective.

With a consistent use of the categories called Inflectional Mutation and Global Mutation, the Golden rule of morphology will always be valid.

The concept of multiple lemma is introduced in a chapter dealing with lemmatization. It describes lemma variants.

We give a detailed description of so-called **compounds**, which incorporate word forms of the type *zač*, *proň*, *koupilas*, *koliks*. The concept of multiple lemma is also used for their lemmatization. According to the word class of their components we divide the compounds into several types. We also deal with the problem of their searching in language corpora.

The second part of the thesis describes a system of patterns for word description. It is divided according to the part of speech. Each pattern has a special set of parameters that allow to grasp a large variability in word formation. We also deal with regular derivations of related words using suffixes.

Obsah

1	Úvod	1
1.1	Základní definice	2
1.2	Zlaté pravidlo morfolgie	6
2	Lemma a lemmatizace	7
2.1	Vícenásobné lemma	7
2.2	Vybrané problémy lemmatizace	8
2.2.1	Lemmatizace sloves	8
2.2.2	Záporná lemmata	12
2.2.3	Slovní tvary „bez lemmat“	13
3	Mutace	15
3.1	Motivace	16
3.2	Rozdělení mutací	17
3.3	Dosavadní pojetí variant v pražském a brněnském systému	18
3.4	Diskuse o hodnotách kategorie Mutace	18
4	Morfologické kategorie	21
4.1	Globální morfologické kategorie	22
4.1.1	Slovní druh (POS)	22
4.1.2	Poddruh (SUB)	24
4.1.3	Funkce (FCE)	30
4.1.4	Vid (ASP)	34
4.1.5	Zkratka (ABR)	34
4.1.6	Globální mutace (GMU)	34
4.2	Flektivní morfologické kategorie	37
4.2.1	Rod (GEN)	38
4.2.2	Číslo (NUM)	38
4.2.3	Duál (DUA)	38
4.2.4	Pád (CAS)	41
4.2.5	Osoba (PER)	41
4.2.6	Stupeň (DEG)	42
4.2.7	Negace (NEG)	42
4.2.8	Slovesný tvar (VRB)	43
4.2.9	Jmenný tvar přídavných jmen (NOM)	46
4.2.10	Stupeň intenzity slovesného děje (INT)	46
4.2.11	Typ složeniny (CMP)	47
4.2.12	Flektivní mutace (FMU)	47
4.3	Morfologická značka	49
4.4	Relevantnost kategorií	51
5	Kondicionál	54

6	Složeníy	56
6.1	Lemma složenin	56
6.2	Relevantní morfologické kategorie složenin	56
6.3	Typy složenin	58
6.3.1	Typy zájmenné ... n, c	58
6.3.2	Typ zájmenně-slovesný ... t	59
6.3.3	Typ zkratkový ... Z	59
6.3.4	Typy slovesné ... N, A, P, C, V, D, T, J, S	60
6.4	Vyhledávání složenin v korpusech	63
7	Morfologický slovník	68
7.1	Vztah morfologického slovníku a morfologických nástrojů	68
7.1.1	Guesser	69
7.2	Struktura slovníku	71
8	Vzory	75
8.1	Stručné porovnání pražského a brněnského systému vzorů	75
8.2	Nové vzory	76
8.2.1	Flektivní vzory	77
8.2.2	Derivační vzory	80
9	Vzory podstatných jmen	82
9.1	Obecné vlastnosti	82
9.2	Neživotné vzory	85
9.2.1	HRAD	85
9.2.2	STROJ	86
9.2.3	Kolísání mezi vzory HRAD a STROJ	88
9.2.4	ŽENA	88
9.2.5	PÍSEŇ	90
9.2.6	KOST	90
9.2.7	Kolísání mezi vzory KOST a PÍSEŇ	91
9.2.8	NŮŠE	92
9.2.9	MĚŠTO	93
9.2.10	MOŘE	94
9.2.11	KUŘE	96
9.2.12	STAVENÍ	96
9.3	Životné vzory	97
9.3.1	PÁN	97
9.3.2	MUŽ	99
9.3.3	Kolísání mezi vzory PÁN a MUŽ	99
9.3.4	PŘEDSEDA	100
9.3.5	SODUCE	100
9.4	Adjektivní vzory	101
10	Vzory přídavných jmen	102
10.1	Skloňování a stupňování	102
10.1.1	Základní část vzoru — skloňování	102
10.1.2	Stupňování	105
10.2	Derivace	108

Obsah

10.2.1	Tvoření jmenného tvaru	108
10.2.2	Tvoření příslovce	109
10.2.3	Tvoření podstatného jména na <i>-ost</i>	110
10.3	Příklady	110
10.4	Adjektivní skloňování dalších slovních druhů	112
11	Vzory pro příslovce	113
12	Slovesné vzory	116
12.1	Flektivní vzor	117
12.1.1	1. pozice — Imperativ	117
12.1.2	2. pozice — Prézens	118
12.1.3	3. pozice — Préteritum	119
12.1.4	4. pozice — Infinitiv	120
12.1.5	Přechodník P	120
12.1.6	Trpný rod T	121
12.2	Derivační vzory	121
12.2.1	Přídavná jména slovesná A	123
12.2.2	Deverbativní příslovce D	123
12.2.3	Podstatná jména slovesná N/O	124
12.2.4	Iterativní sloveso	124
12.3	Sdružené slovesné vzory	124
12.3.1	Příklady	127
13	Vzory zájmen a číslovek	129
13.1	Číslovky	129
13.1.1	Číslovky základní	129
13.1.2	Číslovky řadové a druhové	130
13.1.3	Číslovky úhrnné a souborové	130
13.1.4	Číslovky násobné, opakovací a výčtové	131
13.1.5	Číslovky dílové	131
13.2	Zájmena	131
13.2.1	Zájmena substantivní	131
13.2.2	Zájmena přivlastňovací	132
13.2.3	Zájmena ukazovací a vymezovací	132
13.3	Ostatní zájmena	132
14	Závěr	133
	Literatura	136
	A Přehled kategorií a jejich hodnot	139
	B Kopie účastnického slibu z Konkláve	144
	Rejstřík	145

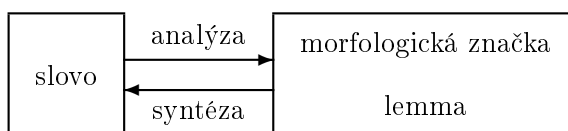
1 Úvod

Základem automatického zpracování jazyka je většinou rozsáhlý morfologický slovník, který popisuje slova daného jazyka. Slovník se využívá v zásadě k řešení dvou duálních úkolů. Prvním je rozpoznávání slov v textu, druhým vytváření slov (do textu).

Rozpoznání slova znamená jeho popis pomocí nějakých vlastností. Jednou ze základních vlastností každého slova je jeho základní tvar, neboli lemma, a slovní druh. Další vlastnosti se u flektivních jazyků potom liší, především podle slovního druhu. Každá vlastnost má nějaké hodnoty, které dohromady vytvářejí tzv. morfologickou značku. Rozpoznání slova tedy znamená určení lemmatu a morfologické značky, která toto slovo popisuje. Tomuto procesu se říká **morfologická analýza**¹. Morfologická analýza je základem jakéhokoli dalšího (automatického) zpracování textu. Bez ní, a potažmo bez morfologického slovníku, se sotva obejde náročnější jazyková aplikace, ať už je to strojový překlad, rozpoznávání mluvené řeči, dialogové systémy, či další složité úlohy.

Vytváření slov jazyka, neboli **morfologická syntéza** nebo také **generování**, je opačný proces než morfologická analýza. Slova se generují na základě lemmatu a morfologické značky.

Vztah morfologické analýzy a generování ilustruje schéma na obrázku 1.1.



Obrázek 1.1: Dualita morfologické analýzy a generování

Morfologický slovník, o kterém jsme se zmínili hned na začátku, by měl obsahovat popis co nejvíce (nejlépe všech) slov jazyka. Způsob popisu může být různý, většinou se používá nějaký systém vzorů.

Pro český jazyk existuje několik automatických popisů, z nichž nejznámější a zejména v akademickém světě nejpoužívanější jsou systémy pražský a brněnský. Každý z nich používá vlastní morfologický slovník. Existuje snaha oba systémy sjednotit, což znamená nejen spojit oba slovníky, ale především jednoznačně definovat morfologické kategorie a jejich hodnoty. Výsledkem by měl být takový popis české morfologie, který bude co nejúplnější², ale zároveň nebude zbytečně přegenerovávat, tzn. nebude obsahovat slova v češtině neexistující.

Přesto, že oba zmíněné systémy již léta slouží lingvistům (zejména v podobě automaticky označovaných jazykových korpusů) i informatikům (jako základ

¹Někdy se za morfologickou analýzu považuje jen přiřazení morfologických značek.

²Vzhledem k obecné povaze jazyka (a nejen českého) nemůže být úplný nikdy. Otázka, co do jazyka patří a co už nebo ještě ne, je velmi subjektivní.

nejrůznějších aplikací), stále je co vylepšovat. Tato práce je pokusem o takové vylepšení.

Naším cílem tedy je vytvořit rámec pro popis slovních tvarů českého jazyka, tzn. přesně definovat kategorie, které se k popisu používají, a stanovit pravidla, podle kterých se slovním tvarům přiřazuje jejich základní tvar, čili lemma. Zdůrazňujeme, že cílem je popis, nikoli vysvětlování ani zdůvodňování jazykových jevů. Tam, kde to je možné, používáme samozřejmě zavedených a odůvodněných lingvistických popisů, občas jsme však v zájmu jednoduchosti a jednoznačnosti popisu odhlédli od lingvistických hledisek a použili lingvisty neoblíbené „technické“ řešení.

Nezabýváme se tedy konkrétními „slovíčky“. Navrhujeme systém, který nejen umožní přesný popis pravidelných morfologických jevů, ale bude schopen konzistentně pojmut i výjimky. Na některé upozornujeme v textu, ale zcela jistě ne na všechny.

Při popisu jednotlivých morfologických kategorií budeme vycházet ze závěrů tzv. Morfologického konkláve (dále jen Konkláve), které se sešlo 21.–23. října roku 2005, aby definovalo jednotlivé morfologické kategorie (viz kopie účastnického slibu v příloze B na straně 144). Jeho závěry však nikdy nebyly dovedeny ke zdárnému konci. Složení Morfologického konkláve bylo (podle abecedy, bez titulů): Jan Hajič, Jaroslava Hlaváčová, Karel Oliva, Klára Osolsobě, Karel Pala a Vladimír Petkevič. Dokument, který na jeho podkladě vznikl, nebyl nikdy publikován. Budeme se snažit pečlivě oddělit výsledky Konkláve od vlastních řešení, ať už pouze doplňujících, nebo zcela odlišných.

Jsme si vědomi toho, že některé kategorie, které používáme, mají spornou definici. Někdy záměrně směšujeme více jevů do jedné kategorie (např. mutace), jindy naopak rozdělujeme zavedené kategorie a jejich hodnoty (např. kategorie Číslo a kategorie Duál). Primárním důvodem je vytvořit systém kategorií takový, aby se pomocí jejich hodnot daly popsat slovní tvary alespoň v takovém rozsahu, v jakém se dají popsat pomocí zmiňovaných systémů nyní.³ Navíc se snažíme odstranit nedostatky, které oba systémy mají a jichž si jsou jejich uživatelé, ale i tvůrci, vědomi. Snažíme se veškeré návrhy dobře popsat a zdůvodnit. K tomu využíváme příkladů z korpusu, z internetu i z vlastních pozorování. Příklady jsou průběžně číslovány. Není-li uvedeno jinak, jsou z korpusu SYN.

1.1 Základní definice

Klíčovými termíny pro popis systému morfologie jsou: slovní tvar, lemma, morfologická značka. Ač jsou tyto koncepty na první pohled jednoduché, při podrobné práci s nimi se dostáváme do situací, které nelze řešit bez pečlivě rozmyšlených definic a pravidel.

Slovo je řetězec písmen, který je na začátku i na konci ohraničen oddělovačem.

Oddělovač je řetězec nealfanumerických znaků. Většinou jde o mezery, o jiné tzv. „bílé znaky“ (white spaces) a znaky z množiny $\{(),,;':!?,\langle\rangle\}$, případně

³Výjimkou jsou pražské kategorie *Přivlastňovací číslo* a *Přivlastňovací rod*, které se patrně nevyužívají, a hodnota „Rodina“ kategorie *Rod* používaná v Brně, která by si zasloužila podrobnější rozbor.

další znaky. Oddělovače nejsou součástí slova, řadíme je do **interpunkce**. S interpunkcí se v českých korpusech zachází jako se zvláštním typem slovního tvaru. V pražském systému tvoří jednu z hodnot kategorie **Slovní druh**, podle níž ji lze vyhledávat. Interpunkce je zpracována dostatečně, není proto důvod se jí znovu zabývat⁴.

Často se uvádí, že slovo musí mít nějaký význam. Nepovažujeme tento požadavek za zásadní, a to ze dvou důvodů:

1. V textu se můžeme setkat s řetězci písmen, jejichž významu nerozumíme, buď proto, že žádný význam nemají, nebo proto, že ho nechápeme. Může jít třeba jen o překlepy (známá chyba vzniklá posunutím ruky písáčky na klávesnici), které učiní dané slovo nesrozumitelným. Chceme-li o takovém řetězci mluvit, používáme i v tomto případě termínu slovo. Slova jsou i záměrně zašifrované řetězce.
2. Existují věty, i správně syntakticky utvořené, které nemají žádný význam, např. *Rychlý strom respektuje rozpustnosti.*, abychom nepoužívali stále stejnou Chomského větu o zelených myšlenkách. Podle Patricka Hanksa a jeho teorie „Corpus Pattern Analysis“ (viz (Hanks – Pustejovsky, 2004)) nemají jednotlivá slova význam, ten dostávají až v kontextu, ve kterém se ocitají (viz též lexikální sémantika Cruse v (Cruse, 1986)). Slova v nesmyslném kontextu tedy mají nesmyslný význam, což se dá také říci tak, že význam nemají žádný.

Poznámka k významu slova „význam“ Už v úvodu, při první definici, jsme se dostali k termínu „význam“, který vůbec není jednoduchý. Mnoho lingvistů, filosofů, matematiků i jiných odborníků se pokoušelo význam definovat. Hned na začátek předesíláme, že se významem slova „význam“ nebudeme zabývat. To ovšem neznamená, že ho nebudeme používat. Naopak, používat ho musíme, neboť bez něj bychom nebyli schopni definovat další pojmy. Termín význam budeme tedy chápat velmi intuitivně takto: Slovo má význam, jestliže existuje kontext, ve kterém něco označuje. Slova z předchozí poznámky 2 tedy přece jen v tomto smyslu význam mít mohou.

Slovní tvar je slovo, které má význam.

Rozlišujeme tedy slovo jako řetězec písmen a slovní tvar jako řetězec písmen s významem. Tím se lišíme např. od definice Havránka a Jedličky (Havránek – Jedlička, 1981): „Slovo je skupina hlásek, která má zřejmý význam.“ Za slovní tvary nepovažujeme interpunkční znaménka. Ta patří mezi oddělovače.

Množinu slovních tvarů jazyka budeme značit S .

Příklady: *nějakou, bývala, stolečku*.

Lemma je základní slovní tvar. Ve slovnících se používá jako slovníkové heslo.

Množinu lemmat jazyka budeme značit L .

Příklady z předchozího odstavce mají lemmata *nějaký, bývat, stoleček*.

⁴Uvažujeme o vytvoření klasifikace funkcí jednotlivých interpunkčních znamének.

Všechny slovní tvary, které lze vytvořit z jednoho lemmatu pomocí skloňování, časování nebo stupňování (obecně ohýbání), tvoří tzv. **paradigma**⁵. Můžeme také říci, že paradigma je množina slovních tvarů, které náležejí danému lemmatu. Na rozdíl od většiny klasických mluvnic zahrnujeme do paradigmatu i nespisovné (nekodifikované) slovní tvary.

Morfologická kategorie je vlastnost slovních tvarů.

Každá morfologická kategorie má předem definovanou konečnou množinu hodnot, kterých může nabývat. Jestliže pro nějaký slovní tvar daná morfologická kategorie nenabývá žádné hodnoty, řekneme, že tato kategorie není pro tento slovní tvar relevantní. Můžeme také říci, že tato kategorie není relevantní pro celé paradigma nebo pro dané lemma jako jeho reprezentanta. A konečně, relevantnost morfologických kategorií se může týkat celých tříd lemmat se společnými vlastnostmi.

Morfologickou kategorii, která popisuje množinu slovních tvarů, budeme nazývat **relevantní morfologickou kategorií** této množiny. Množina může být jednoprvková, tedy jeden konkrétní slovní tvar, nebo víceprvková. V tom případě jde většinou o celé lemma, množinu lemmat stejného slovního druhu, případně i poddruhu. Obecně může jít o množinu libovolnou.

Hodnota kategorie, která není relevantní pro daný slovní tvar, lemma nebo třídu lemmat, má nedefinovanou hodnotu (UNDEF), tedy např. stupeň podstatného jména má hodnotu UNDEF.

Jako příklad morfologické kategorie uveďme slovesný vid, který má tři hodnoty: dokonavý, nedokonavý, obouvidý. Tato kategorie je relevantní pro slovesa a deverbativa. Není relevantní např. pro předložky.

V našem návrhu budeme hodnotám morfologických kategorií přiřazovat kódy, a to tak, aby byly co nejvíce v souladu s kódy používanými v současných morfologických systémech, ne vždy však bude možné shodu dodržet.

K popisu slovního tvaru je třeba většinou více morfologických kategorií. Jejich kódy potom vytvářejí morfologickou značku. Mluvíme-li tedy o morfologické značce, máme na mysli hodnoty relevantních morfologických kategorií daného slovního tvaru.

Hodnoty kategorií se určují v závislosti na ostatních kategoriích, především na kategorii slovního druhu a poddruhu, viz tabulku 4.7 na str. 52.

Za morfologické značky však považujeme jen takové řetězce, které kódují hodnoty všech relevantních morfologických kategorií pro daný slovní tvar (např. rod, číslo a pád pro podstatná jména). Takto chápaná morfologická značka vlastně popisuje slovní tvar obecně (např. podstatné jméno rodu ženského ve třetím pádě jednotného čísla), ve spojení s lemmatem potom popisuje konkrétní slovní tvar. Jinými slovy, máme-li lemma a morfologickou značku, můžeme vytvořit jednoznačně slovní tvar (tzv. Zlaté pravidlo morfologie, viz oddíl 1.2).

Značky, které jsou podspecifikované, tedy nemají vyplněné všechny relevantní kategorie, nepovažujeme za morfologické značky. Podle takové značky bychom totiž nebyli schopni pro dané lemma vygenerovat jednoznačný slovní tvar. Např. hodnota morfologické kategorie **Číslo** pro lemma *jarní* vygeneruje

⁵Termín paradigma se někdy používá ve významu „vzor“. My tyto dva termíny rozlišujeme (paradigma a vzor).

množství tvarů, které mají různé rody, pády, stupně. Zakódování této hodnoty samostatně tedy pro nás není morfologickou značkou⁶.

Takto tedy vypadá definice morfologické značky:

Morfologická značka je řetězec znaků, který kóduje hodnoty všech relevantních morfologických kategorií pro nějaký slovní tvar nějakého lemmatu.

Množinu morfologických značek budeme značit M . Této množině se říká (morfologický) **tagset**, my se ale budeme snažit tomuto převzatému termínu vyhnout.

Hodnoty morfologických kategorií popisují slovní tvary, proto můžeme i o morfologické značce říci totéž.

Dva nepoužívanější české systémy morfologických značek jsou systém pražský (viz (Hajič, 2004)) a brněnský (viz (Sedláček, 1999)). O jejich přednostech, záporech i rozdílech mezi nimi se už hodně mluvilo. Nejpřehlednější porovnání provedla Klára Osolobě (Osolobě). Budeme se o nich zmiňovat jen tehdy, jestliže bude nutné poukázat na nějaké rozdíly mezi řešením navrhovaným a již existujícím.

Pražský morfologický systém používá tzv. pozičního systému značek, kde každá pozice kóduje určitou kategorii, i když je pro daný slovní tvar nerelevantní, brněnský systém používá kompaktní značky uvádějící kódy jen relevantních morfologických kategorií. Oba systémy však mohou být ekvivalentní (bohužel nejsou — to ale není chyba kódování).

Při popisu morfologických kategorií a jejich hodnot zavádíme kódy především proto, abychom s nimi mohli v této práci dále pracovat, zejména pomocí nich vytvářet dotazy.

V kapitole 4.3 navrhujeme způsob, jak vytvořit morfologickou značku. Většinou však v celé práci pracujeme jen s hodnotami jednotlivých kategorií, protože konkrétní tvar morfologické značky není podstatný. Podstatné je pouze to, aby obsahovala všechny relevantní hodnoty.

Lemmatizací rozumíme zobrazení, které každému slovnímu tvaru přiřadí množinu jeho lemmat (viz (Hajič, 2004)):

$$\lambda: S \rightarrow 2^L$$

kde S je množina slovních tvarů a L množina lemmat.

Obvykle se lemmatizací rozumí zobrazení, přiřazující slovnímu tvaru jeho (jedno) lemma. Existují však lemmata, která jsou navzájem ortografickými variantami (*univerzita* — *universita*). V takových případech sdružujeme varianty pod společné lemma. Dalším důvodem k zavedení vícenásobných lemmat jsou tzv. složeniny⁷, které zavádíme v kapitole 4 jako zvláštní slovní druh. Složeniny nemají jednoduché lemma, protože jsou složeny z více slov, každé s jiným lemmatem. Určit jednoslovné lemma slovních tvarů jako *zač*, *každémus* nebo *proňš*

⁶Přesto se částečné kódování jen některých kategorií může pro určité aplikace hodit. Často používané je značkování pouze podle slovního druhu.

⁷Tímto termínem nemyslíme kompozita (slova složená), ta lemmatizujeme jako jedno slovo, neboť to nepřináší žádné problémy, např. *spolupřadatel*, *koupěschopný*. Termín se může zdát nevhodný, lepší jsme však ani po konzultacích s odborníky nevymysleli. Konkláve používalo termín ještě nevhodnější.

není jednoduché. Zde dobře poslouží množina lemmat jednotlivých složek složeniny, tedy vícenásobné lemma. O složeninách více v kapitole 6, o lemmatizaci v kapitole 2.

Morfologická analýza je zobrazení, které každému slovnímu tvaru přiřadí množinu dvojic ⟨lemma, morfologická značka⟩:

$$\mu: S \rightarrow 2^{L \times M}.$$

Lemmatizace je tedy součástí morfologické analýzy.

U homonymních slovních tvarů dostáváme dvě i více lemmat, např.

$$\lambda(\text{pekla}) = \{\text{peklo}, \text{péci}\}.$$

Každé z lemmat navíc může (a v uvedeném příkladě to tak opravdu je) být členem více než jedné dvojice ⟨lemma, morfologická značka⟩.

1.2 Zlaté pravidlo morfologie

Fakt, že jednomu slovnímu tvaru přiřadí zobrazení μ více různých hodnot, nevádí. Mnohem více vadí v jistém smyslu duální skutečnost, že jedné dvojici ⟨lemma, morfologická značka⟩ může odpovídat více než jeden slovní tvar. Např. 6. pád podstatného jména *hrad* v jednotném čísle může být slovní tvar *hradu* i *hradě*.

Bylo by výhodné, kdyby každá dvojice ⟨lemma, morfologická značka⟩ jednoznačně popisovala nejvýše jeden slovní tvar. Tomuto požadavku říkáme Zlaté pravidlo morfologie.

Při různých automatických aplikacích, které využívají generování slovních tvarů, je totiž těžké rozhodování, která z variant se má vybrat. Příkladem takové aplikace je třeba strojový překlad do češtiny, který v určité fázi musí vybírat v cílovém jazyce správný slovní tvar. Jestliže lemma i morfologická značka jsou pro dva tvary stejné, zodpovědný výběr je prakticky nemožný.

Je tedy třeba popis variantních slovních tvarů rozšířit tak, aby dvojice ⟨lemma, morfologická značka⟩ byla pro každý slovní tvar jednoznačná.

To lze udělat několika způsoby. Je možné zahrnout informaci o variantách do lemmatu nebo do morfologické značky, nebo vyčlenit tuto kategorii jako další atribut slovního tvaru. Poslední řešení jsme zvolili my. Z důvodů, které uvedeme dále v kapitole 3, jsme tuto kategorii nazvali **Mutace**.

Zlaté pravidlo morfologie tedy vypadá schematicky takto:

lemma + morfologická značka + mutace = jednoznačný slovní tvar
--

2 Lemma a lemmatizace

Základní jednotkou morfologického slovníku je lemma, které zastupuje celé paradigma slovních tvarů.

Lemmatizaci chápeme jako zobrazení λ z množiny slovních tvarů do množiny lemmat. Již v úvodu jsme naznačili, že zobrazení λ obecně nepřirazuje jediné lemma, ale množinu lemmat: $\lambda: S \rightarrow 2^L$, kde 2^L označuje množinu podmnožin množiny L .

Každému slovnímu tvaru přiřazuje zobrazení λ alespoň jedno lemma (např. $\lambda(\text{okna}) = \{\text{okno}\}$). Případy, kdy toto zobrazení není jednoznačné, nejsou v češtině řídké (např. $\lambda(\text{pekla}) = \{\text{peklo}, \text{péci}\}$). Je to způsobeno vysokou slovnědruhovou a morfologickou homonymií českého jazyka. Homonymie se bez kontextu (a někdy ani s ním) zbavit nelze.

2.1 Vícenásobné lemma

Kromě homonymie však existuje ještě jeden problém, který s lemmatizací souvisí. Jsou to varianty. Vezměme si např. slovní tvary *diskuze* a *diskuse*. Máme je analyzovat jako dvě různá lemmata, nebo varianty lemmatu jednoho?

Tato otázka má závažné praktické pozadí. Jestliže bude např. uživatel korpuse vyhledávat slovní tvary lemmatu *diskuze*, mají se zobrazovat jen slovní tvary se *-z-*, nebo i ty se *-s-*? S problémem se potýká i syntéza. Podle čeho se má z více vygenerovaných slovních tvarů se stejnými charakteristikami vybrat jeden?

Dosavadní morfologické slovníky se tímto problémem příliš nezabývají, a tak můžeme nalézt varianty na úrovni lemmatu, které jsou v pražském morfologickém slovníku zahrnuty pod jedno společné lemma (např. varianty *diskuze* i *diskuse* mají jediné společné lemma *diskuse*), i takové varianty, které jsou rozlišeny jako dvě různá lemmata (*citron* a *citrón*).

Ideální by bylo, kdyby lemma vždy odpovídalo slovnímu tvaru, ale kdyby se zároveň všechny varianty jednoho lemmatu sdružily, aby se daly například snadno vyhledat v korpusech. Toho lze dosáhnout zavedením konceptu vícenásobného lemmatu. Vícenásobným lemmatem z našeho příkladu jsou tedy dvouprvkové množiny $\{\text{diskuze}, \text{diskuse}\}$ a $\{\text{citron}, \text{citrón}\}$. Prvkům této množiny budeme říkat **variantní lemmata**.¹

Vícenásobné lemma zavádíme i pro taková variantní lemmata, která jsou nespisovná, zastaralá nebo jinak příznaková. Máme tedy např. i vícenásobné

¹Vícenásobnými lemmaty se zabýval také Karel Kučera (viz (Kučera, 2007)), ovšem z diachronního hlediska. Vzhledem ke změnám pravopisu slov v průběhu dějin potřeboval sdružit slova v různých etapách vývoje jejich zápisu. Na rozdíl od našeho řešení však zvolil tzv. hyperlemma jakožto zástupce množiny (historických) lemmat se stejným významem. Jeho hyperlemma je jediné a vybírá se ze současné slovní zásoby (pokud takové příslušné lemma existuje). Naše vícenásobné lemma je pohled na stejnou problematiku z hlediska synchronního.

lemma {*otevřít, vteřít*}, {*okénko, okýnko, vokýnko*}, {*blůza, bluza, blusa, blůsa, blůza*}.

Množina lemmat přiřazená homonymním slovním tvarům však vícenásobné lemma není, přestože je to také výsledek zobrazení λ .

Vícenásobné lemma je množina lemmat se stejným významem lišících se pouze zápisem (ortografické varianty).

Pro vícenásobné lemma definujeme ještě tzv. **rozšířené paradigma** jako sjednocení paradigmat jednotlivých variantních lemmat.

Kdykoli budeme v následujícím textu mluvit o lemmatu, budeme mít na mysli i případné vícenásobné lemma, pokud neuvedeme jinak. Pojem lemma tedy rozšíříme i na vícenásobná lemmata. Je to v souladu s naší definicí lemmatizace, tedy zobrazení λ , které nepřisuzuje slovním tvarům jednotlivá lemmata, ale množiny lemmat (i jednoprvkové).

Přesto však v případě, že lemma není vícenásobné, nebudeme nadále pro jeho vyjádření používat množinový zápis. Tedy např. $\lambda(\textit{polévkou}) = \textit{polévka}$. Stejně tak upustíme od množinového zápisu v případě lemmatu sloves, které, bráno zcela striktně, je také vždy vícenásobné, neboť infinitiv má vždy dva tvary, jak ukazuje příklad s vícenásobným lemmatem {*péci, péct*}.

Vícenásobné lemma jakožto množina lemmat přiřazená jednomu slovnímu tvaru poslouží i v případě složenin, u nichž není možné přirozeně jednoznačně a jednoduše zavést základní slovní tvar. Blíže se budeme lemmatizaci složenin věnovat v kapitole 6 o složeninách.

Pro složeniny tedy definujeme vícenásobné lemma poněkud odlišně:

Vícenásobné lemma složeniny je množina lemmat jednotlivých složek složeniny.

2.2 Vybrané problémy lemmatizace

2.2.1 Lemmatizace sloves

Lemmatem slovesa je jeho infinitiv. Toto jednoduché tvrzení je třeba podrobit detailnějším zkoumání pro některé speciální jevy.

2.2.1.1 Zvratná slovesa

Spory vzbuzuje dosavadní praxe přiřazovat infinitiv vždy bez zvratné částice, a to i v případě, že se jedná o reflexivum tantum. Námitka, že *smát* je nesmyslné slovo, je jistě správná, lemma by mělo být *smát se*, tedy dvě slova.

Taková lemmatizace by však mohla znamenat komplikace. V tom případě by totiž bylo logické, aby se stejné lemma přiřadilo nejen slovním tvarům paradigmatu *smát*, ale i vlastní zvratné částici. Pomiňme nyní nesmírně obtížné rozpoznávání zvratné částice, která ke slovesu patří, ve složitějších kontextech, kdy mohou být obě části zvratného slovesa od sebe vzdáleny, a to libovolným počtem slovních tvarů, dokonce na obě strany. Pomiňme i z toho vyplývající fakt, že takováto lemmatizace by nebyla možná v okamžiku morfologické analýzy, ale až po nějaké formě desambiguace, což by v důsledcích pravděpodobně znamenalo potřebu radikálně změnit zavedený postup automatického zpracování textů.

Ani samo přiřazení slovesného lemmatu zvrtné částici se nám však nejeví jako rozumné z hlediska automatického zpracování a posléze využívání lematizovaných korpusů. Kdybychom totiž důsledně zařazovali zvrtnou částici do lemmatu příslušného slovesa, bylo by potom asi také třeba obě části sdruženého lemmatu jednotně značkovat, tedy přiřadit i zvrtné částici hodnoty slovesných kategorií osoby, čísla, vidu, slovesného tvaru a dalších (viz kap. 4). Slovní poddruh zvrtné částice by potom ztratil význam a neznačkoval by se, což by mohlo přinést problémy při jejich vyhledávání v korpusech. Navíc takové značkování není obvyklé. Samozřejmě není třeba držet se zaběhaných postupů, jestliže zjistíme, že už nevyhovují. V tomto případě však takový krok není nutný.

Je tu i možnost, že by se zvrtná částice značkovala nezávisle na zbytku zvrtného slovesa, přičemž toto sloveso by se lematizovalo i se zvrtnou částicí, tedy např. $\lambda(směje) = \{smát\ se\}$ a $\lambda(se) = \{se\}$. Tato alternativa nám připadá nekonzistentní, a navíc zbytečná. Tím, že do lemmatu přidáme další slovo, totiž zvrtnou částici, pouze upozorňujeme na fakt, že jde o reflexivum tantum (jiná reflexiva by se, vzhledem k jejich nejasným vymezením, jako reflexiva značit neměla). Vyhledávání lemmatu *smát se* v korpusech by stejně mohlo přinést jako výsledek pouze tvary slovesa *smát*, bez zvrtné částice, a to ani kdyby se nacházela v těsné blízkosti. Takové řešení nám připadá nelogické a zbytečně složité.

Mnohem jednodušší je považovat lemma za technickou entitu, která reálně v jazyce nemusí existovat, ale která slouží k identifikaci celého paradigmatu. Přestože tedy samostatný tvar *smát* neexistuje, je možno pracovat s tímto tvarem jako s lemmatem na morfologické rovině. Teprve další roviny jazykového popisu sdruží jednotlivé slovní tvary ve správné spojení *smát se*.

2.2.1.2 „Stupňování“ sloves

Mnoho nedokonavých (avšak ne iterativních) sloves má schopnost spojovat se s některými speciálními předponami a se zvrtnou částicí *se* nebo *si*, a tím vytváří celé paradigma nových slovních tvarů s poměrně přesně definovaným významem. Předpony, k nim příslušející zvrtné částice a význam celého prefixovaného slovesa ukazuje tabulka 2.1.

Předpona	Sloveso	zvrtná částice	Význam
<i>roz-</i>	X	<i>se</i>	začít X
<i>po-</i>	X	<i>si/se*</i>	X v klidu, většinou příjemně
<i>za-</i>	X	<i>si/se*</i>	X po delší dobu a užít si to
<i>na-</i>	X	<i>se</i>	hodně X
<i>vy-</i>	X	<i>se</i>	hodně X a být s tím spokojen
<i>u-</i>	X	<i>se</i>	X až do vyčerpání

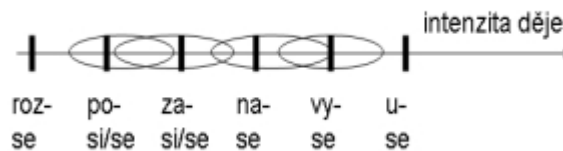
Tabulka 2.1: Přehled stupňovacích slovesných předpon. Hvězdička (*) v tabulce znamená, že pokud jde reflexivum tantum, zůstává i po přidání těchto prefixů zvrtná částice *se*.

Dosadíme-li v tabulce místo X např. sloveso *mávat*, dostaneme sadu nových paradigmat.

Uvádíme vymyšlené příklady uspořádané do malého příběhu, ze kterého by měl být zřejmý význam:

Včera se na nádraží pořádal kompars na nový film. Měli jsme za úkol (1) mávat na odjíždějící vlak. Když dal režisér pokyn, rozmávali jsme se. Nejdřív to vypadalo, že si pomáváme a půjdeme domů. Scéna s máváním se však mnohokrát opakovala, takže jsme si zamávali víc, než se nám líbilo. Namávali jsme se opravdu hodně, vymávali jsme se do sytosti. Měli jsme strach, že se umáváme k smrti.

Tímto způsobem jsme vytvořili sadu paradigmat s povinnou zvratnou částicí. Jednotlivá prefigovaná zvratná slovesa je možno s určitou tolerancí uspořádat podle intenzity děje (viz obr. 2.1).



Obrázek 2.1: Osa s naznačenou posloupností předpon podle stupňující se intenzity děje (zleva doprava).

Krajní body tvoří předpony *roz-* a *u-*, uprostřed je podle intenzity posloupnost *po-*, *za-*, *na-* a *vy-* s vágním až překrývajícím se rozsahem. Z tohoto důvodu nazýváme tento způsob tvoření s jistou nadsázkou pracovní „stupňování“ intenzity slovesného děje.

O všech těchto případech píše Kopečný v (Kopečný, 1962b), v poslední kapitole o českých slovesných předponách, není zde však učiněn závěr o „stupňování“. Uvádíme stručné charakteristiky z tohoto textu, u každého několik příkladů z korpusu SYN na dokreslení. Zmínuje se o nich i Šmilauer v (Šmilauer, 1971), který některá tato slovesa řadí do skupiny sloves vyjadřujících velkou „míru děje“.

2.2.1.2.1 roz- O jednom z významů předpony *roz-* Kopečný píše: „rozprouďení činnosti, obyčejně až po dosažení náležité míry“.

Někde v tom mlází se znova rozřukal datel (2)

Promnul si prsty, samým vzrušením se mu rozbrněly. (3)

2.2.1.2.2 po- Krátce si všímá i námi uvedeného významu předpony *po-*, o kterém píše, že může znamenat „i velkou míru děje“, což je jen částečně ve shodě s naším pozorováním.

místečko ve stínu, kde by si každý druhý pejsek rád pěkně pochrupkal (4)

pohrál si se startovacími klapkami na křídlech (5)

2.2.1.2.3 za- Kopečný se zmiňuje i o předponě *za-*, již přiřazuje „význam vzplanutí děje, jeho začetí“ a spojuje ji „s pocitem malé míry“, což nám nepřípadá zcela přesné, viz příklad ve vymyšleném příběhu. Ani ostatní uvedené příklady ze současné češtiny (6)–(8) tomuto hodnocení nenasvědčují.

o jejich tématech si budou moci rovnou i zachatovat. (6)

Poté zašel do tělocvičny, aby si podle svého rituálu ještě zaposiloval. (7)

Na druhou stranu : trochu si čas od času zašítet v takovéto nevinné záležitosti patří v této nespravedlností a stresem napěchované době skoro k léčebným procedurám. (8)

2.2.1.2.4 na- Výrazům s předponou *na-* (*namodlit se, nasmát se* apod.) říká Kopečný „augmentativnost“, případně „také intenzitivnost“.

Člověk se hrozně naběhá. (9)

matinka zatím doma vykládala, co se nastará a naběhá,... (10)

Co jsem se jen natančila to léto (11)

2.2.1.2.5 vy- U předpony *vy-* uvádí ještě případ zdvojené předpony *vyna-*, která však nezapadá zcela do naší pomyslné škály.

vyplačte se do sytosti (12)

V pohodě se tam celé dny můžeme do sytosti vyjezdít (13)

2.2.1.2.6 u- Předponu *u-* hodnotíme shodně jako Kopečný, když říká: „Reflexivní typ *upracovat se* je téměř paradigmatický.“ My považujeme za paradigmatické všechny právě vyjmenované typy.

To by se asi norští fanoušci uslavili k smrti. (14)

jak bylo zjištěno, unudit se nikdo nemůže (15)

málem se uštěkal (16)

Přestože se pomocí prefixu a reflexiva vytvoří nové sloveso s celým paradigmatickým, umísťujeme tato slova do paradigmatického neprefigovaného (základního) slovesa. Pokládáme zde prefixaci za tvoření slovesného tvaru, nikoli za slovo-tvorbu. Jinými slovy: tvrdíme, že všechny takto utvořené tvary mají společné lemma, a to X, tedy v našem příkladě *mávat*.

Důvodů je hned několik.

Předně je to velká produktivita. Není sice pravda, že takto lze vytvářet celou sadu od každého nedokonavého slovesa (protipříkladem budiž třeba tvar **zadotýkat se*), přesto lze takto vytvořit velké množství „nových“ slov. Navíc předpony v těchto slovech mají VŽDY stejný význam. Tento význam je naznačen v tabulce 2.1 a vyplývá též z příkladu (1).

Kdybychom tyto slovní tvary lemmatizovali jako samostatné předponové sloveso, museli bychom pro každé takové slovo zavést nové lemma, někdy homonymní s lemmatem již existujícím.

Některá takto vytvořená slova již ve slovní zásobě existují, ale mají jiný význam. V našem příkladě je to tvar *zamávat*, ovšem bez zvrtné částice. Příkladem úplné homonymie, včetně zvrtné částice, je *vysmát se*. Běžný význam je zřejmý z příkladu (17):

Budu na něj hodná a on se mi pak vysměje. (17)

Ale vyskytují se i příklady ve významu, který popisujeme zde:

Stavil jsem se tady jen proto, že se tady člověk může v klidu vysmát. (18)

I z těchto případů je zřejmý rozdíl významu. V první větě jde o *výsměch*, zatímco ve druhé o *smích*. Další rozdíl spočívá ve valenci. Zatímco v prvním příkladě jde o sloveso s akuzativní valencí, druhý příklad je intranzitivní. To však není pravidlem v jiných případech.

Podobné je *usmát se*, které může vypovídat buď o *úsměvu*, nebo opět o *smíchu*. I zde je rozdíl ve valenci (*usmát se* na koho — *usmát se* (bez valence)).

Výskyt stupňované intenzity není zpravidla vysoký, ale najdou se výjimky. Některá takto vytvořená slova jsou naopak velmi běžná, i s uvedeným významem, např. *rozesmát se*. V těchto případech je ovšem rozumné předponovou odvozeninu přímo zahrnout do slovníku.

Rozhodnutí, která prefigovaná (stupňovaná) slovesa jsou běžná, lexikalizovaná, a kde jde jen o okazionalismy, je samozřejmě velmi obtížné. Pravděpodobně v tom nebude panovat shoda, navíc se názory budou měnit v čase. Pro současný morfologický slovník je nejspíš nejlepší konzervativní řešení, tedy ponechat ve slovníku ta paradigmata, která tam jsou, včetně zavedené lemmatizace, ale nepřidávat globálně nová. Vycházíme z toho, že současný slovník již naprostou většinu běžných slov obsahuje.

Rozhodně do morfologického slovníku patří ta slovesa, která nejsou zvrtná nebo jsou tranzitivní, a dále potom tzv. odvozená reflexiva. Např. *utancovat se* je odvozené reflexivum od nezvrtného tranzitivního slovesa *utancovat* (koho). Lemma *utancovat* tedy do morfologického slovníku zahrnujeme.

Na závěr této kapitoly je třeba dodat, že naznačené uchopení zvrtných sloves s předponami *roz-*, *po-*, *za-*, *na-*, *vy-* a *u-* je třeba zpracovat ještě mnohem důkladněji. Tato kapitola nechtě slouží jen jako upozornění na velmi pravidelný, byť poměrně řídký úkaz ve vyjadřování expresivity sloves.

2.2.2 Záporná lemmata

Tato kapitolka se týká sloves, přídavných jmen a příslovcí. Začneme slovesy.

Předpona záporu *ne-* se obvykle při lemmatizaci odstraňuje, takže lemma je kladné, např. $\lambda(\textit{nenechal}) = \textit{nechat}$. Existují však slovesa, která se v kladném smyslu vůbec nepoužívají. Mohlo by se zdát, že v těchto případech by mělo smysl ponechat lemma záporné. Podívejme se však na tyto případy blíže.

Vezměme si příklad slovesa *nedutat*, které je lemmatizováno jako *dutat*. Tvrdíme, že tato lemmatizace je správná, přestože se sloveso *dutat* v kladném smyslu téměř nevyskytuje. Pro to, zda lemmatizovat sloveso se zápornou či nikoli, navrhuje test pomocí opisného budoucího času²:

²Tento test je možný jen u nedokonavých sloves, která tvoří budoucí čas s pomocným slovesem *být* (*bude...*).

V případě, že zápor budoucího času daného slovesa převede zápornku *ne-* k pomocnému slovesu, je třeba lematizovat sloveso kladně. V opačném případě tvrdíme, že kladný význam skutečně neexistuje.

Sloveso *nedutat*: říkáme *nebudu dutat*, nikoli **budu nedutat*. Zde je lemmatem *dutat* bez zápornky.

Komu se nelení, tomu se zelení. (19)

Lemma z příkladu (19) je *lenit*, nebo *nelenit*? Budoucí čas: *nebude lenit*, nikoli **bude nelenit*. Lemma je tedy *lenit*, a to přesto, že se v kladném smyslu sloveso *lenit* téměř neužívá.

Opačným příkladem je sloveso *nenávidět*, kde nelze říci **nebudu návidět*, vždy jen *budu nenávidět*. Zde je lemmatem *nenávidět* se zápornkou. Podobné je sloveso *nedoslýchchat*, kde sice máme sloveso *doslýchchat se*, ovšem s odlišným významem, navíc zvrátané.

Stejně jako existují záporná slovesa, existují i záporná přídavná jména a příslovce, v menší míře i záporná podstatná jména. U nich neumíme vytvořit podobně jednoduchý test, jako u sloves test s budoucím časem. Pouze přídavná jména, podstatná jména a příslovce od sloves odvozená převezmou záporné nebo nezáporné lemma podle (ne)zápornosti základního slovesa.

V případě, že lemma se zápornkou má odlišný význam než zápor lemmatu bez zápornky, je lemmatem tvar se zápornkou. Např. lemma *nesmyslný* není záporem lemmatu *smyslný*, jde tedy o dvě lemmata — *smyslný* a *nesmyslný*. První z nich teoreticky může přibrat zápornku a popisovat někoho, kdo není smyslný, druhé však už se zápornkou možné není (**nemesmyslný*).

Bohužel však existují slova, u nichž nepanuje shoda. Např. příslovce *kale*, které se běžně používá ve východních Čechách (viz příklad (20)), mnozí mluví jako nezáporné vůbec neznají a tvrdí, že jediné přípustné použití je se zápornkou, tedy *nekale* (příklad (21)).

Kale tomu nerozuměl. (odposlechnuto) (20)

Obě lupičské tlupy si nekale konkurovaly... (21)

Sporné je např. i podstatné jméno *nekřtěňátko*. V kladné podobě, tedy jako *křtěňátko*, je v korpusu SYN obsaženo jen jednou z 22 vyskytů. V pražském systému je lemma bez zápornky. Ve všech sporných případech musí rozhodnout tvůrce slovníku, kterou alternativu zvolit.

U rčení typu *hlava nehlava* lematizujeme obě složky stejně, a to jako nezáporná lemmata (v našem příkladě tedy *hlava*). Podobných rčení lze nalézt v korpusech celou řadu, jak ukazují příklady (22) a (23) z korpusu SYN.

režisér musí tvrdě zakročit, hvězda nehvězda (22)

Sochařka nesochařka, nejspíš bydlí v Greenwich Village. (23)

2.2.3 Slovní tvary „bez lemmat“

Na závěr se ještě dotkneme problematických slovních tvarů, které nemají jasně definované lemma.

Takovým slovním tvarům sice lze nějaké lemma přiřadit, ale stává se, že v tomto tvaru lemma vůbec neexistuje.

Např. slovo *krážem* ze spojení *křížem krážem* je sporné. V analogii s křížem by se mohl nabízet tvar **kráž*, avšak takový slovní tvar (současná) čeština nemá. V pražských korpusech je *krážem* lemmatizováno jako příslovce *krážem*, což je logické, ovšem předchozí *křížem* má lemma *kříž*. V tomto konkrétním případě by i slovní tvar *křížem* měl být lemmatizován nejspíš jako příslovce s lemmatem *křížem*. Samozřejmě by bylo nejlepší v tomto případě tato slůvka vůbec od sebe neoddělovat a lemmatizovat je dohromady. Taková spojení však zásadně necháváme na zpracování ve vyšších rovinách lingvistického popisu textu. Navíc jsme v korpusu SYN našli i příklad (24), i když zřejmě velmi netypický, kdy by takovéto spojení nebylo vhodné.

prošed křížem a krážem město (24)

Dalším takovým slovním tvarem je *bycha* ze spojení *pozdě bycha honit*. Lemmatem zřejmě nemůže být *bych* jako podstatné jméno, neboť takové neexistuje. Pražské řešení $\lambda(\textit{bycha})=\textit{bycha}$ nám připadá rozumné, včetně přiřazení akuzativu namísto pro lemma obvyklého nominativu.

Takových slovních tvarů je celá řada. Jejich lemmatizaci je třeba řešit individuálně. Považujeme za rozumné použít jako lemma sám slovní tvar (jako v příkladě *bycha*), i když nesplňuje podmínku obvykle na lemmata kladenou, totiž že se má jednat o infinitiv, resp. nominativ. My tuto podmínku od lemmat nevyžadujeme (viz definice lemmatu na str. 3).

3 Mutace

Již v úvodu jsme narazili na problémy týkající se variant. V kapitole 2 o lemmatech a lemmatizaci jsme zavedli pojem vícenásobné lemma s variantními lemmaty.

Problém variant je však mnohem širší, dokonce ani sám termín varianta není jednoduchý. Jeho pojetí v lingvistické literatuře totiž zdaleka není jednotné, navíc se používá ještě termín dubleta, s podobným významem. Přehled různých chápání obou termínů je stručně podán v (Tušková, 2006). Zde si autorka vzala za východisko Mluvnici češtiny (viz (Akademická mluvnice)), která dělí varianty na rovnocenné a diferencované. Rovnocenné jsou ty, které jsou „rovnocenné významově, funkčně i stylově, a jsou navzájem volně zaměnitelné“, zatímco diferencované nelze volně zaměňovat kvůli stylovému zabarvení, dobové vázanosti, frekvenci nebo různému významu. Jak je vidět, toto dělení není zdaleka jednoznačné ani objektivní.

Patrně by nebylo těžké zavést nějakou formální definici založenou na ortografických rozdílech v zápisu, ale je jisté, že s takovou definicí by nikdy nebyli spokojeni všichni. Ukazuje se totiž, že varianta je pojem velmi různorodý a většinou je spojen s nějakým typem hodnocení — stylu, časového zařazení, dialektu, a podobně.

Při formálním morfologickém popisu jednotlivých slovních tvarů nás však tato hodnocení nezajímají. Naopak bychom se jim chtěli vyhnout, protože často nemají jednoznačná kritéria. Někakou kategorii, která rozliší slovní tvary se stejnými hodnotami morfologických kategorií, však potřebujeme, abychom mohli vždy zachovat platnost Zlatého pravidla morfologie. Protože však termín „varianta“ má již v lingvistice svůj význam (i když ne zcela přesně vymezený), navrhuje pro naše účely termín jiný, a to mutace. Jeho vymezení je čistě technické:

Mutace jsou takové dvojice slovních tvarů, které mají stejné lemma a které nelze rozlišit hodnotou žádné jiné morfologické kategorie. Jinými slovy jsou to takové dvojice slovních tvarů, pro které mají všechny morfologické kategorie stejnou hodnotu.

Pojem mutace je širší než varianta, mezi mutace řadíme totiž nejen varianty (v obvyklém významu), můžeme mezi ně zařadit např. i dvojici vokalizované a nevokalizované předložky, které se za varianty nepovažují. Také rozdílné tvary osobních zájmen, např. *jeho*, *ho*, *něho*, *něj*, *jej* nejsou pravými variantami, přestože hodnoty všech jejich klasických morfologických kategorií jsou stejné. V takových případech bychom sice mohli zavést nové kategorie se speciální sadou hodnot, které uvedené tvary rozliší, ale zavedení kategorie **Mutace** umožní vyřešit tento problém pro všechny podobné případy najednou.

Zavádíme tedy novou kategorii **Mutace**, která svými hodnotami rozliší mutace, jak jsme je vymezili v předchozích odstavcích.

Nejsnadnější je vymezení tzv. flektivních mutací, které se liší v zakončení (např. *hradu* — *hradě*). Zahrnujeme mezi ně i nespisovné varianty.

U ostatních (globálních) mutací, které se projevují v celém paradigmatu, je situace složitější. Někdy jde o pouhé ortografické varianty (*atomismus* — *atomizmus*), někdy o různou výslovnost (*citron* — *citrón*), případně ovlivněnou obecnou češtinou (*mýdlo* — *mejdlo*), jindy o různé způsoby tvoření (*brzy* — *brzo*).

V této kapitole popíšeme právě zavedenou kategorii **Mutace** v souvislosti s nevyhovujícím řešením variant v současných morfologických systémech.

Hlavním myšlenka spočívá v rozdělení mutací na dvě skupiny, globální a flektivní, podle toho, zda se týkají celého paradigmatu, nebo jen některých jeho slovních tvarů.

3.1 Motivace

Představme si, že bychom chtěli vyhledat v nějakém korpusu všechny slovní tvary náležející lemmatu *okénko*. V lemmatizovaném korpusu to uděláme snadno dotazem např. [lemma="okénko"]. Mělo by nás však také napadnout, že lemma *okénko* má variantu *okýnko*, která, ač ji můžeme chápat jako „míň spisovnou“, je uváděna jako rovnocenná (SSJC). A pokud bychom chtěli být opravdu důslední, měli bychom zahrnout i nespisovnou variantu *vokýnko*. Výsledný dotaz na korpus by potom mohl vypadat např. takto: [lemma="v?ok[éy]nko"]. Přitom předpokládáme, že slovník, podle kterého se náš korpus lemmatizoval a značkoval, obsahuje lemmata *okénko*, *okýnko* a *vokýnko* (zadaný dotaz připouští i možnost **vokénko*, která však zřejmě neexistuje).

Vzhledem k tomu, že se jedná o varianty jednoho lemmatu, není asi rozumné vytvářet tři různá lemmata. Mnohem přirozenější by bylo, kdyby všechny tvary všech tří variant měly jedno společné lemma. Rozhodněme tedy, že společným lemmatem uvedených tří variant bude např. *okénko*. Budeme-li nyní chtít vyhledat v korpusu např. 7. pád množného čísla lemmatu *okénko*, vyhledají se nám tyto tvary:

- | | |
|-------------------|---------------------|
| 1. <i>okénky</i> | 4. <i>okénkama</i> |
| 2. <i>okýnky</i> | 5. <i>okýnkama</i> |
| 3. <i>vokýnky</i> | 6. <i>vokýnkama</i> |

První dvě varianty jsou spisovné a víceméně rovnocenné, zbylé jsou nespisovné, přičemž varianta 3 (*?vokýnky*) je podivná, neboť má nespisovný kmen a spisovnou koncovku. Mnohem přirozenější je *vokýnkama* (6). S takovou vyhledávkou budeme spokojeni.

Co však s opačnou úlohou, kdybychom chtěli 7. pád množného čísla lemmatu *okénko* vygenerovat? Dostaneme 6 různých slovních tvarů. Mohli bychom je rozlišit na úrovni stylových příznaků, avšak vidíme, že spisovných i nespisovných tvarů máme více, a to hned ze dvou důvodů. V příkladech 3 a 6 jde o nespisovnost způsobenou protetickým *v-*, v příkladech 4, 5 a 6 je nespisovnost způsobena hovorovou koncovkou *-ama*. My bychom ale chtěli, aby každá relevantní kombinace hodnot morfologických kategorií pro dané lemma vygenerovala jednoznačný slovní tvar (Zlaté pravidlo morfologie — viz 1.2). Jinými

slovy, aby se popis uvedených šesti slovních tvarů v něčem lišil: v lemmatu nebo v hodnotě nějaké morfologické kategorie. Tuto kategorii jsme nazvali Mutace.

3.2 Rozdělení mutací

Jak už jsme si všimli v předchozím oddíle, ne všechny mutace jsou stejného typu. Můžeme rozlišovat mutace na úrovni jednotlivých slovních tvarů, vyjadřujících určitou kombinaci hodnot morfologických kategorií, např. *okýnky* a *okýnkama* pro 7pl, a na úrovni celých paradigmat, např. dvojice *okýnko* — *vo-kýnko*, *okýnka* — *vokýnka* atd. V prvním případě jde o spisovnou a nespisovnou koncovku, které mohou být připojeny k libovolnému podstatnému jménu rodu středního, které se skloňuje podle vzoru město (*městy* — *městama*), zatímco ve druhém případě jde o nespisovnost týkající se libovolného tvaru většiny lemmat začínajících na *o*¹.

Ve slovním tvaru *vokýnkama* se potom spojují oba druhy nespisovnosti, tento tvar je tedy nespisovný „na druhou“.

Některé slovní tvary však vykazují mutace, které nejsou systematické vůbec, týkají se pouze a jenom příslušného lemmatu. Snad nejznámějším případem je dvojice *brzy* a *brzo*.

Mutace tedy mohou být dvojího typu, přičemž jeden nevylučuje druhý:

1. mutace týkající se celého paradigmatu, tj. všech slovních tvarů,
2. mutace týkající se jen některých tvarů daného paradigmatu.

První typ budeme nazývat mutací globální, druhý typ mutací flektivní.

Flektivní mutace jsou takové mutace, které se projevují jen v některých tvarech paradigmatu.

Globální mutace jsou takové mutace, které se projevují ve všech tvarech paradigmatu, a to všude stejně.

Obrázek 3.1 ukazuje, jak se mohou globální i flektivní mutace kombinovat.

	globálně	
flektivně		
spisovné	<i>okny</i>	<i>vokny</i>
nespisovné	<i>oknama</i>	<i>voknama</i>

Obrázek 3.1: Příklad kombinace globálních a flektivních mutací

¹Zdaleka ne všechny však mohou protetické *v*-přijímat. Protipříklady: **voficiální*, **vovšem*, **vondatra*, **vovoce*, **votec*, ale *vocet* ano.

3.3 Dosavadní pojetí variant v pražském a brněnském systému

Současné systémy se problému variant (mutací) věnují jen okrajově, a proto v nich také často nebývá splněno Zlaté pravidlo morfologie. V pražském systému je variantám vyhrazena 15. pozice v morfologické značce (viz (Hajič, 2004)). Pro ni je vytvořen číselník, který kóduje stylové příznaky. Globální a flektivní mutace zde nejsou rozlišeny, což vede k velké nekonzistenci v popisu.

Většinou se rozlišují mutace flektivní, ale ne vždy. Flektivní mutace, které jsou systematické, bývají obvykle zahrnuty ve vzoru, podle kterého se generují všechny tvary daného paradigmatu, ale značí se i mutace nesystematické, např. mutace *dýchá* se značkou VB-S---3P-AA--- / *dýše* se značkou VB-S---3P-AA--2.

Některé globální mutace se popisují samostatnými lemmaty, např. *busola* — *buzola*, jiné se sice sdružují pod společné lemma, ale potom nastává situace, kdy se k danému lemmatu a morfologické značce vygenerují dva různé slovní tvary, např. lemma *klauzule*, ke kterému existuje (zastaralá) mutace *klausule*. Ve slovníku jsou zahrnuty obě mutace pod společným lemmatem *klauzule*, ale nejsou jako mutace označeny. Každá morfologická značka tedy vytvoří pro toto lemma dva slovní tvary, jeden se -s- a druhý se -z-, čímž porušuje Zlaté pravidlo morfologie.

Brněnský systém se mutacemi zabývá ještě méně. Různé globální mutace jednoho lemmatu se považují za lemmata dvě. Např. *busola* i *buzola* mají stejné morfologické značky ve všech příslušných dvojicích svých tvarů, ale lemmata jsou různá. Svým způsobem je toto řešení lepší než pražské, protože se tím neporuší Zlaté pravidlo morfologie.

Zato flektivní mutace v brněnském systému Zlaté pravidlo porušují často. Existuje sice gramatická kategorie „Stylistický příznak tvaru“, ale v korpusu DESAM je využívána minimálně. Např. slovní tvary *soudci* i *soudcové* mají stejnou značku k1gMnPc1, zatímco v pražském systému dostávají *soudci* značku NNMP1-----A---- a *soudcové* NNMP1-----A---1, jsou tedy rozlišeny na 15. pozici jako rovnocenné varianty.

V přehledu morfologických kategorií (viz kap. 4) jsme umístili flektivní mutaci (FMU) mezi flektivní kategorie a globální mutaci (GMU) mezi kategorie globální. Důvod je zřejmý. Globální mutace popisuje celé paradigma, zatímco flektivní mutace se týká jen některých jeho slovních tvarů. Oba typy mutací se mohou navíc kombinovat, a na to by jediná kategorie nestačila.

Zbývá nyní vyřešit otázku hodnot kategorií FMU a GMU.

3.4 Diskuse o hodnotách kategorie Mutace

Vzhledem k tomu, že každý slovní tvar může vykazovat jen malé množství mutací, ať už globálních, nebo flektivních, stačilo by oba typy mutací v rámci jednoho slovního tvaru jednoduše číslovat. Neznáme případ, kdy by slovní tvar měl více než 10 mutací (platí i pro globální a flektivní mutace zvlášť), takže by k popisu stačily číslice 0 až 9.

Nejjednodušší by bylo použít číslice pouze k formálnímu odlišení mutací. Na druhou stranu jsou ale uživatelé korpusů zvyklí na to, že hodnota mutace, když

už je někde uvedena, také o něčem vypovídá. V současných systémech se její hodnota využívá k odlišení stylových příznaků mutací.

V příspěvcích (Hlaváčová, 2008) a (Hlaváčová – Lopatková, 2008) jsme proto navrhli — ve shodě s dosavadní praxí — zakódovat do hodnot FMU i GMU stylový příznak. Tato vlastnost je velmi subjektivní, proměnlivá v čase i místně, takže je velmi obtížné ji stanovit. Pokud by záleželo skutečně jen na rozlišení jednotlivých mutací kvůli jednoznačné morfologické značce, na konkrétní hodnotě (číslu) by nemuselo záležet. Přesto jsme navrhli kódovat mutace pomocí pomyslné škály, kdy 0 by ležela uprostřed jako nejběžnější synchronní mutace. Číslice 1, 2, ... by vyjadřovaly mutace rovnocenné, přičemž řada by směřovala „do budoucnosti“. To znamená, že tyto hodnoty by dostaly mutace nespisovné, případně hovorové či obecně české (budoucnost je tu míněna jako potenciální čas, kdy by nespisovné mutace eventuálně mohly být uznány jako spisovné). Číslice 9, 8, 7, ... by potom směřovaly „do minulosti“, což znamená, že by popisovaly mutace už nepoužívané, zastaralé, archaické, případně nářeční.

Je jasné, že taková škála má velmi mnoho nevýhod, a je tedy snadno napadnutelná. Je velmi abstraktní, časové vztahy jsou, zvláště co se týče budoucnosti, často jen hypotetické. Některé mutace, např. nářeční, na takovou škálu ani nepatří, potřebovaly by svou samostatnou škálu.

Další námitkou proti tomuto řešení je fakt, že stylové příznaky nejsou přijímány celou lingvistickou komunitou jednoznačně. Spíše by měly být (a také jsou) cílem výzkumu (viz např. (Křístek, 2002)). Morfologický slovník by měl být na subjektivních názorech jednotlivých badatelů nezávislý. Z toho důvodu bychom se neměli snažit hodnoty morfologických kategorií, tedy ani mutací, jakkoliv hodnotit.

Z výše uvedených důvodů od takového návrhu upouštíme a přiřazujeme kategoriím **Flektivní mutace** a **Globální mutace** nezávislou sadu hodnot, prostou jakéhokoli hodnocení. Pomocí těchto hodnot se snažíme vyjádřit obecné vlastnosti mutací. Nejběžnější dvojice hodnot obou kategorií uvádí tabulka 3.1. Z uvedených příkladů vidíme, že tyto hodnoty se mohou týkat nejrůznějších typů slovních tvarů (v případě FMU) i celých paradigmát (v případě GMU).

Dvojice hodnot	Vysvětlení	Příklad
D — K	delší — kratší (počet písmen)	<i>pustější — pustší</i> <i>skáčeme — skáčem</i> <i>vracejí — vrací</i>
d — k	dlouhá — krátká (samohláska)	<i>musím — musím</i> <i>zavříno — zavřeno</i> <i>tráv — trav</i> <i>salón — salon</i>
t — m	tvrdá — měkká	<i>vlaštovka — vlastovka</i> <i>student — študent</i> <i>mazám — mažu</i>

Tabulka 3.1: Základní hodnoty kategorií **Flektivní mutace** a **Globální mutace**

Již v úvodní kapitole jsme naznačili, že mutace, tzn. kategorie FMU ani GMU, nebudeme zahrnovat do morfologické značky. Vyčleňujeme je jako další, speciální atributy popisu a souhrnně jim říkáme **Mutace**. V rámci této kate-

3 Mutace

gorie lze libovolně kombinovat jednotlivé typy mutací, globálních i flektivních. Každý slovní tvar je tedy jednoznačně popsán hodnotami svého lemmatu, morfologické značky (bez mutací) a hodnotou kategorie **Mutace**.

Toto řešení nám připadá velmi slibné, neboť jím lze zachytit neomezené množství kombinací jednotlivých mutací. Vyhneme se tím také jednomu z obvyklých požadavků, totiž stanovení „základní mutace“, která by měla být nejběžnější, což je ovšem většinou velmi těžké stanovit. Kód, který zaznamená hodnoty kategorie **Mutace**, zapíšeme pomocí regulárního výrazu takto:

$$\text{MUT} = \text{F.}+\text{G.}+$$

Znaky zapsané na místě tečky za písmenem F kódují typy flektivní mutace, znaky za G typy mutace globální.

Kromě všeobecných hodnot z tabulky 3.1 existuje celá řada mutací typických pouze pro některé kombinace hodnot morfologických kategorií. Množinu hodnot typických pro kategorii **Globální mutace** uvádíme v tabulce 4.4 na str. 34. O kategorii **Flektivní mutace** pojednáme blíže v oddíle 4.2.12, o jednotlivých hodnotách potom v kapitolách 9 až 13, kde popisujeme vzory konkrétních slovních druhů.

4 Morfologické kategorie

Nejprve vytvoříme množinu všech morfologických kategorií¹, které se k popisu slovních tvarů používají. Množinu hned zpočátku rozdělíme na dvě podmnožiny. První bude obsahovat tzv. globální kategorie, druhá kategorie flektivní.

Globální morfologická kategorie je taková kategorie, jejíž hodnota je stejná pro celé paradigma. Příkladem globální morfologické kategorie je slovní druh.

Flektivní morfologická kategorie je taková kategorie, jejíž hodnoty se pro jednotlivé slovní tvary jednoho paradigmatu liší. Příkladem flektivní morfologické kategorie je pád. Také kategorie rod je flektivní, ovšem jen v rámci přídavných jmen a některých zájmen a číslovek.

U každé kategorie uvedeme:

1. hodnoty, kterých může nabývat;
2. pro jaké druhy slovních tvarů je relevantní; nemá např. smysl určovat rod příslovcí, nebo stupeň přivlastňovacích přídavných jmen. Pro jednoduché formulace v následujícím textu si situaci zjednodušíme tak, že uvažujeme každou kategorii pro každý slovní tvar s tím, že jedna z hodnot může být UNDEF.

Každé kategorii přidělíme zkratkové jméno.

Každé hodnotě přiřadíme jednoduchý symbol. Jak už bylo uvedeno, tento symbol může být v různých systémech různý. Zde ho uvádíme proto, aby bylo možno jednoduše tvořit dotazy a nezaplést se do zbytečně zdlouhavých popisů. Hodnoty morfologických kategorií vycházejí převážně ze symbolů pro tytéž nebo podobné hodnoty v pražském systému (viz (Hajič, 2004)).

Hodnoty jednotlivých kategorií jsou v ideálním případě ekvivalencí na množině všech slovních tvarů. Tato ekvivalence rozděluje slovní tvary na navzájem disjunktní třídy, které celou množinu pokrývají. To mimo jiné znamená, že hodnoty každé kategorie jsou „vyčerpávající“, tzn. že neexistuje slovní tvar, který by neměl přiřazenou hodnotu každé kategorie, počítáme-li i hodnotu UNDEF.

V praxi ovšem často dochází k případům, že jeden slovní tvar lze popsat více hodnotami jedné kategorie (např. kategorie Pád pro slovní tvar *stavení*). Při postupném přechodu do dalších rovin lingvistického popisu se však množina hodnot jednotlivých kategorií obvykle zmenšuje. Už syntaktický rozbor věty vybere z množiny všech možných hodnot dané kategorie většinou hodnotu jedinou. V některých případech je třeba pro desambiguaci zkoumat další roviny, sémantickou, případně pragmatickou, které berou v úvahu kontext širší než jedna věta, nebo i mimojazykové znalosti. Některé věty zůstávají víceznačné i po podrobných rozborech.

¹Mluvíme o morfologických kategoriích, přestože některé z nich ze striktního pohledu čistě morfologické nejsou. Už sama nejdůležitější kategorie, tedy slovní druh, nevyjadřuje jen morfologické vlastnosti.

Na úvod ještě poznamenejme, že většina kategorií a jejich hodnot je tradiční, některé jsou však nové, vytvořené speciálně pro snadnější automatické zpracování textů. Tradiční lingvisté se mohou nad některými návrhy pozastavovat. Budeme se však snažit nové přístupy vždy důkladně odůvodnit.

Ve zbytku kapitoly představíme všechny morfologické kategorie, které se používají k popisu českých slovních tvarů. U nových kategorií nebo hodnot, které vyžadují podrobnější vysvětlení, necháváme rozbor na zvláštní kapitoly.

4.1 Globální morfologické kategorie

Připomeňme si definici globální kategorie:

Globální morfologická kategorie je taková kategorie, jejíž hodnota je stejná pro celé paradigma.

Zde je výčet všech globálních kategorií

- Slovní druh — POS
- Poddruh — SUB
- Funkce — FCE
- Vid — ASP
- Zkratka — ABR
- Globální mutace — GMU

4.1.1 Slovní druh (POS)

Tato kategorie je základní v tom smyslu, že má definovanou konkrétní hodnotu pro každý slovní tvar. Jinými slovy, v této kategorii nezavádíme hodnotu UNDEF. Tento požadavek se dá chápat také tak, že kategorie slovní druh rozkládá množinu všech slovních tvarů na podmnožiny, které celou množinu pokrývají. Ideální by bylo, kdyby tyto množiny mohly být disjunktní, takže by se o každém slovu dalo jednoznačně říci, jaký má slovní druh. To bohužel neplatí, a to z několika důvodů:

- Homonymie
Homonymní slovní tvary mají více slovnědruhových interpretací, např. *drát* jako podstatné jméno a *drát* jako sloveso. Tento problém se snadno vyřeší tím, že se víceznačná lemmata nějakým způsobem odliší. V našem příkladě budeme mít potom dvě různá lemmata *drát*: *drát-1* a *drát-2*.
- S-formy
Název jsme převzali od brněnských kolegů. Jde o skupinu slov, která nejsou jednoznačně zařaditelná. Může to být proto, že v různých kontextech nabývá slovní druh různých hodnot, přestože můžeme říci, že jejich význam je stále týž. Často se ani renomovaní lingvisté neshodnou, jaká hodnota by to v daném kontextu měla být. Většinou jde o slova synsemantická (funkční). I tento problém lze vyřešit zmnožením jednotlivých

lemmat. V tomto případě je však podstatné povolit při morfologické analýze více možných hodnot, potažmo více morfologických značek. Příklady s-forem: *dokud* (příslovce, spojka), *tedy* (spojka, částice), *jak* (příslovce, spojka, částice i podstatné jméno), *jednou* (číslovka, příslovce, částice).

Hodnoty kategorie Slovní druh

- N: podstatné jméno
- A: přídavné jméno
- P: zájmeno
- C: číslovka
- V: sloveso
- D: příslovce
- R: předložka
- J: spojka
- I: citoslovce
- T: částice
- F: cizí slovo (K)
- G: prefixový segment (K)
- S: složenina (K)
- X: neznámé slovo

Slovní druhy označené znakem (K) v závorce byly nově zavedeny na zasedání Konkláve. Znak X označuje slovní druh „neznámé slovo“, které sice k tradičním slovním druhům nepatří, ale lingvisté se s ním již dávno setkávají při práci s jazykovými korpusy. Ostatní slovní druhy jsou tradiční, budeme je nadále považovat za známé a jejich definice uvádět nebudeme.

Uveďme si charakteristiky nových slovních druhů.

4.1.1.1 Cizí slovo

je slovo, které nepodléhá české flexi a nemá v češtině vlastní význam.

Příklady cizích slov: *the, you, der, di, du, to, company, ...*

Nepatří sem však neskloňná slova, která jsou součástí české slovní zásoby, jako např. *kupé², lila* (barva). Ta sice české flexi také nepodléhají, ale mají v češtině jasný význam. Může se stát, že cizí slova, zvláště ta krátká, jsou homonymní se slovy českými (z uvedených příkladů jsou to slova *der* — imperativ slovesa *drát*, *di* a *du* jako nespisovné tvary slovesa *jít*, i slovní tvar *to* homonymní s českým ukazovacím zájmenem). Lemma cizího slova je vždy stejné jako slovo samo.

Zavedení slovního druhu cizí slovo velmi usnadní automatické zpracování.

²Děkuji vedoucímu práce za upozornění na existenci neobvyklého tvaru *kupém*, který však spíše potvrzuje, že toto slovo nepatří mezi slova cizí, neboť jakési české flexi zdá se podléhá.

4.1.1.2 Prefixový segment

je začátek slova, který stojí samostatně, a teprve někde dál v textu je doplněn na plnovýznamové slovo. Většinou to bývá předpona. Jde-li o jiný řetězec, vždy se v daném kontextu jako předpona chová.

Příklady: *česko a rusko - německý, tři až čtyřprocentní*. Spojovník, který případně může stát bezprostředně za prefixovým segmentem, ať už s mezerou, nebo bez ní, není součástí ani slova, ani lemmatu.

Lemma prefixového segmentu se shoduje se slovním tvarem. V uvedených třech příkladech jsou tedy lemmaty slova *česko, rusko* a *tři*.

4.1.1.3 Složenina

popisuje slovní tvar, který vznikl z více slov (většinou různých slovních druhů) a určení jeho slovního druhu je problematické. Jejich vymezení a podrobný popis jsou uvedeny v kapitole 6.

Příklady: *zač, oň, byls*.

Pod termínem „složenina“ chápeme tedy něco jiného než např. Šmilauer v (Šmilauer, 1971), který pokládá termíny „složenina“ a „složené slovo“ za synonyma. Složená slova (kompozita), sice také vznikla z více slov, avšak jejich zařazení mezi slovní druhy nečiní potíže (např. *středomořský, vlastizrádce, spolumvíník*).

4.1.1.4 Neznámé slovo

je takové slovo, jehož slovní druh neumíme určit. Jsou to slova, která morfologická analýza nerozpozná.

Tato hodnota kategorie **Slovní druh** už v pražském systému existuje a uživatelé s ní pracují.

Při ruční anotaci se neznámým může stát pouze takové slovo, které anotátor nezná a nemůže určit. Může to být například nějaká šifra nebo naprosto nerosozumitelný překlep (např. poslední slovo před touto závorkou se posunutím pravé ruky na klávesnici může stát neznámým slovem *ořejkeoz*).

Neznámá slova nejsou součástí morfologického slovníku. Při automatické morfologické analýze je tedy tato hodnota přiřazena těm slovním tvarům, které nejsou rozpoznány. Z neznámého slova se může stát známé buď přidáním do slovníku, nebo dalšími metodami (guesser — viz oddíl 7.1.1).

4.1.2 Poddruh (SUB)

Kategorie **Poddruh** je relevantní pro všechny slovní druhy kromě citoslovcí, předložek, cizích slov, prefixových segmentů a neznámých slov.

Hodnoty kategorie **Poddruh** jsou závislé na slovním druhu. Některé slovní druhy mají dva poddruhy (zájmena, číslovky, příslovce). V tomto případě nazveme druhou kategorií poddruhu jinak (**Funkce**, viz 4.1.3), kvůli možné kombinovatelnosti.

Podívejme se tedy na jednotlivé slovní druhy. U každého slovního druhu uveďme, jaké poddruhy se u něj rozlišují, a přiřadíme jim kromě názvu i jednopísmenný kód. Tento kód je pokud možno v souladu s kódem 2. pozice v pražském

systemu morfologických značek, která se shodou okolností také nazývá Poddruh. My však tuto kategorii chápeme zcela odlišně. Poddruh v dosavadním pražském systému je směs hodnot popisující nejrůznější slovní tvary. Některé hodnoty se týkají jednotlivých slovních tvarů, jiné celých lemmat, není to tedy ani flektivní, ani globální kategorie. Náš Poddruh je kategorie striktně globální, je tedy relevantní vždy pro celé paradigma.

Oba poddruhy, náš globální, i dosavadní pražský, mají jednu stejnou podstatnou vlastnost: jejich kód je v rámci celé kategorie jedinečný. Znamená to, že z hodnoty kategorie Poddruh je možno určit slovní druh. Jak bylo řečeno výše, nezáleží na konkrétním kódu hodnoty, ale je důležité, aby tyto kódy byly různé.

Tento požadavek, zavedený Hajičem (viz (Hajič, 2004)), není na první pohled příliš důležitý. Kategorie poddruh je napříč slovními druhy velmi různorodá, u každého slovního druhu popisuje jiné vlastnosti, a tak by se mohlo zdát, že nezáleží na jednoznačnosti jejího kódu. Podstatně to však zjednoduší vyhledávání v korpusech podle této kategorie. Usnadní to posléze i kódování kategorie Poddruh v rámci složenin. Jednoznačnost kódu má tři výjimky:

1. Každý slovní druh, pro který je kategorie Poddruh relevantní, musí být hodnotami této kategorie zcela pokryt. U některých slovních druhů proto zavádíme hodnotu „Ostatní“, která je vždy kódována hodnotou 0 (nula), je tedy stejná pro více slovních druhů. Konkrétně jsou to podstatná jména, přídavná jména, slovesa, příslovce a částice. Tato hodnota vlastně vyjadřuje jen tu informaci, že dané slovo nemá žádnou z vlastností sledovaných v kategorii Poddruh. Proto není nutné tuto hodnotu dělit podle slovních druhů.
2. Poddruh „přivlastňovací“ sdílí zájmena a přídavná jména. Jde totiž o velmi podobnou vlastnost u obou slovních druhů.
3. Poddruh „deverbativní“ sdílí podstatná jména, přídavná jména a příslovce. I zde jde o vyjádření podobné vlastnosti, totiž blízké příbuznosti se slovesem.

4.1.2.1 Poddruh podstatných jmen

- S: deverbativní typu *věznění, pokrytí,...*
- 0: ostatní

Hodnota S náleží podstatným jménům odvozeným ze sloves, vyjadřující slovesný děj a končící na *-ní* nebo *-tí*. Tato podstatná jména se mohou chovat ve větě jinak než ostatní, primární podstatná jména. Jedině ona totiž mohou mít zvrtnou částici a adverbialní rozvití adverbii pravidelně odvozenými od přídavných jmen (*zpívání falešně*). Viz též příklad (25).

Mužské bádání, v ženském jazyce nazývané též štítění se práce, dokáže (25) jít rozumnému člověku (míním tím ženu) na nervy.

Kopečný je v (Kopečný, 1962a) považuje za slovesa pro jejich „naprostou totožnost významovou (obsahovou) a pro paradigmaticnost, s níž se od slovesných základů derivují.“ Syntakticky i morfologicky ovšem patří mezi podstatná jména, proto je mezi nimi ponecháme. K odlišení stačí vyjádření poddruhu.

Je třeba ovšem rozlišovat lexikalizovaná deverbativní podstatná jména, která nevyjadřují slovesný děj. Kopečný uvádí příklad *vázání ječmene — lyžařské vázání*. V tomto případě patří do slovníku dva záznamy, a to *vázání-1* s poddruhem 0 a *vázání-2* s poddruhem S. Dalším příkladem může být *krmení zvířat — krmení pro zvířata*. Najdou se i případy, kdy tato distinkce bude sporná.

4.1.2.2 Poddruh přídavných jmen

- U: přivlastňovací (*matčín, otcův,...*)
- G: od přechodníku přít. (*mající, sedící, beroucí,...*)
- M: od přechodníku min. (*ušedší, nakupovavší,...*)
- S: ostatní deverbativní (*namazaný, zemřelý, nakousnutý, namazán, nakousnut,...*)
- 0: ostatní (jarní, starý,...

Přídavná jména utvořená od přechodníků nejsou opět podle Kopečného „normální adjektiva“, neboť neplní všechny své funkce. Konkrétně nemohou být použita ve funkci jmenného přísudku (**žák je sedící*). Tato námitka neplatí obecně, neboť některá tato přídavná jména se už lexikalizovala. Příkladem může být lemma *vynikající* v příkladě (26).

Ten čaj je vynikající. (26)

Vzhledem k jejich adjektivní flexi je považujeme za přídavná jména.

Rozhodnutí, která přídavná jména jsou deverbativní (SUB=S) a která ostatní (SUB=0), je někdy těžké. Nejednoznačné případy (např. *šílený, vleklý*) bude třeba řešit jednotlivě.

4.1.2.3 Poddruh zájmen

U zájmen je třeba rozlišovat na této pozici dvě kategorie; jedné necháváme název Poddruh, druhou nazýváme Funkce. Vzhledem k tomu, že Funkce je kategorie společná pro více slovních druhů, totiž pro zájmena, číslovky a příslovce, pojednáme o ní jako o zvláštní morfologické kategorii dále v oddíle 4.1.3.

Toto dělení jsme zavedli proto, že běžné dělení zájmen nebere v úvahu dvojí podstatu užívaných hodnot. Např. zájmeno *něčí* je současně přivlastňovací i neurčité. Dvojí dělení zájmen bylo prvně použito v brněnském systému a přejalo ho i Konkláve.

Hodnoty obou kategorií, Poddruh i Funkce, se samozřejmě mohou kombinovat (proto byly zavedeny), ovšem ne zcela libovolně. Možné kombinace jsou uvedeny v tabulce 4.1 na str. 32, která by měla obsahovat úplný výčet všech českých zájmen.

Hodnoty kategorie **Poddruh** pro zájmena jsou:

- Z: substantivní (*já, kdo, nikdo, oni,...*)
- U: přivlastňovací (*můj, čí,...*)
- D: ukazovací (*ten, takový,...*)
- V: vymežovací (*každý, všechen, týž, sám*)
- 0: ostatní

Až na zájmena substantivní jde o klasické hodnoty.

Zájmena substantivní nahrazují osobní zájmena, ale navíc ještě některá zájmena neurčitá (*někdo, kdosi,...*), tázací (*kdo, co*) a záporná (*nikdo, nic*), která mají jednak substantivní skloňování a jednak podobné syntaktické postavení ve větě jako zájmena osobní (určitá).

Substantivní zvrtná zájmena *sebou, sobě, sebe, se, si* nemají první pád, jejich lemma tedy položíme rovno slovnímu tvaru, jmenovitě *sebou, sobě, sebe, se, si*. Viz též oddíl 2.2.3, kde jsme se zamýšleli nad lemmatizací slovních tvarů, které nemají nominativ, resp. infinitiv.

Slovní tvary *se, si* jsou homonymní. Jde buď o zájmena, nebo o částice, podle kontextu. Tvar *si* je navíc obecně českou variantou 2. osoby jednotného čísla přítomného času lemmatu *být*.

Slovní tvar *toť*, který se tradičně řadí mezi ukazovací zájmena, je homonymní se složeninou. Z korpusů jsme totiž zjistili, že jako klasické zájmeno se vyskytuje jen ve velmi omezeném typu vět. Většinou spíše zastupuje frázi *to je*, případně *to jsou*, a je tudíž složeninou. Viz příklady (27) až (29):

Toť se ví, že... (zájmeno ukazovací) (27)

Toť vše. (složenina v jednotném čísle) (28)

... *toť pomyslné kóty* (složenina v množném čísle) (29)

4.1.2.4 Poddruh číslovek

Slovní druh Číslovky je velmi různorodý. Viz např. nejnovější pojednání o číslovkách a kvantifikátorech (Šimandl, 2007) a (Jiranová, 2008).

Hodnoty kategorie **Poddruh**, které zde uvádíme, jsou ty, na kterých se shodlo Konkláve. Kromě nich lze porůznu v literatuře najít ještě další druhy číslovek. Jde např. o číslovky velikostní (*tisícový, osmitisícový*), které ovšem považujeme za přídavná jména.

Na druhou stranu číslovky úhrnné, souborové a druhové mají nejasně rozlišené skloňování (viz též tabulku 13.3 a zmínku v kapitole 13 na straně 130), takže by možná bylo rozumné je sdružit do poddruhu jediného. Tato otázka by zasloužila ještě podrobnější zkoumání.

Slova, která se občas označují jako dílové číslovky (např. *polovina, osmina, setina*) chápeme jako podstatná jména. Výjimku tvoří pouze tři dílové číslovky uvedené v tabulce 4.1. Stejně tak tradiční číslovky skupinové (např. *dvojice, pětka, tisícovka*) považujeme za podstatná jména. Je to proto, že se ve větě vždy syntakticky projevují jako klasická podstatná jména.

Slova typu *kopa, tucet, hromada, spousta, kus, kousek* jsou podstatná jména. Jejich kvantitativní vlastnost se v případě potřeby musí zachytit jinak, podle aplikace.

Uvědomujeme si, že s takovým rozdělením číslovek leckdo nemusí souhlasit. Číslovky jsou bohužel velmi těžko vymezitelná kategorie. Pro potřeby automatického zpracování však je potřeba vymezit číslovky jednoznačně a co nejjednodušeji.

Hodnoty kategorie **Poddruh** pro číslovky tedy jsou:

- 1: základní (*jedna, sto,...*)
- r: řadové (*druhý, pátý,...*)
- u: úhrnné (*dvé, patero,...*)
- s: souborové (*dvoje, paterý,...*)
- d: druhové (*dvojí, paterý,...*)
- n: násobné (*dvakrát, pětkrát,...*)
- o: opakovací (*podruhé, popáté,...*)
- v: výčtové (*zadruhé, zapáté,...*, ale i *druhé* z dvojice *za druhé,...*)
- p: dílové (*půl, čtvrt, třet*)

4.1.2.5 Poddruh sloves

- m: modální (*moci/moct, mít, mívat, muset, musívat, smět, chtít, hodlat, dát se, dávat se, dovést, umět*)
- f: fázová (*začít, začínat, přestat, přestávat, zahájit, skončit,...*)
- b: pomocná (*být, bývat, mít, dostat*)
- 0: ostatní (*navštívit, koupat se,...*)

Podle závěrů Konkláve se u pomocných sloves *být* a *bývat* jako pomocné značí jen tvoření minulého času, podmiňovacího způsobu a budoucího času. V tom případě jsou pro něj relevantní kategorie Osoba, Číslo i Slovesný tvar (indikativ *přišel* **jsí**, příčestí činné **byl** *bys* **přišel** nebo budoucí čas **budeš** *skákat*). Totéž pro *bývat*, kde však může být pouze příčestí činné (*byl* *bys* **býval** *přišel*). Slovní tvary lemmatu *bývat* v indikativu přítomnosti nejsou nikdy pomocným slovesem a jednoduchý budoucí čas *bývat* nemá.

4.1.2.6 Poddruh příslovcí

- P: místní (*kudy, tudy, odkud, nikudy, nikam; daleko, nedaleko,...*)
- T: časová (*kdy, nikdy; včera, odpoledne,...*)
- J: způsobová (*jak, všelijak; krásně, velmi, široce,...*)
- R: predikativní (*jasno, možno, teplo, volno,...*)

První skupina příkladů (před středníkem) místních, časových a způsobových příslovce se často řadí mezi tzv. zájmenná příslovce. U nich je možno určovat i kategorii **Funkce** s hodnotami určitá, neurčitá, záporná jako u číslovek a zájmen (viz dále, zejména tabulku 4.3). Pomocí kategorie **Funkce** lze též zájmenná příslovce vymezit. Jsou to ta příslovce, pro něž je kategorie **Funkce** relevantní, tedy POS=D a FCE≠UNDEF.

4.1.2.7 Poddruh spojek

- ^ (stříška): souřadící (*a, ale, nebo,...*)
- , (čárka): podřadící (*protože, když, až, -li,...*)
- * (hvězdička): matematické operace (*plus, minus/mínus, krát, děleno* — neplést s *děleno* jako jmenný tvar přídavného jména *dělený*, případně rod trpný od slovesa *dělit*)

Slovní tvary *abys, abyste, abychom, kdybych, kdybyste, kdybychom*, ale i nepisovné *abysem, abysme, kdybysem a kdybysme* jsou také spojky. Podrobněji o nich pojednáme v kapitole 5 o kondicionálu.

4.1.2.8 Poddruh částic

- 7: zvrtné (*se, si*)
- c: kondicionálová (pouze *by*)
- 0: ostatní (*ba, ano, bože, nechť, ať,...*)

Bylo by třeba ještě zvážit, zda mezi zvrtné částice nepočítat i tvar *sebou* ve spojení např. *hodit sebou*.

Částice by se mohly členit ještě více, např. podle (Akademická mluvnice). Částicemi se intenzivně zabývá i práce (Čermák, 2007). Inspirací může být také klasifikace slovenských částic (viz např. (Slovenská morfolgie) a (Šimková, 2001)).

Oba tyto problémy necháváme zatím otevřené.

Ostatní slovní druhy, tedy předložky³, citoslovce, cizí slovo⁴, prefixový segment a neznámé slovo, poddruh nemají.

4.1.2.9 Poddruh složenin

Složeniny „dědí“ poddruh některé své složky. Složeninami včetně jejich poddruhů, se budeme zabývat podrobně v kapitole 6.

³Na Konkláve jsme místo Poddruhu zavedli pro předložky kategorii Vokalizace s hodnotami +/-, nový návrh ale toto rozlišení řeší pomocí kategorie Mutace (viz kap. 3).

⁴Výhledově by se mohlo uvažovat o přiřazení jazyka, ze kterého cizí slovo pochází, případně o jeho slovním druhu v daném jazyce. Pro zpracování českého textu ale zřejmě tyto údaje nemají význam.

4.1.3 Funkce (FCE)

Tato kategorie je relevantní pro zájmena, číslovky a příslovce, i když ne všechny vyjmenované slovní druhy mohou nabývat všech jejích hodnot.

Uvědomujeme si, že termín „funkce“ je již obsazen mnoha dalšími významy, přesto si myslíme, že tento název nejlépe vystihuje podstatu této kategorie. Na Konkláve byla tato kategorie nazvána „Neurčitost“ s tím, že jde pouze o pracovní název.

Hodnoty tázací a vztažná jsou v praxi velmi obtížně rozlišitelné (viz např. (Ševčíková, 2008)). Jedním z důvodů, proč je rozlišovat na morfologické rovině, spatřujeme v existenci slovních tvarů *jenž*, *kdož*, *jakž*⁵ apod., které mohou být pouze vztažné, nikoli tázací. Podrobná diskuse o vhodnosti či nevhodnosti rozlišování těchto dvou hodnot je však mimo rámec této práce, proto obě hodnoty zachováváme, přičemž by nebyl problém je v konkrétní implementaci sloučit do hodnoty jediné.

- U: určitá (všechna osobní zájmena, určité číslovky, tady, teď,...)
- N: neurčitá (*někdo*, *čísi*, *několik*, *někdy*,...)
- Z: záporná (*nikdo*, *ničí*, *nijak*,...)
- T: tázací (*kdo*, *čí*, *kolik*, *kde*,...)
- V: vztažná (*kdo*, *čí*, *jenž*, *kdy*,...)
- S: zvrtná (*se*, *si*, *sobě*, *sebe*, *sebou*)

Neurčitá funkce (FCE=N) se projevuje většinou pomocí speciálních předpon a přípon, které se připojují ke slovním tvarům s funkcí tázací (FCE=T). Kombinovatelnost ukazuje obrázek 4.1. Předpony a přípony z levého rámečku lze kombinovat se slovy v pravém rámečku, čímž vznikne slovo s neurčitou funkcí. Některé kombinace možné nejsou, např. nelze **něodkud*, dále se netvoří neurčité číslovky pomocí přípon, pouze pomocí předpon. Některé kombinace jsou velmi neobvyklé, ale netroufneme si tvrdit, že zcela nemožné, např. ?*kamžkolivěk*.

<i>lec-</i> , <i>leda-</i> , <i>všeli-</i> , <i>ně-</i> , <i>bůhví-</i> , <i>čertví-</i> , <i>pánbůhví-</i> , <i>pámbuví-</i> , <i>kdoví-</i> , <i>nevím-</i> , <i>kde-</i> , <i>lec-</i> , <i>leda-</i> , <i>ledas-</i> , <i>málo-</i> , <i>všeli-</i> , <i>všelis-</i> , <i>zřídka-</i> , <i>sotva-</i> , <i>-koli</i> , <i>-koliv</i> , <i>-žkolí</i> , <i>-žkoliv</i> , <i>-kolivěk</i> , <i>-si</i>	<i>kdo</i> , <i>co</i> , <i>čí</i> , <i>jaký</i> , <i> který</i> , <i>kde</i> , <i>kam</i> , <i>kudy</i> , <i>odkud</i> , <i>kdy</i> , <i>jak</i> , všechny tázací číslovky
---	--

Obrázek 4.1: Kombinovatelnost předpon s tázacími zájmeny, číslovkami a příslovci při tvoření slovních tvarů s funkcí neurčitou (FCE=N)

Uvádíme tři tabulky, ze kterých je vidět ortogonálnost kategorie **Poddruh** a **Funkce** pro zájmena (tabulka 4.1), číslovky (tabulka 4.2) a příslovce (tabulka 4.3). Seznam z pravého rámečku schématu 4.1 lze v tabulkách 4.2 a 4.3 doplnit

⁵Jde o příklady typu ... *uhlazujete povrch štětkou*, *jakž je zvykem lepičů plakátů*., nikoli o výrazy *jakž takž*.

na místa neurčitých číslovek resp. příslovčí. V tabulce 4.1, která popisuje funkce zájmen, jsou možnosti vypsány explicitně.

Tabulku pro zájmena vytvořila autorka této práce, modifikoval Karel Oliva, doplnila opět autorka. Ostatní tabulky jsou autorčiny.

Tabulka 4.1 by měla být vyčerpávající, to znamená, že jsme se snažili, aby obsahovala veškerá zájmena. Jako zdroj posloužily dostupné slovníky a korpusy SYN2000 a SYN2005.

Prázdné závorky () je možno nahradit slovy ze stejné řádky, která jsou ve sloupci tázacích zájmen, konkrétně v řádku substantivních zájmen jsou to *kdo*, *co*, v řádku přivlastňovacích zájmen *čí* a v řádku ostatních zájmen *jaký* a *který*.

Hodnoty kategorie **Funkce** pro číslovky jsou stejné jako pro zájmena, ale ne všechny se pro určování číslovek uplatní. Zejména hodnota **Funkce** zvrtná je pro číslovky nepoužitelná.

Tabulka 4.2 je analogická tabulce 4.1 pro zájmena, liší se však v jednom podstatném bodě, není totiž úplná. Přesněji, sloupec s funkcí určitá není úplný až na číslovky dílové a potom také políčko s číslovkami základními neurčitými. Ostatní neurčité číslovky získáme ze schématu 4.1 na str. 30, takže tento sloupec, i sloupec s číslovkami tázacími za úplný považovat můžeme. Neúplnost je naznačena třemi tečkami v příslušných políčkách.

Kategorii **Funkce** využíváme i pro specifikaci zájmenných příslovčí. Tabulka 4.3 ukazuje kombinovatelnost zájmenných příslovčí s různými hodnotami kategorie **Funkce**.

Hodnoty určitá, neurčitá a tázací je možné využít též pro bližší určení speciálních přídavných jmen s číselnou předponou (např. *pětičlenný*, *několikaletý*, *kolikawattový*). Dokonce můžeme určit tyto hodnoty kategorie **Funkce** i u podstatného jména *-násobek* s číselnými předponami (*pětinásobek*, *několikanásobek*, *kolikanásobek*). Nejzajímavější je u těchto slovních tvarů hodnota tázací, neboť může uvozovat vedlejší větu, což u „normálních“ podstatných a přídavných jmen není možné, viz příklady (30) a (31)⁶.

Dopita ještě nemá jasno, kolikaletý kontrakt podepíše. (30)

... koeficienty, jež určují, kolikanásobek této základny budou ústavní činitelé dostávat. (31)

⁶Za upozornění na tyto typy přídavných a podstatných jmen děkuji Karlu Olivovi.

4 Morfologické kategorie

	určitá	neurčitá	záporná	tázací	vztaž- ná	refle- xivní
substantivní	<i>já, ty, on, ona, ono, my, vy, oni, ony.</i>	<i>leccos, ledacos, všelicos, ně(), bůhví(), pánbůhví(), pámbuví(), čertví(), kdoví(), nevím(), kde(), lec(), leda(), ledas(), málo(), všeli(), všelis(), zřídka(), sotva(), ()koli, ()koliv, ()žkoli, ()žkoliv, ()kolivěk, ()s, ()si.</i>	<i>nikdo, nic, pranic.</i>	<i>kdo, co, ()pak, ()že, ()ž.</i>	<i>jaký, který, jenž, kdo, co, an.</i>	<i>si, se, sebou, sobě, sebe.</i>
přívlastňovací	<i>můj, tvůj, jeho, její, náš, váš, jejich.</i>	<i>něčí, čisi, čikoli, čikoliv, čikolivěk, bůhvíčí, pánbůhvíčí, pámbuvíčí, čertvíčí, ledačí, ledasčí, leccí, kdovíčí, nevímčí, kdečí.</i>	<i>ničí.</i>	<i>čí, čípak, číze.</i>	<i>jehož, jejíž, jejichž, čí.</i>	<i>svůj.</i>
ukazovací	<i>ten, tento, takový, tehle, onen, onaký, týž, tentýž, takovýto, tenhleten, tamten, tuten, taký, tamhleten, tuhleten, tadyten, toť.</i>	--	--	--	--	--
vymezovací	<i>každý, sám, všechen, týž, tentýž, tatáž, totéž, titíž, tytéž, samý, všecek, všecken, všeliký, veškery.</i>	--	--	--	--	--
ostatní	--	<i>ně(), bůhví(), čertví(), kdoví(), nevím(), kde(), lec(), leda(), ledas(), málo(), všelijaký, zřídka(), ()koli, ()koliv, ()žkoli, ()kolivěk, ()si, ()s.</i>	<i>nijaký, žádný, nižádný, pražádný.</i>	<i>jaký, který, ký, ()pak, ()že, ()ž.</i>	--	--

Tabulka 4.1: Kombinovatelnost hodnot kategorií Podruh a Funkce u zájmen

4 Morfologické kategorie

	určitá	neurčitá	tázací
základní	<i>jedna, raz, dva, dvě, pět, sto, ...</i>	<i>několik, hodně, málo, poskrovnu, víc, dost ...</i>	<i>kolik</i>
úhrnné	<i>dvě, patero, tisícero, obé, ...</i>	<i>několikero</i>	<i>kolikero</i>
souborové	<i>dvoje, patery, tisícery, oboje, ...</i>	<i>několikery</i>	<i>kolikery</i>
druhové	<i>dvoji, paterý, tisícery, obojí, ...</i>	<i>několikerý</i>	<i>kolikerý</i>
násobné	<i>dvakrát, pětkrát, (po)obakrát, ...</i>	<i>několikrát</i>	<i>kolikrát</i>
řadové	<i>druhý, pátý, ...</i>	<i>několikátý</i>	<i>kolikátý</i>
opakovací	<i>podruhé, popáté, ...</i>	<i>poněkolikáté</i>	<i>pokolikáté</i>
výčtové	<i>zaprvé, zadruhé, zapáté, ...</i>	<i>zaněkolikáté</i>	<i>zakolikáté</i>
dílové	<i>půl, čtvrt, třet'</i>	--	--

Tabulka 4.2: Kombinovatelnost hodnot kategorií Poddruh a Funkce u číslovek

	určitá	neurčitá	tázací a vztažná	záporná
místní	<i>tady, zde, tam, tudy, tamtudy, ...</i>	<i>někam, někudy, odněkud, ...</i>	<i>kam, kudy, odkud, ...</i>	<i>nikde, nikam, odnikud, ...</i>
časová	<i>ted', nyní, ...</i>	<i>někdy</i>	<i>kdy</i>	<i>nikdy</i>
způsobová	<i>tak, takto, ...</i>	<i>nějak</i>	<i>jak</i>	<i>nijak</i>

Tabulka 4.3: Kombinovatelnost hodnot kategorií Poddruh a Funkce u zájmených příslovcí

4.1.4 Vid (ASP)

- D: dokonavý (*koupit, napsat, doručit, narodit se,...*)
- N: nedokonavý (*kupovat, psát, doručovat, chodívat,...*)
- O: obouvidý (*referovat, absolvovat, izolovat,...*)

Stejně jako dosavadní morfologické systémy, nezavádíme další hodnoty kategorie Vid pro iterativní slovesa, i když je, pokud je to možné, pravidelně vytváříme pomocí slovesných derivačních vzorů (viz 12.2.4). V případě potřeby by nebyl problém hodnoty doplnit.

V valenčním slovníku VALLEX (viz (Lopatková et al., 2006)) jsou vidové dvojice zpracovány jako jedno slovníkové heslo. V morfologickém slovníku to neděláme. Členy vidové dvojice považujeme za dvě různá slova. Bylo by však záhodno jejich slovníkové záznamy propojit zvláštním typem odkazu (viz také kapitulu 7).

4.1.5 Zkratka (ABR)

Tato kategorie je relevantní pro všechny slovní druhy.

Kategorie zkratka má pouze dvě hodnoty, a to:

- + : ano
- UNDEF

Kladnou hodnotu globální kategorie Zkratka dostávají zkratky, ostatní slovní tvary nemají tuto hodnotu definovanou. V brněnském systému je zkratka samostatným slovním druhem se značkou **KA**, v pražském systému se zkratky značkují pomocí 2. nebo 15. pozice.

Zkratka ve větě často zastupuje konkrétní slovní druh, a její prohlášení za samostatný slovní druh tak může komplikovat další zpracování. Pražské řešení nám proto připadá lepší, není však uplatňováno konzistentně, což je obecně slabá stránka obou zmiňovaných pozic v pražské morfologické značce.

Zkratka jako samostatná globální kategorie podle našeho názoru nejlépe postihne všechny možnosti, jak může zkratka vypadat.

Zkratka tedy může být libovolný slovní druh. V případě, že zkratka zastupuje jeden konkrétní slovní tvar, je slovním druhem této zkratky slovní druh tvaru, které zkratka zastupuje. Např. *č.* jako zkratka slovního tvaru *číslo* je podstatné jméno a jako takové jsou pro ni relevantní všechny kategorie relevantní pro ostatní podstatná jména. Zkratky, které nezastupují jedno slovo, ale celé slovní spojení, je třeba hodnotit individuálně. Zastupuje-li např. zkratka jmenovou frázi (*USA, ODS*, apod.), je POS=N. U takových zkratk je možno stanovit i rod, číslo a pád podle toho, co zkratka zkracuje, nebo případně jak se používá. Je možno též využít hodnoty X (sdružená hodnota). Zastupuje-li zkratka složitější frázi (např. *atd., např.*), je POS=S (více viz zkratkové složeniny v kapitole 6).

4.1.6 Globální mutace (GMU)

Tato kategorie je relevantní pro všechny slovní druhy.

Podle naší definice se globální mutace projevují ve všech slovních tvarech paradigmatu. Z toho mimo jiné vyplývá, že globální jsou všechny mutace neohebných lemmat. Jde především o mutace příslovcí, která se nestupňují, tedy např. *zítra* — *zejtra*.

Mezi globální zahrnujeme i vokalizované mutace předložek (*od* — *ode*, *k* — *ke* — *ku*). Jsme si vědomi toho, že toto je zcela jiný typ mutací než např. mutace ortografické. Opět je třeba připomenout, že nám jde o co nejjednodušší popis, takže kategorie mutace, zde konkrétně mutace globální (GMU) využíváme k rozlišení slovních tvarů, jejichž ostatní relevantní kategorie nabývají hodnot totožných.

Zajímavější jsou slova ohebná. Mnoho globálních mutací je specifických, týkajících se jednoho konkrétního lemmatu. I zde však existuje několik systematických typů variantních dvojic, které se mohou uplatnit u mnoha, někdy dokonce všech lemmat určitých vlastností. Asi nejznámější jsou mutace cizích slov přejatých do češtiny, kde se v původním jazyce (většinou latina) píše *s*, ale v češtině vyslovuje *z*. Podle doporučení posledního vydání Pravidel českého pravopisu (viz (Pravidla)) existuje několik pravidel a mnoho výjimek, jak taková slova správně psát, ale uživatelé jazyka si s tím často hlavu nelámou a píšou různě. Tvary se *z* se většinou považují za spisovné, mutace se *s* za knižní nebo zastaralé, ale v textech se setkáme s oběma. Všechny takové dvojice musíme považovat za mutace. Nejznámější mutace se týkají přípon *-ismus* — *-izmus*. Zde však připouštějí Pravidla možnosti obě a používá se více mutace *-ismus*⁷, a to i tehdy, když se vyskytuje ve slově s více možnými *s/z*, např. *izomorfismus* (15) — *isomorfismus* (2). Čísla v závorce udávají frekvenci v korpusu SYN2005.

Další systematické mutace vznikají již zmíněným přidáním protetického *v* před lemmata začínající na *o*-.

Důležitou vlastností globálních mutací je to, že obě (všechny) mohou být použity při vytváření odvozenin. Např. z mutací podstatných jmen *okno* — *vokno* lze utvořit mutace přídavného jména *okenní* — *vokenní*. Tyto mutace jsou opět globální.

Z toho, co bylo řečeno o globálních mutacích, je zřejmé, že základní tvary globálních mutací jsou vždy variantními lemmaty vícenásobného lemmatu (viz též kap. 2). Neplatí to však naopak — existují i taková vícenásobná lemmata, jejichž prvky jsou flektivní mutace (např. vícenásobné lemma {*myslit*, *myslet*}), nebo nejsou mutacemi vůbec, což je případ všech složenin (viz kap. 6).

Tabulka 4.4 ukazuje hlavní typy českých globálních mutací, bez ohledu na jejich klasifikaci, to znamená, že nedělá rozdíl mezi kodifikovanými a nekodifikovanými mutacemi. Poslední sloupec tabulky uvádí kódy pro hodnoty kategorie **Globální mutace**. Kódy v horní části tabulky vyjadřují hláskovou změnu v mutacích. Kód *d* zastupuje „dlouhé“ mutace, *k* „krátké“. Podobně *m* znamená „měkké“, *t* „tvrdé“. Mutace, které se vymykají běžným typům, se označují čísly, jak ukazuje poslední řádek tabulky.

⁷Nepoměr mezi oběma mutacemi nás velmi překvapil: Nejrozsáhlejší obecně dostupný český korpus SYN obsahuje pouze necelých 2100 výskytů slov s lemmatem zakončeným sufixem *-ismus* oproti téměř 133 tisícům s lemmatem na *-ismus*.

4 Morfologické kategorie

Typ	Příklad	Hodnoty GMU
o — vo	<i>okno — vokno</i>	0 — v
ý — ej	<i>mýdlo — mejdlo</i>	0 — j
z — s	<i>klauzule — klausule</i>	z — s
t — th	<i>tema — thema</i>	0 — h
é — í	<i>kolébka — kolíbka</i>	e — i
é — ý	<i>okénko — okýnko</i>	e — y
á — e	<i>originální — originelní</i>	a — e
á — a	<i>Abrahám — Abraham</i>	d — k
é — e	<i>acetylén — acetylen</i>	
ó — o	<i>salón — salon</i>	
ý — y	<i>apetýt — apetyt</i>	
í — i	<i>alexandrín — alexandrin</i>	
ů — u	<i>přezůvky — přezuvky</i>	
ú — u	<i>Plútarchos — Plutarchos</i>	
s — š	<i>student — študent</i>	t — m
t — ř	<i>vlaštovka — vlašřovka</i>	
n — ň	<i>šňůra — šňůra</i>	
d — ř	<i>dolík — řolík</i>	
e — ě	<i>Bardejov — Bardějov</i>	
z — ž	<i>zbrzdování — zbržďování</i>	
jiné	<i>Afganistan — Afghanistan</i>	0 — 1

Tabulka 4.4: Přehled nejčastějších typů globálních mutací s příklady

Hromadění typů mutací v jednom lemmatu se vyjádří vícero hodnotami, viz tabulka 4.5. Jsou zde naznačeny možné kombinace tří typů globálních mutací: s — z a d — k, přičemž dlouhá mutace se v tomto případě rozpadá na další dvě možnosti, a to ú — ů. Tento poslední typ není zcela typický, proto jsme ho nezahrnuli do tabulky 4.4. Mutaci s -ů- jsme tedy nechali jako „dlouhou“ s hodnotou d, mutaci s -ú- jsme označili číslicí 1.

Lemma	Hodnota GMU
<i>bluza</i>	kz
<i>blůza</i>	dz
<i>blůza</i>	1z
<i>blusa</i>	ks
<i>blůsa</i>	ds

Tabulka 4.5: Příklad vícera hodnot kategorie Globální mutace

Pro lemmata, která se vyskytují v mnoha mutacích (většinou jde o cizí vlastní jména) je nejvýhodnější označit mutace pomocí čísel, přestože by někdy bylo možné i v těchto jménech vystopovat uvedené typy. Může totiž dojít k tomu, že se v jednom lemmatu uplatní jeden typ vícekrát. V tom případě by označování mutací mohlo být krkolomné.

Příkladem takových mutací je množina různých zápisů země *Afghánistán*, kde se projevuje typ „jiný“ a dva typy „a—á“. V korpusu SYN se vyskytuje ve všech osmi možných mutacích: *Afghánistán*, *Afgánistán*, *Afganistán*, *Afgha-*

nistán, Afghanistan, Afganistan, Afghánistan, Afgánistan.

Kdybychom chtěli i v takových případech rozlišovat typy globálních mutací, bylo by třeba jejich hodnoty udávat i s místem v konkrétním lemmatu, kde k rozlišení typu dochází. Z uvedeného příkladu je ale zřejmé, že takto podrobný popis globálních mutací by byl pravděpodobně zbytečný.⁸

4.2 Flektivní morfologické kategorie

Flektivní morfologická kategorie je taková kategorie, jejíž hodnoty se pro jednotlivé slovní tvary jednoho paradigmatu liší.

Následuje seznam flektivních morfologických kategorií:

- Rod — GEN
- Číslo — NUM
- Duál — DUA
- Pád — CAS
- Osoba — PER
- Stupeň — DEG
- Negace — NEG
- Slovesný tvar — VRB
- Jmenný tvar přídavných jmen — NOM
- Stupeň intenzity slovesného děje — INT
- Typ složeniny — CMP
- Flektivní mutace — FMU

Probereme je nyní jednu po druhé. Některé jsou tradiční a mají i tradiční hodnoty, jiné jsme zavedli nově, někde jsme dokonce pozměnili tradiční hodnoty. Vše je podřízeno snadnější využitelnosti při automatickém zpracování češtiny při zachování veškeré lingvistické informace, kterou hodnoty kategorií nesou.

Ve výčtu hodnot následujících kategorií bude na posledním místě občas vystupovat hodnota nazvaná „sružená hodnota“. Tím se myslí libovolná z hodnot předcházejícího seznamu. Podobná hodnota již v pražském systému existuje. Tam se však bere jako proměnná, která případně může být nahrazena jednou z konkrétních hodnot (např. rod cizího slova podle přívlastkového rozvití). Při některých experimentech se s ní takto dokonce pracovalo — morfologická značka obsahující tuto hodnotu se rozepsala, aby v ní byly jen konkrétní hodnoty, čímž se rozpadla do množství jednoznačných značek. My ale chápeme sruženou hodnotu odlišně. V našem pojetí znamená to, že hodnota dané

⁸Dokonce by se v podobných případech (tedy u nejednotného zápisu cizích vlastních jmen) dalo uvažovat o tom, že nebude splněno Zlaté pravidlo morfologie.

kategorie se nedá a nikdy nepůjde rozlišit. Případy, kdy může daná kategorie nabývat více hodnot, řešíme možností přiřazení více morfologických značek danému slovu. Tím se nahradí částečně sdružené hodnoty u kategorií rod, číslo a pád, které obsahuje pražský systém.

4.2.1 Rod (GEN)

Kategorie je relevantní pro podstatná jména, přídavná jména, některá zájmena, některé číslovky, přechodníky a slovesa v přičestí činném.

- M: mužský životný
- I: mužský neživotný
- F: ženský
- N: střední
- X: sdružená hodnota

Kromě tradičních hodnot zachováváme i již zavedenou korpusovou praxi z pražského i brněnského systému a rozlišujeme dva mužské rody. Pravděpodobně přirozenější by bylo mít jen jeden mužský rod a životnost vyjádřit pomocí další kategorie. Nechceme však zavádět nové kategorie zbytečně, když se zdá, že současné pojetí nečiní problémy.

4.2.2 Číslo (NUM)

Kategorie je relevantní pro podstatná jména, přídavná jména, některá zájmena, některé číslovky, slovesa a kondicionálové částice a spojky.

- S: jednotné
- P: množné
- X: sdružená hodnota

Tradiční hodnota duál, která byla zachována i na Konkláve, se v našem návrhu stává samostatnou kategorií.

4.2.3 Duál (DUA)

Kategorie je relevantní pro duálová podstatná jména (viz dále) a slova s adjektivním skloňováním.

- +
- UNDEF

V mluvnících českého jazyka se o duálu mluví jako o třetí hodnotě kategorie **Číslo**. Anž bychom chtěli jakkoli napadat tento fakt, navrhneme na duál jiný pohled, který usnadní automatické zpracování českého jazyka. Zdůrazňujeme, že nám jde o zjednodušení automatických analýz, nikoli o zpochybňování lingvistických tradic.

Motivace

Hodnota duál kategorie **Číslo** má nepříjemnou vlastnost, způsobuje totiž v některých případech neplatnost jedné ze základních syntaktických vlastností češtiny, a to shody v čísle.

Např. ve větách (32) a (33) z korpusu (P v závorce znamená plurál, D duál):

zájem však vzbudil svými(P) dvěma(D) knihami(P) (32)

hleděla na ni upřenýma(D), nevidoucíma(D) očima(D), které(P) ani nezamrkaly(P) (33)

Číslovka *dvěma* z prvního příkladu je v duálu (D), přestože zájmeno *svými* i substantivum *knihami* jsou v čísle množném (P). Vzhledem k tomu, že číslovka *dva/dvě* ve skutečnosti žádné množné číslo nemá, stalo se zvykem považovat všechny duálové tvary této číslovky za tvary množného čísla. Pokud přijmeme tuto tezi, shoda tu samozřejmě je.

Ve druhém případě bychom očekávali shodu v čísle mezi podstatným jménem *očima* a vztažným zájmenem *keré*, uvozujícím vedlejší větu, podobně, jako ve větě (34):

Podívala se na něho brýlemi(P), které(P) v tu chvíli byly(P) velmi jasně. (34)

Shody bychom mohli dosáhnout dvěma způsoby:

1. buď připustíme, že slovní tvar *keré* může být kromě množného čísla i v duálu,
2. nebo prohlásíme všechny duálové tvary za množné číslo.

První varianta by ovšem znamenala, že bychom museli pro všechna vztažná zájmena ve všech pádech připustit dvojí hodnotu — množné číslo a duál. Další důsledek by byl ještě revolučnější — kvůli shodě podmětu s přísudkem bychom museli totéž udělat se všemi slovesy v množném čísle. Tato varianta se tedy zdá nepřijatelná.

Druhé řešení nám připadá rozumnější, musíme ho ale rozšířit, protože úplná ztráta hodnoty duál by nám přinesla zase jiné problémy. Především by se ztratil rozdíl mezi tvary přídavných jmen, zájmen, číslovek a některých podstatných jmen s koncovkami *-mi* a *-ma* v 7pl. Zavádíme tedy novou kategorii Duál s jedinou hodnotou +.

Tuto hodnotu mají

- podstatná jména, která tvoří duální tvary v 7pl, a to ve všech pádech množného čísla. Jde o lemmata (v závorce je uveden tvar pro 7pl): *oko (očima)*, *ucho (ušima)*, *ruka (rukama)*, *noha (nohama)*, *očičko (očičkama)*, *očko (očkama)*, *ouško (ouškama)*, *ručička (ručičkama)*, *ručka (ručkama)*, *nožička (nožičkama)*, *nožka (nožkama)*. Nadále jim budeme říkat duálová slova. Mezi duálová slova nezařazujeme *koleno*, *rameno* ani *prso*, přestože také tvoří duálové tvary, a to v 2pl a 6pl. Ty však nezpůsobují problémy se shodou. Viz též dále.

- přídavná jména, zájmena a číslovky s adjektivním skloňováním včetně základních číslovek *dva, oba, tři a čtyři*, ale jen v 7. pádě množného čísla. Tyto slovní tvary mají v 7pl buď hodnotu kategorie DUA=UNDEF (např. pro koncovku *-mi (krásnými)*), nebo DUA=+ (např. pro koncovku *-ma (krásnými)*). Dvě možnosti duálové hodnoty, totiž + a UNDEF, mají i tvary vztažných zájmen *jimiž, jejichž* a *jehož* v množném čísle, přestože se slovní tvary duálu neliší od prostého množného čísla. Viz příklady dále. Kladnou hodnotu duálu může mít v 7. pádě množného čísla i číslovka *tři*, tvar *třema*, přestože ji žádná nám známá mluvnice jako možnost neuvádí, viz např. 2. díl Mluvnice češtiny ((Akademická mluvnice)) na str. 405). Podle této mluvnice by tedy bylo správné např. poněkud absurdní spojení *třemi nohama*.

4.2.3.1 Odbočka ke shodě

I když naše práce pojednává o morfologii, považujeme za vhodné se na tomto místě dotknout otázky shody, abychom podrobněji vysvětlili právě navržené řešení duálu. Shoda v čísle se nyní rozpadá na shodu v čísle a shodu v duálu.⁹

Shoda v čísle se naším návrhem výrazně zjednoduší. Např. ve větě

Oči (P) se dívaly (P) (35)

zůstává zachovaná shoda v čísle a nevádí, že *oči* mají hodnotu Duálu + (nadále budeme značit D+) a *dívaly* nemají hodnotu Duálu žádnou (v následujících několika příkladech budeme tento stav značit D-).

V našich příkladech dostáváme:

zájem však vzbudil svými (P D-) dvěma (P D+) knihami (P D-) (36)

hleděla na ni upřenýma (P D+), nevidoucíma (P D+) očima (P D+), které (P D-) ani nezamrkaly (P D-) (37)

Shoda v čísle je všude jednoduše zachována.

Shoda v duálu se vyžaduje pouze u výše vyjmenovaných slov, a to jen v 7pl.

tmavé (P D-) oči (P D+), ale tmavýma (P D+) očima (P D+) (38)

očima (P D+), kterýma (P D+) (39)

oči (P D+), kterýma (P D+) (40)

oči (P D+), kterými (P D-).. špatně, není shoda v duálu, přestože je to 7.pád. (41)

očima (P D+), které (P D-).. dobře, není 7.pád (42)

očima (P D+), jimiž (P D+) (43)

jejichž (P D+) očima (P D+) (44)

⁹Nezabýváme se shodou v rodě, která nečiní potíže, i když v našich příkladech vystupují *oči*, které jsou v jednotném čísle rodu středního a v množném rodu ženského. Podobnou vlastnost mají i *uši*.

V posledních dvou příkladech (43) a (44) mohou mít vztažná zájmena samozřejmě i hodnotu kategorie Duál nedefinovanou, ale v těchto konkrétních kontextech je jejich hodnota kladná.

V příkladě (36) je sice tvar *dvěma* v duálu, ale shoda se nevyžaduje, neboť podstatné jméno *knih*a nepatří mezi duálová slova.

Pro úplnost je třeba zde dodat, že duál se projevuje ještě v 2pl a 6pl následujících tří lemmat středního rodu: *rameno* (*ramenou*), *koleno* (*kolenou*), *prso* (*prsou*). Vzhledem k tomu, že se vyskytují i nespisovné tvary *ramenech*, *kolenech* a alternativní *prsech* (zde ale dochází k nerozlišitelné homonymii s podstatným jménem *prs*), stanovujeme i v těchto případech pro rozlišení kladnou hodnotu kategorie Duál, která však nemá pro shodu praktický význam. Shodu v duálu stačí požadovat skutečně jen u duálových slov v 7pl.

Příklady:

na svých (P D-) *kolenou* (P D+)
se svými (P D-) *koleny* (P D-)

4.2.4 Pád (CAS)

Kategorie je relevantní pro podstatná jména, přídavná jména, zájmena, číslovky a předložky. U předložek je význam kategorie Pád poněkud odlišný — vyjadřuje rekcí. V tomto případě to sice není morfologická kategorie, avšak její přiřazení předložkám je velmi užitečné.

- 1 až 7
- X: sdružená hodnota

Hodnoty pádu jsou tradiční, není třeba se jimi dále zabývat.

4.2.5 Osoba (PER)

Kategorie je relevantní pro zájmena, slovesa, kondicionálovou částici *by* a kondicionálové spojky *aby*, *kdyby*.

- 1 až 3
- v: vykání

Hodnota *v* není tradiční. Ukazuje se, že při analýze textů činí značný problém rozpor mezi množným číslem zájmena *vy* (často psáno jako *Vy*) a jednotným číslem minulého přičestí slovesa (*vy jste mluvil*, *vy byste mluvil*), pasiva (*vy jste chycen*), případně přídavného jména po sponě (*vy jste sám*, *vy jste hezká*). Rosen a Saloni v (Rosen – Saloni, 2006) navrhují zahrnout tato tzv. honorifika do slovesného paradigmatu, kde se z nejasných důvodů tradičně neuvádějí. Stejný názor je vyjádřen i v práci Panevové ((Panevová, 2008)).

Navrhujeme tedy zavést kromě tří klasických hodnot kategorie *Osoba* ještě hodnotu čtvrtou, tzv. vykání. Přiřazujeme jí kód *v*.

Kromě hodnoty PER=2 mají tedy určitá zájmena *vy* a *váš* také hodnotu *v* (PER=*v*). Jde zde o pravidelnou homonymii. Stejná homonymie nastává i v prezentu sloves (*vy rádi zpíváte* — *vy rád zpíváte*).

Panevová se ve své práci (Panevová, 2008) zabývá ještě otázkou, zda může být honorativ také v množném čísle, jestliže vykáme celé skupině osob, případně jen někomu ze skupiny. Nenachází však na ni jednoznačnou odpověď.

Zájmena *vy* a *váš* jako honorifika (PER=v), tedy mohou být jak v jednotném, tak v množném čísle. Jsou-li v množném čísle, homonymii nelze vyřešit jinak než ze sémantiky kontextu. V příkladě (45) mají tedy tvary *vy* a *jste* jednoznačné hodnoty kategorie PER, neboť NUM=S, zatímco v příkladě (46) jsou možné hodnoty kategorie PER dvě.

vy (PER=v NUM=S) *jste* (PER=v NUM=S) *mluvil* (NUM=S) (45)

vy (PER=v/2 NUM=P) *jste* (PER=v/2 NUM=P) *mluvili* (NUM=P) (46)

4.2.6 Stupeň (DEG)

Kategorie je relevantní pro přídavná jména a příslovce.

- 1: pozitiv
- 2: komparativ
- 3: superlativ
- s: typ sebe + komparativ

U stupňování se vedou spory, zda patří do morfologie, nebo spíše do slovo-tvorby, viz např. obšírné pojednání (Karlík – Hladká, 2004) s argumenty pro obě zařazení. Pro morfologickou analýzu i syntézu je důležité, aby bylo možno zachytit všechny slovní tvary a přiřadit jim rozumné hodnoty. Zařazujeme tedy kategorii **Stupeň** mezi morfologické, a to flektivní kategorie.

První tři hodnoty, vyjádřené čísly, nepotřebují komentář.

Poslední hodnota, tedy s, však není běžná. Nikde se do stupňování nepočítá. Týká se slov typu *sebekrásnější*, *sebekrásněji*. Tvoření přídavných jmen a příslovcí tímto způsobem je však velmi pravidelné a týká se všech stupňovatelných lemmat. Je tedy přirozené zachytit takto vytvořené tvary s lemmatem společným i pro stupňované tvary.

V dosavadní praxi se tyto typy přídavných jmen a příslovcí lemmatizují jako samostatné jednotky, např. $\lambda(\textit{sebemenší}) = \textit{sebemenší}$. Mnoho takových slov např. v pražském morfologickém slovníku není, a tak zůstávají nerozpoznána. Vzhledem k naprosté pravidelnosti jejich tvoření, vysoké produktivitě a zjevné příslušnosti k pozitivu přídavného jména či příslovce je přirozené je začlenit do paradigmatu příslušného pozitivu. Kategorie **Stupeň** nám pro takový popis připadá nejvhodnější. Pokračovat v číselné stupňovací řadě, jak by se na první pohled mohlo zdát přirozené, nám však nepřipadá vhodné, neboť slova typu *sebekrásnější* nezapadají logicky do řady 1., 2. a 3. stupně. Mají totiž odlišný význam. Proto jsme zvolili pro tento typ kód mimo číselnou řadu¹⁰.

¹⁰ Ale jak už bylo několikrát zdůrazněno, není to podstatné. Jde samozřejmě o pouhý kód.

4.2.7 Negace (NEG)

Kategorie je relevantní pro slovesa, přídavná jména, příslovce a (v omezené míře i) pro podstatná jména (*hlava nehlava*).

Má tyto hodnoty:

- N: pro záporné slovní tvary, které začínají záporkou *ne-*
- A: pro ostatní slovní tvary

V případě negace dochází u některých ajdektiv, substantiv a sloves ke sporům, zda existují pozitivní tvary či nikoliv. Souvisí to i s lemmatizací — má být lemmatem záporných tvarů základní tvar v pozitivním, nebo negativním tvaru? Tento problém se týká i lemmatizace, zabývali jsme se jím tedy i v kapitole 2.

Uveďme několik dalších příkladů:

Substantivum *nepřítel* je sice opakem slova *přítel*, ale většinou se již takto nechápe, proto by mělo být lemmatizováno jako *nepřítel*, s hodnotou NEG=A.

Nemoc není opakem *moci*, i zde jde o dvě lemmata, obě s hodnotou NEG=A.

Podobně na tom je přídavné jméno *nesmyslný*. I zde jsou dvě lemmata — *smyslný* i *nesmyslný*, obě s NEG=A.

Jsou ale i slova, která jsou sporná, např. adverbium *nekale* (viz též oddíl 2.2.2).

Nejjednodušší by bylo technické řešení, kdy by se prohlásilo, že všechna slova s předponou *ne-* mají hodnotu kategorie NEG=N a lemma se rovná základnímu slovnímu tvaru bez předpony *ne-*. Uvědomujeme si však, že takové řešení by se setkalo s velkou nevolí na straně většiny uživatelů korpusů, proto necháme rozhodnutí na správci konkrétního slovníku.

4.2.8 Slovesný tvar (VRB)

Kategorie je relevantní pro slovesa, přídavná jména (pasivum), částice a spojky (kondicionál).

- P: indikativ přítomného času (*kolíbá*)
- B: budoucí čas (*ponese, bude*)
- F: infinitiv (*otevřít*)
- I: imperativ (*peč*)
- L: příčestí činné (*strouhal*)
- T: příčestí trpné (*zavřen*)
- K: kondicionál (*aby, kdyby, by*)
- p: přechodník přítomný (*starajíc*)
- m: přechodník minulý (*vtoupiv*)

Podobnou množinu hodnot má i brněnský systém pro kategorii Mód. I Konkláve se rozhodlo, že nebude kódovat tradiční kategorie sloves jako je Čas a Slovesný rod, protože existuje jen několik málo smysluplných kombinací hodnot těchto kategorií. V pražském systému to byly kombinace hodnot tří kategorií, a to Detailní určení slovního druhu (SUBPOS — pozice 2), Čas (TENSE — pozice 9) a Aktivum/pasivum (VOICE — pozice 12). Jediné kombinace, které se mohly vyskytnout, uvádí tabulka 4.6.

	Číslo pozice			Stručná vysvětlivka
	2	9	12	
★	B	P	A	přítomný čas
★	B	F	A	budoucí čas
★	f	-	-	infinitiv
★	i	-	-	imperativ
★	p	R	A	příčestí činné (včetně přidaného -s)
	s	H	P	pasivní příčestí se zakončením -s
★	s	X	P	pasivní příčestí
★	e	-	-	přechodník přítomný
★	m	-	-	přechodník minulý
	c	-	-	kondicionál slovesa <i>být</i>
	q	R	A	min. čas archaický (<i>vstalt</i>)
	t	F	A	archaický budoucí čas s - <i>ť</i> (<i>budut</i>)
	t	P	A	archaický přítomný čas s - <i>ť</i> (<i>dávámť</i> , ale i <i>poradímť</i>)

Tabulka 4.6: Jediné možné kombinace 2., 9. a 12. pozice pražského systému.

Kombinace hodnot, u nichž je v tabulce 4.6 uvedena hvězdička (★), jsme sdružili do jedné hodnoty nové kategorie, kterou nazýváme **Slovesný tvar**.

Řádky tabulky bez hvězdičky probereme postupně:

Řádka s názvem „pasivní příčestí se zakončením -s“ zahrnuje slova typu *rytas*, *propuštěnas*. V korpusu SYN se 600 miliony slov jsme takových slov našli 121, avšak všechny jsou buď špatně označovány, nebo jde o překlepy. Kdyby se někdy takový tvar vyskytl, bude se jednat o slovesnou složeninu (viz kapitulu o složeninách 6).

Řádka, nazvaná „kondicionál slovesa *být*“, označuje slovní tvary *by*, *bys*, *bych*, *bychom*, *byste*. Tyto slovní tvary pojednáváme ve zvláštní kapitole 5 spolu s dalšími dvěma kondicionálovými lemmaty *aby* a *kdyby*. Na tomto místě je pouze třeba upozornit, že nově zavedená hodnota Kondicionál (K) není relevantní pro slovesa, ale pro spojky (*aby*, *kdyby*) a částice (*by*).

Poslední tři řádky, kde je na druhé pozici q nebo t, jsou archaické, s -*ť* na konci slovního tvaru. Zde není třeba zvláštní značky. Rozdíl oproti slovnímu tvaru bez -*ť* zachycujeme pomocí kategorie Flektivní mutace FMU (viz kap. o mutacích 3).

Místo uvedených tří kategorií (pozic pražského systému) jsme se tedy rozhodli zavést jednu, jejíž hodnoty budou odpovídat všem jejích možným kombinacím. Kupodivu je jich jen 8 — ty, co jsou v tabulce 4.6 označeny hvězdičkou. Tyto kombinace se tedy staly hodnotami nové kategorie nazvané **Slovesný tvar**. K nim přidáváme devátou hodnotu Kondicionál, která je však odlišná od kondicionálu, uvedeného v tabulce 4.6 na 4. řádku zdola.

Na chvíli se zastavme u některých hodnot kategorie **Slovesný tvar**.

Příčestí trpné Tato hodnota kategorie Slovesný tvar se neurčuje u slovního druhu sloveso. V našem novém návrhu jsme totiž všechna příčestí trpná zařadili mezi jmenné tvary přídavných jmen. Formálně se tak chovají a často je velmi obtížné rozlišit, zda jde o příčestí trpné slovesa, nebo o jmenný tvar přídavného jména. Pro přídavná jména hovoří i fakt, že u příčestí trpného se může měnit, byť velmi omezeně, pád. Ve všech rodech i číslech může vystupovat v akuzativu, jak ukazují příklady (47) až (54):

Za hodinu jsme měli připravenu hromadu klestí (47)

... měl najatu restauraci (48)

Hranol má hranu podstavy rovnu $a=24$ cm (49)

... máme... hotovu dokumentaci (50)

Psychotesty již máme hotovy (51)

Základní návrh chceme mít hotov v březnu. (52)

Musíme mít připraven mírový plán (53)

Budeme mít připravena i vodní děla. (54)

Přesto, že příčestí trpné od slovesa vzniklo, tvoří protiklad k příčestí činnému a většinou vystupuje ve větě v její přísudkové části, formálně je možné ho vždy nahradit jmenným tvarem přídavného jména. Vzhledem k obtížné rozlišitelnosti je jednodušší, když ho vždy považujeme za přídavné jméno.

Kvůli zachování těsné vazby ke slovesu však vyplňujeme u těch jmenných tvarů přídavných jmen, která jsou zároveň příčestím trpným, i tuto hodnotu kategorie Slovesný tvar. Například tedy slovní tvar *otevřen* má tyto morfologické hodnoty:

POS=A, (Slovní druh: přídavné jméno)

SUB=S, (Poddruh: deverbativní)

GEN=M/I, (Rod: mužský životný nebo neživotný)

NUM=S, (Číslo: jednotné)

CAS=1, (Pád: 1)

VRB=T, (Slovesný tvar: příčestí trpné)

NEG=A, (Negace: pozitiv)

lemma=*otevřený*

Bude-li třeba vytvořit dotaz na vyhledání všech sloves, včetně příčestí trpných, lze to učinit takto: $POS=V \vee (POS=A \wedge VRB=T)$ ¹¹.

Složené slovesné tvary Složené slovesné tvary neurčujeme jako celek. Jsme si vědomi, že by bylo vítané, kdybychom rozpoznali ve větě složené slovesné tvary, ale tato úloha přesahuje rámec morfologie. Jak už jsme zdůraznili na začátku, zabýváme se jednotlivými slovy, nikoli jejich kombinacemi. Z toho důvodu také nemohou mít nedokonavá slovesa nikdy hodnotu kategorie VRB=B (budoucí čas).

¹¹Vyhledají se všechna slovesa (POS=V) a ta přídavná jména (POS=A), která jsou současně slovesným trpným rodem (VRB=T).

Budoucí čas Hodnota budoucí čas se týká jen tvarů budoucího času slovesa *být* (tedy *budu*, *budeš*, *bude*, *budeme*, *budem*, *budete*, *budou*) a dokonavých sloves tvořících budoucí čas pomocí předpony *po-*.

Všechna ostatní dokonavá slovesa, byť sémanticky vyjadřující budoucí čas, mají hodnotu indikativ přítomnosti (VRB=P). Tedy: *nesu* (VRB=P), *ponesu* (VRB=B), ale *přinesu* (VRB=P).

Nedokonavá slovesa tvoří budoucí čas pouze ve složených slovesných tvarech, ta tedy nemají VRB=B nikdy (viz předchozí odstavec).

Kondicionál Tato hodnota popisuje tvary částice *by* a spojek *kdyby* a *aby*. Podrobně o kondicionálu pojedáváme v kapitole 5.

Přechodníky Přechodník přítomný určujeme pouze u nedokonavých sloves, přechodník minulý pouze u sloves dokonavých. U obouvidých sloves se mohou vyskytovat oba typy přechodníků. Podobně přechodníky pomocného slovesa *být* mohou být buď přítomné (*jsa*, *jsouc*, *jsouce* mají VRB=p), nebo minulé (*byv*, *byvši*, *byvše* mají VRB=m).

Oproti současnému pražskému systému vypouštíme možnost tvoření přechodníku přítomného pro dokonavá slovesa, jakožto zastaralý a už dávno neuzívaný slovesný tvar. Při morfologické analýze jsme schopni takový tvar rozpoznat a správně určit pomocí guessru.

4.2.9 Jmenný tvar přídavných jmen (NOM)

Kategorie je relevantní pro přídavná jména.

Možné hodnoty:

- J: jmenný tvar
- UNDEF

Jmenné tvary přídavných jmen jsou v pražském systému popsány pomocí 2. pozice SUBPOS. Při zběžném pohledu by se mohlo zdát, že jmenný rod je poddruhem přídavného jména. Toto zařazení bylo přijato i na Konkláve. V tom případě bychom však jmenný rod nemohli zahrnout pod společné lemma dlouhého tvaru, neboť poddruh je kategorie globální. Jestliže chceme, aby např. $\lambda(\textit{sláb}) = \lambda(\textit{slabý}) = \{\textit{slabý}\}$, musíme tuto kategorii vyčlenit zvlášť.

Mezi jmenné tvary přídavných jmen počítáme i tvary přičestí trpného sloves, neboť je často velmi obtížné je od sebe rozlišit. Slovní tvar *ukryt* mohl být odvozen jak z přídavného jména *ukrytý*, tak ze slovesa *ukrýt*. Podobně slovní tvar *spokojen* můžeme chápat jako odvozeninu od *spokojit* i od *spokojený*. Vzhledem k tomu, že se ve větě chovají tyto tvary stejně jako jmenné tvary ostatních přídavných jmen jmenných (např. *mlád*), mají dokonce omezenou flexi (4. pád všech rodů i čísel), zařazujeme je do této kategorie. Z toho speciálně vyplývá, že jejich slovní druh není sloveso, ale přídavné jméno.

Abychom však zachytili v popisu jejich slovesný charakter, je pro ně relevantní flektivní morfologická kategorie **Slovesný tvar**, a to s hodnotu T (přičestí trpné). Uvědomujeme si, že toto řešení není zcela v souladu s běžným chápáním slovesného trpného rodu, ale podařilo se nám tak jednoznačně

popsat sporné případy, kdy není jasné, zda jde o přídavné jméno či o sloveso, aniž by se tím ztratila jakákoliv informace. Viz též 4.2.8.

4.2.10 Stupeň intenzity slovesného děje (INT)

Kategorie je relevantní pro slovesa¹².

V oddíle 2.2.1.2 jsme upozornili na pravidelné tvoření zvrtných podob nedokonavých sloves pomocí určitých předpon. Ukázali jsme, že s takto vytvořenými tvary by se nemělo zacházet jako se samostatnými slovesy, a zařadili jsme je pod lemma jejich neprefigovaného základního slovesa. Pomocí kategorie INT odlišíme tvary základního slovesa od z něho odvozených zvrtných sloves prefigovaných.

Hodnoty jsou:

- r: pro předponu *roz-*
- p: pro předponu *po-*
- z: pro předponu *za-*
- n: pro předponu *na-*
- v: pro předponu *vy-*
- u: pro předponu *u-*

4.2.11 Typ složeniny (CMP)

Tato kategorie je relevantní jen pro složeniny a pojednáme o ní v kapitole 6 o složeninách.

4.2.12 Flektivní mutace (FMU)

Flektivní mutace se projevují především pomocí koncovek a jsou většinou systematické. Podle definice se nikdy netýkají celého paradigmatu, vždy jen některých kombinací hodnot gramatických kategorií. Systematické hodnoty kategorie FMU jsou z velké části zahrnuty přímo do ohýbacích vzorů, pojednáme o nich tedy v kapitolách o vzorech. Mezi flektivní mutace zahrnujeme i nekodifikované koncovky, které se však běžně používají, takže by se do systému paradigmat měly zahrnout jako varianty (mutace) koncovek spisovných. Na tento fakt poukázali už v roce 1992 Sgall a Hronek (viz (Sgall – Hronek, 1992)). Příkladem je koncovka *-ma* v 7pl všech skloňovaných slov, nebo používání *-ej-* u tvrdého adjektivního skloňování. V současném pražském systému nespisovné varianty většinou zahrnuty jsou a my je zachováváme. Lišíme se jen v jejich značení.

K mutacím dochází navíc u všech slov, jejichž skloňování kolísá mezi dvěma vzory stejného rodu, např. *stroj* a *hrad*, *muž* a *pán*, *kost* a *píseň*. Některé kombinace morfologických kategorií vytvoří podle obou vzorů stejný slovní tvar, u jiných je tvar odlišný. Právě tehdy je kategorie FMU relevantní.

¹²Dalo by se uvažovat i o slovesných odvozeninách — deverbativech.

Toto však není případ lemmat s kolísajícím rodem (např. *kredenc*), viz pojednání (Brabcová, 2004). Přestože se z lingvistického hlediska může jednat o varianty, není třeba zde tvary odlišovat pomocí kategorie FMU, neboť jsou rozlišeny hodnotou kategorie Rod.

Systematické případy variant (mutací) jsou uvedeny v kapitolách o vzorech. Nechceme tvrdit, že jsme na žádnou flektivní mutaci nezapomněli. Vyčerpávající přehled české morfologie ani nebyl cílem této práce. Nabízíme však způsob, jak konzistentně flektivní mutace zachytit. Není problém neuvedený typ flektivní mutace do seznamu zahrnout. Při kódování jejího typu stačí dodržet podmínku jednoznačné hodnoty v rámci konkrétní kombinace morfologických kategorií, jichž se mutace týkají.

Nesystematické mutace frekventovaných slov mohou mít své vlastní, specifické hodnoty, ostatní navrhuje značit pomocí číslic. Existuje-li např. lemma X , u něhož se vyskytnou dva různé slovní tvary X_1 a X_2 se stejnými hodnotami všech relevantních morfologických kategorií, můžeme položit $FMU=1$ pro X_1 a $FMU=2$ pro X_2 , jestliže typ mutace je neobvyklý a není pokryt standardním číselníkem.

Pro ilustraci se na tomto místě se ještě zmiňme o několika nepravidelných, leč četných mutacích.

Jde o tvary slovesa *jít*, které v přítomném a budoucím čase a v imperativu ztrácejí počáteční *j-*, např. *jdu — du, půjdu — pudu, jděte — děte*. Přiřazujeme jim $FMU=g$. V první osobě množného čísla se zde dokonce kombinují dva druhy mutace, oba flektivní. Jeden je pravidelný, tedy ztráta koncového *-e*, druhý nepravidelný, ztráta počátečního *j-*: *jdeme, deme, dem, jdem*. Tato kombinace hodnot gramatických kategorií lemmatu *jít* má tedy 4 různé flektivní mutace, jednu spisovnou, ostatní nespisovné.

Jiné nepravidelné flektivní mutace jsou např. *míň* a *méně* pro lemma *málo* a *líp*, *lépe* pro lemma *dobře*. Zde můžeme využít hodnot $FMU=K$ pro kratší tvary *míň*, *líp* a $FMU=D$ pro delší tvary *méně*, *lépe*.

Mezi flektivní mutace zařazujeme i mutace skloňování osobních zájmen. Jde o krátké a dlouhé tvary (*tebe — tě, mne — mě, mně — mi, jeho — jej*). Odlišujeme je opět hodnotami $FMU=D$ a $FMU=K$.

Polský morfologický systém, vytvořený pro morfologickou anotaci korpusu IPI PAN (viz (Przepiórkowski, 2004)) zavádí pro zachycení této variability zvláštní morfologickou kategorii „Accentability“. My jsme se rozhodli využít kategorie **Flektivní mutace**. Důvodem je především určitá šetrnost — zavedení nové kategorie nám připadá zbytečné, když je možno využít v tomto případě jinak nevyužitou kategorii FMU, která je i tak velmi nesourodá.

Osobní zájmena pro 3. osobu jednotného i množného čísla, tedy *on, ona, ono, oni, ony* také mají mutace. Předchází-li těmto zájmenům předložka, mění se ve všech pádech počáteční *j-* na *n-*, případně *je-* na *ně-*. Dostáváme tedy dvojice *jeho — něho, jej — něj, ji — ni, jimi — nimi* atd. s hodnotami $FMU=J$ a $FMU=N$.

I tyto mutace řeší v Polsku pomocí zvláštní kategorie „Post-prepositionality“. Důvod, proč jsme zvolili řešení pomocí kategorie **Flektivní mutace**, je stejný jako v předchozím případě.

Pro 4. pád jednotného čísla lemmatu *on* tak máme dokonce 5 různých slov-

ních tvarů: *ho, jeho, něho, jej, něj*¹³, s hodnotami po řadě FMU=K, Dj, Dn, Kj, Kn).

Při označování hodnot kategorie **Flektivní mutace** většinou nepřičítáme hodnotu mutacím, které se užívají v psaném textu nejběžněji. Považujeme je za tzv. nulové mutace. Alternativní přístup by mohl všem takovým mutacím přiřadit FMU=0, což ale považujeme za zbytečné.

Námítka, že právě popsáný způsob označování mutací je složitý, je správná. Vzhledem k tomu, že jde o problém velmi rozsáhlý a mnohotvárný, domníváme se, že jednoduché řešení ani neexistuje. Naším cílem bylo navrhnout jednoznačné rozlišení mutací slovních tvarů a lemmat, aby vždy mohlo být splněno Zlaté pravidlo morfologie. Tohoto cíle jsme zřejmě dosáhli.

Pokud by si uživatelé slovníku, potažmo korpusů pomocí slovníku anotovaných, přáli mít hodnoty mutací jiné, je možno navržené hodnoty ohodnotit podle nějakého kritéria (např. podle kritizovaného stylového příznaku) a toto hodnocení vložit jako hodnoty do speciální nové kategorie. Tím by se různé hodnoty mutací sdružily do několika tříd, podle přání uživatelů.

Poznámka: Syntaktický slovní druh (SYN)

Na Konkláve se hovořilo též o kategorii **Syntaktický slovní druh**. Byla vymezena poměrně vágně:

Syntaktický slovní druh je kategorie, která vyjadřuje, jak se dané slovo obvykle chová v rovině povrchové syntaxe.

Tato kategorie měla usnadnit práci na pravidlové desambiguaci,

V novém návrhu od zavedení této kategorie v rámci morfologie upouštíme, neboť ve skutečnosti nejde o kategorii morfologickou. Neslouží totiž k popisu slovních tvarů, jedná se o kategorii syntaktickou. Je sice pravda, že některé kategorie, které jsme již zavedli, také nejsou zcela morfologické (např. kategorie **Poddruh**), ale pomáhají při popisu jednotlivých slovních tvarů. Kromě toho také rozlišují relevantnost dalších kategorií. Kategorie **Syntaktický slovní druh**, která měla být relevantní pro všechny slovní druhy, však žádné rozlišení na morfologické rovině nepřináší.

4.3 Morfologická značka

Morfologická značka je kód, pomocí něhož lze jednoznačně určit hodnoty všech relevantních morfologických kategorií pro daný slovní tvar.

V úvodní kapitole jsme prohlásili, že se konkrétní podobou morfologické značky zabývat nechceme. Kdybychom chtěli být důslední, tento oddíl o morfologické značce bychom vůbec do své práce nezařazovali. Morfologické kategorie a jejich hodnoty, jak jsme je zavedli v předchozích oddílech, je možné použít k vytvoření jakéhokoli kódu. Abychom si zjednodušili práci s vyjmenováváním kategorií a jejich hodnot v následujících kapitolách, přece jen morfologickou značku zavedeme.

Následující popis morfologické značky je tedy třeba brát jako příklad, jak také je možno morfologické kategorie kódovat. Pro konkrétní aplikace je možno použít jen některé morfologické kategorie, např. jen kategorie POS. Takovým

¹³Tvary typu *doň, zaň* sem nepočítáme, to jsou podle nového návrhu složeniny s jiným lemmatem.

kódům už nebudeme říkat morfologická značka, protože nekódují všechny relevantní morfologické kategorie, ale ve speciálních případech mohou být i takové kódy užitečné.

Morfologická značka musí být jednoznačná, aby bylo možno podle ní rozlišit rozdílné hodnoty morfologických kategorií. Teoreticky by bylo možno všechny kombinace hodnot morfologických kategorií nějakým způsobem očíslovat a jako kódu používat čísel. Takové kódování je samozřejmě nežádoucí, neboť není ani trošku intuitivní. Předpokládáme, že hodnoty morfologických kategorií se využijí k sestavení takového kódu, ze kterého půjdou jednoduše vyčíst.

Základní typy morfologických značek používaných pro češtinu jsou:

- poziční
- kompaktní
- hodnotový

Příkladem poziční značky je dosavadní morfologická značka pražská. Má jednotnou délku 15 pozic a každá její pozice kóduje jednu konkrétní kategorii (viz (Hajič, 2004)). Je zřejmé, že poziční značky „plýtávají místem“, protože neexistuje slovní tvar, pro který by bylo všech 15 kategorií relevantních. Na druhou stranu se s ním dobře pracuje, neboť každá kategorie má ve značce své neměnné místo.

Pražský systém používá i tzv. kompaktní značky (viz též (Hajič, 2004)), které jsou utvořeny tak, aby obsahovaly pouze hodnoty relevantních morfologických kategorií, a přitom zůstaly jednoznačné. Tento typ značek šetří místem, je však méně přehledný.

Příkladem hodnotového typu je značka brněnská (viz (Sedláček, 1999)). Její délka je proměnlivá, avšak vždy sudá, neboť obsahuje vždy dvojici (název kategorie, její hodnota) pro každou morfologickou kategorii relevantní pro daný slovní tvar.

V našem návrhu morfologické značky použijeme poziční typ a typ hodnotový. Poziční typ nám připadá přehlednější pro zachycení hodnot morfologických kategorií. Hodnotový typ má proti pozičnímu tu výhodu, že je snadné ho rozšířit tak, aby mohly mít kategorie více hodnot. Proto ho využijeme pro kódování hodnot mutací, a to jak globálních (kategorie GMU), tak flektivních (kategorie FMU).

Morfologickou značku rozdělíme na dvě části, globální a flektivní. Definovali jsme celkem 6 globálních morfologických kategorií a 12 flektivních. Z tohoto celkového počtu vyjmemme kategorie **Globální mutace** a **Flektivní mutace**, protože ty kódujeme jinak (viz kap. 3). Zbylých 5 globálních a 11 flektivních kategorií zakódujeme do poziční značky takto: 1. až 5. pozice tvoří globální část morfologické značky a obsahuje hodnoty těchto kategorií:

1. POS
2. SUB
3. FCE
4. ASP
5. ABR

Druhá, flektivní část značky, tedy pozice 6 až 16, obsahuje po řadě hodnoty těchto kategorií:

6. GEN
7. NUM
8. CAS
9. DUA
10. PER
11. DEG
12. NEG
13. VRB
14. NOM
15. CMP
16. INT

4.4 Relevantnost kategorií

V kapitole 4 jsme vyjmenovali všechny kategorie, které popisují v našem pojetí české slovní tvary. Ne všechny kategorie jsou však relevantní pro všechny slovní tvary. Relevantnost kategorie závisí především na kategorii **Slovní druh**, často ještě na kategorii **Poddruh** a na kategorii **Slovesný tvar**.

Množiny kategorií relevantních pro jednotlivé dvojice hodnot POS a SUB jsou přehledně zachyceny v tabulce 4.7. Číslování v první řádce se odkazuje k pozicím v návrhu morfologické značky. Relevantní kategorie jsou vyplněny znakem \oplus . Znak \ominus v některých buňkách tabulky znamená, že tato kategorie je relevantní pouze někdy:

U zájmen záleží nejen na konkrétní kombinaci hodnot kategorií **Poddruh** a **Funkce**, ale i na konkrétním lemmatu, která kategorie je relevantní. Např. u substantivního určitého lemmatu *já* není zvykem určovat rod, zatímco u lemmatu *on*, které je rovněž substantivní určité, ano. Pouze kategorie **Pád** je relevantní pro všechna zájmena.

Podobné je to u číslovek. Číslovky 1 až 4 mohou vyjadřovat duál, ostatní substantivní určité číslovky ne. Číslovky 1 a 2 navíc vyjadřují i rod.

Číslovka řadová *první* má jako jediná relevantní kategorii **Stupeň** (viz příklady (55) a (56)¹⁴), což je naznačeno znakem \ominus v příslušné buňce tabulky.

dle volebních preferencí čím dál prvnější (55)

jeho otec byl nejprvnějším dělníkem v přístavě (56)

Dále je to kategorie **Slovesný tvar** v případě spojek. Ta je relevantní pouze pro spojky *aby* a *kdyby*.

Pro slovní druh **Složeniny** jsme použili znak \ominus u všech kategorií kromě **Typ složeniny**, která je naopak relevantní právě pouze pro složeniny. Relevantnost ostatních kategorií závisí právě na ní. Podrobná tabulka týkající se složenin je uvedena v kapitole o složeninách na str. 67.

¹⁴V pražském slovníku jsou tyto tvary klasifikovány jako přídavná jména.

4 Morfologické kategorie

1 POS	2 SUB	3 FCE	4 ASP	5 ABR	6 GEN	7 NUM	8 CAS	9 DUA	10 PER	11 DEG	12 NEG	13 VRB	14 NOM	15 CMP	16 INT
N	S		⊕	⊕	⊕	⊕	⊕	⊕			⊕				
	0			⊕	⊕	⊕	⊕	⊕			⊕				
	U			⊕	⊕	⊕	⊕	⊕							
A	GM		⊕	⊕	⊕	⊕	⊕	⊕			⊕		⊕		
	S			⊕	⊕	⊕	⊕	⊕		⊕	⊕	⊕			
	ostatní			⊕	⊕	⊕	⊕	⊕		⊕	⊕		⊕		
P				⊕	⊖	⊖	⊕	⊖	⊖						
	1		U	⊕	⊖	⊖	⊕	⊖							
	1	ostatní		⊕			⊕								
C	usd			⊕		⊕	⊕	⊕							
	r			⊕	⊕	⊕	⊕	⊕		⊖					
	p			⊕	⊕	⊕	⊕								
	nov			⊕											
				⊕	⊕	⊕	⊕				⊕	Lpm			⊕
V			⊕	⊕		⊕			⊕		⊕	PBI			⊕
				⊕							⊕	F			⊕
	zájmenná		⊕	⊕						⊕					
D	ostatní			⊕						⊕	⊕				
				⊕			⊕								
R				⊕											
J				⊕								⊖			
I															
	c				⊕										
	7								2						
T	ostatní														
S	⊖		⊖	⊖	⊖	⊖	⊖	⊖	⊖	⊖	⊖	⊖	⊖	⊕	⊖
F				⊕	⊕	⊕	⊕								
G				⊕											

Tabulka 4.7: Relevantní kategorie v závislosti na kategoriích Slovní druh a Poddruh

4 Morfologické kategorie

Příklady Několik příkladů morfologické značky je uvedeno v tabulce 4.8.

Slovní tvar	Morf. značka	Mutace	Lemma
<i>okénko</i>	NO---NS1---A----		{ <i>okénko</i> , <i>okýnko</i> , <i>vokýnko</i> }
<i>vokýnko</i>	NO---NS1---A----	Gvy	
<i>okénky</i>	NO---NP7---A----		
<i>vokýnkama</i>	NO---NP7---A----	GvyFa	
<i>pánové</i>	NO---MP5---A----	Fv	<i>pán</i>
<i>páni</i>	NO---MP5---A----	Fi	
<i>neutřen</i>	AS-D-IS1---NTJ--		<i>utřený</i>
<i>nes</i>	VO-N--S--2-AI---		{ <i>nést</i> , <i>nésti</i> }
<i>pones</i>	VO-N--S--2-AI---	Fp	
<i>doběhla</i>	VO-D--S--3-AL---		<i>doběhnout</i>
<i>doběhnula</i>	VO-D--S--3-AL---	F _n	
<i>něj</i>	PZU--MS2-3-----	FK _n	<i>on</i>
<i>jeho</i>	PZU--MS2-3-----	FD _j	
<i>jeho</i>	PUU--XXX-3-----		

Tabulka 4.8: Příklady morfologického popisu slovních tvarů

5 Kondicionál

Kondicionál se v češtině vyjadřuje pomocí „slůvek“ *by*, *aby* a *kdyby*. Už samo jejich zařazení mezi slovní druhy je problematické a vedou se o něm stále spory. Neděláme si zde nárok na jejich vyřešení, ale pro automatické zpracování češtiny je třeba tato slova nějak klasifikovat.

Mohlo by se zdát, že nemá smysl zavádět speciální kategorii pro pouhá tři lemmata, když je lze zadat výčtem. Ze strany uživatelů korpusů a ostatních lingvistů se však ukazuje potřeba mít jejich kondicionálovou povahu zvlášť vyznačenu.

V brněnské morfologii pro ně byl vytvořen nový slovní druh se značkou **kY**. Nemá speciální název.

Pražská morfologie řadí *by* (historicky) mezi slovesa, *kdyby* a *aby* (tradičně) mezi spojky. Od ostatních sloves se *by* odlišuje zvláštní hodnotou na 2. pozici (poddruh). Tvary *aby* a *kdyby* se odlišují od ostatních spojek pouze ve svých časovaných tvarech.

Ve slovenštině řeší podobný problém zavedením kategorie „kondicionálnost“, která je relevantní pro slovní druh spojka a částice.

Slovenské řešení nám připadá nejvhodnější, ale rozhodli jsme se pro jeho odlišné zpracování.

Kondicionál se projevuje výhradně ve spojení se slovesem, pomocí něj se tvoří podmiňovací způsob. Proto ho zařazujeme jako novou hodnotu kategorie **Slovesný tvar**. Slovnědruhové zařazení však neměníme, takže se kategorie **Slovesný tvar** stává (trochu překvapivě) relevantní pro slovní druhy spojek a částic. Hodnota kondicionál se vlastně netýká žádného slovesa, ale právě jen spojek a částic, a to jen oněch tří vyjmenovaných v úvodu této kapitoly.

Zvláštní povaha těchto tří kondicionálových slov se projevuje už tím, že ač řazena mezi neohebné slovní druhy (částice a spojky), mají vyjádřenu osobu a číslo.

V tabulce 5.1 jsou přehledně zpracovány hodnoty všech relevantních morfologických kategorií pro paradigmata všech tří kondicionálových slov.

Dvojí hodnotu ve sloupci kategorie **Číslo** a **Osoba** u slovních tvarů *aby*, *kdyby*, *by* nelze libovolně kombinovat. Možné kombinace čísla a osoby jsou pouze S2 (2. osoba singuláru, viz příklady (57) až (59)), S3 (3. osoba singuláru, viz (60) až (62)) a P3 (3. osoba plurálu, příklady (63) až (65)). Kombinace P2 (2. osoba plurálu) možná není.

záleží na tom, aby sis udržel chladnou hlavu (57)

Chováš se, jako kdyby ses narodila v jiném století (58)

Myslím, že by ses měl připravit na menší šok (59)

Přemýšlí o muškách, které by si dala k snídani (60)

A kdyby se za něj provdala, (61)

Nechci, aby to někdo věděl (62)

5 Kondicionál

Lemma	Slovní tvar	POS	SUB	VRB	NUM	PER	FMU
<i>aby</i>	<i>aby</i>	J	,	K	SP	23	—
	<i>abych</i>	J	,	K	S	1	—
	<i>abys</i>	J	,	K	S	2	—
	<i>abychom</i>	J	,	K	P	1	0
	<i>abysme</i>	J	,	K	P	1	1
	<i>abyste</i>	J	,	K	P	2	—
<i>kdyby</i>	<i>kdyby</i>	J	,	K	SP	23	—
	<i>kdybych</i>	J	,	K	S	1	—
	<i>kdybys</i>	J	,	K	S	2	—
	<i>kdybychom</i>	J	,	K	P	1	0
	<i>kdybysme</i>	J	,	K	P	1	1
	<i>kdybyste</i>	J	,	K	P	2	—
<i>by</i>	<i>by</i>	T	c	K	SP	23	—
	<i>bych</i>	T	c	K	S	1	—
	<i>bys</i>	T	c	K	S	2	—
	<i>bychom</i>	T	c	K	P	1	0
	<i>bysme</i>	T	c	K	P	1	1
	<i>byste</i>	T	c	K	P	2	—

Tabulka 5.1: Relevantní kategorie a jejich hodnoty kondicionálových slov

A to mě přivádí k myšlence, jaké by to bylo, kdyby v zoologických byly taky ukázky lidí. (63)

jako by tomu Jihoafričané nechtěli věřit (64)

pokřikovali na muzikanty, aby zase hráli (65)

6 Složeniny

Složenina popisuje slovní tvar, který zastupuje dva nebo více slovních tvarů (**složek** složeniny) a většinou mu není možné přiřadit jednoduše slovní druh. Příklady: *naň = na něj, byls = byl jsi*.

Složenina vzniká spojením těchto slovních tvarů, nejde však o prosté zřetězení. Většinou lze složeninu ve větě původními slovními tvary nahradit, aniž by se změnil smysl věty.

Složenině nelze přiřadit žádný z klasických slovních druhů, až na slovesnou složeninu typu V (viz dále). To byl důvod, proč jsme složeninu zavedli jako samostatný slovní druh. Toto řešení bylo přijato už na Konkláve, ač pod jiným názvem. Rozdělení do typů a jejich přesné vymezení provedla autorka této práce.

Složeninou v našem pojetí však nejsou slova vzniklá jedním z tradičních způsobů slovo tvorby, která mají svůj vlastní význam a již se zařadila jednoznačně do některého slovního druhu. Složeninami například nejsou slova *černobílý, novotvar, spoluporaďat*, přestože také vznikla složením různých slov.

6.1 Lemma složenin

Pro lemma složenin využijeme nově zavedeného konceptu vícenásobného lemmatu (viz kap. 2). Lemma složenin je tedy vícenásobné a jeho prvky jsou lemmata jednotlivých složek složeniny. To je výhodné pro vyhledávání jednotlivých lemmat v korpusech. Vícenásobné lemma zajistí, že složenina bude ve výsledku vyhledávání podle lemmatu libovolné své složky. To neplatí pro konkrétní slovní tvary, např. dotaz na slovní tvar *jemu* nezahrne slovesnou složeninu *jemus*, což je ale správné, protože se skutečně jedná o dva různé slovní tvary. Dotaz na lemma *on* s hodnotou kategorie CAS=3 již ale oba slovní tvary najde, neboť lemma *on* je součástí vícenásobného lemmatu slovního tvaru *jemus*.

6.2 Relevantní morfologické kategorie složenin

Podívejme se, jaké morfologické kategorie jsou pro složeniny relevantní a jakých mohou nabývat hodnot.

Možnost, že složenina bude mít tolik morfologických značek, kolik má složek, nepřichází v úvahu, neboť složenina sama o sobě je samostatným slovním tvarem a jako taková musí být popsána morfologickou značkou jako celek. Zavedení vícenásobné morfologické značky (podobně, jako jsme zavedli vícenásobné lemma) by bylo navíc neekonomické, protože u většiny složenin nedochází ke konfliktu mezi hodnotami kategorií relevantních pro jednotlivé složky. Složenina tedy může být popsána stejně jako každé jiné slovo jedinou morfologickou značkou.

Označme K_i množinu kategorií, která je relevantní pro i -tou složku složeniny.

Množina K kategorií relevantních pro složeninu vznikne sjednocením množin kategorií relevantních pro její jednotlivé složky:

$$K = K_1 \cup \dots \cup K_n.$$

Označme dále průnik těchto množin:

$$G = K_1 \cap \dots \cap K_n.$$

Hodnoty těch kategorií, které neleží v průniku těchto množin (tedy kategorie z množiny $K - G$), se přenáší i na výslednou složeninu.

Hodnoty kategorií z množiny G však musíme vyřešit pro jednotlivé typy složenin zvlášť. V případě, že se jejich hodnoty shodují u všech složek složeniny, měla by i výsledná hodnota být stejná. Jestliže se však hodnoty kategorií z množiny G pro jednotlivé složky liší, je třeba tento konflikt vyřešit tak, aby výsledkem byla jednoznačná hodnota, samozřejmě bez ztráty informací o jednotlivých složkách složeniny.

Příklad:

Složenina *zač* má 1. složku *za*, s množinou relevantních kategorií

$$K_1 = \{\text{POS}, \text{CAS}\},$$

a 2. složku *co*, s množinou relevantních kategorií

$$K_2 = \{\text{POS}, \text{SUB}, \text{FCE}, \text{CAS}, \text{GEN}, \text{NUM}\}.$$

Tabulka 6.1 ukazuje jejich hodnoty:

	<i>za</i>	<i>co</i>	<i>zač</i>
POS	R	P	S
SUB	-	Z	Z
FCE	-	T	T
CAS	4	4	4
GEN	-	N	N
NUM	-	S	S

Tabulka 6.1: Hodnoty relevantních kategorií složek složeniny *zač*.

Množinu K tvoří kategorie v levém sloupci tabulky, v množině G jsou všechny kategorie, které jsou relevantní pro obě složky složeniny, tedy ty, které mají vyplněnou hodnotu ve 2. a 3. sloupci tabulky 6.1:

$$G = \{\text{POS}, \text{CAS}\}.$$

Hodnota kategorie POS z množiny G má pro výsledný slovní tvar *zač* hodnotu složenina (viz dále), hodnota kategorie CAS není u obou složek v konfliktu, výsledná hodnota může být tudíž stejná, tedy 4. pád. Ostatní kategorie (z množiny $K - G$) přebírá složenina od své složky *co*. Výsledné hodnoty všech kategorií jsou v posledním sloupci tabulky 6.1.

Slovní druh složenin je vždy „složenina“, nezávisle na tom, jaké hodnoty kategorie **Slovní druh** mají jednotlivé složky. Kategorie **Slovní druh** je totiž až na slovesné složeniny typu V (viz dále) vždy konfliktní, neboť $\text{POS} \in G$.

Slovní druhy, které tvoří složky složeniny, je možné odvodit z typu složeniny. Také je možné snadno upravit dotazy využívající hodnot kategorie **Slovní druh** tak, aby zahrnuly i složky složenin, jak ukážeme v oddíle 6.4 o vyhledávání složenin v morfologicky anotovaných korpusech.

Množina G se liší podle typu složeniny, proto teď tyto typy probereme jednotlivě. Podíváme se také na problematiku vyhledávání složenin. Uživatel hledající v korpusu zřejmě bude chtít, aby se složeniny zahrnuly mezi odpovědi na obecné dotazy podle příslušných morfologických kategorií. Proto bude třeba modifikovat některé jednoduché dotazy, aby se mezi výsledky dostaly i příslušné složeniny.

Relevantní kategorie pro jednotlivé typy složenin ukazuje tabulka 6.5 na straně 67.

6.3 Typy složenin

Složeniny rozdělíme do několika základních skupin a v jejich rámci potom na typy. Typ složeniny (CMP) je flektivní morfologickou kategorií, která je relevantní pouze pro složeniny a zastupuje vlastně kategorii poddruh, která je u složenin vyhrazena k popisu jedné ze složek (viz dále).

6.3.1 Typy zájmenné ... n, c

Zájmenné složeniny jsou slovní tvary vzniklé spojením předložky s akuzativní rekcí a tázacího zájmena *co* nebo substantivního určitého (osobního) zájmena *on*. Podle toho rozlišujeme dva typy, které jsou tvořeny uzavřenou množinou tvarů. Můžeme je tedy zadat výčtem.

Typ c

Výčet všech složenin:

1. *zač* = *za co*, lemma: {*za*, *co*}
2. *nač* = *na co*, lemma: {*na*, *co*}
3. *oč* = *o co*, lemma: {*o*, *co*}
4. *več* = *v co*, lemma: {*v*, *co*}
5. *začpak* = *za copak*, lemma: {*za*, *copak*}
6. *načpak* = *na copak*, lemma: {*na*, *copak*}
7. *očpak* = *o copak*, lemma: {*o*, *copak*}
8. *?večpak* = *v copak*, lemma: {*v*, *copak*}
9. *Xzač* = *za Xco*, lemma: {*za*, *Xco*}
10. *Xnač* = *na Xco*, lemma: {*na*, *Xco*}
11. *Xoč* = *o Xco*, lemma: {*o*, *Xco*}
12. *?Xveč* = *v Xco*, lemma: {*v*, *Xco*}

Za X v posledních řádcích tabulky je možné dosadit jednu z neurčitých předpon, které se podílejí na tvorbě neurčitých zájmen, číslovek a příslovcí, totiž *kdoví*, *bůhví*, a další. Jejich seznam je uveden v obr. 4.1 na str. 30.

Mezi složeniny nepatří *proč*, protože jde už o lexikalizovanou spřežku, kterou nelze zaměnit za *pro co*. Tvary označené otazníkem jsou možné zřejmě jen hypoteticky.

Typ n

Výčet všech složenin:

1. *zaň* = *za něho*, lemma: {*za*, *on*}
2. *naň* = *na něho*, lemma: {*na*, *on*}
3. *oň* = *o něho*, lemma: {*o*, *on*}
4. *proň* = *pro něho*, lemma: {*pro*, *on*}
5. *doň* = *do něho*, lemma: {*do*, *on*}
6. *?veň* = *v něho*, lemma: {*v*, *on*}

Neobvyklý archaický tvar *veň* jsme skutečně našli v Ottově slovníku naučném:

... na krátko byv zbaven svého úřadu v Irsku a znova veň uveden jako (66)
tajemník...

6.3.2 Typ zájmenně-slovesný ... t

Tento typ je zastoupen pouze jediným tvarem: *toť*. Jak již bylo zmíněno v oddíle o poddruzích zájmen (4.1.2.3), je toto slovo homonymní s ukazovacím zájmenem a s částicí (viz tabulka 6.2 a příklady (67) až (70)). I jako složenina je homonymní. Může totiž zastupovat dvě různé dvojice složek, lišící se v kategorii Číslo, pokaždé však se stejným vícenásobným lemmatem.

Slovní tvar	Složky složeniny	Lemma	Slovní druh
<i>toť</i>	<i>to je</i>	{ <i>to</i> , <i>být</i> }	složenina
<i>toť</i>	<i>to jsou</i>	{ <i>to</i> , <i>být</i> }	složenina
<i>toť</i>	—	<i>toť</i>	zájmeno
<i>toť</i>	—	<i>toť</i>	částice

Tabulka 6.2: Slovní tvar *toť* a jeho možné interpretace

Příklady:

Kniha, toť kouzelná brána k dobrodružství... (složenina v jednotném čísle) (67)

... *co opravdu rád slyším, toť hlasy ptactva a noční vítr.* (složenina v množném čísle) (68)

Inu, inu, toť, toť... (částice) (69)

Toť se ví. (ukazovací zájmeno) (70)

Další příklady ((27) až (29)) byly uvedeny na str. 27.

6.3.3 Typ zkratkový ... Z

Složeninami jsou i zkratky víceslovných frází, např. *atd*, *atp*, *např.* Jejich lemmata však nerozebíráme na jednotlivé složky, neboť by to často vedlo ke složitým vícenásobným lemmatům. Lemma zkratkového typu složeniny je konkrétní slovní tvar.

Ze stejného důvodu neurčujeme ani ostatní morfologické kategorie.

Jediné relevantní morfologické kategorie složenin typu Z tedy jsou POS=S a CMP=Z.

6.3.4 Typy slovesné ... N, A, P, C, V, D, T, J, S

Slovesné složeniny může tvořit slovo téměř libovolného slovního druhu s přidaným zakončením *-s*, které zastupuje slovní tvar *jsi* (*převrátils, jehos*). Typ slovesné složeniny vyjadřujeme znakem, který kóduje příslušný slovní druh.

Slovesné složeniny se netvoří z předložek, citoslovcí, prefixových segmentů a zřejmě ani z cizích slov.

Slovesné složeniny nemůžeme zadat výčtem, jak jsme učinili u předchozích, zájmených typů. Jde o otevřenou třídu. Přesto můžeme slovesné složeniny rozdělit do několika typů, podle slovního druhu jejich první složky.

Přes vysokou produktivnost tvoření slovesných složenin však jejich výskyt není příliš častý. Do morfologického slovníku je nezařazujeme, rozpoznávají se pomocí guessru. Jejich rozpoznání je velmi snadné — odtržení koncového *-s* u neznámých slov ponechá rozpoznatelný slovní tvar. Pokud jsou splněny podmínky (nečetné) uvedené dále, jde o slovesnou složeninu.

Než vyjmenujeme a popíšeme jednotlivé typy, poznamenejme, že slovesné složeniny bez ohledu na typ obvykle nevznikají ze slov, u nichž by přidané *-s* činilo potíže s výslovností: **vlass, *pařezs*.

Pro rozpoznávání slovesných složenin v průběhu morfologické analýzy může být důležité ještě jedno zjištění, a to pořadí slovesné složeniny v klauzi. Slovní tvar *jsi*, který je implicitně ve slovesné složenině přítomen, má většinou funkci pomocného slovesa, a jako takový je příklonkou. Slovesná složenina tedy většinou stojí ve větě na jejím začátku, aby bylo splněno Wackernaglovo pravidlo.

Toto pravidlo však není absolutně spolehlivé, jak ukazují příklady (71) až (76) z korpusu SYN:

- To není tak samozřejmé, jak říkals.* (71)
někam jsem za Tebou prostě šla, kde docela určitě a úplně volně stáls. (72)
Seňko, Seňuško, ty jedinej mně zůstals. (73)
zpívej, jako jsi zpíval, ještě než zešedivěls. (74)
doznals, co nečinils (75)
Blesku, český Blesku, nevím, kolik lidí již svou září osvítils. (76)

Podle slovního druhu první složky složeniny rozeznáváme typy slovesných složenin. Tyto typy označujeme podle kódu pro příslušný slovní druh.

Typ N

První složku slovesné složeniny typu N tvoří podstatné jméno.

Příklady, které uvádíme pod čísla (77) až (81), jsou vymyšlené, neboť se tento typ příliš často nevyskytuje, a vzhledem k tomu, že dosud nebyl v korpusech značkován, není snadné ho cíleně vyhledat. Pouze příklad (82) je z korpusu SYN, ale byl nalezený víceméně náhodou:

- Bez oknas/okenas nemohl vidět ven.* (77)
Oknus/Oknúms přidělal okenice. (78)
To oknos/Ta oknas rozbil ty. (79)
O okněs/oknechs nemluvil. (80)

Oknems/Oknys viděl dobře. (81)

Z latinys měl reparát loni. (82)

Typ A

První složku slovesné složeniny typu A tvoří přídavné jméno.

Ani tento typ složenin se příliš často nevyskytuje, podařilo se nám však nalézt v korpusu SYN příklad s první složkou v 1. pádě (příklad (83)), z čehož zejména vyplývá, že v tomto případě nevystupuje implicitní *jsi* ve funkci pomocného slovesa:

Věrnýs jak kůň, jak býk všaks vášnivý. (83)

Příklad (84) ukazuje slovesnou složeninu typu A utvořenou ze jmenného tvaru přídavného jména:

Salome, podobnas úponku (z písně Karla Kryla) (84)

Typ P

První složku slovesné složeniny typu P tvoří zájmeno.

Příklady na slovesnou složeninu typu P již tak řídké nejsou. Např. tvarů [tT]ys se v korpusu SYN2000 vyskytuje 1068, převážná většina pochopitelně v beletristické části. Ve slovesných složeninách typu P se však vyskytují i jiná zájmena, jak ukazují příklady (85) až (87).

Copaks to musel řešit zrovna takhle? (85)

Tos řekl ty, já ne. (86)

Všechno, o čems mluvil... (87)

Typ C

První složku slovesné složeniny typu C tvoří číslovka.

Nejčastější případy jsou zřejmě číslovky tázací, ale je možné jsou i jiné číslovky, jak ukazuje příklad (90).

Koliks jich koupila? (88)

Kolikráts to viděl? (Internet) (89)

Pěts jich nemohl porazit. (vymyšleno) (90)

Typ V

První složku slovesné složeniny typu V tvoří sloveso.

Slovesné složeniny se slovesy vyžadují sloveso v příčestí minulém činném (VRB=L)¹, jednotném čísle (NUM=S) a ve 2. osobě (PER=2). Nelze tedy např. *kupujs, *kupuješs, *kupujs ani kupovalis.

¹Složeniny s pasivem jsou řazeny mezi složeniny typu A.

Implicitně přítomné sloveso *jsi* ve složenině je v tomto případě vždy pomocné, protože není možné, aby v jedné klauzi byla dvě finitní slovesa.

Příklady:

... *má milá ženo, bylas tak statečná...* (91)

Koupils ho Ireně. (92)

Typ D

První složku slovesné složeniny typu D tvoří příslovce.

Vytváření slovesných složenin typu D nemá zdá se žádná omezení, jak je vidět z příkladů (93) až (98):

Včeras měl narozeniny. (93)

Posledněs říkala, že... (94)

A určitěs to ztratila? (95)

A ještěs mě nikdy neodměnil. (96)

Nikdys nechtěla vařit. (97)

Jaks k tomu došla? (98)

Typ T

První složku slovesné složeniny typu T tvoří zvratná částice.

Tento typ slovesných složenin lze zadat výčtem. Jsou to tyto slovní tvary: *sis, ses*.

Mezi složeniny typu T by se mohly počítat i tvary *bych, bys, bychom, bysme, byste*, protože i u nich jde o spojení částice (*by*) a tvaru slovesa *být*. My je však řadíme mezi částice s jednoduchým lemmatem *by*.

Typ J

První složku slovesné složeniny typu J tvoří spojka.

Některé spojky slovesné složeniny patrně netvoří. Jsou to zejména **as, *is, ?ales, ?čis*. Z většiny ostatních spojek slovesné složeniny utvořit lze, např. *nebos, protos, nebofs, protožes, zdas*. Jak je vidět, nezáleží to na jejich souřadnosti nebo pořadnosti.

Příklady (99) až (101) pocházejí z korpusu SYN:

... kdyžs teda říkal... (99)

... neměls už čas se zase stejnou cestou vrátit, nebos na to zapomněl. (100)

Nevím, jestlís ho vůbec znal. (101)

Mezi složeniny typu J nepočítáme spojky *abys, kdybys*, ani jejich ostatní tvary *abych, kdybych, abychom, kdybychom, abyste, kdybyste, abysme, kdybysme*, i když bychom je za složeniny považovat mohli. Podobně jako kondicionálová částice *by*, i v tomto případě volíme tradiční řešení a všechny uvedené tvary řadíme mezi spojky, s lemmaty *aby, kdyby*.

Typ S

První složku slovesné složeniny typu S tvoří složenina. Vícenásobné lemma těchto složenin má tři prvky: dvouprvkové vícenásobné lemma první složeniny a lemma *být* slovního tvaru *jsi*.

Jde o tvary vytvořené ze zájmenných složenin, tedy např. *oňš, začs*. Že nejde o pouhý teoretický případ, dokazuje úryvek z textu písně Hany Zagorové:

... o čem sníl jsi ty, načs přísahal... (102)

Běžně se však tento typ složenin opravdu nevyskytuje.

6.4 Vyhledávání složenin v korpusech

V tomto oddíle se zamyslíme nad způsobem, jak zahrnout složky různých typů složenin do vyhledávacích dotazů.

Otázka souvisí s množinou relevantních kategorií pro jednotlivé typy složenin, viz 6.5 na str. 67.

Typy zájmenné a zájmenně-slovesné

U zájmenných typů jde o předložku s akuzativní rekcí a zájmeno v akuzativu. V množině G je kromě kategorie **Slovní druh** pouze kategorie **Pád**, který je ovšem pro obě složky shodný — akuzativ. Předložka další relevantní kategorie nemá, pro výslednou složeninu jsou tedy relevantní všechny kategorie, které jsou relevantní i pro příslušné zájmeno. Jsou to **Rod** (střední, mužský životný nebo neživotný pro typ n a střední pro typ c), **Číslo** (jednotné), **Osoba** (3), **Poddruh** (substantivní), **Funkce** (určitá pro typ n , tázací nebo vztažná pro typ c).

Složeniny tohoto typu by se měly, pokud si to uživatel přeje, zahrnout do výsledku na dotaz požadující všechny výskyty předložek. K tomu je ovšem třeba dosud jednoduchý dotaz $POS=R$ modifikovat, a to tak, aby se našly nejen předložky ($POS=R$), ale i složeniny ($POS=S$) typu n a c ($CMP=[nc]$). Výsledný dotaz lze zapsat např. takto (před šípkou je jednoduchý dotaz, za šípkou dotaz modifikovaný):

$$POS = R \rightarrow (POS = R) \vee (POS = S \wedge CMP = [nc])$$

Podobně musíme modifikovat i dotaz na zájmena. Sem je třeba navíc zahrnout zájmenně-slovesnou složeninu typu t a slovesnou složeninu typu P ($CMP=[nctP]$):

$$POS = P \rightarrow (POS = P) \vee (POS = S \wedge CMP = [nctP])$$

Dotazy na ostatní morfologické kategorie mohou zůstat beze změny, neboť jejich hodnoty se stávají hodnotami příslušných kategorií zájmenné složeniny.

Složenina *toť* má 1. složku *to*, s množinou relevantních kategorií

$$K_1 = \{POS, SUB, GEN, FCE, CAS, NUM\}.$$

6 Složeniny

a 2. složku *je/jsou*, s množinou relevantních kategorií

$$K_2 = \{\text{POS, SUB, PER, NUM, VRB, NEG}\},$$

Tabulka 6.3 ukazuje jejich hodnoty.

	<i>to</i>	<i>je/jsou</i>	<i>toť</i>
POS	P	V	S
SUB	D	0	D
FCE	U	—	U
CAS	1	—	1
GEN	N	—	N
NUM	S/P	S/P	S/P
PER	—	3	3
VRB	—	P	P
NEG	—	A	A

Tabulka 6.3: Hodnoty relevantních kategorií složek složeniny *toť*.

Tabulka 6.4 uvádí přehled hodnot relevantních kategorií pro zájmenné složeniny.

Složenina	SUB	FCE	GEN	NUM	CAS	FMU	Příklad
typ c	Z	T/V	N	S	4	—	<i>oč</i>
typ n	Z	U	M/I	S	4	n	<i>oň</i>

Tabulka 6.4: Zájmenné složeniny a jejich relevantní kategorie

Typy slovesné

Přidaná koncovka *-s*, zastupující slovní tvar *jsi* slovesa *být*, nese tyto hodnoty relevantních morfologických kategorií:

$$\begin{aligned} \text{PER} &= 2 & \text{SUB} &= \text{b}/0 \\ \text{NUM} &= \text{S} & \text{NEG} &= \text{A} \\ \text{VRB} &= \text{P} \end{aligned}$$

Všechny tyto kategorie mohou (ale nemusí) ležet v množině G .

V případě konfliktu hodnot kategorií PER, SUB, NUM, NEG bude hodnota výsledné složeniny rovna hodnotě náležející první složce složeniny, tedy např. pro tvar *židlemís* bude NUM=P, pro slovo *nevysokýs* bude NEG=N, pro slovo *měs* bude PER=1 a pro slovo *tomus* bude SUB=D.

Kategorie VRB leží v množině G pouze pro složeniny typu V a dále pak pro ty složeniny typu A, jejichž první složka je trpným rodem slovesa (např. *ukrytas*, *podobnas*), i když ta se vyskytují opravdu zřídka. V obou případech je hodnota této kategorie v konfliktu s hodnotou první složky, která je v prvním případě VRB=L, ve druhém VRB=T. Tyto hodnoty převádíme na hodnoty kategorie Slovesný tvar složeniny. Z typu složeniny je možné odvodit i hodnotu kategorie Slovesný tvar druhé složky (*jsi*) a vytvořit podle toho dotaz.

U ostatních slovesných složenin je VRB=UNDEF, i když tam není tato kategorie v konfliktu s kategoriemi první složky ($\text{VRB} \notin G$). Kdybychom položili

VRB=P podle hodnoty slovního tvaru *jsi*, nebylo by to konzistentní s předchozím rozhodnutím o hodnotě této kategorie u složenin typu V a A v pasivu. Slovesné složeniny zahrneme v případě potřeby do dotazu na všechna slovesa v přítomnosti jednoduše vyloučením neslovesných složenin ($CMP \neq [ncZ]$):

$$VRB = P \rightarrow (VRB = P) \vee (POS = S \wedge CMP \neq [ncZ])$$

Podobně řešíme i kategorii *Osoba*, která je sice konfliktní pouze pro složeniny typu $CMP=V$ a $CMP=P$, ale v zájmu konzistence ji u ostatních typů nepovažujeme za relevantní ($PER=UNDEF$). Následující dotaz vyhledá všechna slova s hodnotou kategorie $PER=2$ a slovesné složeniny kromě zájmenných (typ n a c), zájmenně slovesných (typ t) a zkratkových (typ Z).

$$PER = 2 \rightarrow (PER = 2) \vee (POS = S \wedge CMP \neq [ncZt])$$

Ve všech slovesných složeninách kromě typu t je vždy druhá složka v jednotném čísle. Z dotazu na jednotné číslo tak musíme vyloučit složeniny typu t, kde může být druhá složka v čísle množném, jednotné číslo je však vždy pokryto přímo hodnotou kategorie $NUM=S$:

$$NUM = S \rightarrow (NUM = S) \vee (POS = S \wedge CMP \neq [Zt])$$

Podobný dotaz lze vytvořit i pro negaci. Slovní tvar *jsi* je ve slovesných složeninách vždy s hodnotou $NEG=A$, dokonce i když je první složka negativní (*nebyls*). Je tedy otázkou, zda vůbec má dotaz na $NEG=A$ v takovém případě smysl. Pro úplnost ho ale uvádíme:

$$NEG = A \rightarrow (NEG = A) \vee (POS = S \wedge CMP \neq [ncZ])$$

Dotaz na poddruh pomocné sloveso je komplikovanější. Ve většině případů je *jsi* ve slovesné složenině pomocné, ale není tomu tak vždy. Výjimkou mohou (ale nemusí) být složeniny se jménem v 1. pádě, viz příklady (103), (104) a také (83) a (84).

Drahoušku, ale tys moje žena! (103)

tys nejen blázen, ale ke všemu ještě pitomec! (104)

Ve větě *Tys blázen* má implicitní *jsi* $SUB=0$, zatímco ve větě *Tys byl blázen* má $SUB=b$. Hodnota celé složeniny *tys* je $SUB=Z$ (substantivní zájmeno). V dotazu je tedy třeba vyloučit slovní tvary v 1. pádě. V případě, že v takové složenině je implicitní *jsi* ve funkci slovesa pomocného, je tato alternativa pokryta první možností v disjunkci nového dotazu. Vyhledají se jen ty případy, kdy jde skutečně o pomocné sloveso².

Kromě toho je třeba v dotazu vyloučit složeniny neslovesné, tj. zájmenné, zájmenně-slovesné a zkratkové.

Výsledný dotaz tedy bude vypadat takto:

$$SUB = b \rightarrow (SUB = b) \vee (POS = S \wedge CMP \neq [ncZt] \wedge CAS \neq 1)$$

²Za předpokladu, že je kategorie *Poddruh* ve složenině správně určena.

Ve složenině *toť* typu t se sice vyskytuje tvar slovesa *být*, ale zde nikdy nevystupuje v roli slovesa pomocného.

Je třeba ještě doplnit dotazy na slovní druhy jednotlivých složek složeniny:

$$\text{POS} = x \rightarrow (\text{POS} = x) \vee (\text{POS} = S \wedge \text{CMP} = x),$$

kde x zastupuje [NADCJT], tedy jeden ze slovních druhů, které mohou být první složkou slovesné složeniny. Dotaz je přímočarý — hledáme slovní druh x a složeniny typu x, což jsou právě ty, jejichž první složkou je tvar s hodnotou POS=x.

Vynechali jsme slovní druh sloveso a zájmeno. Dotaz na zájmeno musí zahrnout i složeniny zájmenné a zájmenně-slovesné:

$$\text{POS} = P \rightarrow (\text{POS} = P) \vee (\text{POS} = S \wedge \text{CMP} = [\text{PncS}])$$

Sloveso je přítomno ve všech slovesných složeninách a ve složenině zájmenně-slovesné, proto bude v dotazu jednodušší vyloučit složeniny ostatní (neslovesné):

$$\text{POS} = V \rightarrow (\text{POS} = V) \vee (\text{POS} = S \wedge \text{CMP} \neq [\text{ncZ}])$$

Na první pohled mohou modifikované dotazy vypadat složitě. Většina korpusových vyhledávačů (manažerů) však umožňuje, aby si uživatel definoval a pojmenoval určité dotazy, aby je mohl neustále využívat bez složitého vytváření. Toto jsou případy, kdy by se takové definice hodily.

Kromě toho je třeba poznamenat, že dosud takové vyhledávky nebyly možné vůbec. Uživatel, který nebude chtít složeniny do svých dotazů zahrnout, může i nadále používat dotazy, na které je zvyklý.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	
POS	SUB	FCE	ASP	ABR	GEN	NUM	CAS	DUA	PER	DEG	NEG	VRB	NOM	CMP	INT	
S	Z	U			MI	S	4		3					n		
	Z	TV			N	S	4							c		
	D	U				SP	1							t		
				⊕											Z	
	⊕				⊕	⊕	⊕	⊕	⊕			⊕			N	
	⊕				⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	(T)	⊕	A	
	⊕	⊕	⊕		⊕	⊕	⊕	⊕	⊕	⊕					P	
	⊕	⊕	⊕		⊕	⊕	⊕	⊕	⊕						C	
	⊕						S			2		⊕	L		V	⊕
	⊕	⊕	⊕				S				⊕	⊕			D	
	⊕						S								T	
	⊕						S								J	
	Z	UTV			MIN	S	S			3					S	

Tabulka 6.5: Přehled relevantních kategorií pro jednotlivé typy složeniny. Relevantní kategorie jsou pro jednotlivé typy vyplněny buď kódem (případně více možnými kódy) svých hodnot, nebo znakem \oplus , který zastupuje libovolnou hodnotu z množiny hodnot své kategorie.

7 Morfologický slovník

7.1 Vztah morfologického slovníku a morfologických nástrojů

Morfologický slovník v ideálním případě obsahuje „všechny“ slovní tvary českého jazyka spolu s jejich morfologickým popisem. Slůvko všechny je v uvozovkách z pochopitelných důvodů — nikdy totiž nebudeme schopni vytvořit takový slovník, který by obsahoval úplnou slovní zásobu nějakého živého jazyka.

Existující morfologické slovníky češtiny, o kterých jsme se již zmiňovali, totiž pražský a brněnský, jsou však dostatečně velké na to, aby se pomocí nich mohly rozpoznat běžně užívané slovní tvary jak standardní, tak i obecné češtiny. Oba slovníky jsou neustále obohacovány o nové slovní tvary, které tam dosud nebyly ať už proto, že se na ně „zapomnělo“, nebo proto, že nově vznikly (neologismy).

Přidávání nových a nových slovních tvarů do slovníku však nemusí být vždy jen pozitivní, jak by se na první pohled mohlo zdát. Nebezpečné je především přidávání cizích vlastních jmen, která mohou být silně homonymní. Jedno slovo tak může znamenat např. cizí příjmení, ale současně i název firmy. V každém z těchto případů má jiný rod, a pokud přebírá českou flexi, skloňuje se pokaždé podle jiného vzoru. Od českého slovníku samozřejmě nelze očekávat, že bude obsahovat všechna cizí slova. V okamžiku, kdy se v textu vyskytne takové slovo v jiném významu, než bylo již zachyceno ve slovníku, dochází k chybě. Podle slovníku se totiž dané slovo rozpozná, ale špatně. V takových případech by paradoxně bylo lepší, kdyby se nerozpoznalo. Jiné nástroje, využívající pravidel slovtvorby, by si s jeho rozpoznáním poradily lépe. Pro další zpracování, zejména desambiguaci, je totiž většinou výhodnější hodnoty morfologických kategorií raději podspecifikovat, tedy přiřadit jim sdruženou hodnotu, než je určit špatně.

Při návrhu struktury a obsahu morfologického slovníku je třeba mít na paměti, jakým způsobem se slovník bude využívat. Jistě by bylo hezké, kdyby morfologický slovník obsahoval veškeré morfologické informace o všech slovních tvarech daného jazyka, a to nezávisle na programových nástrojích, které ho využívají.

Na druhou stranu však je výhodné některé jevy nepopisovat pomocí slovníku, ale nechat to na nástrojích morfologické analýzy nebo syntézy. Jde zejména o tvary negace. Až na několik výjimek u slovesných infinitivů se tvoří zcela pravidelně pomocí předpony *ne-*. Toto pravidlo je velmi jednoduché, téměř univerzální, a zahrnuje navíc obrovské množství slovních tvarů.

Zdá se tedy rozumné nevkládat pravidelně negované slovní tvary do slovníku, ale nechat jejich zpracování na nástrojích. Slovníkové heslo by mělo pouze obsahovat informaci o tom, zda lze z daného slovního tvaru utvořit negaci, či nikoli, ale ne samotný negovaný tvar. Stejně lze zacházet se stupňováním pří-

davných jmen a příslovcí. Negaci, stejně jako stupňování, musí umět rozpoznat analýza a vytvořit syntéza.

Tato symetrie morfologické analýzy a syntézy však není nutná. Např. cizí slova, o kterých byla řeč výše, stačí umět jen rozpoznat, a na jejich rozpoznání není morfologický slovník třeba. V naprosté většině případů jde totiž o podstatná jména. Jako podstatná jména lze analyzovat i součásti víceslovných cizojazyčných celků, např. názvů písniček, měst, různé slogany, a to i tehdy, když ve svém původním jazyce patří k jiným slovním druhům. Podle našeho názoru tedy z cizích slov do morfologického slovníku patří jen ta nejběžnější. Je však třeba umět v textu rozpoznat i ta ostatní a přiřadit jim správné hodnoty morfologických kategorií, třeba podspecifikované, tzn. se sdruženou hodnotou některých kategorií. K rozpoznání neznámých slov se používá guesser.

7.1.1 Guesser

Morfologická analýza nerozpozná úplně všechna slova v neznámém textu. Z experimentů s Českým národním korpusem vyplývá, že 2 až 3 procenta slov zůstávají nerozpoznána (viz (Hlaváčová, 2001)). Mnoho slov, která nejsou obsažena ve slovníku, však lze rozpoznat pomocí tzv. guessru. Guesser využívá ortografických a morfologických pravidel k analýze neznámých slov. Jeho výsledkem je odhad relevantních kategorií neznámého slova a jejich hodnot a také rekonstrukce pravděpodobného lemmatu.

Guesser tedy přiřazuje hodnoty morfologických kategorií neznámým slovním tvarům. V ideálním případě přiřadí jedno (správné) lemma a jednu (správnou) morfologickou značku. Většinou je však možností více. Guesser by však měl být navržen tak, aby jich nebylo příliš, protože to potom ztěžuje následná zpracování, především desambiguaci. Pochopitelně by však mezi výsledky měla být i správná dvojice lemmatu a morfologické značky.

V českém jazyce lze využít k sestavení guessru jednak slovních zakončení, tedy přípon a koncovek, a jednak předpon.

7.1.1.1 Prefixový guesser

Využití předpon pro práci guessru je přímočaré. Mějme nerozpoznané slovo z , které lze rozdělit na dvě části: $z = p \cdot w$, kde p je známá předpona a w slovní tvar rozpoznatelný na základě morfologického slovníku. Je-li $\lambda(w) = \bar{w}$, potom $\lambda(z) = p \cdot \bar{w}$. Navíc obě slova, w i z , mají stejné relevantní morfologické kategorie se stejnými hodnotami, tedy:

Je-li $\mu(w) = \langle \bar{w}, m \rangle$, potom $\mu(z) = \langle p \cdot \bar{w}, m \rangle$, kde m je morfologická značka.

Vyjádřeno slovy — lemma neznámého slovního tvaru, který se skládá ze známé předpony a známého slovního tvaru, lze zrekonstruovat přidáním předpony k lemmatu známého slovního tvaru. Hodnoty morfologických kategorií jsou stejné (zde mohou nastat výjimky ve změně vidu u sloves).

Příklad:

Předpokládejme, že slovní tvar $z = \textit{eurooken}$ není obsažen v morfologickém slovníku. Můžeme ho však rozdělit na předponu $p = \textit{euro-}$ a slovní tvar $w = \textit{oken}$, který ve slovníku je a my ho umíme rozpoznat. Budeme-li mít \textit{euro} v seznamu

předpon, můžeme původní tvar *eurooken* lemmatizovat jako *eurookno* a přiřadit mu tyto hodnoty: POS=N, GEN=N, CAS=2, NUM=P, které jsou stejné jako hodnoty známého slovního tvaru *oken*.

Seznam předpon jsme vytvořili pomocí statistických metod ze slovních tvarů korpusu SYN2000. Popis metod i jejich výsledky lze nalézt v příspěvcích (Hlaváčová – Hrušecký, 2008) a (Urrea – Hlaváčová, 2005).

U každé předpony je možno uvést, s jakým slovním druhem se může kombinovat. Např. uvedená předpona *euro-* se nemůže připojit ke slovesu, ale jen k podstatnému, případně ještě přídavnému jménu. Většinu předpon, které se pojí se jmény, by možná někteří lingvisté nepovažovali za klasické předpony, ale spíše za složky slova vzniklého skládáním. Pro naše účely je důležité, že lze pomocí jejich odtržení velmi spolehlivě analyzovat neznámá slova.

Seznam předpon získaný automaticky pomocí nástrojů zmíněných výše jsme použili při nové implementaci morfologické analýzy, viz (Hlaváčová – Kolovratník, 2008), která pracuje s dosavadním pražským morfologickým slovníkem.

7.1.1.2 Postfixový guesser

Postfixový guesser má k dispozici seznam zakončení slovních tvarů s množinou možných morfologických značek a pravidel na vytvoření k nim příslušejících lemmat. Jestliže neznámé slovo končí jedním nebo více z těchto zakončení, guesser mu přiřadí příslušnou množinu morfologických značek a podle pravidel odvodí pravděpodobné tvary lemmat.

Některá slovní zakončení určují své morfologické kategorie jednoznačně, jiná nabízejí možností mnoho. Záleží zejména na délce zakončení, podle kterého se snažíme odhad provést. Při implementaci guessru je tedy třeba se rozhodnout, s jak dlouhými zakončeními chceme pracovat. Čím delší zakončení, tím méně možností guesser nabídne, ale tím více jich je zapotřebí. Krátkých zakončení není třeba tolik, jsou však většinou neúnosně mnohoznačná. Je třeba zvolit vhodný kompromis. Guesser použitý pro korpus SYN2000 (viz (Hlaváčová, 2001)) pracoval s postfixy délky 4, která vycházela z výsledků ještě starší práce na projektu MOZAIKA (viz (Kirschner, 1983))¹.

¹Uvedený postfixový guesser by bylo možno rozšířit následujícím způsobem i na běžnější slovní tvary, ale v tom případě by bylo třeba zajistit i opačný postup jejich možného generování, tedy by to vlastně přestal být guesser.

Existuje např. řada známých zakončení, která nemusí splňovat požadavky na automatické vyhledávání postfixů, konkrétně na jejich délku. Je vhodné seznam známých postfixů přidat k seznamu vytvořenému automaticky. Máme teď na mysli zejména slova s číselnými předponami, které se sice rozpoznají pomocí prefixového guessru (viz 7.1.1.1), ale po jejich odtržení nevznikne analyzovatelný slovní tvar, např. *letý, nohý, hlavý, dílný, místný, ramenný* po odtržení předpon *dvou-, čtyř-, sedmi-,...* Není rozumné ani proveditelné zahrnout všechny možné kombinace těchto postfixů i prefixů do slovníku, zvláště u prefixů číslovkových. Seznam takových postfixů však tak rozsáhlý není a mohl by posloužit podobně jako už dnes slouží guesser prefixový. Navíc tyto ruční postfixy většinou „zaberou“ stoprocentně, tedy určují morfologické kategorie i lemma daného slovního tvaru jednoznačně.

Dále do tohoto seznamu můžeme přiřadit postfixy *-koli, -si, -pak* a další, pomocí kterých vznikají zájmena neurčitá odvozená od zájmen tázacích (*kdosi, jakýsi, čísi, kdokoli, jakýkoli, číkoli, kdopak, jakýpak, čípak,...*) a další. Jejich výčet je uveden v obrázku 4.1 a v tabulce 4.1.

Jednotlivé slovní tvary z paradigmatu těchto zájmen jsou zatím vyjmenovány jednotlivě v pražském morfologickém slovníku, ale jejich zahrnutí do takového postfixového guessru by bylo přehlednější a elegantnější. Zahrneme-li navíc požadavek na slovní druh, případně i poddruh řetězce vzniklého po odtržení postfixu, byly by odhady guessru i jednoznačné.

Např. postfix *-si* je příliš krátký na to, abychom mohli spoléhat na jednoznačnost morfolo-

Např. zakončení *-kyní* patří téměř na 100 % podstatnému jménu rodu ženského v 7sg nebo 2pl a s lemmatem *-kyně*. S touto znalostí můžeme úspěšně analyzovat slovní tvary, které neobsahuje slovník, viz příklady (105) až (108) z korpusu SYN. Příklad (108) obsahuje slovní tvar s překlepem, přesto je možno ho v textu správně rozpoznat.

byla ve výuce dospělých ostřílenou borkyní (105)

*Falstaff zatroubil na ústup před tou čarodějnící, před tou ďáblovou
spřeženkyní* (106)

s jedinou nadšenkyní, která se k nim přidala (107)

v rozhovoru s nejvyšší státní zastupkyní (108)

Zkrátíme-li postfix na *-yní*, stále ještě můžeme s velkou pravděpodobností hádat na stejný typ podstatného jména, ovšem vyloučíme-li velmi časté slovo *nyní*.

Zkrácení postfixu jen na *-ní* už nám příliš nepomůže, neboť slov s tímto zakončením je mnoho — nejčastější jsou podstatná jména rodu středního (vzor *stavení*) s velkou homonymií ve většině kombinací čísla a pádu, přídavná jména měkká (*jarní*) s homonymií ještě větší, přídavná jména tvrdá v 1pl rodu mužského životného (*úspěšní muži*), slovesa 4. třídy (*špiní*).

7.2 Struktura slovníku

Morfologický slovník by měl v ideálním případě obsahovat právě všechny slovní tvary daného jazyka, opatřené potřebnou morfologickou informací, tedy lemmatem a hodnotami všech relevantních morfologických kategorií vyjádřených pomocí morfologické značky a mutace.

Celý slovník si můžeme velmi zjednodušeně pro začátek představit jako tabulku s právě vyjmenovanými údaji, viz příklad v tabulce 7.1.

Slovní tvar	Lemma	Morf. značka	Mutace
okny	{ <i>okno, vokno</i> }	N-----NNP7-----A-----	-
oknama	{ <i>okno, vokno</i> }	N-----NNP7-----A-----	Fa
voknama	{ <i>okno, vokno</i> }	N-----NNP7-----A-----	FaGv

Tabulka 7.1: Příklad několika položek jednoduchého modelu morfologického slovníku

Značení hodnot kategorie *Mutace* je popsáno v oddíle 3.4.

Právě naznačený jednoduchý model slovníku jako tabulky však nezachycuje vztahy mezi jednotlivými lemmaty, které bychom rádi také ve slovníku měli. Jde nám především o vztahy derivační, které se ukazují být klíčovými pro praktické využití morfologického slovníku při automatických překladech. Jednotlivá paradigmatata proto sdružujeme do záznamů.

Základní jednotkou morfologického slovníku je tedy záznam reprezentovaný lemmatem. Lemma může být vícenásobné. Každý záznam obsahuje informace

gických vlastností slovních tvarů s tímto zakončením. Přidáme-li však požadavek, že řetězec, vzniklý po odtržení tohoto postfixu, je nejen rozpoznatelný morfologickým analyzátořem, ale je to navíc zájmeno, je jeho analýza již jednoznačná — v tom smyslu, že má všechny kategorie až na FCE=N stejné jako základní zájmeno.

o jednom paradigmatu. Toto paradigma je tímto záznamem zcela a jednoznačně morfologicky popsáno.

Vztahy mezi lemmaty jsou vyjádřeny pomocí odkazů. My zde popíšeme odkazy vyjadřující derivační vztahy, ale obecně je možno definovat i jiné typy odkazů, např. synonymické, hypo- a hyperonymické, odkazy mezi vidovými dvojicemi. Teoreticky si může uživatel slovníku přidat libovolné typy odkazů, podle aplikace, ke které slovník využívá.

Reprezentantem paradigmatu je lemma, které slouží jako slovníkové heslo. Kdyby neexistovala v češtině homonymní lemmata, mohli bychom každé lemma považovat za klíč, podle kterého se ve slovníku vyhledává. Kvůli nejednoznačnosti to bohužel nejde. Homonymní lemmata rozlišujeme přidáním přirozených čísel. Klíčem záznamu v morfologickém slovníku je tedy buď lemma, jestliže k němu neexistuje homonymní lemma, nebo lemma s číselným sufixem, který rozlišuje homonymní lemmata (např. *žít-1* — *žiji*, *žít-2* — *žnu*).

Problém polysémie (víceznačnosti) v morfologickém slovníku neřešíme. Máme-li tedy dvě lemmata s odlišným významem, ale shodnými paradigmaty, budou zastoupena jediným záznamem (např. *kolej*).²

Rozšířené paradigma globálních mutací je reprezentováno vícenásobným lemmatem, které jsme definovali jako množinu. Množina nemůže být klíčem záznamu, proto ji vyjadřujeme jako jeden řetězec pomocí regulárních výrazů. Vícenásobné lemma {*okno*, *okno*} z předchozího příkladu se pomocí regulárního výrazu zapíše jako *v?okno*.

Slovníkový záznam morfologického slovníku tedy obsahuje celé paradigma jednoho lemmatu, společně s morfologickými značkami jednotlivých slovních tvarů a označením typu jejich mutací. Tabulka 7.2 ukazuje příklad jednoho slovníkového záznamu.

Lemma	Slovní tvar	Morf. značka	Mutace	
město	<i>město</i>	N-----NNS1-----A----		
		N-----NNS4-----A----		
		N-----NNS5-----A----		
	<i>města</i>	N-----NNS2-----A----		
		N-----NNP1-----A----		
		N-----NNP4-----A----		
	městu	městu	N-----NNS3-----A----	
			N-----NNS6-----A----	Fu
		<i>městě</i>	N-----NNS6-----A----	Fe
		<i>městem</i>	N-----NNS7-----A----	
		<i>měst</i>	N-----NNP2-----A----	
		<i>městům</i>	N-----NNP3-----A----	
		<i>městech</i>	N-----NNP6-----A----	
		<i>městý</i>	N-----NNP7-----A----	
<i>městama</i>		N-----NNP7-----A----	Fa	

Tabulka 7.2: Příklad jednoho záznamu morfologického slovníku

Vzhledem k velké pravidelnosti českého gramatického systému je výhodné

²Morfologický slovník by se mohl v budoucnu rozšířit o další elementy popisu, včetně popisu sémantického.

do slovníku vkládat nikoli jednotlivé slovní tvary, ale využít pravidla, podle kterých se slovní tvary dají vytvořit, tedy vzory. Vzory výrazně zmenší velikost slovníku. Tabulka 7.3 ukazuje příklad z tabulky 7.2 zapsaný pomocí vzoru. Veškeré tvary, značky i kódy mutací jsou součástí vzoru.

Lemma	Kofix	Vzor
<i>město</i>	<i>měst</i>	mt

Tabulka 7.3: Příklad jednoho záznamu morfologického slovníku zapsaného pomocí vzoru

O vzorech obecně pojednáváme v kapitole 8, kapitoly 9 až 13 se věnují vzorům jednotlivých slovních druhů.

Kromě popisu paradigmatu obsahuje slovníkový záznam ještě derivační odkazy na jiné slovníkové záznamy, které popisují paradigmatu příbuzných lemmat. Odkazy si můžeme představit jako šipky mezi jednotlivými záznamy. Zásadně používáme šipky oboustranné. To znamená, že jestliže vede odkaz od záznamu A k záznamu B, vede odkaz i od záznamu B k záznamu A. Z tohoto důvodu není označení derivační odkaz zcela přesné, protože nepopisuje vždy skutečné odvození jednoho lemmatu z druhého. Kvůli nejružnějším aplikacím je však výhodné mít odkazy obousměrné. Vyhýbáme se tím občas i obtížnému rozhodování, které lemma je původní a které odvozené. Z hlediska aplikací je tato otázka podružná. Lemmata záznamů propojených derivačními odkazy ale raději nazýváme příbuzná lemmata.

Bylo třeba rozhodnout, jak derivační odkazy propojovat v případech, kdy je jich více. Zvolili jsme „hvězdicové uspořádání“. Znamená to, že jeden záznam je zvolen za základní a od něj vedou derivační odkazy k lemmatům příbuzným. Příklady jsou uvedeny na obrázcích 7.1 a 7.2. Na obrázcích jsou zachyceny jen ty derivace, které lze odvodit ze vzoru — pro zobrazené záznamy ze slovesných vzorů (viz kap. 12). Ostatní příbuzná slova, která do „hvězdy“ také patří (k obr. 7.1 by to bylo např. lemma *skok*), je třeba zadat do slovníku „ručně“. Rovněž ručně se musí propojit vidové dvojice, např. lemmata *skočit* a *skákat* z uvedených obrázků. Samozřejmě lze pomocí odkazů propojit i slova odvozená pomocí derivací prefixových.

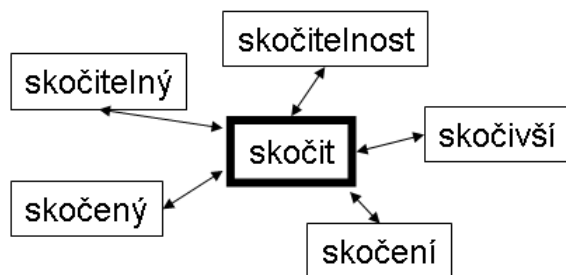
V zásadě nezáleží na tom, které lemma zvolíme za základní, důležitá je existence derivačních odkazů, které je propojují. Přesto jsme se pokusili o systematické řešení:

Základním lemmatem (středem hvězdy) je v případě slovesných derivací sloveso. V případě, že je mezi příbuznými lemmaty více sloves, je základním sloveso nejkratší. Tímto požadavkem upozadujeme slovesa iterativní.

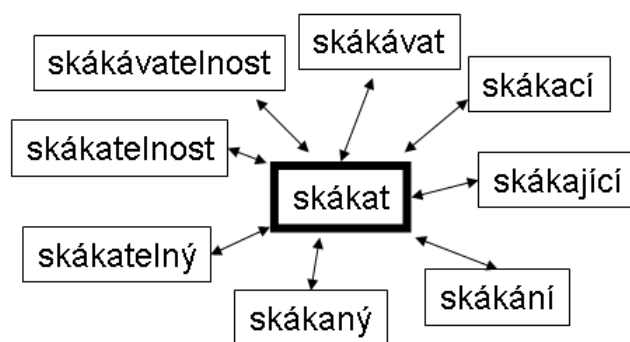
Není-li mezi příbuznými lemmaty sloveso, je základním lemmatem podstatné jméno. Není-li tam ani podstatné jméno, je základním lemmatem přídavné jméno.

Naskýtá se otázka, co všechno považovat za příbuzná lemmata. V prvním plánu jsme mezi příbuzná lemmata zahrnuli jen ta lemmata, která lze odvodit pravidelně pomocí přípon. Uvědomujeme si, že toto řešení je alibistické a má mnohá úskalí. Při implementaci konkrétního slovníku však nic nebrání propojení více lemmat. V rámci navrženého formátu je možné dokonce vytvářet i jiné typy odkazů. Důležitým typem odkazů bude např. propojení vidových

dvojice sloves.



Obrázek 7.1: Ukázka derivačních odkazů odvozených pravidelně ze základního slovesa *skočit*



Obrázek 7.2: Ukázka derivačních odkazů odvozených pravidelně ze základního slovesa *skákat*

Jelikož je většina derivačních vztahů velmi pravidelná, přidáváme pravidla pro vytvoření odvozených lemmat přímo do jednotlivých záznamů jako součást vzoru. Slovníku s tímto typem informace o derivacích budeme říkat kompaktní slovník. Pro uchování kompaktního slovníku použijeme formátu PML, který byl vytvořen jako jednotný datový formát pro ukládání lingvistických dat.³

Pro automatické použití kompaktního slovníku je třeba ho nejprve „rozgenerovat“ podle derivačních pravidel, čímž se vytvoří záznamy pro odvozená lemmata a oboustranné odkazy mezi nimi. Tyto odkazy se zakódují tak, aby je bylo možno využít v aplikacích, ke kterým se slovník používá. Teprve v další fázi je možno vytvářet jednotlivé slovní tvary.

³http://ufal.mff.cuni.cz/jazz/PML/doc/pml_doc.html

8 Vzory

V první části práce jsme se zabývali vymezením kategorií, které charakterizují české slovní tvary, a jejich hodnotám. Toto vymezení je postačující pro vytvoření morfologického slovníku, který obsahuje všechny slovní tvary spolu s jejich popisem pomocí uvedených kategorií. Vzhledem k pravidelnostem v českém gramatickém systému však je výhodné neukládat do slovníku jednotlivé slovní tvary každý zvlášť, ale zprostředkovaně, ve formě vzorů, které umožňují podstatně zmenšit objem slovníku. Vzory představují jakási pravidla, jak vytvářet a rozpoznávat slovní tvary jazyka. Potřebujeme je pouze kvůli popisu paradigmát.

Systém vzorů pro ohýbání českých slov se učí děti na základní škole. Pro automatické zpracování jazyka je však třeba systém školních vzorů zjemnit. Stejně jako máme v současné době dva hlavní systémy morfologických značek, tak i systém vzorů byl pro češtinu vytvořen dvakrát, v Praze a v Brně.

V následujícím oddíle stručně popíšeme oba zmíněné systémy vzorů a vysvětlíme, proč jsme se rozhodli pro vytvoření systému nového. V dalších oddílech potom podrobně popíšeme nový systém ohýbacích vzorů.

8.1 Stručné porovnání pražského a brněnského systému vzorů

Oba systémy, tedy pražský i brněnský, vycházejí ze základních vzorů, které se učí děti na základní škole. Vzhledem ke způsobu, jakým se vzory pracuje, bylo třeba v obou systémech repertoár vzorů podstatně rozšířit.

Klára Osolsobě ve své disertační práci (Osolsobě, 1996), na jejímž základě je brněnský systém postaven, o vzorech píše: „Pod pojmem *vzor* rozumíme ... konkrétní slovo reprezentující množinu všech slov, která tvoří ohebné tvary pomocí identického inventáře koncovek, jejichž společným rysem dále je, že tvoří paradigmaticky odvoditelné tvary podle příslušného slovního druhu, a u kterých dochází ke stejným změnám finální skupiny kmene.“

Pražský systém (viz (Hajič, 2004)) chápe vzor obdobně, až na to, že místo konkrétního slova jako reprezentanta množiny slov se stejnými koncovkami definuje abstraktní množiny „koncovek“. Poslední slovo je v uvozovkách proto, že se ve skutečnosti nejedná o klasické koncovky, ale obecnější řetězce. O nich pohovoříme za chvíli. Pražský vzor tedy není konkrétní slovo, ale řetězec znaků, který do jisté míry, zejména pro podstatná jména, využívá zkratk klasických vzorů (např. *hd* pro *hrad*, *zn* pro *žena*, *kr* pro *kuře*). Jejich jemnější členění se ale vyjadřuje pomocí dalších číslic i písmen (*hd1*, *hd1x*, *zn15e* atd.). Vzory pro ohýbání dalších slovních druhů vykazují také prvky logického systému, avšak již nepřipomínají klasické školní vzory.

Odlíšné je řešení výjimek. Brněnský systém pokrývá vzory celou slovní zásobu, tedy i výjimky, které jsou jedinečné. Podle takového vzoru se potom

ohýbá právě jedno lemma, které je v tom případě shodné s názvem vzoru. Pražský systém řeší výjimky jinak — doplněním výjimečných tvarů s jejich konkrétními značkami ke vzorům, nebo, v případě, že nelze použít vzor ani pro část paradigmatu, uvedením všech slovních tvarů s konkrétními morfologickými značkami.

V principu nezáleží na tom, jak se vzor pojmenuje, ani jakým způsobem se zachází s výjimkami. Důležité je, aby systém pokryl veškeré možnosti, které se v jazyce vyskytnou. Oba systémy tuto podmínku splňují, ne však beze zbytku. Za hlavní nevýhody považujeme:

- Brněnské názvy vzorů nemají systém (není to kritika vzorů, ale jejich názvů!). Ztěžuje to práci se slovníkem. Ten, kdo slovník obhospodařuje, se musí naučit všechny vzory nazpaměť.
- Pražské vzory systém mají, ale zdaleka ne dokonalý. I zde je obtížné udržovat přehled o názvech vzorů.
- Ani jeden ze systémů nepokrývá veškeré flektivní varianty (mutace).
- Ve vzorech se kódují dohromady globální i flektivní kategorie.
- Není dobře vyřešeno odvozování příbuzných lemmat. Pražský systém se derivacemi zabývá, ale silně přegenerovává, takže lze odvodit i slovní tvary, které v jazyce nejsou.

Proto jsme se rozhodli navrhnout nový systém vzorů. Netvrdíme, že je zcela dokonalý, ale uvedené nedostatky odstraní. Dalším důvodem je změna systému kategorií, včetně jejich hodnot, popsána v předchozích kapitolách. Právě s nimi by morfologie měla pracovat.

8.2 Nové vzory

Při návrhu vzorů vycházíme z pražského systému. Ten není přímo založen na morfologické stavbě slova, ani na fonologických pravidlech. Jeho hlavní myšlenkou je práce s řetězci¹.

Hlavní rozdíl našeho pojetí vzorů oproti oběma zmiňovaným systémům je ten, že kóduje pouze flektivní kategorie. Globální kategorie jsou totiž stejné pro celé paradigma (tak byly definovány), není tedy třeba je kódovat pro každý slovní tvar zvlášť. Můžeme tak použít stejný vzor pro více slovních druhů, např. adjektivní vzory pro podstatná jména, číslovky a zájmena.

Další změna spočívá ve striktním oddělení flektivní části vzoru a části derivační. Toto oddělení se sice v pražském systému používá, je však pro každý vzor neměnné. Každý vzor má napevno přiřazenu množinu derivací, což na jednu stranu velmi zjednodušuje popis², na druhou stranu však vede k vytváření neexistujících lemmat (např. od slovesa *pospíchat* se pravidelnou derivací utvoří příslovce *pospíchaně*.)

¹Lingvisticky přijatelnější popis nabízí např. dvouúrovňová morfologie, kterou navrhl Kimmo Koskeniemi v r. 1983 (Koskeniemi, 1983). Pro češtinu implementovala dvouúrovňový popis Hana Skoumalová (Skoumalová, 1997).

²Autor tohoto systému Jan Hajič měl dobré důvody k tomu, aby byl popis co nejjednodušší, protože ve své době potřeboval, aby zabíral v paměti počítače co nejméně místa. Jeho kritici si toto zhusta neuvědomují.

Vzory, které budeme popisovat, se tedy skládají ze dvou částí — flektivní a derivační. Přesto, že obě části zapisujeme do jednoho řetězce, budeme jim říkat flektivní vzor a derivační vzor. **Flektivní vzor** popisuje slovní tvary, **derivační vzor** kóduje pravidla, podle kterých se tvoří nová lemmata.

8.2.1 Flektivní vzory

Flektivní vzor je množina řetězců a k nim náležejících flektivních částí morfologických značek (nadále budeme psát zkráceně jen o morfologické značce) a flektivních mutací, jde tedy o množinu trojic ⟨řetězec, morfologická značka, flektivní mutace⟩. Připojí-li se řetězec z flektivního vzoru ke správnému začátku slova (jinému řetězci), vyjde platné české slovo (slovní tvar), jemuž náleží morfologická značka a flektivní mutace z druhých dvou členů příslušné trojice vzoru. Zmíněným řetězcům ze vzoru říkáme **zakončení**. Obecně nemá zakončení žádný konkrétní lingvistický význam. Může to být koncovka, ale nemusí.

Formálně:

Flektivní vzor Ω je množina trojic ⟨ s , M , F ⟩, kde s je řetězec, M platná morfologická značka a F flektivní mutace. Nulová flektivní mutace se může z trojice vypustit.

Řekneme, že slovní tvar w s morfologickou značkou M a flektivní mutací F byl vytvořen podle flektivního vzoru Ω , jestliže existuje trojice ⟨ s , M , F ⟩ $\in \Omega$ taková, že $w = p \cdot s$ pro nějaký řetězec p . Tečka \cdot je zde i dále znak pro operaci konkatenace (zřetězení).

Jestliže existuje řetězec p a flektivní vzor $\Omega = \{\langle s_i, M_i, F_i \rangle; i = 1 \dots n\}$, pro který množina $\{p \cdot s_i; i = 1 \dots n\}$ tvoří celé paradigma nějakého lemmatu \bar{w} , potom říkáme, že lemma \bar{w} se ohýbá (skloňuje, časuje nebo stupňuje) podle vzoru Ω . Flektivní vzor však nemusí popisovat celé paradigma, ale jen jeho podmnožinu.

Řetězec p nemusí mít žádný gramatický význam, i když ho často má. Může být kmenem, jeho částí, může obsahovat i celou nebo jen část přípony, a samozřejmě může obsahovat i předpony. Z toho důvodu pro něj nemáme žádné přesné lingvistické pojmenování. V dalším textu ho budeme nazývat kofix. Tento termín vymyslel student MFF UK David Kolovratník jako zkratkové slovo pro nejasné spojení řetězců tvořících kmen, prefix i sufix nebo jejich částí, které dohromady tvoří začátek slovního tvaru.

Kofix je nejdelší počáteční řetězec, který sdílí všechny slovní tvary popsané jedním vzorem.

Stejně tak řetězce s_i , které vystupují v nějakém vzoru, nejsou (gramatické) koncovky. Právě proto jim říkáme vágně **zakončení**.

Pro skloňování stačí, až na několik výjimek, jeden flektivní vzor. Například paradigma lemmatu *hrádek* má kofix *hrád* a množinu zakončení $\{ek, ku, kem, ky, ků, kům, kách, cích\}$.

U slovesných flektivních vzorů je situace komplikovanější. Slovesné tvary jednoho lemmatu se často tvoří podle několika různých vzorů, neboť v nich častěji dochází k hláskovým změnám. V takovém případě má paradigma více kofixů i více množin zakončení. Kdybychom chtěli zachovat zásadu, že k jednomu lemmatu přísluší jediný flektivní vzor, museli bychom pracovat s velmi krátkými kofixy a počet vzorů by musel být velký. Proto k vytvoření vzorů využíváme

společných vlastností jednotlivých slovesných tvarů a hláskové změny řešíme pomocí vícera kofixů. Jedno slovesné paradigma se potom časuje podle několika flektivních vzorů s několika kofixy.

Flektivní vzory jsou navrženy tak, aby pokrývaly i systematické flektivní mutace. Ty jsou uvedeny jako třetí člen trojic ve vzoru.

V této práci navrhuje nový systém flektivních vzorů, ale základní myšlenka pro vytváření slovních tvarů pomocí vzorů a kofixů zůstává stejná jako v pražském systému. Jen místo dvojic ⟨kofix, vzor⟩, které jsou přiřazeny lemmatům v současném systému, používáme trojice ⟨kofix, vzor, mutace⟩. Mutace z této trojice může být flektivní, globální i smíšená. Jde o hodnotu kategorie *Mutace*, která byla zavedena v oddíle 3.4.

Lemmata, která se ohýbají nepravidelně, nemají flektivní vzor. Místo toho záznam slovníku obsahuje všechny slovní tvary jejich paradigmatu s příslušnými morfologickými značkami. Alternativně by se daly nepravidelnosti řešit pomocí individuálních vzorů, které by platily vždy jen pro jediné lemma. Tento způsob je použit v systému brněnském. My se držíme praxe pražské, která výjimky popisuje pomocí jednotlivých slovních tvarů. Vyhneme se tak příliš vysokému počtu flektivních vzorů. Navíc je možno jednoduše přidávat nepravidelné tvary bez ohledu na vzory. Např. lemma *otevřít* má kromě pravidelného časování dva archaické tvary *otevru* a *otevrou*, které jsou flektivními mutacemi k tvarům *otevřu* a *otevřou*. Podobné mutace vykazují ještě lemmata *zavřít* a *uzavřít*. Tyto mutace nejsou systematické, nepopisujeme je tedy pomocí vzoru, nýbrž je coby výjimky přidáváme do paradigmatu zvlášť.

Paradigma popisujeme tedy dvojím způsobem:

1. pomocí kofixů a flektivního vzorů,
2. pomocí slovních tvarů a jejich morfologických značek.

První způsob lze chápat jako zkrácený zápis způsobu druhého.

Oba typy mohou být v záznamu přítomny současně. Je-li třeba, musí být navíc doplněny informací o mutaci, a to jak globální, tak flektivní, aby vždy platilo Zlaté pravidlo morfologie.

Součástí slovníkového záznamu jsou tedy dva typy trojic:

1. ⟨kofix, flektivní vzor, mutace⟩ nebo
2. ⟨slovní tvar, morfologická značka, mutace⟩.

V případě, že mutace je pro daný kofix nebo slovní tvar nulová, nemusí se uvádět. Množina takových trojic definuje celé paradigma lemmatu, které je klíčem daného záznamu. Slovníkový záznam může obsahovat jednu nebo více trojic jednoho nebo obou typů.

Všechny takové trojice přiřazené jednomu lemmatu tvoří předpis pro vytvoření celého paradigmatu tohoto lemmatu.

Flektivní vzory, které popíšeme v následujících oddílech, jsou rozděleny podle slovních druhů. Místo velkého množství vzorů, které popisují slovní zásobu, jsme zavedli pro každý slovní druh vzorů jen několik, zato parametrizovatelných. Potřebné množství vzorů nahrazujeme parametry, které jsou však pro každý vzor jiné. Jsme přesvědčeni, že tento systém flektivních vzorů je vhodnější než dosavadní vzory obou zmiňovaných systémů. Jeho hlavní předností je to, že lze parametry snadno kombinovat tak, aby popsaly různé možnosti

flexe vyskytující se v českých paradigmatech. Parametry jsou zejména vhodné pro popis flektivních mutací slovních tvarů, což je zřejmě největší slabinou současných systémů.

Stávající vzory, jak brněnské, tak i pražské, lze většinou pomocí vhodné zvolených parametrů převést na vzory nové. U pražského systému se tímto způsobem můžeme zbavit poměrně velkého množství výjimek, které doplňují paradigmata některých lemmat, protože stávající množina vzorů je nezahrnuje. Nelze převést některé brněnské vzory, a to ty, které popisují nepravidelné alternace v kmeni. Jak již bylo řečeno, brněnský systém se vypořádává s výjimkami zavedením vzorů i pro velmi malé množiny lemmat, v nejzazším případě i pro lemma jediné. Pražský systém výjimky popisuje pomocí přiřazení značek jednotlivým tvarům. Tento způsob jsme převzali i my.

Názvy vzorů vytváříme výhradně pomocí znaků anglické abecedy. Jednotlivé znaky tedy kódujeme bez háčeků a bez čárek. Znak *ě* kódujeme systematicky pomocí znaku *j*.

8.2.1.1 Stupňování

Vzory pro skloňování přídavných jmen obsahují i předpis pro stupňování³. Pomocí tohoto předpisu se vytvářejí pouze tvary druhého stupně (komparativu). Třetí stupeň a nově zavedený stupeň *s* (*sebe-* + komparativ) se tvoří pravidelně ze stupně druhého, a to pomocí prefixů *nej-* a *sebe-*. Stejně pravidelné je i stupňování negovaných adjektiv (*nebezpečnější, nejnebezpečnější, sebenebezpečnější*).

Pražský systém vzorů obsahuje i předpis na tvoření superlativu, a to pomocí značky *+* u zakončení. Naše nové vzory to nedělají, neboť tvoření je naprosto pravidelné a může se tvořit od každého přídavného jména, které má druhý stupeň. Totéž se týká i stupňování příslovcí. Je-li tedy ve vzoru uveden předpis na vytvoření 2. stupně, automaticky to znamená, že se vytvoří i stupeň 3 a stupeň *s*. Starost o vytváření tvarů 3. stupně a stupně *s* u přídavných jmen a příslovcí necháváme tedy na nástrojích morfologické analýzy a syntézy.

8.2.1.2 Negace

Podobně zacházíme i s negací. Předpona *ne-*, která slouží k vytvoření negovaného tvaru přídavného jména, příslovce a slovesa, řidčeji i podstatného jména, nebývá obsažena ve slovníku, neboť je pravidelná. Proto ani u negace nepoužíváme praxe z pražského systému, který možnost negace vyznačuje speciálním znakem ve všech zakončeních uvedených ve vzoru.

Zde však existuje několik výjimek. Infinitiv některých sloves se totiž nevytváří prostým připojením *ne-* na začátek slova, ale dochází zde ke krácení v kmeni. Jako příklad uveďme infinitivy *brát* — *nebrat*, *spát* — *nespat*. Tyto tvary je třeba do slovníku zadat explicitně a současně zabránit, aby se tvořily nesprávné dlouhé tvary automaticky z afirmativního tvaru. Tomu zabráníme zavedením **parametru negace** přímo do názvu vzoru. Tento parametr vystupuje ve vzorech ve formě nepovinného prefixu (proto mu také říkáme **prefix**

³O stupňování se lingvisté přou, zda patří do morfologie nebo do slovtvorby, viz např. (Karlík – Hladká, 2004). Držíme se zavedené pražské i brněnské praxe a popisujeme stupňování v rámci morfologie.

negace), který může mít hodnotu A nebo N. Hodnota A znamená, že vzor lze použít pouze ke generování afirmativních tvarů, hodnota N připouští pouze tvary negativní. Absence prefixu neklade na negaci žádná omezení, tedy povoluje tvary jak afirmativní, tak negativní. Parametr negace může vystupovat jen u vzorů přídavných jmen, sloves, příslovcí a podstatných jmen. Vzhledem k velmi volným pravidlům při tvoření negace je lepší omezení na negaci ve vzorech nepoužívat, nemění-li negace kmen slovního tvaru. Zabránili bychom tak rozpoznání neobvyklých tvarů, jako v příkladu (109) z korpusu SYN a v příkladu (110) nalezeném na internetu.

Blbec neblbec, ale hlavně, že jsem zdravý (109)

V tradiční (nerychlé) restauraci však můžeme vidět také velké rozdíly (110)

Právě popsaný přístup k negaci a stupňování již byl částečně implementován, prozatím s původními pražskými značkami, viz (Hlaváčová – Kolovratník, 2008). S prefixem negace se zde zachází podobně jako s jinými prefixy s tím rozdílem, že se při rozpoznání změní hodnota kategorie NEG z NEG=A na NEG=N. Prefixy používané ke stupňování zase mění hodnotu kategorie DEG z DEG=2 na DEG=3 nebo DEG=s.

8.2.2 Derivační vzory

Velké množství českých slov vzniká odvozením z jiných slov. Na otázku, zda derivace patří do morfologie či nikoli, není jednoznačná odpověď. Vzhledem ke způsobu použití morfologického slovníku však derivační vztahy do slovníku zahrnujeme. Např. v automatickém překladu je občas nutné použít místo slovesa slovesné přídavné jméno (*Udělal jsem to. — Mám to udělané / uděláno.*). Jestliže budeme mít ve slovníku odkazy mezi slovesem a příslušnými deverbativy, můžeme tyto vztahy využít při nejrůznějších konstrukcích.

Jsou v zásadě dvě možnosti, jak to udělat:

1. pomocí pravidel,
2. pomocí odkazů mezi jednotlivými hesly.

Nejoperativnější je kombinovaný přístup. V případech, kdy se derivace tvoří naprosto pravidelně, postačí do slovníkového záznamu vložit pravidlo, jak z daného lemmatu vytvořit lemma odvozené.

Derivační vzor však nepřirazujeme přímo k lemmatu, ale ke kofixu. Pomocí tohoto vzoru se odvodí z daného kofixu nové lemma se svým flektivním vzorem (včetně přesně nastavených hodnot všech parametrů). Z toho vyplývá, že derivační pravidla většinou nepoužíváme k odvozování nepravidelných lemmat, které nemají jeden pravidelný flektivní vzor. Derivační vzor se také využije pro vytvoření derivačního odkazu mezi oběma lemmaty — tedy toho, v jehož záznamu je uvedeno, a toho, které se pomocí něho odvodí.

Na rozdíl od ohýbání tedy tvary vzniklé na základě derivačních vzorů nejsou součástí paradigmatu lemmatu, jehož tvary vzor popisuje. Derivační vzory vždy vytvářejí nová lemmata s vlastními paradigmaty.

Použití vzorů pro tvoření derivačních odkazů má několik výhod. Jednak se tím může výrazně zmenšit velikost slovníku. Víme-li např., jak z nedokonavých sloves páté třídy (vzor *dělá*) vytvořit iterativní slovesa (*dělávat*), není třeba

je všechny zahrnovat do slovníku, vygenerují se samy pomocí jednoduchého pravidla.

Druhá výhoda spočívá v implicitní možnosti propojit v rámci slovníku příbuzná slova. Příbuzná slova lze propojit ručně, avšak automatické propojení pomocí derivačních vzorů je elegantnější a méně náchylné k chybám. V těch případech, kdy není možné odvozené slovo utvořit podle jednoduchého pravidla, je stále možnost v rámci slovníku pomocí odkazů příbuzná lemmata propojit ručně. Odvozovací pravidla však tento proces zjednodušují.

Další výhodou je snížení pravděpodobnosti chyby při údržbě slovníku. Jakýkoli ruční zásah do obsahu slovníku s sebou nese riziko chyby, proto je třeba se snažit o co největší využití pravidelných operací.

Některá odvozovací pravidla jsou velmi produktivní a pokrývají velké množství příbuzných dvojic. Existují však i derivace nepravidelné. Jistě je možno vytvořit pravidlo pro každé odvození. Ovšem derivační vzor, který by platil jen pro malý počet, případně dokonce jen jedno odvození, by neúměrně zvyšoval objem celého systému. Přijímáme stejnou zásadu jako v případě ohýbacích vzorů, které nevytváříme pro nepravidelná paradigmata. Ani zde nebudeme vytvářet vzory pro nepravidelná odvození.

Místo toho zařadíme odvozené slovo do slovníku jako plnohodnotné heslo a jeho souvislost s původním slovem zajistíme pomocí derivačního odkazu.

Vzhledem k tomu, že často není ani lingvistům jasné, které lemma je původní a které odvozené, nebudeme se snažit o zachycení směru tohoto vztahu a derivační odkazy umístíme do slovníku oboustranné bez jakýchkoli preferencí. Přesto budeme tyto odkazy i nadále nazývat derivační.

Derivační vztah v našem pojetí tedy propojuje dva záznamy bez ohledu na skutečné (nebo aspoň všeobecně přijímané) časové vztahy (co bylo dřív a co potom).

9 Vzory podstatných jmen

9.1 Obecné vlastnosti

Pro skloňování podstatných jmen používáme základní vzory, které se učí děti na základní škole. Zapisujeme je zkratkou složenou ze dvou písmen. Tato zkratka je většinou shodná se zkratkou vzorů pražského systému a mnemotechnicky vyjadřuje klasický základní vzor.

Zakončení z jednotlivých vzorů podstatných jmen jsou pravidelná, až na některé kombinace pádu a čísla, ve kterých dochází k více možnostem tvoření příslušného slovního tvaru. Je-li tato kombinace přípustná pro celou množinu lemmat ohýbaných podle daného vzoru, je hodnota kategorie **Flektivní mutace** přímo zahrnuta do příslušné trojice ⟨řetězec, morfologická značka, flektivní mutace⟩ vzoru. Příkladem může být dvojitý tvar podstatných jmen v 7pl (*hrady* — *hradama*). Jestliže nějaká kombinace připouští více možností zakončení, ale ne pro všechna lemmata, říkáme jí **kritická kombinace**. U každého lemmatu musíme stanovit, které zakončení je pro něj správné, a přiřadit mu i hodnotu kategorie **Flektivní mutace** (např. 2sg *lesa* s FMU=a, ale *hradu* s FMU=u). U každého vzoru jsou kritické kombinace jiné, některé vzory kritické kombinace nemají, protože jsou zcela pravidelné.

Celý vzor podstatného jména se skládá z následujících částí:

1. z dvoupísmenné části kódující základní vzor,
2. z části kódující koncové znaky kofixu (většinou přímo slovního kmene),
3. z části kódující možná zakončení v kritických kombinacích,
4. z části definující derivace,
5. z nepovinné části kódující různá omezení na tvoření některých slovních tvarů.

Části 2 až 4 jsou závislé na základním vzoru, část 5 je nepovinná a může být přidána k libovolnému substantivnímu vzoru.

Část 2 říká, jak vypadá posledních několik znaků kmene, neboť podle nich se řídí množina zakončení. Často jsou to přímo koncové znaky, a stávají se tak součástí zakončení ze vzoru.

Možnosti zakončení v kritických kombinacích jsou zakódovány ve 3. části vzoru. Obecně budeme kritické kombinace zapisovat ve tvaru Ksg pro jednotné číslo nebo Kpl pro číslo množné, kde K je číslo pádu (hodnota kategorie **Pád**). V případě, že je ve vzoru uvedeno více hodnot pro některou kritickou kombinaci, vygenerují se podle nich příslušné hodnoty kategorie **Flektivní mutace**. Pořadí, ve kterém se mutace uvedou do vzoru, není významné, neboť mutace nejsou umístěny na žádnou škálu. Při práci na přiřazování vzorů se však snažíme zadávat jako první možnost tu mutaci, která je spisovná nebo, v případě obou spisovných nebo obou nespisovných, běžnější. Kdyby se totiž někdo rozhodl mutace číslovat, mohlo by mu takové řazení usnadnit práci.

Vzor může také obsahovat předpis na vytvoření derivací. V současné práci popisujeme pouze odvození přivlastňovacích přídavných jmen a u životných vzorů také odvozená feminina. Přivlastňovací přídavná jména je možno tvořit příponou *-ův* pro hodnoty GEN=M nebo příponou *-in* pro hodnotu GEN=F. Tvoření přivlastňovacího přídavného jména je pro jednotlivé vzory pravidelné, stačí tedy zadat, zda se tvoří, nebo netvoří. Příznak V tvoření přídavného jména přivlastňovacího umísťujeme do části 4, ale explicitně ho uvádíme jen u ženských vzorů, neboť u mužských životných vzorů lze přídavné jméno přivlastňovací utvořit zřejmě vždy.

O tvoření odvozených feminin pojednáme v oddíle 9.3 o životných vzorech.

Každý vzor navíc může obsahovat na konci znak P, znamenající, že se generují jen tvary plurálu, nebo S pro tvoření pouze singuláru. Není-li přítomen ani jeden z těchto znaků, generují se všechny tvary plurálu i singuláru.

Další vlastnost společná všem vzorům může být označení skutečnosti, že vzor se nemá použít pro vytvoření tvaru pro 2pl. Existuje totiž poměrně značné množství paradigmat, ve kterých právě v tomto pádě dochází k alternaci ve kmeni, přičemž všechny ostatní tvary jsou pravidelné, nebo aspoň pravidelné tvary tvoří jednu z více mutací (např. *bouda — bud, dílo — děl, kráva — krav, chvíle — chvíl*). Jak je vidět z příkladů, je tato vlastnost společná více vzorům, i když zdaleka ne všem. Přesto je výhodnější připustit tuto možnost obecně pro substantivní vzory, než ji vypisovat u každého zvlášť. Fakt, že vzor platí pro všechny tvary kromě 2pl, se vyznačí uvedením znaku 2 na konec názvu vzoru.

Obě omezení, na číslo i na (ne)pravidelné tvoření 2pl, jsou obsahem části 5.

Substantivní vzory tedy mohou končit těmito řetězci, které omezují tvoření slovních tvarů:

- S: jen tvary jednotného čísla
- P: jen tvary množného čísla
- 2: všechny tvary kromě 2pl
- 2P nebo P2: všechny tvary množného čísla kromě 2pl (průnik dvou předchozích omezení)

Omezení negace se vyjádří pomocí parametru negace, který se však umísťuje systematicky u vzorů všech slovních druhů jako prefix před základním vzorem. Viz oddíl 8.2.1.2 v obecné kapitole o vzorech.

Všechny vzory podstatných jmen mají pravidelné flektivní mutace FMU=a se zakončením *-ma*, které jsou možné vždy v 7pl: *pány — pánama, stroji — strojema, ženami — ženama, kuřaty — kuřatama* atp. Tyto mutace jsou též součástí všech vzorů, kde o nich též pojednáme.

Nejběžnější mutace podstatných jmen jsme umístili do tabulky 9.1, která je uvedena na str. 84.

Následuje seznam vzorů, jejich popis a příklady.

V názvu vzoru je vždy tučně uveden dvoupísmenný kód názvu vzoru. Podtržené znaky jsou parametry, jejichž možné hodnoty jsou uvedeny ve výčtu, který vždy následuje. Hodnota ϵ znamená prázdný znak a ve skutečnosti se nezapíše. Část vzoru 5 popisující možnost obecných omezení již u jednotlivých

9 Vzory podstatných jmen

POS	GEN	CAS	NUM	Mutace	Její kód	Mutace	Její kód	Mutace	Její kód
					Příklad		Příklad		Příklad
N	I	2	S	a	a	u	u		
					kouta		koutu		
N	I	3	S	u	u	i	i		
					kořenu		kořeni		
N	I	6	S	u	u	e/ě	e	i	i
					obchodu/ lesu/ kořenu		lese/obchodě		kořeni
N	IMN	6	P	ích	i	ách	a	ech	e
					domcích/ ptáčích/ ramíncích/hotelích		domkách/ ptákách/ ramínkách		hotelech/spisovatelech
N	M	36	P	u/í	K	ovi	D		
					pánu/muži/soudci		pánovi/mužovi/soudcovi		
N	M	15	P	i	i	é	e	ové	v
					invalidi		invalidé		invalidové
N	MIN	3	P	ům	d	um	k		
					domům/pánům/městům		domum/pánum/městum		
N	F	1245	P	i	i	e/ě	e		
					lodi/noci		noce/lodě		
N	F	3	P	ím	i	em/ěm	e		
					lodím/nocím/Dejvicím		nocem/loděm		
				ům	u	um	uk		
					Dejvicům		Dejvicum		
N	FM	6	P	ích	i	ech	e		
					nocích/obyvatelích		nocech/obyvatelech		
N	F	7	P	mi	K	emi/ěmi	o	ema/ěma	ea
					lodmi/nocmi		nocemi/loděmi		nocema/loděma
				ma	ka	ima	ia		
					lod'ma/nocma		lodima		
N	F	1	S	e	e	a	a		
					Marie		Maria		
N	F	36	S	e	e	i	i		
					Saše		Saši		
N	F	145	P	e	e	i	i		
					Saše		Saši		
N	N	2	S	í	o	ího	h		
					stavení		staveního		
N	N	6	S	u	u	e/ě	e		
					městu/mléku		mléce/městě		
N	F	2	P	í	i	-	o		
					jeskyní		jeskyň		
N	M	1-7	SP	vložené e	e	bez e	E		
					Bergerovi		Bergrovi		
N	N	2	P	vložené e	e	bez e	E		
					stanovisek		stanovisk		
N	M	7	P	i	o	ama	a	ema	ea
					obyvateli/muži		obyvatelama		obyvatelema
N	MI	kromě 1S		[oeus]s.*	s	*	o		
					Kolumbusovi		Kolumbovi		
NA	.	7	P	mi	o	ma	a		
					ženami/ drahými/ pány		ženama/ drahýma/ pánama		

Tabulka 9.1: Nejběžnější mutace podstatných jmen

vzorů nezmiňujeme, neboť je obecně použitelná u libovolného substantivního vzoru.

9.2 Neživotné vzory

9.2.1 HRAD

hd_x-2sg-6sg-6pl

hd_x-2sg-6sgS

hd_x-2pl-6plP

Vzor *hrad* má tři varianty. První varianta je nejčastější. Podle ní se vytvářejí všechny tvary jednotného i množného čísla. Kritické kombinace jsou, jak je vidět ze vzoru, 2sg, 6sg a 6pl.

Druhá varianta popisuje *singularia tantum*. V tomto případě jsou kritické kombinace pouze 2sg a 6sg a název celého vzoru musí být zakončen znakem S.

Třetí varianta popisuje *pluralia tantum*. Z kritických kombinací pochopitelně vypadnou 2sg a 6sg, naopak k 6pl přibude 2pl, neboť existuje velká skupina substantiv, která má prázdné zakončení v 2pl (např. *Dukovany*, 2pl *Dukovan*, nikoli **Dukovanů*, jak by odpovídalo vzoru *hrad*). V tomto případě musí být na konci vzoru P.

Druhou a třetí variantu lze použít i pro popis takových paradigmat, jejichž množné a jednotné číslo mají různý kofix.

Význam jednotlivých částí:

x je zakončení kofixu. Většinou se jedná o poslední jedno až dvě písmena lemmatu.

Možné hodnoty parametru x:

- r
- ch
- h
- g
- ek (ve 2pl se vypouští -e-)
- k
- en (ve 2pl se vypouští -e-)
- e1 (ve 2pl se vypouští -e-)
- et (ve 2pl se vypouští -e-)
- us
- ky
- y
- ε — ostatní možná zakončení

2sg je zakončení 2. pádu jednotného čísla. Možné hodnoty parametru 2sg jsou současně skutečná zakončení. U kódů uvádíme též hodnotu kategorie **Flektivní mutace**.

- a, FMU=a
- u, FMU=u

Jestliže jsou možné obě hodnoty parametru 2sg, použijí se nejen k vytvoření slovních tvarů, ale i k přiřazení hodnoty **Flektivní mutace**.

6sg je zakončení 6. pádu jednotného čísla. Možné hodnoty parametru 6sg a kategorie **Flektivní mutace**:

- u, FMU=u
- e, FMU=e
- j pro zakončení -ě, FMU=e

Kategorie FMU je stejná u dvou zakončení proto, že tyto dvě alternativy se nikdy nemohou vyskytnout u stejné kombinace pádu a čísla jednoho lemmatu. To vyplývá z historického vývoje českého skloňování. Zakončení -ě nastává po hláskách b, f, m, p, v, d, t, n, zatímco -e po hláskách l, s, z. Se stejným typem mutace se setkáme ještě několikrát.

6pl je zakončení 6. pádu množného čísla. Možné hodnoty parametru 6pl:

- a pro zakončení -ách, FMU=a
- e pro zakončení -ech, FMU=e
- i pro zakončení -ích s automatickým změkčením předchozí souhlásky, FMU=i

2pl je zakončení 2. pádu množného čísla. Možné hodnoty parametru 2pl:

- u pro zakončení -ů, FMU=u
- 0 pro prázdné zakončení, FMU=0

Příklady různých lemmat se vzorem **hrad** jsou v tabulce 9.2. V posledním sloupci je uveden současný vzor z pražského systému.

9.2.2 STROJ

sj_x

Vzor **stroj** je velmi pravidelný, nemá žádné kritické kombinace. Jediným parametrem je zakončení kofixu **x**.

Možné hodnoty parametru **x** jsou:

- e1
- ec
- en pro zakončení -eň
- ε — ostatní možná zakončení

Příklady různých lemmat se vzorem **stroj** jsou v tabulce 9.3. V posledním sloupci je uveden současný vzor z pražského systému.

9 Vzory podstatných jmen

Lemma	Kofix	Vzor	Současný pražský vzor
<i>problém</i>	<i>problém</i>	hd-u-u-e	hd1
<i>hotel</i>	<i>hotel</i>	hd-u-u-ei	hd1xx
<i>cirkus</i>	<i>cirkus</i>	hd-u-u-e	hd1
<i>virus</i>	<i>vir</i>	hdus-u-u-e	hdus
<i>kurs</i>	<i>kurs</i>	hd-u-eu-e	hd2
<i>oceán</i>	<i>oceán</i>	hd-u-ju-e	hd2x
<i>javor</i>	<i>javor</i>	hd-au-u-e	hd4
<i>Motol</i>	<i>Motol</i>	hd-a-euS	hd4
<i>Vlkov</i>	<i>Vlkov</i>	hd-au-ju-e	hd5x
<i>skřek</i>	<i>skře</i>	hdk-u-u-ia	hd1k
<i>hrádek</i>	<i>hrád</i>	hdek-u-u-ia	hd1ek
<i>srpen</i>	<i>srp</i>	hden-a-u-e	hd1ena
<i>svícen</i>	<i>svíc</i>	hden-u-u-e	hd1en
<i>tucet</i>	<i>tuc</i>	hdet-u-u-e	hd1et
<i>jazyk</i>	<i>jazy</i>	hdk-au-eu-ia	hdka
<i>gong</i>	<i>gon</i>	hdg-u-u-ia	hd1g
<i>běh</i>	<i>bě</i>	hdh-u-u-ia	hd1h
<i>poslech</i>	<i>posle</i>	hdch-u-u-ia	hd1ch
<i>Slapy</i>	<i>Slap</i>	hdy-u-eP	hdpy
<i>Dukovany</i>	<i>Dukovan</i>	hdy-0-eP	hdpy0
<i>Jeseníky</i>	<i>Jesení</i>	hdky-u-aiP	hdpsy

Tabulka 9.2: Příklady lemmat ohýbaných podle vzoru *hrad*

Lemma	Kofix	Vzor	Současný pražský vzor
<i>jetel</i>	<i>jetel</i>	sj	sj1
<i>řetězec</i>	<i>řetěz</i>	sjec	sj1ec
<i>obratel</i>	<i>obrat</i>	sjel	sj1el
<i>oheň</i>	<i>oh</i>	sjen	--

Tabulka 9.3: Příklady lemmat ohýbaných podle vzoru *stroj*

9.2.3 Kolísání mezi vzory HRAD a STROJ

hs

Název vzoru je utvořen z počátečních písmen obou vzorů. Podstatná jména, která kolísají mezi vzory **hrad** a **stroj**, jsou většinou pravidelná, proto v jejich vzoru nejsou žádné parametry. Pravděpodobně jedinými nepravidelnými výjimkami jsou lemmata *den* a *týden* (viz (Osolsobě, 1996) na str. 90).

V případě, že se tvary podle obou vzorů neshodují, jsou rozlišeny hodnotou kategorie **Flektivní mutace**. Kombinace kategorií jsou uvedeny v tabulce 9.4.

Pád	Zakončení (hrad)	FMU	Zakončení (stroj)	FMU	Příklad
2sg	-u	u	-e	e	<i>bobelu/bobele</i>
3sg/6sg	-u	u	-i	i	<i>bobelu/bobeli</i>
5sg	-e	e	-i	i	<i>bobele/bobeli</i>
1pl/4pl	-y	y	-e	e	<i>bobely/bobele</i>
6pl	-ech	e	-ích	i	<i>bobelech/bobelích</i>
7pl	-y	y	-i	i	<i>bobely/bobeli</i>

Tabulka 9.4: Tabulka mutací při kolísání mezi vzory **hrad** a **stroj**

9.2.4 ŽENA

zn \underline{x} **Význam jednotlivých parametrů:**

U vzoru **žena** je parametr **x** povinný (není nikdy ϵ), a pokud není uvedeno jinak, znamená poslední souhláskovou skupinu před koncovým *a* v 1sg. Hodnoty jsou uvedeny v tabulce 9.5, ze které jsou též vidět klíčové kombinace pádu a čísla, kvůli nimž se jednotlivé hodnoty musí rozlišovat. V posledním sloupci tabulky je zakončení lemmatu přídavného jména přivlastňovacího, pokud se tvoří.

V posledních třech řádcích tabulky jsou klíčová zakončení uvedena v hranatých závorkách ([x], [y], [j]) proto, že zastupují více znaků. Konkrétně to jsou:

x : b, d, f, m, n, p, t, v**y** : l, s, z**j** : ž, š, č, ř, c, j... 1pl není *y*, ale *i*

Příklady uvádíme v tabulce 9.6. V posledním sloupci je uveden současný vzor z pražského systému.

9 Vzory podstatných jmen

Par.	1sg	[36]sg	2pl	V	Par.	1sg	[36]sg	2pl	V
k	-ka	-ce	-k	-čín	re	-ra	-ře	-er	-řin
r	-ra	-ře	-r	-řin	ke	-ka	-ce	-ek	-čin
g	-ga	-ze	-g	-zin	dj	-ďa	-dě	-ď	-din
h	-ha	-ze	-h	-žin	tj	-ťa	-tě	-ť	-tin
c	-cha	-še	-ch	-šin	nj	-ňa	-ně	-ň	-nin
e	-ea	-ey/eje	-í	—	nk	-ňka	-ňce	-něk	-ňčin
ve	-va	-vě	-ev	-vin	tk	-tka	-tce	-těk	-tčin
ne	-na	-ně	-en	-nin	dk	-dka	-dce	-děk	-dčin
le	-la	-le	-el	-lin	x	-[x]a	-[x]ě	-[x]	-[x]in
be	-ba	-bě	-eb	-bin	y	-[y]a	-[y]e	-[y]	-[y]in
te	-ta	-tě	-et	-tin	j	-[j]a	-[j]e	-[j]	-[j]in
me	-ma	-mě	-em	-min					

Tabulka 9.5: Možná zakončení vzoru žena

Lemma	Kofix	Vzor	Současný pražský vzor
<i>půda</i>	<i>půd</i>	znx	zn1
<i>teta</i>	<i>tet</i>	znxV	zn1n
<i>bouda</i>	<i>boud</i>	znx2	zn1x
<i>tráva</i>	<i>tráv</i>	znx2	zn1x
<i>dělba</i>	<i>děl</i>	znbe	zn2e
<i>tóga</i>	<i>tó</i>	zng	zn3
<i>duha</i>	<i>du</i>	znh	zn4
<i>archa</i>	<i>ar</i>	znc	zn5
<i>informatika</i>	<i>informati</i>	znk	zn6
<i>Maruška</i>	<i>Maruš</i>	znkeV	zn6en
<i>námítka</i>	<i>námit</i>	znke	zn6e
<i>koza</i>	<i>koz</i>	zny	zn7
<i>jehla</i>	<i>jeh</i>	znle	zn7e
<i>farma</i>	<i>far</i>	znme	zn8e
<i>pryčna</i>	<i>pryč</i>	znne	zn9e
<i>lodka</i>	<i>lo</i>	zndk	zn10e
<i>laňka</i>	<i>la</i>	znnk	zn12e
<i>mezera</i>	<i>meze</i>	znr	zn13
<i>hra</i>	<i>h</i>	znre	zn13e
<i>jachta</i>	<i>jach</i>	znte	zn14e
<i>jizva</i>	<i>jiz</i>	znve	zn15e
<i>Korea</i>	<i>Kore</i>	zne	zn19
<i>Saša</i>	<i>Saš</i>	znjV	zn20
<i>dožínky</i>	<i>dožín</i>	znkeP	znp6ge
<i>Karpaty</i>	<i>Karpat</i>	znxP	zn26
<i>Káťa</i>	<i>Ká</i>	zntjV	--
<i>Tatry</i>	<i>Tat</i>	znreP	--

Tabulka 9.6: Příklady lemmat ohýbaných podle vzoru žena

9.2.5 PÍSEŇ

ps_xe

Význam jednotlivých částí:

x je kód posledního znaku lemmatu:

- d pro zakončení *-d'*
- t pro zakončení *-t'*
- n pro zakončení *-ň*
- v pro zakončení *-v*
- c pro zakončení *-č*
- ε pro ostatní zakončení

Parametr e může mít tyto hodnoty:

- e když se v 2sg vypouští e
- ε v ostatních případech

Vzory ps_v ani ps_c neexistují, protože v případech, kdy se před koncové *-v*, *-č* nekládá ve 2pl *-e-*, se použije obecnější vzor ps (*kleč*, *křeč*, a další). Příklady jsou v tabulce 9.7.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>třešeň</i>	<i>třeš</i>	psne	ps1e
<i>Výtoň</i>	<i>Výto</i>	psn	ps1
<i>brukev</i>	<i>bruk</i>	psve	ps2
<i>Ohaveč</i>	<i>Ohav</i>	psceS	--

Tabulka 9.7: Příklady lemmat ohýbaných podle vzoru píseň

9.2.6 KOST

kt

Vzor kost je zcela pravidelný, nemá tedy žádné parametry. Příklady jsou v tabulce 9.8.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>kost</i>	<i>kost</i>	kt	kt1
<i>dveře</i>	<i>dveř</i>	ktP	--

Tabulka 9.8: Příklady lemmat ohýbaných podle vzoru kost

Existuje však množství lemmat, jejichž skloňování kolísá mezi vzory kost a píseň.

9.2.7 Kolísání mezi vzory KOST a PÍSEŇ

kp_x-2sg-36pl-7pl

x je kód posledního znaku lemmatu:

- d pro zakončení *-d'*
- t pro zakončení *-t'*
- n pro zakončení *-ň*
- ε pro ostatní zakončení

Kritické kombinace můžeme rozdělit do tří skupin. Jsou to:

1. 2sg, 1pl, 4pl a 5pl, kódováno parametrem 2sg
2. 3pl a 6pl, kódováno parametrem 36pl
3. 7pl, kódováno parametrem 7pl

V ostatních kombinacích čísla a pádu jsou oba vzory totožné.

2sg je kód zakončení tvaru v 2sg, 1pl, 4pl a 5pl:

- i pro zakončení *-i* — FMU=i
- e pro zakončení *-e* — FMU=e
- j pro zakončení *-ě* — FMU=e

36pl je kód zakončení tvarů v 3pl a 6pl:

- i pro zakončení *-ím* v 3pl a *-ích* v 6pl — FMU=i
- e pro zakončení *-em* v 3pl a *-ech* v 6pl — FMU=e

7pl je kód zakončení tvaru v 7pl. Flektivní mutaci *-ma* ve výčtu pro 7pl neuvádíme, ta je pravidelná pro všechna podstatná jména. Může být k uvedeným mutacím přidána (např. tedy zakončení *-ema* má FMU=ea).

- e pro zakončení *-emi* — FMU=e
- j pro zakončení *-ěmi* — FMU=e
- 0 pro zakončení *-mi* — nulová flektivní mutace

Příklady jsme většinou převzali ze seznamu vzorů skupiny **kost** z disertační práce Kláry Osolsobě (Osolsobě, 1996). Jsou v tabulce 9.9.

Současné pražské vzory, u kterých je uveden znak *, nepopisují celé paradigma příslušného lemmatu. Tvary, které jsou v paradigmatu navíc, jsou uvedeny ve slovníku jako výjimky.

9 Vzory podstatných jmen

Lemma	Kofix	Vzor	Současný pražský vzor
<i>mast</i>	<i>mast</i>	kp-ij-ie-0j	kt1
<i>moc</i>	<i>moc</i>	kp-ie-ie-e	ps18
<i>myš</i>	<i>myš</i>	kp-i-i-0	--
<i>noc</i>	<i>noc</i>	kp-i-i-e	--
<i>žluč</i>	<i>žluč</i>	kp-ie-i-e	ps5 *
<i>ocel</i>	<i>ocel</i>	kp-ie-i-e0	ps5 *
<i>oběť</i>	<i>oběť</i>	kpt-i-ei-0 ¹	--
<i>hrud'</i>	<i>hru</i>	kpd-ij-i-j	ps5 *
<i>odpověď</i>	<i>odpověť</i>	kpd-i-ei-j0	--
<i>choť</i>	<i>choť</i>	kpt-ij-i-j	ps7 *
<i>lod'</i>	<i>lo</i>	kpd-ij-i-0j	ps6

Tabulka 9.9: Příklady lemmat, která kolísají mezi vzory *píseň* a *kost*

9.2.8 NŮŠE

nsx-2pl-1sg

Kritickými kombinacemi vzoru *nůše* je 2pl, který může mít zakončení *-í* nebo prázdnou koncovku (např. *žákyně* — *žákyně* i *žákyní*).

U mutací typu *Marie* — *Maria* dochází k nejednoznačnosti i v 1sg.

Význam jednotlivých částí:

x je kód předposledního znaku lemmatu, před koncovým *-e/ě*:

- **n** pro zakončení *-ně*, je-li v 2pl zakončení *-ň*
- **ε** pro ostatní zakončení

2pl je kód zakončení 2pl. V případě, že zde dochází k alternaci ve kmeni, která znemožňuje pravidelné tvoření, může být tento parametr zcela vypuštěn. V tom případě je však povinné označení 2 na konci názvu vzoru. Parametr 2pl může mít tyto hodnoty:.

- **i** pro zakončení *-í*, FMU=D
- **0** pro prázdné zakončení, FMU=K

1sg je kód zakončení tvaru v 1sg:

- **e** pro zakončení *-e*, FMU=e
- **j** pro zakončení *-ě*, FMU=e
- **a** pro zakončení *-a*, FMU=a

Příklady uvádí tabulka 9.10.

¹Mutaci *obětní* (FMU=t) je třeba zachytit jako výjimku.

9 Vzory podstatných jmen

Lemma	Kofix	Vzor	Současný pražský vzor
<i>nůše</i>	<i>nůš</i>	ns-i-e	ns1
<i>košile</i>	<i>košil</i>	ns-0i-e	ns1
<i>země</i>	<i>zem</i>	ns-i-j	ns2
<i>tradice</i>	<i>tradic</i>	ns-0i-e	ns10
<i>justice</i>	<i>justic</i>	ns-i0-e	ns3x
<i>učnice</i>	<i>učnic</i>	ns-0-e	ns10
<i>žákyně</i>	<i>žáky</i>	nsn-i0-jV	ns3n
<i>Natalie</i>	<i>Natali</i>	ns-i-eaV	ns1
<i>brýle</i>	<i>brýl</i>	ns-i-eP	ns7
<i>Bohunice</i>	<i>Bohunic</i>	ns-eS	ns07
<i>saně</i>	<i>san</i>	ns-i-jP	--

Tabulka 9.10: Příklady lemmat ohýbaných podle vzoru *nůše*

9.2.9 MĚSTO

mtx-6sg-2pl

mtx-6sgS

mtx-2plP

Vzor *město* má tři varianty. První varianta je nejběžnější. Podle ní se vytvářejí všechny tvary jednotného i množného čísla. Kritické kombinace jsou, jak je vidět ze vzoru, 6sg a 2pl.

Druhá varianta popisuje singularia tantum. V tomto případě je kritická kombinace pouze 6sg a název celého vzoru musí být zakončen znakem S.

Třetí varianta popisuje pluralia tantum. Z kritických kombinací pochopitelně vypadne 6sg, zůstane jen 2pl. V tomto případě musí být na konci vzoru P.

Mutace *-kách* a *-cích* (*políčko* — *políčkách* i *políčkách* v 6pl se týkají jen lemmat končících na *-ko* a jsou zcela pravidelné, není tedy třeba je zahrnovat mezi kritické.²

Pod vzor *město* spadají i neutra z latiny a řečtiny, končící na *-um* a *-on*.

²V práci (Osolobě, 1996) jsou zmíněny i tvary *tanzích* (lemma *tango*), *tiších* (lemma *ticho*) a *vlazích* (lemma *vlaho*), které zde nepopisujeme, neboť, ač teoreticky přípustné, se nevyskytují. Na internetu jsme našli jediný výskyt tvaru *tiších*, a to v básni Otokara Březiny.

Význam jednotlivých částí:

x je většinou kód předposledního znaku lemmatu, před koncovým *-o*:

- **r** pro zakončení *-ro*
- **k** pro zakončení *-ko*
- **v** pro zakončení *-vo*, je-li v 2pl vložené *e*
- **l** pro zakončení *-lo*, je-li v 2pl vložené *e*
- **m** pro zakončení *-mo*, je-li v 2pl vložené *e*
- **n** pro zakončení *-no*, je-li v 2pl vložené *e*
- **j** pro zakončení s měkkou souhláskou *-čo* nebo *-jo*
- **kum** pro zakončení *-kum*, je-li v 2pl vložené *e*
- **rum** pro zakončení *-rum*, je-li v 2pl vložené *e*
- **um** pro zakončení *-um*, není-li v 2pl vložené *e*
- **on** pro zakončení *-on*
- **ε** pro ostatní zakončení

6sg je kód zakončení tvaru v 6sg:

- **u** pro zakončení *-u*, FMU=**u**
- **e** pro zakončení *-e*, FMU=**e**
- **j** pro zakončení *-ě*, FMU=**e**
- **i** pro zakončení *-i*, FMU=**i**

2pl je kód zakončení tvaru v 2pl:

- **0** pro prázdné zakončení
- **e** pro vložené *-e*, FMU=**E**
- **i** pro zakončení *-í* (*stadium* — *stadií*)³

Příklady jsou v tabulce 9.11.

9.2.10 MOŘE**mr_x**

Všechna slova vzoru **moře** mají zakončení *-e* nebo *-tě*.

Pod tento vzor zahrnujeme i převzatá slova se zakončením *-e*, která jsou částečně nesklonná, např. *finále*, *promile*, *faksimile*. Někdy se vyskytují s českými koncovkami, jindy se neskloňují, jak ukazují příklady (111) a (112) z korpusu SYN. Mají tak v kombinacích 3sg, 6sg, 7sg, 2pl, 3pl, 6pl a 7pl dvě hodnoty kategorie Mutace — MUT=0 pro klasické zakončení vzoru **moře**, MUT=**e** pro zakončení *-e*.

byl s finále dvouhry spokojen (111)

budu spokojena s finálem (112)

³Tento parametr určuje i tvar v 7sg: *-i* místo běžného *-y*. Jde vlastně o kolísání mezi vzory **město** a **moře**.

9 Vzory podstatných jmen

Lemma	Kofix	Vzor	Současný pražský vzor
<i>slovo</i>	<i>slov</i>	mt-ju-0	mt1x
<i>maso</i>	<i>mas</i>	mt-eu-0	mt1e
<i>kolo</i>	<i>kol</i>	mt-e-0	mt1e
<i>teplo</i>	<i>tep</i>	mt-euS	mt8e
<i>Náchodsko</i>	<i>Náchods</i>	mtk-uS	mts
<i>kakao</i>	<i>kaka</i>	mt-u-i	--
<i>zastupitelstvo</i>	<i>zastupitelst</i>	mtv-uj-e	mt4
<i>lečo</i>	<i>leč</i>	mtx-eu-0	mt1i
<i>vojsko</i>	<i>vojs</i>	mtk-u-0	mt7
<i>ložisko</i>	<i>ložis</i>	mtk-u-e0	mt7y0
<i>sluníčko</i>	<i>sluníč</i>	mtk-u-e	mt7e
<i>patro</i>	<i>pat</i>	mtr-eu-e	mt3e
<i>divadlo</i>	<i>divad</i>	mtl-e-e	mt8e
<i>pásmo</i>	<i>pás</i>	mtm-uj-e	--
<i>specifikum</i>	<i>specifi</i>	mtkum-u-0	mt12
<i>neutrum</i>	<i>neut</i>	mtrum-u-e	mt12r
<i>stadium</i>	<i>stadi</i>	mtum-u-i	mt12i
<i>album</i>	<i>alb</i>	mtum-u-0	mt12
<i>epiteton</i>	<i>epitet</i>	mtou-u-0	--
<i>sympozion</i>	<i>sympozi</i>	mtou-u-i	--
<i>vrata</i>	<i>vrat</i>	mt-0P	mtp
<i>vrátka</i>	<i>vrát</i>	mtk-eP	mtp7
<i>kamna</i>	<i>kam</i>	mtn-eP	--
<i>stehno</i>	<i>steh</i>	mtn-ju-e	mt9e

Tabulka 9.11: Příklady lemmat ohýbaných podle vzoru *město*

Parametr **x** vzoru *moře* může mít tyto hodnoty:

- **t** pro zakončení *-tě*
- **e** pro zakončení *-e* — typ *finále* FMU=e
- **ε** pro ostatní zakončení

Příklady jsou v tabulce 9.12.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>moře</i>	<i>moř</i>	mr	mr1
<i>staveniště</i>	<i>staveniš</i>	mrt	mr10
<i>finále</i>	<i>finál</i>	mre	--

Tabulka 9.12: Příklady lemmat ohýbaných podle vzoru *moře*

9.2.11 KUŘE

kr_x

Do skupiny slov ohýbaných podle vzoru **kuře** zařazujeme, stejně jako v pražském systému, i skupinu přejatých slov se zakončením *-a* (*revma*, *klíma*, *kóma*). **x** může mít tyto hodnoty:

- **t** pro zakončení *-tě*
- **d** pro zakončení *-dě*
- **n** pro zakončení *-ně*
- **e** pro zakončení *-e*
- **j** pro zakončení *-ě*
- **a** pro zakončení *-a*

Příklady jsou v tabulce 9.13.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>kuře</i>	<i>kuř</i>	kre	kr1
<i>nemluvně</i>	<i>nemluv</i>	krn	kr4
<i>koště</i>	<i>koš</i>	krt	--
<i>mládě</i>	<i>mlá</i>	krd	--
<i>poupě</i>	<i>poup</i>	krj	kr1x
<i>dítě</i>	<i>dít</i>	krjS	--
<i>revma</i>	<i>revma</i>	kra	kr5

Tabulka 9.13: Příklady lemmat ohýbaných podle vzoru **kuře**

9.2.12 STAVENÍ

st

Všechna substantiva ohýbající se podle vzoru **stavení** mají jednotný systém zakončení, není třeba žádných dalších parametrů, viz tabulka 9.14.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>stavení</i>	<i>stavení</i>	st	st
<i>listí</i>	<i>listí</i>	stS	st

Tabulka 9.14: Příklady lemmat ohýbaných podle vzoru **stavení**

9.3 Životné vzory

Součástí všech vzorů podstatných jmen rodu mužského životného je i informace o tom, jestli se z nich dá vytvořit femininum. Tato informace se kóduje znakem F v derivační části vzoru. Za ním musí následovat způsob, jakým se ženský rod tvoří. To s sebou nese i informaci o jeho skloňování a tvoření přídavného jména přivlastňovacího, viz tabulka 9.15.

Zakončení feminina	Hodnota parametru	Ženský vzor
-ka	Fa	znkeV
-kyně	Fy	nsn-i0-jV
-ice	Fi	ns-0-eV
-ová	Fo	NFy ⁴

Tabulka 9.15: Hodnoty parametru popisujícího možnosti tvoření feminina od životných substantivních vzorů

Tyto hodnoty pokrývají jen nejběžnější případy. Z mnoha dalších maskulin se sice ženský rod tvoří, ale dochází k alternaci v kmeni, což z praktického hlediska znamená, že se mění kofixy. Pro tyto případy by bylo třeba navýšit počet vzorů, což nám nepřipadá rozumné. Jednodušší je zachovat pouze tato nejběžnější tvoření a ostatní případy vyřešit pomocí odkazů ve slovníku (viz kap. 7). Jako příklad můžeme uvést slovo *tulák* s ženským protějškem *tulačka*. Stejně řešení, tedy slovtvorný odkaz ve slovníku, budou mít i méně častá tvoření, např. *žid* — *židovka*.

Z podstatného jména rodu mužského životného se tvoří pravidelně přivlastňovací přídavné jméno. Příznak V pro jeho tvoření, jak byl zaveden na str. 83, tedy není třeba explicitně zadávat. Přídavné jméno přivlastňovací se tvoří vždy.

Všechny mužské životné vzory mají společné kritické kombinace 1pl a 5pl:

- i pro zakončení -i, FMU=i
- e pro zakončení -é, FMU=e
- v pro zakončení -ové, FMU=v

Často jsou možné dvě mutace, občas i všechny tři. V obou kombinacích přichází vždy v úvahu stejná množina zakončení, tzn. že není možné, aby zakončení 1pl bylo odlišné od 5pl. Tato zakončení 1pl a 5pl jsou stejná pro všechny životné mužské vzory. Otázkou kodifikovanosti se ani zde nezabýváme.

Vzory se vytvářejí stejně jako dosud, tedy ze 2 znaků připomínajících školní vzory.

9.3.1 PÁN

pn_x

Hodnoty parametru **x** jsou uvedeny v tabulce 9.16, ze které jsou též vidět klíčové pády, kvůli nimž se jednotlivé hodnoty musí rozlišovat. Stejná hlásková

³Adjektivní vzor, viz kap. 10

9 Vzory podstatných jmen

změna, která je uvedena ve sloupci 1pl, se projevuje i v 6pl v mutaci *i*, pokud připadá u konkrétního lemmatu v úvahu, např. *archeolog* — *archeoložích*, *vlk* — *vlcích*.

Parametr	1sg	5sg	1pl	
k	<i>-k</i>	<i>-ku</i>	<i>-ci/kové</i>	
ek	<i>-ek</i>	<i>-ku</i>	<i>-ci/kové</i>	
nk	<i>-něk</i>	<i>-ňku</i>	<i>-ňci/ňkové</i>	
ik	<i>-ík</i>	<i>-íku</i>	<i>-íci/íkové</i>	odvozené fem. <i>-ice</i>
r	<i>-r</i>	<i>-re</i>	<i>-ři/rové</i>	
er	<i>-er</i>	<i>-re</i>	<i>-ěři/erové/ři/rové</i>	
t	<i>-r</i>	<i>-ře</i>	<i>-ři/rové</i>	
h	<i>-h</i>	<i>-hu</i>	<i>-zi/hové</i>	
g	<i>-g</i>	<i>-gu</i>	<i>-zi/gové</i>	
c	<i>-ch</i>	<i>-chu</i>	<i>-ši/chové</i>	
es	<i>-es</i>	<i>-e</i>	<i>-i/ové</i>	
us	<i>-us</i>	<i>-e</i>	<i>-i/ové</i>	
ε	ostatní zakončení			

Tabulka 9.16: Možná zakončení vzoru *pán*

V současném pražském systému vzorů se několik alternativ vzoru *pán* liší pouze v pořadí mutací 1pl (včetně jejich hodnocení z hlediska spisovnosti), což však není patrné z názvu vzoru. Z našeho vzoru je pořadí vidět v názvu vzoru, i když, jak už bylo zmíněno výše, mutace nikterak nehodnotíme.

Příklady jsou v tabulce 9.17.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>páv</i>	<i>páv</i>	pnivFi	pn1
<i>žid</i>	<i>žid</i>	pnie	pn5
<i>machr</i>	<i>mach</i>	pntiv	pn6vv
<i>autor</i>	<i>auto</i>	pnrFa	pn6f
<i>Šuster</i>	<i>Šust</i>	pnrFo	pn6
<i>zloduch</i>	<i>zlodu</i>	pnciv	pn22
<i>odborník</i>	<i>odborn</i>	pnikiFi	pn16ik
<i>jézéďák</i>	<i>jézéďá</i>	pnki	pn16
<i>daněk</i>	<i>da</i>	pnnkiv	-
<i>sok</i>	<i>so</i>	pnkivFy	pn18
<i>oslík</i>	<i>osl</i>	pnikivFi	pn16ik
<i>sládek</i>	<i>slád</i>	pnekivFo	pn17
<i>chirurg</i>	<i>chirur</i>	pngv	pn19

Tabulka 9.17: Příklady lemmat ohýbaných podle vzoru *pán*

9.3.2 MUŽ

mz_x

x může mít tyto hodnoty:

- l pro zakončení *-el*
- c pro zakončení *-ec*, odvozené femininum: *-kyně*
- ε pro ostatní zakončení

Příklady jsou v tabulce 9.18.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>plynař</i>	<i>plynař</i>	mziFa	mz1z
<i>vědec</i>	<i>věd</i>	mzciFy	mz5
<i>obyvatel</i>	<i>obyvat</i>	mzleFa	mz1e
<i>tchoř</i>	<i>tchoř</i>	mzi	mz1
<i>vítěz</i>	<i>vítěz</i>	mzviFa	mz1z

Tabulka 9.18: Příklady lemmat ohýbaných podle vzoru muž

9.3.3 Kolísání mezi vzory PÁN a MUŽ

pm_x

Stejně jako u neživotných vzorů *hrad* a *stroj*, i mezi měkkým a tvrdým vzorem životným není ostrá hranice. Existuje několik lemmat, jejichž skloňování mezi těmito dvěma vzory kolísá. Aby bylo zachováno Zlaté pravidlo morfolgie, je třeba pro tuto kategorii vymezit samostatný vzor, který se postará o správné přiřazení hodnot kategorií *Flektivní mutace*. Kombinace kategorií jsou uvedeny v tabulce 9.19.

Pád	Zakončení (pán)	FMU	Zakončení (muž)	FMU	Příklad
2sg/4sg	<i>-a</i>	a	<i>-e</i>	e	<i>markýza/markýze</i>
3sg/6sg	<i>-u</i>	u	<i>-i</i>	i	<i>markýzu/markýzi</i>
5sg	<i>-e</i>	e	<i>-i</i>	i	<i>markýze/markýzi</i>
4pl	<i>-y</i>	y	<i>-e</i>	e	<i>markýzy/markýze</i>
6pl	<i>-ech</i>	e	<i>-ích</i>	i	<i>markýzech/markýzích</i>
7pl	<i>-y</i>	y	<i>-i</i>	i	<i>markýzy/markýzi</i>

Tabulka 9.19: Tabulka mutací při kolísání mezi vzory pán a muž

9 Vzory podstatných jmen

Příklad ukazuje tabulka 9.20.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>vyvrhel</i>	<i>vyvrhel</i>	pmv	pn3j mz1

Tabulka 9.20: Příklad lemmatu s kolísavým skloňováním podle vzorů pán a muž

9.3.4 PŘESEDÁ

pd_x

x může mít tyto hodnoty:

- k pro zakončení *-ka*
- j pro měkká zakončení
- ϵ pro ostatní zakončení

Příklady jsou v tabulce 9.21.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>obejda</i>	<i>obejd</i>	pdv	pd1
<i>předseda</i>	<i>předsed</i>	pdvFy	pd1
<i>Sašenka</i>	<i>Sašenk</i>	pdkv	pd5
<i>turista</i>	<i>turist</i>	pdieFa	pd2f
<i>Míša</i>	<i>Miš</i>	pdjv	pd1

Tabulka 9.21: Příklady lemmat ohýbaných podle vzoru předseda

9.3.5 SOUDCE

sc

Všechna lemmata skloňovaná podle vzoru **soudce** mají zakončení *-ce*. Do konce jsou u tohoto vzoru vždy možná zakončení 1pl *-ci* a *-cové*, stejně tak lze vždy tvořit ženský vzor pomocí zakončení *-kyně*. Pro jednotnost s ostatními životnými vzory ale tyto alternativy vždy vypisujeme. Vzor **soudce** tedy vypadá takto: **sci**vFy.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>soudce</i>	<i>soud</i>	scivFy	sc1o

Tabulka 9.22: Příklady lemmat ohýbaných podle vzoru *soudce*

9.4 Adjektivní vzory

Některá podstatná jména mají adjektivní skloňování. Často jde o vlastní jména osob, ale mohou to být i jména obcí (*Vřesová*), nebo apelativa (*kapesné, vrátný*). V takovém případě pro skloňování přebíráme vzory přídavných jmen, ovšem s poznámkou, že jde o jména podstatná. To se vyznačí pomocí prefixu N před vlastním adjektivním vzorem.

Kromě toho je třeba ještě dodat informaci o rodu, protože adjektivní vzory obecně vytvářejí všechny rody. Rod se ve vzoru zachycuje hned za úvodním prefixem N, a to pomocí běžných kódů pro rod v rámci systému značek. Není třeba zaznamenávat možnost tvoření negace, ani stupňování, protože u podstatných jmen to není relevantní. Ze zpodstatnělých přídavných jmen se také netvoří žádné automatické derivace.

Příklady adjektivně skloňovaných substantiv jsou uvedeny v kapitole o adjektivních vzorech, v oddíle 10.4.

10 Vzory přídavných jmen

Podle adjektivních vzorů se ohýbají především přídavná jména, ale využívají se i k vytváření slovních tvarů některých zájmen (např. *který*), číslovek (*první*) a podstatných jmen (*kapesné*). Pražský systém vzorů má pro každý slovní druh zvláštní vzor, ale to není třeba. Je pouze nutno do vzoru zaznamenat, jakého slovního druhu se týká, aby se správně přiřadila hodnota kategorie **Slovní druh**.

Názvy adjektivních vzorů tedy mají povinný prefix, označující slovní druh, jehož paradigma je vzorem vytvářeno. Prvním znakem adjektivního vzoru je tedy kód slovního druhu: N pro podstatná jména, C pro číslovky, P pro zájmena a A pro přídavná jména.

Následuje nepovinný prefix negace (viz 8.2.1.2).

Dále se adjektivní vzor skládá z následujících částí:

1. z části kódující základní skloňování ve všech pádech, číslech a rodech,
2. z části kódující stupňování,
3. z části kódující možné derivace.

Jednotlivé části vzoru je možno téměř libovolně kombinovat, přičemž každá část generuje jen určitou podmnožinu slovních tvarů paradigmatu, třetí část potom podává předpis, jak vytvořit odvozené lemma s pevně daným vzorem pro ohýbání. V případě některých derivací, kde se liší kofix pro derivace od kofixu pro ohýbání, je možno první dvě části vzoru vynechat. Např. jmenný tvar *mlád* se ohýbá podle vzoru AJ, kde J patří až do třetí části vzoru — viz 10.2.1.

První část můžeme chápat jako vlastní název vzoru. Aby byl název co nejintuitivnější, je jím v případě tvrdého skloňování (zkrácená) poslední slabika lemmatu. Podle ní se totiž řídí skloňování i případné stupňování. V případě skloňování měkkého stačí jeden znak (volíme i).

Další části vzoru můžeme chápat jako parametry, které rozšiřují nebo upravují vlastnosti základní části vzoru.

10.1 Skloňování a stupňování

10.1.1 Základní část vzoru — skloňování

Běžné dělení adjektivních paradigmat na tvrdá a měkká pro automatické zpracování nestačí, neboť u tvrdých přídavných jmen dochází v 1pl a při stupňování ke změně v kmeni. Klasický tvrdý vzor je tedy třeba rozepsat.

Tabulka 10.1 vyjmenovává základní části adjektivních vzorů. Druhý sloupec tabulky stanoví zakončení lemmatu, jehož se vzor týká. Další sloupce potom ukazují změny v zakončení slovních tvarů pro zvláštní hodnoty morfologických kategorií. Z tabulky je vidět, čím se jednotlivé vzory od sebe liší a proč jsou od sebe odlišeny. V posledních dvou řádcích tabulky jsou uvedeny vzory pro

10 Vzory přídavných jmen

přídavná jména přivlastňovací.

Hvězdička * v posledních dvou sloupcích znamená, že se tvar řídí podle dalších částí vzoru, pomlčka – znamená, že vzor příslušné tvary nevytváří.

Vzor	Zakončení lemmatu	1pl živ.	Stupň.	Odvozené D
sky	-ský	-ští	-štější	-sky
cky	-cký	-čtí	-čtější	-cky
hy	-hý	-zí	-ž-	* (-ze/-ho)
chy	-chý	-ší	-š-	* (-še/-cho)
ry	-rý	-ří	* (-rší/-řejší)	* (-ře/-ro)
ky	-ký ¹	-cí	* (-čí/-ší)	* (-ce/-ko)
ny	-ný ²	-ní	-n-	* (-ně/-no)
y	-ý ³	-í	*	*
yy ⁴	-ý	–	*	*
yí ⁵	–	-í	*	*
i ⁶	-í	-í	*	*
y23 ⁷	–	–	-ší	*
uv	-ův	-ovi	–	–
in	-in	-ini	–	–

Tabulka 10.1: Tabulka adjektivních vzorů

Vzor y23 se použije v tom případě, kdy ve 2. a 3. stupni dochází ke změně v kmeni, a je tedy nutné použít jiný kofix. Podle tohoto vzoru se samozřejmě skloňují jen stupňovaná přídavná jména.

Každým z těchto vzorů se definuje skloňování 1. stupně (pozitivu) pro všechny rody a obě čísla. Vzory sky a cky díky své pravidelnosti umožňují i jednoznačně určit zakončení 2. stupně, totiž -štější, -čtější. Vzor y23 naopak slouží pouze pro stupňování, a to vždy se zakončením -ší. Tomu je třeba přizpůsobit i kofix. Je-li uveden vzor y23, je třeba ke kofixu přidat kód K (krátký) nebo D (dlouhý) kategorie **Flektivní mutace**. U ostatních vzorů to není třeba, protože typ je dán kódem pro stupňování, který uvedeme v následující části vzoru — viz oddíl 10.1.2.

Stupňují se jen ta přídavná jména, která to nemají „zakázáno“ — viz alternativa 0 jako pokračování vzoru.

Přídavná jména se zakončením na -ný mohou mít buď vzor ny, nebo y, jestliže netvoří jmenný tvar v maskulinu, nebo jejichž jmenný tvar nemá epentetické e (*plný* — *pln*). Samozřejmě je třeba vzoru přizpůsobit kofix.

Vzory y23, sky a cky nemohou pokračovat žádným z následujících způsobů, které popisují stupňování, neboť jejich stupňování je jednoznačné a vždy stejné.

¹ale ne na -ský ani -cký

²jmenný tvar v maskulinu s epentetickým e (*schopný* — *schopen*)

³ostatní tvrdá, všechny tvary

⁴ostatní tvrdá kromě 1pl masc. živ.

⁵tvrdá v 1pl masc. živ.

⁶všechna měkká

⁷skloňování pouze 2. a 3. stupně

10.1.1.1 Pravidelné mutace ve skloňování přídavných jmen

U všech přídavných jmen dochází k mutacím v 7pl ($[-\acute{y}\acute{i}]mi$, $[-\acute{y}\acute{i}]ma$, u tvrdých navíc i $-ejma$).

U přídavných jmen tvrdých je pravidelná mutace $-ej$ spisovné mutace $-\acute{y}$, a to ve všech slovních tvarech, kde je $-\acute{y}$ přítomno ve spisovné formě, tedy nejen *nový* — *novej*, ale i *nových* — *novejch*, *novými* — *novejma*, apod.

Další systematické mutace $-í$ a $-\acute{y}$ jsou v 1pl rodu mužského životného, podobně $-é$ a $-\acute{y}$ (*velké/velký domy*).

Přehled systematických mutací adjektivního skloňování je uveden v tabulce 10.2 i s kódy.

Kód	Vysvětlení	Příklad
y	\acute{y} místo \acute{e} nebo \acute{i}	<i>velký (domy, páni), velkého</i>
k	krátká mutace	<i>velkým, otcovým</i>
a	zakončení $-ma$ v 7pl	<i>velkýma</i>
j	$-ej-$ místo $-\acute{y}-$	<i>velkej, velkejch</i>
ja	příklad kombinace	<i>velkejma</i>

Tabulka 10.2: Systematické flektivní mutace v adjektivních vzorech

Přídavná jména přivlastňovací mají několik specifických mutací. Mimo jiné se zde vyskytují (nekodifikované) mutace $-ovo$ a $-ic$ pro libovolný rod, číslo i pád, viz příklady (113) až (119) z korpusu SYN.

<i>Zemřel v <u>Havlíčkovo</u> Brodě</i>	(113)
<i>šerifovo pomocníku vůbec nenapadlo</i>	(114)
<i>Na jeviště nastoupil <u>Dismanovo</u> soubor</i>	(115)
<i>táhnul to později s <u>kostelníkovo</u> ženou</i>	(116)
<i>sleduje <u>sousedovic</u> slepice</i>	(117)
<i>se <u>sousedovic</u> ohařem</i>	(118)
<i>od <u>sousedovic</u> trávníku</i>	(119)

Přehled těchto i dalších systematických flektivních mutací specifických pro přídavná jména přivlastňovací obsahuje tabulka 10.3.

Kód	Vysvětlení	Příklad
o	nesklonná mutace <i>-ovo</i>	<i>Novákovo (domy)</i>
c	nesklonná mutace <i>-ovic</i>	<i>Novákovíc (domy)</i>
e	<i>-ině/-ově</i>	<i>matčině, otcově</i>
u	<i>-inu/-ovu</i>	<i>matčinu, otcovu</i>
k	krátká mutace	<i>otcovym, otcuv</i>
E	vložené <i>-e-</i> před r	<i>Bergerův</i>
R	bez vloženého <i>-e-</i> před r	<i>Bergrův</i>
Eja	příklad kombinace	<i>Bergerovejma</i>

Tabulka 10.3: Systematické flektivní mutace vzorů přivlastňovacích přídavných jmen

Tabulka 10.4 na str. 106 uvádí příklad vzoru **Aky0**. Liché sloupce obsahují zakončení, sudé sloupce potom značku, která se skládá z pěti pozic. První pozice kóduje kategorii **Slovní druh**, který je všude stejný, totiž přídavné jméno. Další pozice kódují kategorie **Rod**, **Číslo**, **Pád** a **Stupeň**, který je opět všude stejný. Značka není úplná, protože v ní chybí zakódování kategorie **Negace**. Je zřejmé, že příslušné slovní tvary mohou být jak negativní (**NEG=N**), tak pozitivní (**NEG=A**). Není to tedy morfologická značka podle přesné definice, která vyžaduje přítomnost všech relevantních morfologických kategorií. Kód za pomlčkou u některých zakončení kóduje flektivní mutaci.

V tabulce 10.5 na str. 107 uvádíme druhou část definice vzoru **Akys**, která obsahuje předpis na stupňování. První část vzoru, která popisuje skloňování pozitivu, je vlastně popis vzoru **Aky0**, tedy tabulka 10.4.

10.1.2 Stupňování

Druhá část vzoru, která popisuje stupňování, používá následujících kódů:

- e pro zakončení *-ejší*, FMU=D
- j pro zakončení *-ější*, FMU=D
- s pro zakončení *-ší*, FMU=K
- c pro zakončení *-čí*, FMU=K
- 0 pro nestupňovatelné

V prvních dvou případech se všem tvarům přiřadí dlouhá flektivní mutace (FMU=D), ve druhých dvou případech krátká flektivní mutace (FMU=K), a to i tehdy, když se u konkrétního lemmatu mutace netvoří.

Výjimky jsou, jak již bylo uvedeno, vzory **y23**, který slouží pouze pro stupňování, a vzory **sky** a **cky**, kde je zakončení stupňování pravidelné, takže není třeba ho zvlášť specifikovat. Jediné přípustné možnosti jsou **sky0** a **cky0** pro zákaz stupňování. Ostatní základní vzory nemají zakončení stupňování jednoznačně dáno (např. *divoký* — *divočejší*, ale *lehký* — *lehčí*), tudíž je třeba ho specifikovat touto druhou částí vzoru.

10.1.2.1 Mutace ve stupňování přídavných jmen

Několik přídavných jmen má krátkou a dlouhou mutaci pro stupňování. Dlouhá mutace končí na *-[e]jší*, krátká je bez *-[e]j-*. Jako příklady mohou sloužit dvojice *snadnější* — *snazší*, *hrubější* — *hrubší*. Není jich mnoho, jde o uzavřenou skupinu. Kódování ukazují tabulka 10.6.

10 Vzory přídavných jmen

ká	ANP41	kejch	AFP61-Fj	ký	ANS41-Fy	kých	ANP61
ká	AFS51	kejch	AMP61-Fj	ký	AMP41-Fy	kých	AMP61
ká	ANP51	kejch	AIP21-Fj	ký	AFP41-Fy	kých	AIP61
ká	ANP11	kejch	ANP61-Fj	ký	ANP11-Fy	kých	AMP21
ká	AFS11	kejch	AFP21-Fj	ký	AIS51	kých	AFP21
ké	ANS11	kejch	ANP21-Fj	ký	ANS11-Fy	kých	AFP61
ké	AFS31	kejch	AIP61-Fj	ký	AFS61-Fy	kých	AIP21
ké	AIP11	kejch	AMP21-Fj	ký	AFP11-Fy	kých	ANP21
ké	AIP41	kejm	ANP31-Fj	ký	ANS51-Fy	kym	AIS61-Fky
ké	AFS21	kejm	AMP31-Fj	ký	AIS41	kym	ANS71-Fky
ké	ANS51	kejm	AIP31-Fj	ký	ANP41-Fy	kym	AMS61-Fky
ké	AMP41	kejm	AFP31-Fj	ký	AIS11	kym	AMS71-Fky
ké	ANS41	kejma	AFD71-Fj	ký	ANP51-Fy	kym	AIS71-Fky
ké	AFP11	kejma	ANP71-Fj	ký	AMP11-Fy	kym	ANS61-Fky
ké	AFS61	kejma	AIP71-Fj	ký	AMP51-Fy	kým	AMP31
ké	AIP51	kejma	AMP71-Fj	ký	AFP51-Fy	kým	AIS61-Fy
ké	AFP41	kejma	AFP71-Fj	ký	AFS21-Fy	kým	AMS71
ké	AFP51	kém	ANS61	ký	AMS51	kým	AMS61-Fy
kého	AMS21	kém	AIS61	ký	AIP41-Fy	kým	AFP31
kého	ANS21	kém	AMS61	ký	AMS11	kým	ANS61-Fy
kého	AIS21	kému	AMS31	ký	AIP51-Fy	kým	AIP31
kého	AMS41	kému	AIS31	ký	AFS31-Fy	kým	ANP31
cí	AMP51	kému	ANS31	ký	AIP11-Fy	kým	AIS71
cí	AMP11	kou	AFS41	kýho	AIS21-Fy	kým	ANS71
kej	AIS41-Fj	kou	AFS71	kýho	AMS41-Fy	kýma	AFD71
kej	AMS11-Fj	kými	AFP71	kýho	AMS21-Fy	kýma	ANP71-Fy
kej	AIS11-Fj	kými	AMP71	kýho	ANS21-Fy	kýma	AMP71-Fy
kej	AIS51-Fj	kými	AIP71	kýmu	AMS31-Fy	kýma	AFP71-Fy
kej	AMS51-Fj	kými	ANP71	kýmu	AIS31-Fy	kýma	AIP71-Fy
				kýmu	ANS31-Fy		

Tabulka 10.4: Příklad vzoru Aky0

10 Vzory přídatných jmen

ší	ANP1[23s]	ším	AIP3[23s]-Fk	šíma	AFD7[23s]
ší	AIS1[23s]	ším	AMS7[23s]-Fk	šíma	AMP7[23s]-Fa
ší	AMS5[23s]	ším	ANP3[23s]-Fk	šíma	AIP7[23s]-Fa
ší	AFS7[23s]	ším	ANS6[23s]-Fk	šíma	ANP7[23s]-Fa
ší	AFP5[23s]	ším	ANS7[23s]-Fk	šíma	AFP7[23s]-Fa
ší	AIP5[23s]	ším	AIS6[23s]-Fk	šími	AFP7[23s]
ší	AIP1[23s]	ším	AIS7[23s]-Fk	šími	ANP7[23s]
ší	AFS3[23s]	ším	AMP3[23s]-Fk	šími	AIP7[23s]
ší	AMP4[23s]	ším	AMS6[23s]-Fk	šími	AMP7[23s]
ší	AFS1[23s]	ším	AFP3[23s]-Fk	šímu	ANS3[23s]
ší	AFS2[23s]	ším	AMS6[23s]	šímu	AIS3[23s]
ší	AMS1[23s]	ším	AMS7[23s]	šímu	AMS3[23s]
ší	AMP1[23s]	ším	AIS6[23s]	šího	AIS2[23s]
ší	AMP5[23s]	ším	AIS7[23s]	šího	ANS2[23s]
ší	AFS6[23s]	ším	AMP3[23s]	šího	AMS2[23s]
ší	AIS4[23s]	ším	ANS6[23s]	šího	AMS4[23s]
ší	AFS5[23s]	ším	ANS7[23s]	ších	ANP6[23s]
ší	ANS1[23s]	ším	ANP3[23s]	ších	AMP6[23s]
ší	AFP1[23s]	ším	AFP3[23s]	ších	AFP2[23s]
ší	ANS5[23s]	ším	AIP3[23s]	ších	AIP2[23s]
ší	ANP4[23s]			ších	AMP2[23s]
ší	AIP4[23s]			ších	ANP2[23s]
ší	AFP4[23s]			ších	AIP6[23s]
ší	ANS4[23s]			ších	AFP6[23s]
ší	AIS5[23s]				
ší	AFS4[23s]				
ší	ANP5[23s]				

Tabulka 10.5: Příklad: část vzoru *Akys*, který kóduje stupňování se zakončením *-ší* (např. *těžší*). Značky jsou zapsány pomocí regulárního výrazu — ve všech případech tedy jde o zakončení stupně 2, 3 nebo *s*.

Kód	Vysvětlení	Příklad
FD	dlouhé stupňování <i>-[eě]jší</i>	<i>hrubější, trpčejší</i>
FK	krátké stupňování <i>-ší/-čí</i>	<i>hrubší, trpčí</i>

Tabulka 10.6: Systematické flektivní mutace ve stupňování přídavných jmen

Tyto mutace se pochopitelně kombinují s mutacemi ve skloňování, takže např. pro komparativ 7pl lemmatu *hrubý* dostáváme 4 variantní tvary: *hrubějšími, hrubšími, hrubějšíma, hrubšíma*, které mají po řadě hodnoty FMU=D, FMU=K, FMU=Da, FMU=Ka.

Existují ještě dvě lemmata s mutacemi ve stupňování, které však nezapadají do paradigmatického vzorce z předchozího odstavce, protože mutace nejsou v zakončení, nýbrž ve kmeni. Jsou to lemmata *bílý* a *svatý*. U obou dochází při stupňování k alternaci v kmeni, čímž vznikne mutace. Tyto mutace se týkají všech slovních tvarů 2. i 3. stupně, i nově zavedeného stupně S (typ *sebekrásnější*).

Máme totiž *svatější* i *světější*, jakož i *bílejší* i *bělejší*, i když *bílejší* je výrazně méně časté. Jelikož se nejedná o systematický jev, používáme v tomto případě pro odlišení mutací obyčejné číslování. Na konkrétní hodnotě nezáleží, podstatné je mutace pouze odlišit.

10.2 Derivace

Uvedeme zde nejběžnější derivace z přídavných jmen. Výjimky je třeba zadat do slovníku zvlášť. Uvedené derivace popisují, jakým způsobem se z přídavných jmen tvoří odvozená lemmata. U každého typu derivace tedy ještě uvádíme, jakým způsobem se odvozená lemmata ohýbají.

Do derivační části umísťujeme i pravidlo na tvoření jmenného tvaru, přestože to derivace není. Se skutečnými derivacemi má však společné to, že se netvoří vždy, ale jen od některých lemmat. O derivaci v pravém slova smyslu však nejde, jmenný tvar patří do paradigmatu dlouhého tvaru přídavného jména, které je jeho lemmatem.

10.2.1 Tvoření jmenného tvaru

J následováno buď

- ničím (*jistý — jíst*), nebo
- e pro vložení epentetického *-e-* u mužského rodu (*schopný — schopen*), nebo
- - pro vypuštění mužského singuláru (např. adjektivum *snadný* netvoří jmenný rod pro rod mužský, ale ostatní rody možné jsou).

Ohýbání jmenných tvarů spočívá ve vytvoření tvarů pro 1. a 4. pád všech rodů a čísel (v případě J- bez 1sg mužského rodu), i když 4. pád lze většinou považovat za archaismy (zejména u mužského rodu). Důležité je, že všechna tato tvoření jsou pravidelná, a proto pro ně není třeba specifikovat zvláštní vzory.

Výjimečné tvoření jmenného tvaru se vyskytuje mj. u těchto adjektiv: *mlád*, *sláb*, *stár*, *zdráv*. Zde jde o změnu délky v kmeni, takže vždy bude třeba použít dvojí kofix a vzor:

- zdrav Ajj pro skloňování a stupňování adjektiva *zdravý*,
- zdráv AJ pro odvození jmenných tvarů *zdráv*, *zdráva*, *zdrávo*, *zdrávi*, *zdrávy*, *zdrávu*.

Vzory Ay[ejsc0]Je, Ayy[ejsc0]Je, ani Ayi[ejsc0]Je nemají smysl, jsou tedy nepřijatelné. Epentetické *-e-* se totiž vkládá vždy před poslední souhlásku před zakončením *-ý*, takže vzor tuto souhlásku musí kódovat.

Měkká adjektiva jmenné tvary netvoří nikdy, vzory AiJ[e-]? tedy také nemají smysl.

10.2.2 Tvoření příslovce

Používáme znak D následovaný:

- e pro zakončení *-e*
- j pro zakončení *-ě*
- u pro zakončení *-u*⁸
- y pro zakončení *-y*
- o pro zakončení *-o*

Poznámka:

Je-li u vzoru ky derivace Dy, adverbium má zakončení *-ky*, je-li derivace De, zakončení je *-ce*.

Odvozená příslovce se mohou stupňovat v závislosti na typu přídavného jména, ze kterého byly odvozeny (viz též tabulka 10.1).

sky, cky ... stupňuje se pravidelně (*-štěji*, *-čtěji*),

ky ... stupňuje se pravidelně (*-čeji*), ale je hodně výjimek (*hluboký — hluboce — hlouběji*, *úzký — úzce — úžeji*, *krátce — kratčeji*, *řídce — řídčeji*, *vysoký — vysoko*, *vysoce — výše*) — u těch je třeba derivaci uvést zvlášť,

Do ... nestupňuje se

De ... stupňuje se *-ji* (*milý — mile — mileji*), tedy pravidelně

Du ... stupňuje se *-eji* (*pomalý — pomalu — pomaleji*), tedy pravidelně

Dj ... stupňuje se *-ji* (*čistý — čistě — čistěji*), tedy pravidelně

Je-li původní přídavné jméno nestupňovatelné (má ve vzoru 0), nestupňuje se ani příslovce.

⁸Zdá se, že se zakončením *-u* existuje jen jediné příslovce tvořené z přídavného jména, a to *pomalu*.

10.2.3 Tvoření podstatného jména na *-ost*

Takto derivované podstatné jméno se vždy skloňuje podle vzoru *kost*, stačí tedy možnost této derivace zaznamenat jediným znakem, a to N. U vzoru *ky* se vytváří se zakončením *-kost*.

Tvoření jiných typů podstatných jmen není tak pravidelné, nezavádíme pro ně proto žádné pravidlo.

Jako příklad může sloužit kofix *strnul* se vzorem *AyeN*, který vytvoří všechny tvary přídavného jména *strnulý*, s pravidelným stupňováním *strnulejší* a odvozeným podstatným jménem *strnulost*.

10.3 Příklady

V tomto oddíle přinášíme příklady přídavných jmen spolu s kofixy a vzory, pomocí kterých se tvoří celé paradigma. Jako příklady jsme vybrali téměř všechny adjektivní vzory uvedené v disertační práci Kláry Osolsobě (viz (Osolsobě, 1996)). Jsou uvedeny v tabulce 10.7. Z příkladů je patrna parametrizovatelnost adjektivních vzorů. Pro názornost jsou v posledním sloupci uvedeny 4 tvary odvozené z příslušného kofixu a vzoru. Jsou to po řadě: komparativ, jmenný tvar, derivované podstatné jméno a derivované příslovce. V případě, že se některý tvar netvoří, je místo něj uvedena pomlčka (–). Pomlčka tedy v tomto sloupci neslouží jako oddělovač.

Poznámky pod čarou na této stránce se vztahují k odkazům v tabulce 10.7 na str. 111.

⁸Jmenný tvar v rodě mužském *mrtev* pomocí vzorů nemůžeme popsat. Do slovníku je zadán explicitně se svou morfologickou značkou. Vzor jsme nevytvořili proto, že jde o výjimku (ještě se vyskytuje *polomrtev*). Většina jmenných tvarů s vloženým *-e-* v rodě mužském končí na *-en*, což je ve vzorech zachyceno. Další výjimky jsou ještě *blízek*, *vesel* a *dalek*, které ale zavedené vzory popsat mohou.

⁹Príslovce *těžko/těžce* se zde odvodit nedají kvůli jinému kofixu při stupňování (*tíže/tížeji*).

10 Vzory přídavných jmen

Lemma	Kofix	Vzor	
<i>čilý</i>	<i>čil</i>	AyeNDe	<i>čilejší – čilost čile</i>
<i>kovový</i>	<i>kovov</i>	AyODj	– – – <i>kovově</i>
<i>mrtvý</i>	<i>mrtv</i>	Ay0J-Dj	– <i>mrtva</i> – <i>mrtvě</i> ⁹
<i>jiný</i>	<i>jin</i>	Ay0	– – – –
<i>křepký</i>	<i>křep</i>	AkycNDe	<i>křepčí – křepkost křepce</i>
<i>blizoučkový</i>	<i>blizouč</i>	AkyODoe	– – – <i>blizoučko/ce</i>
<i>sladký</i>	<i>slad</i>	AkysNDoe	<i>sladší – sladkost sladko/ce</i>
<i>těžký</i>	<i>těž</i>	AkysN	<i>těžší – těžkost</i> – ¹⁰
<i>daleký</i>	<i>dale</i>	Aky0JNDoe	– <i>dalek dalekost daleko/ece</i>
<i>řídový</i>	<i>říd</i>	AkyONDoe	– – <i>řídlost řídko/ce</i>
	<i>řid</i> FMU=K	Ay23	<i>řidší – – –</i>
<i>hluboký</i>	<i>hlubo</i>	AkyONDoe	– – <i>hlubokost hluboko/ce</i>
	<i>hlub</i> FMU=K	Ay23	<i>hlubší – – –</i>
<i>strohý</i>	<i>stro</i>	AhyONDe	– – <i>strohost stroze</i>
<i>vetšný</i>	<i>vet</i>	AchyeNDe	<i>vetšejší – vetčnost vetše</i>
<i>moudrý</i>	<i>moud</i>	AryeNDe	<i>moudřejší – moudrost moudře</i>
<i>bílý</i>	<i>bíl</i>	AyeDe	<i>bílejší – – bíle</i>
	<i>bělej</i> FMU=1	Ay23	<i>bělejší – – –</i>
<i>nový</i>	<i>nov</i>	AyeNDj	<i>novější – novost nově</i>
<i>pustý</i>	<i>pust</i>	AyseJDjo	<i>pustší/pustější pust – pustě/o</i>
<i>chudý</i>	<i>chud</i>	AysjJNDjo	<i>chudší/chudější chud chudost chudě/o</i>
<i>společenský</i>	<i>společen</i>	AskyD	<i>společenšější – – společensky</i>
<i>pražský</i>	<i>praž</i>	AskyOD	– – – <i>pražsky</i>
<i>otrocký</i>	<i>otroc</i>	AckyD	<i>otročtější – – otrocky</i>
<i>řecký</i>	<i>řec</i>	AckyOD	– – – <i>řecky</i>
<i>starý</i>	<i>star</i>	Ayys	<i>starší – – –</i>
	<i>stař</i>	AyiDe	– – – <i>staře, též staří</i>
	<i>stár</i>	AJ	– <i>stár</i> – –
<i>drahý</i>	<i>dra</i>	AhysNDe	<i>dražší – drahost draze</i>
<i>snadný</i>	<i>snad</i>	AnyjJ-NDoj	<i>snadnější snadna snadnost snadno/ně</i>
	<i>snaz</i> FMU=K	Ay23	<i>snazší – – –</i>
<i>cizí</i>	<i>ciz</i>	AieNDe	<i>cizejší – cizost cize</i>
<i>jarní</i>	<i>jarn</i>	AijDj	<i>jarnější – – jarně</i>
<i>zadní</i>	<i>zad</i>	Aij	<i>zadnější – – –</i>
	<i>zaz</i> FMU=K	Ay23	<i>zazší – – –</i>
<i>ostatní</i>	<i>ostatn</i>	AiODj	– – – <i>ostatně</i>
<i>vlčí</i>	<i>vlč</i>	Ai0	– – – –

Tabulka 10.7: Příklady adjektivních vzorů pro přídavná jména

10.4 Adjektivní skloňování dalších slovních druhů

Podle uvedených vzorů se skloňují ještě některá lemmata jiných slovních druhů, a to zájmena, číslovky a podstatná jména. Pro správné přiřazení kategorie Slovní druh je třeba před vlastní název vzoru přiřadit kód slovního druhu, a to P pro zájmena, C pro číslovky a N pro podstatná jména. Přídavná jména mají označení A.

V těchto případech používáme totožné vzory pro skloňování, nikoli však už část pro stupňování. Je-li specifikace rodu jiná než A (přídavné jméno), automaticky se nevytváří stupňování. Není tedy třeba explicitně ve vzoru uvádět kód 0 v části pro stupňování. Taktéž se nepoužívá derivační část vzoru.

Znak A pro skloňování přídavných jmen na začátku názvu vzoru by se opticky mohl plést s možným označením prefixu negace, ale vzhledem k tomu, že označení slovního druhu je na začátku názvu adjektivního vzoru povinné, k záměně dojít nemůže. Prefix negace následuje vždy až za označením slovního druhu, tedy kdybychom např. chtěli mít NEG=A pro celé paradigma lemmatu *nekalý*, přiřadíme mu vzor AAy. Tím zabráníme tvoření slov s dvojitým prefixem *nene-* (**nenekalý*).

U podstatných jmen je navíc třeba specifikovat rod. To učiníme připojením kódu rodu na konec vzoru. Příklady uvádí tabulka 10.8.

Lemma	Kofix	Vzor	Současný pražský vzor
<i>který</i>	<i>kte</i>	Pry	rypif
<i>takový</i>	<i>takov</i>	Py	-
<i>třetí</i>	<i>třet</i>	Ci	iccr
<i>druhý</i>	<i>dru</i>	Chy	-
<i>sterý</i>	<i>ste</i>	Cry	-
<i>šéfová</i>	<i>šéfov</i>	NFy	mdf
<i>Vřesová</i>	<i>Vřesov</i>	NFyS	mdn
<i>kapesné</i>	<i>kapesn</i>	NNy	mdn
<i>vrátný</i>	<i>vrátn</i>	NMy	mdm
<i>vrátná</i>	<i>vrátn</i>	NFy	mdf
<i>cestující</i>	<i>cestujíc</i>	NMi	jnm
<i>cestující</i>	<i>cestujíc</i>	NFi	jnm
<i>sudí</i>	<i>sud</i>	NMi	jnm

Tabulka 10.8: Příklady podstatných jmen s adjektivními vzory

11 Vzory pro příslovce

Přestože příslovce patří mezi neohebné slovní druhy, popisujeme v rámci morfologie jejich stupňování a možnost negace.

Příslovce, která se nestupňují, vzor nemají. Ve slovníku jsou zachyceny se svojí morfologickou značkou. V naprosté většině se tato příslovce ani nedají negovat. Pokud by se přece vyskytlo nějaké příslovce, ke kterému existuje negace, a přitom bylo nestupňovatelné, je třeba jeho negativní tvar zachytit ve slovníku samostatně.

Podobně jako u stupňování přídavných jmen, i v případě příslovcí necháváme tvoření superlativu a stupně *s* na morfologických nástrojích. I zde platí, že možnost tvoření 2. stupně implikuje automaticky i možnost pravidelného tvoření ostatních dvou stupňů.

Kromě toho máme, podobně jako v pražském systému, ještě dva sdružené vzory, které umožňují popsat jedním vzorem stupňování příslovcí se zakončením *-sky* a *-cky*. Tato příslovce vznikají výhradně z přídavných jmen, v kompaktně zapsaném slovníku by se tedy samostatně neměla vyskytnout. Možnost stupňování těchto typů příslovcí závisí na možnosti stupňování původních přídavných jmen.

Pro příslovce stačí jediný vzor, dali jsme mu název **adv**:

adv_x

Parametr **x** může nabývat těchto hodnot:

- s** → stupňování příslovcí se zakončením *-sky*
- c** → stupňování příslovcí se zakončením *-cky*
- ε** → pravidelné stupňování pomocí zakončení *-ji*
- 1** → pro daný kofix se tvoří jen pozitiv
- 23** → pro daný kofix se tvoří všechny stupně kromě pozitivu

Tabulka 11.1 ukazuje přehled vzorů pro příslovce.

Pravidelné mutace stupňovaných příslovcí

V tabulce 11.1 jsou též zachyceny hodnoty kategorie **Flektivní mutace** pro různá zakončení stupňovaných tvarů. Kromě zde uvedených hodnot má mnoho příslovcí ve druhém a třetím stupni dva tvary — kratší a delší, který se z kratšího tvoří přidáním *-e* na konec slovního tvaru (*blíž* — *blíže*), s mutací FMU=e pro tvar se zakončením *-e*.

U některých příslovcí dochází ke kombinacím, např. pro lemma *snadně* máme 5 mutací druhého stupně (a stejně i stupně třetího a stupně *s*): *snáz*, *snáze*, *snadněji*, *snadněj* a *snadnějc* s mutacemi po řadě: FMU=K, FMU=Ke, FMU=D, FMU=Dj, FMU=Dc.

11 Vzory pro příslovce

Vzor	Zakončení	Stupeň (DEG)	Flektivní mutace (FMU)
adv	0	1	0
	<i>-ji</i>	[23s]	0
	<i>-j</i>	[23s]	j
	<i>-jc</i>	[23s]	c
advc	<i>-cky</i>	1	0
	<i>-čtěji</i>	[23s]	0
	<i>-čtěj</i>	[23s]	j
	<i>-čtějc</i>	[23s]	c
adv _s	<i>-sky</i>	1	0
	<i>-štěji</i>	[23s]	0
	<i>-štěj</i>	[23s]	j
	<i>-štějc</i>	[23s]	c

Tabulka 11.1: Přehled vzorů příslovčí

K některým příslovcím se zakončením *-e/-ě* existují příslovce, většinou predikativní (SUB=R), se stejným kmenem a zakončením *-o*. Často se dají i stupňovat, čímž dochází k homonymii mezi stupňovanými tvary obou příslovčí.

Příbuznost takových dvojic lze vyjádřit pomocí derivačních odkazů ve slovníku.

Příklady

Lemma	Kofix	Vzor	Tvary 2. stupně
<i>rychle</i>	<i>rychle</i>	adv	<i>rychleji</i>
<i>složitě</i>	<i>složitě</i>	adv	<i>složitěji</i>
<i>pomalů</i>	<i>pomalů</i>	adv1	
	<i>pomale</i>	adv23	<i>pomaleji</i>
<i>rusky</i>	<i>rusky</i>	adv1	
<i>divoce</i>	<i>divoce</i>	adv1	
	<i>divoče</i>	adv23	<i>divočeji</i>
<i>divoko</i>	<i>divoko</i>	adv1	
	<i>divoče</i>	adv23	<i>divočeji</i>
<i>domácky</i>	<i>domácky</i>	adv _c	<i>domáčtěji</i>

Existuje poměrně velká množina příslovčí s nepravidelným stupňováním. Tato příslovce nepopisujeme pomocí vzorů. Každý slovní tvar z jejich paradigmatu má ve slovníku svůj záznam s přesným určením hodnot svých relevantních kategorií. Několik příkladů ukazují tabulky 11.2 a 11.3.

Lemma	Slovní tvar, příp. kofix	FMU	Kategorie, příp. vzor
<i>draze</i>	<i>draze</i>		DEG=1
	<i>dráž</i>		DEG=2
	<i>dráže</i>	FMU=e	DEG=2
	<i>nejdráž</i>		DEG=3
	<i>nejdráže</i>	FMU=e	DEG=3
	<i>sebedráž</i>		DEG=s
	<i>sebedráže</i>	FMU=e	DEG=s
<i>draho</i>	<i>draho</i>		DEG=1
	<i>dráž</i>		DEG=2
	<i>dráže</i>	FMU=e	DEG=2
	<i>nejdráž</i>		DEG=3
	<i>nejdráže</i>	FMU=e	DEG=3
	<i>sebedráž</i>		DEG=s
	<i>sebedráže</i>	FMU=e	DEG=s

Tabulka 11.2: Příklad rozepsání paradigmat příslovcí *draho* a *draze* pomocí jednotlivých slovních tvarů a jejich vlastností

Lemma	Slovní tvar, příp. kofix	FMU	Kategorie, příp. vzor
<i>snadno</i>	<i>snadno</i>		DEG=1
	<i>snadně</i>	FMU=D	adv23
	<i>snáz</i>	FMU=K	DEG=2
	<i>snáze</i>	FMU=Ke	DEG=2
	<i>nejsnáz</i>	FMU=K	DEG=3
	<i>nejsnáze</i>	FMU=Ke	DEG=3
	<i>sebesnáz</i>	FMU=K	DEG=s
	<i>sebesnáze</i>	FMU=Ke	DEG=s
<i>snadně</i>	<i>snadně</i>		DEG=1
	<i>snadně</i>	FMU=D	adv23
	<i>snáz</i>	FMU=K	DEG=2
	<i>snáze</i>	FMU=Ke	DEG=2
	<i>nejsnáz</i>	FMU=K	DEG=3
	<i>nejsnáze</i>	FMU=Ke	DEG=3
	<i>sebesnáz</i>	FMU=K	DEG=s
	<i>sebesnáze</i>	FMU=Ke	DEG=s

Tabulka 11.3: Příklad rozepsání paradigmat příslovcí *snadno* a *snadně* pomocí kombinace vzoru a jednotlivých slovních tvarů s vlastnostmi. Vzor **adv23** popisuje dlouhou mutaci stupňování příslovce, tedy s druhým stupněm *snadněji*.

Mutace, které jsou součástí vzoru **adv23**, se připojí k mutaci uvedené u kofixu.

12 Slovesné vzory

Česká slovesa jsou zřejmě nejsložitějším slovním druhem z hlediska popisu.

U sloves dochází k alternaci v kmenech mnohem častěji než u ostatních slovních druhů. Přitom systém koncovek není zas tak rozsáhlý, viz např. (Osolsobě, 1996).

Máme tedy na jedné straně velkou otevřenou množinu různých kofixů, na straně druhé potom uzavřenou a poměrně malou množinu možných zakončení. Pro vytváření vzorů je výhodné si možná zakončení rozdělit do několika základních podmnožin popisujících konkrétní slovesné tvary, a tyto podmnožiny potom přiřazovat kofixům, se kterými tvoří smysluplné slovní tvary.

Podobným způsobem je systém českých sloves popsán v (Osolsobě, 1996). My však budeme podmnožiny definovat přímo v názvech slovesných vzorů, podobně jako dosud u jmenných vzorů.

Slovesný vzor začíná vždy kódem pro slovní druh sloveso, tedy znakem V. Dále se skládá ze dvou částí — flektivní a derivační. V případě, že kofix lze použít jen pro jednu z obou částí, druhá může chybět.

Inspirací pro návrh flektivní části slovesných vzorů byla práce Simeona Romportla — viz (Romportl, 1970), především jeho formální způsob kódování slovesných tvarů. Romportl se však zabývá slovesnými tvary obecně, to znamená i vícečlennými, my popisujeme jen samostatné slovní tvary. Používáme také jinou množinu kategorií, takže jsme jeho způsob popisu museli poměrně značně přetvořit.

K popisu základního časování vymežíme 6 množin tvarů, jejichž kofixy se mohou lišit. Jsou to:

1. imperativ
2. prézens
3. préteritum
4. infinitiv
5. přechodník
6. pasívum

Podle nich se potom vytvářejí slovní tvary odpovídající ostatním hodnotám relevantních morfologických kategorií.

Prvních 5 kategorií je pro každé lemma povinných. Kvůli přehlednosti je zapisujeme ke každému kofixu jako pěti hodnot. V případě, že daný kofix se některé kategorie netýká, má v pěti hodnotu - (pomlčka).

Uvedenou pěti může předcházet prefix B, který znamená, že se vzor pro imperativ a pro prézens použije též pro vytvoření imperativu a budoucího času pomocí předpony *po-* (tzv. determinovaná slovesa).

Další část vzoru se týká tvoření trpného rodu, který je kódován jinak, neboť se nemusí tvořit vždy. Má tak tvar podobný vzoru derivačnímu.

Za pevnou pěticí se uvádí typ derivace podle slovních druhů následovaný hodnotou, která přesně specifikuje, jakým způsobem se z daného kofixu odvozené slovo vytváří a jaký je její flektivní vzor. Odvozeniny s nepravidelným ohýbáním se automaticky nevytvářejí.

Popíšeme teď jednotlivé množiny slovesného vzoru. U každé uvedeme i nejběžnější pravidelné mutace. Jejich celkový přehled je uveden také v tabulce 12.1.

POS	VRB	PER	NUM	Mutace	Její kód	Mutace	Její kód	Mutace	Její kód
				Příklad		Příklad		Příklad	
V	F			<i>ti, ci</i>	i	<i>ct</i>	t	<i>[jý][sz]t</i>	y
				<i>dělati, říci</i>		<i>říct, ne dělat</i>		<i>vízt, nýst, ne číst</i>	
V	I			krátce	K	dlouze	D		
				<i>plav, osvědč</i>		<i>plavej, osvědči</i>			
				<i>po-</i>	p	<i>ž</i>	z		
				<i>poběž</i>		<i>budíž, řekněmež</i>			
V	ITP			měkce	m	tvrdě	t		
				<i>mraž/mražen</i>		<i>mraz/mrazen</i>			
V	PBL			<i>t'</i>	t				
				<i>vímět, budet', vědělit'</i>					
V	P	1	P	<i>me</i>	o	<i>m</i>	K		
				<i>neseme</i>		<i>nesem</i>			
V	P	1	S	<i>i</i>	i	<i>u</i>	u	<i>ím/ám</i>	d
				<i>maží, kupuji</i>		<i>mažu, kupuju</i>		<i>sedím, vím, mazám</i>	
				<i>im</i>	k	měkce	m	tvrdě	t
				<i>sedím, vim</i>		<i>mažu</i>		<i>mazám</i>	
V	P	3	P	<i>í</i>	i	<i>ou</i>	u		
				<i>maží, kupují, sázejí</i>		<i>mažou, kupujou</i>			
				<i>j</i>	j	<i>i</i>	k		
				<i>sázej, dělaj</i>		<i>vědi</i>			
V	L			<i>l[aoiy]?</i>	l	<i>nul[aoiy]?</i>	n	<i>nul[aoiy]</i>	n
				<i>tiskl</i>		<i>tisknul</i>		<i>usnula, ne usl</i>	

Tabulka 12.1: Přehled nejběžnějších mutací v časování sloves

12.1 Flektivní vzor

12.1.1 1. pozice — Imperativ

V tabulce 12.2 jsou uvedeny kódy první pozice flektivní části slovesného vzoru a množiny zakončení, které jsou jimi popsány.

Hodnota	Množina zakončení	Příklady kofixu
0	0, <i>-me, -te</i>	<i>zruš, sud, ohlas, ohlaš</i>
j	<i>-ej, -ejme, -ejte</i>	<i>děl</i>
w	<i>-i, -ěme, -ěte</i>	<i>tiskn</i>
e	<i>-i, -eme, -ete</i>	<i>pošl</i>

Tabulka 12.2: Tabulka vzorů pro rozkazovací způsob

Pravidelné mutace

Ve slovesném imperativu se vyskytuje poměrně značné množství typů mutací.

Jednak jsou to krátké a dlouhé tvary (končící na *-ej*), např. *plav* — *plavej*, *maž* — *mazej*. Pomocí délky lze rozlišovat i rozkazovací způsoby, jejichž delší mutace končí ve 2. osobě singuláru na *-i*: *osvědč* — *osvědčí*, *přivab* — *přivábí*. Kratší mutace mají FMU=K, delší FMU=D.

Někdy se mutace imperativu liší v tvrdosti (FMU=t) — *rozmraz*, *zaráz* a měkkosti — *rozmraž*, *zaráž* s FMU=m.

Dalším případem je mutace v imperativu determinovaných sloves s předponou *po-*, které přiřazujeme FMU=p. Jde např. o tyto dvojice: *nes* — *pones*, *jdi* — *pojď*, *jeď* — *pojeď*, *běž* — *poběž*, *slyš* — *poslyš*, ale už ne *posekej*, neboť existuje lemma *posekat*, ani *popojeď* od lemmatu *popojet*.

Mutace archaická s FMU=z připojuje k imperativu *-ž* (např. *vymyslemež*, *dejž*, *chraňtež*).

12.1.2 2. pozice — Prézens

Druhá pozice flektivní části slovesného vzoru je popsána v tabulce 12.3.

Hodnota	Množina zakončení	Příklady kofixu
o	<i>-u, -eš, -e, -eme/em, -ete, -ou</i>	<i>kop</i>
i	<i>-u/i, -eš, -e, -eme/em, -ete, -í/ou</i>	<i>kryj, maž</i>
a	<i>-ám, -áš, -á, -áme, -áte, -ají/aj</i>	<i>děl</i>
s	<i>-ím/im, -íš, -í, -íme, -íte, -í/ej/ejí</i>	<i>pros, sáz</i>
p	<i>-ím/im, -íš, -í, -íme, -íte, -í/ěj/ějí</i>	<i>trp, vypráv</i>

Tabulka 12.3: Tabulka vzorů a zakončení tvarů přítomného času

Označení vzorů pro přítomný čas je odvozeno od 3. osoby množného čísla, kromě vzorů *s* a *p*, které je nutné rozlišovat pouze kvůli tvarům *prosejí*, *prosej*, *trpějí*, *trpěj*.

Pravidelné mutace

Z tabulky 12.3 jsou vidět pravidelné mutace, ke kterým v přítomnosti dochází. Nejběžnější je zřejmě mutace 1. osoby plurálu *-eme/-em*. Spisovná koncovka je *-me* (*bereme*), ale u všech sloves první, druhé a třetí třídy, které tvoří první osobu množného čísla pomocí *-eme*, má navíc hovorový tvar *-m* (*berem, kynem, kryjem*). Spisovnou koncovku mutací neoznačujeme, koncovka *-m* má FMU=K (kratší).

V obecné češtině také často dochází ke krácení samohlásky. Nejčastěji jde o 1. osobu jednotného čísla (*musím — musím, platím — platím*, apod.), nebo 3. osobu čísla množného (např. *vědí — vědi*). Rozlišujeme zde hodnotu kategorie **Flektivní mutace** krátkou (FMU=k) a dlouhou (FMU=d).

Mutace 3. osoby plurálu mohou být dvojího typu: $-[aeě]jí/-[aeě]j$ a $-í/-ou$. U prvního typu považujeme spisovnou mutaci s $-jí$ za nulovou, obecně česká mutace $-j$ má hodnotu FMU=j (*krájejí — krájej*). U druhého typu jde o slovesa 1. třídy (*mažou — maží, pečou — pečí*) a 3. třídy (*kupují — kupujou, kryjí — kryjou*) s FMU=i resp. FMU=u. Např. lemma *sázet* má v 3. osobě množného čísla tři mutace: *sázejí* (FMU=0), *sázej* (FMU=j) a *sází* (FMU=i). Podobné jsou i mutace 1. osoby singuláru $-í/-u$ u sloves 1. a 3. třídy (*mažu — maží, kupuji — kupuju*), mají proto stejné hodnoty kategorie **Flektivní mutace**.

Časté je také kolísání časování přítomného času mezi pátou a první třídou (*kopu — kopám*). Mutace zde rozlišujeme podle krátké / dlouhé samohlásky. Zakončení z množiny *o* mají FMU=k, zakončení z množiny *a* FMU=d.

Zde se mohou navíc projevovat následující kmenové změny, které přidávají další hodnotu do kategorie **Flektivní mutace**:

š-s (*češu/češi — česám*)

ž-z (*mažu/maži — mazám*)

č-k (*skáču/skáči — skákám*).

Dohromady s mutací $-í/-u$ sloves 1. třídy mají v 1. osobě jednotného čísla trojí hodnotu FMU: Příklady nalevo mají FMU=kmu/kmi (mutace i/u, krátká, měkká), napravo FMU=dt (mutace dlouhá, tvrdá).

Ještě zmíníme pravidelné archaické mutace s koncovým $-ť$, které se vyskytují i u minulého přičestí (např. *udělalť, udělámť*). V pražském systému se s nimi počítá, mají speciální hodnotu na 2. pozici. Přiřazujeme jim FMU=t.

12.1.3 3. pozice — Préteritum

Minulý čas je velmi pravidelný, má dvě možnosti tvoření. První je jednoduchá, druhá umožňuje stejnému kofixu přiřadit dvě různé množiny zakončení, jak je naznačeno v tabulce 12.4.

Hodnota	Množina zakončení	Příklady kofixu
l	<i>-l, -la, -lo, -lí, -ly, -la</i>	<i>nes</i>
n	<i>-l, -la, -lo, -lí, -ly, -la</i> <i>-nul, -nula, -nulo, -nuli, -nuly, -nula</i>	<i>usch</i>
m	<i>-la, -lo, -lí, -ly, -la</i> <i>-nul, -nula, -nulo, -nuli, -nuly, -nula</i>	<i>us</i>

Tabulka 12.4: Množina zakončení tvarů minulého přičestí

Hodnota **m** se od **n** liší tím, že neobsahuje tvoření mužského singuláru.

Pravidelné mutace přičestí minulého činného

Jde o mutace se zakončením $-l — -nul$ u sloves *trhl — trhnul*. U některých sloves tato mutace nepřípadá v úvahu u mužského rodu, u ostatních rodů však

ano (**usl* — *usnul*, ale *usla* — *usnula*). Tento rozdíl vyjadřuje dvojí hodnota parametru ve vzoru — viz hodnoty *m* a *n* v tabulce 12.4. Tvary s *-n-* mají FMU=*n*, stejně jako u podobného případu v infinitivu.

Další mutace se projevují u přičestí minulého zakončeného na *Kl*, kde *K* je souhláska. Zde se vytvářejí nespisovné mutace bez koncového *-l*, ovšem jen u mužského rodu (*vedl* — *ved*, *pletl* — *plet*, *nesl* — *nes*, *zábl* — *záb*, *kopl* — *kop*, *vezl* — *vez*, *klekl* — *klek*, *pomohl* — *pomoh*, *všiml* — *všim*). Mutace bez koncového *-l* značíme jako kratší — FMU=*K*.

12.1.4 4. pozice — Infinitiv

I infinitivy jsou velmi pravidelné. Množinu zakončení uvádí tabulka 12.5. Mutace infinitivu *odřící* — *odřeknout* je třeba vyjadřovat jako dvě položky s různými kofixy a různým zakončením. Označení mutace se přidává ke kofixu.

Hodnota	Množina zakončení	Příklady kofixu
t	-t, -ti	<i>kopa</i>
c	-ci, -ct	<i>mo</i>

Tabulka 12.5: Zakončení tvarů infinitivu

Pravidelné mutace slovesného infinitivu

Všechny infinitivy mají dvě mutace: *-t* s FMU=0 a *-ti* s FMU=i (*být* — *býti*), nebo *-ci* s FMU=i a *-ct* s FMU=t (*moci* — *moct*). Tento typ mutace je tedy přítomen vždy i u následujících typů a kombinuje se s nimi.

Další mutace infinitivu se týkají jen několika málo sloves, např. *hanit* — *hanět*, *myslit* — *myslet*, *bydlit* — *bydlet*. U těchto sloves se projevuje stejná flektivní mutace ještě v přičestí minulém (*myslil* — *myslel*). Označujeme ji FMU=y resp. FMU=e.

Podobné jsou mutace *-ést* — *-íst*, *-ézt* — *-ízt* a *-éct* — *-[íy]ct*, např. u sloves *vést* — *víst*, *vézt* — *vízt* a *péct* — *píct*, případně *téct* — *týct*. Tyto mutace se však neprojevují v minulém čase. Mají stejné označení, tedy FMU=y resp. FMU=e. U těchto sloves zřejmě nedochází ke kombinované mutaci FMU=iy (**pícti*).

U mutací typu *řící* — *řeknout*, *začít* — *začnout* mají tvary s *-n-* FMU=*n*. Tato mutace se také projevuje v minulém přičestí.

12.1.5 Přechodník P

Přechodníky, ač užívané dnes již zřídka, je třeba také umět rozpoznat, výjimečně i tvořit. Tabulka 12.6 uvádí množiny jejich zakončení.

První tři řádky tabulky popisují přechodník přítomný, poslední řádek pak přechodník minulý. Není třeba uvádět do vzoru, o jaký typ přechodníku se jedná. Jde o hodnotu kategorie Slovesný tvar, která se správně přiřadí příslušným vzorem, tedy VRB=p pro vzory *a*, *e*, *w*, VRB=m pro vzor *v*.

Přechodníky se pravidelně využívají k tvoření přídavného jména slovesného se zakončením *-cí* od přechodníku přítomného nebo *-vší* od přechodníku minulého. Oba tyto typy přídavných jmen se skloňují podle měkkého adjektivního

Hodnota	Množiny zakončení	Příklady kofixu
a	-a, -ouc, -ouce	<i>drhn</i>
e	-e, -íc, -íce	<i>vyprávěj</i>
w	-ě, -íc, -íce	<i>hromad</i>
v	-v, -vši, -vše	<i>nahromadí</i>

Tabulka 12.6: Vzory pro vytváření přechodníku

vzoru $Ai0$, tedy bez možnosti stupňování. Tvoření přechodníku tak automaticky implikuje tvoření těchto přídavných jmen.

12.1.6 Trpný rod T

Hodnoty kódů derivačních vzorů pro trpný rod ukazuje tabulka 12.7.

Hodnota	Množina zakončení	Příklady kofixu
Tn	-n, -na, -no, -ny, -ní, -nu	<i>nese</i>
Ta	-án, -ána, -áno, -ány, -áni, -ánu	<i>ps</i>
Tt	-t, -ta, -to, -ty, -ti, -tu	<i>táhnu</i>

Tabulka 12.7: Derivační vzory pro trpný rod

Hodnoty **n** a **a** jsou rozlišeny kvůli jednoduchému vytvoření lemmatu. Lemmatem přídavných jmen slovesných totiž není sloveso, ale přídavné jméno (viz kap. 10.2.1), konkrétně jeho dlouhý tvar. Takže např.

$\lambda(\textit{nesen}) = \textit{nesený}$,

$\lambda(\textit{táhnut}) = \textit{táhnutý}$,

$\lambda(\textit{psán}) = \textit{psaný}$ (změna v délce samohlásky v kmeni).

Mutace deverbativ

Některá slovesa tvoří dvojí trpný rod, bez měkčení kmenové souhlásky (*mrazen*) a s měkčením (*mražen*). Od obou lze odvozovat další deverbativa — podstatná jména (*mrazení* a *mražení*), přídavná jména (*mrazený*, *mražený*) a případně i příslovce. Tyto mutace se týkají i imperativu (*mraz* — *mraž*) (viz též 12.1.1). Přiřazujeme jim tyto hodnoty kategorie **Flektivní mutace**: FMU=m pro změkčené varianty, FMU=t pro nezměkčené.

12.2 Derivační vzory

Pražský systém slovesných vzorů, který vytvořil Jan Hajič, má velmi bohatý repertoár odvozování příbuzných deverbativ pomocí přípon. Náš systém by měl být schopen odvodit veškeré smysluplné derivace, které odvozuje současný pražský systém. Podívejme se tedy nejprve na typy lemmat, které jsou v pražském systému zachyceny. Lze je přehledně zpracovat do tabulky 12.8.

Tabulka má dva hlavní sloupce: první popisuje zakončení tvaru základního slovesa a jeho odvozenin, druhý zakončení tvaru slovesa iterativního a jeho odvozenin.

Základní		Iterativní	
-t		-[íá]vat I	
-cí Ac → Ai0		-[íá]vací → Ai0	
-cí/ší automaticky z přechodníku → Ai0		-[íá]vající → Ai0	
-ní/tí N[nt] → st		-[íá]vání → st	
-ný/tý/lý A[nt1] [j0] → Ay[j0]	-telný Ae → Ayj	-[áí]vaný → Ay[j0]	-[áí]vatelný → Ayj
-ně/tě/le D[nt1] [j0] → adv1?	-telně De → adv	-[áí]vaně → adv1?	-[áí]vatelně → adv
-nost/tost/lost O[nt1] → kt	-telnost Oe → kt	-[áí]vanost → kt	-[áí]vatelnost → kt

Tabulka 12.8: Pravidelné slovesné derivace popsané v pražském systému vzorů

Ve spodní části buněk jsou kódy, které mají následující význam. Kód před šipkou je součástí derivační části slovesného vzoru v případě, že se příslušné deverbativum tvoří. Budeme mu říkat derivační pravidlo. Proč není toto pravidlo uvedeno ve všech buňkách, vysvětlíme vzápětí. Kód za šipkou je kód již zavedeného vzoru, podle kterého se derivované lemma ohýbá.

Tak např. „Ac → Ai0“ znamená, že pomocí derivačního kódu (pravidla) Ac a příslušného kofixu se vytvoří odvozené lemma se zakončením *-cí* a bude se ohýbat podle adjektivního vzoru Ai0 (měkké skloňování bez možnosti stupňování).

Třetí řádek obsahuje přídavná jména vytvořená pravidelně z přechodníků, proto zde žádné derivační pravidlo není třeba.

Alternativa [nt1] vyjadřuje trojí možné zakončení deverbativ, alternativa [j0] potom možnost nebo nemožnost stupňování odvozeného přídavného jména nebo příslovce. Pouze u přídavných jmen se zakončením lemmatu *-telný* volbu stupňování nepřipouštíme, neboť se domníváme, že zde je stupňování možné vždy.

Celý druhý sloupec se týká jen iterativních sloves, popisuje tedy jen lemmata typu *kupovávat*, *lehávat*, *končívát*.

Druhý sloupec obsahuje derivační pravidlo (I) pouze v první řádce. Vyjadřuje možnost tvoření iterativa od slovesa základního (z levého sloupce). Vzory, podle kterých se odvozené iterativní sloveso časuje, jsou uvedeny dále v oddíle 12.2.4.

Ostatní buňky pravého sloupce jsou pravidelně tvořené z tvaru iterativa, neobsahují tedy derivační pravidla (kódy před šipkou). Domníváme se totiž, že není třeba ve vzorech udržovat explicitně toto velké množství odvozenin od iterativních sloves a příslušná derivační pravidla zahrnujeme přímo do pravidla I.

Jestliže tedy existuje k slovesu příbuzné sloveso iterativní, pravidelně se z něj mohou tvořit všechny tvary naznačené ve druhém sloupci. Uvědomujeme si, že i zde může docházet ke generování velmi nepravděpodobných tvarů. Po několika sondách do korpusů i na internetu jsme však zjistili, že by nebylo rozumné derivační pravidla od iterativních sloves omezovat.

V následujících oddílech podrobně rozebereme jednotlivé buňky tabulky 12.8 a uvedeme příklady.

12.2.1 Přídavná jména slovesná A

Hodnoty kódů derivačních vzorů pro přídavná jména slovesná ukazuje tabulka 12.10.

Hodnota	Zakončení lemmatu	Vzor	Příklady kofixu
An	-ný	AyO	<i>unese</i>
At	-tý	AyO	<i>kopnu</i>
Al	-lý	AyO	<i>zemře</i>
Ae	-telný	AyO	<i>snesi</i>
Ac	-cí	AiO	<i>pozměňova</i>

Tabulka 12.9: Derivační vzory pro přídavná jména slovesná

Všechna přídavná jména slovesná odvozená podle uvedených vzorů mají automaticky přiřazenou hodnotu poddruhu SUB=S (deverbativní).

Kromě právě popsaného odvození přídavného jména slovesného se ještě tvoří přídavná jména od přechodníků. Pro ně není třeba zavádět žádné vzory, neboť se tvoří pravidelně, viz 12.1.5. Skloňují se podle vzoru AiO a mají poddruh SUB=G (od přechodníku přítomného), nebo SUB=M (od přechodníku minulého).

12.2.2 Deverbativní příslovce D

Hodnota	Zakončení lemmatu	Vzor	Příklady kofixu
Dn	-ně	adv	<i>unese</i>
Dt	-tě	adv	<i>rozvinu</i>
Dl	-le	adv	<i>ochraptě</i>
De	-telně	adv	<i>snesi</i>

Tabulka 12.10: Derivační vzory pro deverbativní příslovce

Příslovce lze tvořit od většiny přídavných jmen slovesných, ale ne vždy. Protipříkladem je přídavné jméno *kopnutý* s neexistujícím příslovcem **kopnutě*. Ani přídavná jména se zakončením *-ný* nemají automaticky příslovce — *zahájený*, ale **zahájeně*. Na druhou stranu ale, když příslovce existuje, lze ho odvodit pravidelně vždy záměnou koncové dlouhé samohlásky *-ý* za *-ě* v případech *-ný*, *-tý* a za *-e* v případech *-lý* (záviset — závislý — *závisle*). Ze zakončení *-cí* se příslovce nevytvářejí.

12.2.3 Podstatná jména slovesná N/0

Podstatná jména slovesná jsou dvojího druhu:

1. se zakončením *-í*, kód N, skloňování podle vzoru **st**,
2. se zakončením *-ost*, kód 0, skloňování podle vzoru **kt**.

Podstatná jména slovesná se zakončením *-í* se odvozují z tvarů trpného rodu prostým připojením zakončení *-í* v případě hodnot **n** a **t** (*nesení, táhnutí*), nebo zakončení *-aní* pro hodnotu **a** (*psán — psaní*). Mají poddruh SUB=S (deverbativní). Poddruh druhého typu, se zakončením *-ost*, není deverbativní (SUB=0). Tabulka 12.11 rozepisuje jednotlivé možnosti podle kofixu.

Hodnota	Zakončení lemmatu	Vzor	Příklady kofixu
Nn	<i>-ní</i>	st	<i>unese</i>
Nt	<i>-tí</i>	st	<i>kopnu</i>
On	<i>-nost</i>	kt	<i>připojiště</i>
Ot	<i>-tost</i>	kt	<i>netknu</i>
Ol	<i>-lost</i>	kt	<i>poblouďi</i>
Ne	<i>-telnost</i>	kt	<i>nedotknu</i>

Tabulka 12.11: Derivační vzory pro podstatná jména slovesná

Derivace trpného rodu, slovesného podstatného jména, přídavného jména a příslovce se může zaznamenat do vzoru najednou, jestliže mají stejný kofix a jestliže současně následný parametr, tedy **n**, **t** nebo **l**, je shodný. Pro vytváření deverbativ jsou tedy teoreticky možné např. tyto kombinace (na pořadí nezáleží): AN, AD, AT, ONT, DN, ADN, ADT, ADTN, následované (alespoň) jedním z parametrů **n**, **t**, **l**.

12.2.4 Iterativní sloveso

Časování odvozeného iterativního slovesa, tedy derivační pravidlo I ukazuje tabulka 12.12.

Konec kofixu	Vzor	Příklad
<i>-áv</i> <i>-ív</i>	ja--- (imperativ, prézens)	<i>kupovávej, kupovávám</i> <i>chodívej, chodívám</i>
<i>-áva</i> <i>-íva</i>	--lt- (préteritum, infinitiv)	<i>kupovával, kupovávat</i> <i>chodíval, chodívat</i>
<i>-ávaj</i> <i>-ívaj</i>	----e (přechodník přítomný)	<i>kupovávaje</i> <i>chodívaje</i>

Tabulka 12.12: Časování odvozeného iterativního slovesa

Pro odvozeniny z iterativních sloves se použijí stejná pravidla jako pro odvozeniny uvedené výše.

12.3 Sdružené slovesné vzory

Flektivní část vzoru je poměrně složitá proto, že v řadě sloves dochází ke změnám ve kmeni. Kromě toho však existuje velké množství sloves, která jsou

naopak velmi pravidelná a uvedený popis flexe, ač obecný, je pro ně zbytečně složitý. Ke stejnému závěru dospěl i Hajič, když pro ně ve své disertační práci (Hajič, 1994) zavedl sdružené vzory, které i úspěšně implementoval.

Jeho systém sdružených vzorů tak můžeme převzít, ovšem pouze jeho flektivní část. Pravidelné derivace, které jsou ke vzorům v pražském systému napevno připojeny, nevyužijeme. Místo toho použijeme pro derivace systém derivačních vzorů uvedených v oddíle 12.2, zejména pak v tabulce 12.8. Tato tabulka byla sestavena právě na základě derivačních částí sdružených vzorů. Pomocí nového systému si z ní však můžeme vybírat jen to, co se k danému kofixu hodí.

V tabulce 12.13 uvádíme seznam sdružených flektivních vzorů, které přebíráme z pražského Hajičova systému slovesných vzorů. U každého sdruženého vzoru jsou uvedeny kofixy a nově zavedené flektivní vzory, jejichž spojením je možno vygenerovat celé paradigma.

Do tabulky jsme nezahrnuli vidovou dvojici sdružených vzorů *itxd* a *itxn*, protože tyto vzory se liší od ostatních tím, že neobsahují předpis pro tvoření přechodníku přítomného, imperativu a trpného rodu. Umožňují tak jejich nepravidelné tvoření. V našem systému vzorů to nepotřebujeme.

V posledním sloupci tabulky je pro každý vid uveden příklad s kofixem odděleným pomlčkou od zakončení lemmatu.

Ve vzorech *noutd* a *noutn* je zahrnuta varianta ve tvoření minulého přičestí, a to uvedením hodnoty kategorie **Flektivní mutace** u kofixu.

Ke sdruženým vzorům je možno přidávat derivační vzory stejně jako ke vzorům jednoduchým, a tím vybrat pro konkrétní kofix jen ty derivace, které skutečně existují. Tím se můžeme vyhnout tvoření derivací jako např. **ženěně*, **ženěnost*, **ženívací*, **ženívanost*, a dalších, které se derivují zcela pravidelně z pražského vzoru *nitn*.

12 Slovesné vzory

Lemma	Nedokonavý vzor		Dokonavý vzor		Příklad
	Kofix	Vzor	Kofix	Vzor	
<i>Xat</i>	atn		atd		<i>děl-at / uděl-at</i>
	<i>X</i>	ja---	<i>X</i>	ja---	
	<i>Xa</i>	--lt-	<i>Xa</i>	--ltv	
	<i>Xaj</i>	----e	—	—	
<i>Xovat</i>	ovatn		ovatd		<i>děk-ovat / poděk-ovat</i>
	<i>Xuj</i>	Oi--e	<i>Xuj</i>	Oi---	
	<i>Xova</i>	--lt-	<i>Xova</i>	--ltv	
<i>Xet</i>	etn		etd		<i>kráj-et / nakráj-et</i>
	<i>X</i>	js---	<i>X</i>	js---	
	<i>Xe</i>	--lt-	<i>Xe</i>	--ltv	
<i>Xět</i>	wtm		wtd		<i>reziv-ět / zreziv-ět</i>
	<i>Xě</i>	j-lt-	<i>Xě</i>	j-ltv	
	<i>X</i>	-p---	<i>X</i>	-p---	
<i>Xěj</i>	----e		—	—	
	ditn		ditd		<i>chla-dit / ochla-dit</i>
	<i>Xd'</i>	0----	<i>Xd'</i>	0----	
<i>Xd</i>	-p--w	<i>Xd</i>	-p---		
<i>Xdi</i>	--lt-		<i>Xdi</i>	--ltv	
	titn		titd		<i>čtvr-tit / rozčtvr-tit</i>
	<i>Xt'</i>	0----	<i>Xt'</i>	0----	
<i>Xt</i>	-p--w	<i>Xt</i>	-p---		
<i>Xti</i>	--lt-		<i>Xti</i>	--ltv	
	nitn		nitd		<i>že-nit / ože-nit</i>
	<i>Xň</i>	0----	<i>Xň</i>	0----	
<i>Xn</i>	-p--w	<i>Xn</i>	-p---		
<i>Xni</i>	--lt-		<i>Xni</i>	--ltv	
	iten		ited		<i>zuř-it / rozzuř-it</i>
	<i>X</i>	0s--e	<i>X</i>	0s---	
<i>Xi</i>	--lt-	<i>Xi</i>	--ltv		
<i>Xit</i>	itin		itid		<i>barv-it / obarv-it</i>
	<i>X</i>	wp--w	<i>X</i>	wp---	
	<i>Xi</i>	--lt-	<i>Xi</i>	--ltv	
<i>Xit</i>	itOn		itOd		<i>křiv-it / zkřiv-it</i>
	<i>X</i>	Op--w	<i>X</i>	Op---	
	<i>Xi</i>	--lt-	<i>Xi</i>	--ltv	
<i>Xnout</i>	noutn		noutd		<i>sch-nout / usch-nout</i>
	<i>Xn</i>	wo--a	<i>Xn</i>	wo---	
	<i>X</i>	--n--	<i>X</i>	--n--	
	<i>Xnou</i>	---t-	<i>Xnou</i>	---t-	

Tabulka 12.13: Přepis pražských sdružených slovesných vzorů pomocí nových flektivních vzorů. Znak X v tabulce označuje kofix, ke kterému se sdružené vzory vztahují.

12.3.1 Příklady

Tabulka 12.14 ukazuje několik příkladů, jak se používají slovesné vzory pro popis flexe i derivací. Příklad lemmatu *zavřít* ukazuje, jak se kombinují vzory a konkrétní morfologické značky v záznamu jednoho lemmatu. Poslední příklad ukazuje řešení vícenásobného lemmatu.

Lemma	Kofix/Tvar	Vzor/Značka	Odvozená lemmata
<i>zhasnout</i>	<i>zhas</i> <i>zhasnu</i>	Vnoutd VTANtAe	celé paradigma <i>zhasnut, zhasnutý, zhasnutí, zhasnutelný</i>
<i>hýbat</i>	<i>hýb</i> <i>hýba</i> <i>hýbá</i>	Vatn VAcADMe VNnI	celé paradigma <i>hýbací, hýbatelný, hýbatelně, hýbatelnost</i> <i>hýbání, hýbávat</i>
<i>zavřít</i>	<i>zavř</i> <i>zavře</i> <i>zavří</i> <i>zavřeno</i> FMU=d <i>zavru</i> FMU=t <i>zavrou</i> FMU=t	Viu--- V--1-vTADNnNo V---t- AO-D-NS1-----TJ-- VO-D--S--1----- VO-D--P--3-----	imperativ, přezens <i>zavřel, zavřev, zavřen, zavřený/ně/ní/nost</i> infinitiv flektivní dlouhá mutace jmenného tvaru flektivní tvrdá mutace 1.os. sg flektivní tvrdá mutace 3.os. pl
{ <i>ukrást</i> , <i>ukradnout</i> }	<i>ukrad</i> <i>ukrade</i> <i>ukrás</i> FMU=d <i>ukradnu</i> FMU=n	Vnoutd VTANn V---t- VTANT	celé paradigma <i>ukraden/ný/ní</i> <i>ukrást</i> (mutace s dlouhou samohláskou) <i>ukradnut/tý/tí</i> (flekt. mutace s <i>n</i>)

Tabulka 12.14: Příklady slovesných vzorů

13 Vzory zájmen a číslovek

Číslovky a zájmena se většinou mohou skloňovat podle adjektivních vzorů, proto jsme se o nich zmínili již v kapitole 10, oddíle 10.4. Na tomto místě jen připomeneme, že se používá tvrdý i měkký adjektivní vzor bez možnosti stupňování. Před adjektivní vzor se uvádí kód slovního druhu, tedy C pro číslovky a P pro zájmena.

13.1 Číslovky

13.1.1 Číslovky základní

Skloňování číslovek základních je většinou nepravidelné. Zavádíme vzory jen pro lemmata se zakončením *-t* (*jedenáct* až *dvacet*, *třicet*, atd., a tvary typu *dvaapadesát*).

Vezmeme-li v úvahu i nekodifikované tvary typu *jedenácte*, *padesáte* a typu *dvacíti*, musíme použít dva vzory. Nazveme je podle nejmenší číslovky, která se podle vzoru skloňuje, C11 a C20.

Vzory uvádíme jako množinu trojic ⟨zakončení, morf. charakteristika, FMU⟩. Morfologická charakteristika zahrnuje kvůli přehlednosti místo celé morfologické značky jen hodnoty kategorií Pád a Číslo.

$$\begin{aligned} \text{C11} = \{ & \langle 0, \text{CAS}=[145] \text{ NUM}=\text{S}, 0 \rangle, \\ & \langle -te, \text{CAS}=[145] \text{ NUM}=\text{S}, \text{D} \rangle, \\ & \langle -ti, \text{CAS}=[2367] \text{ NUM}=\text{P}, 0 \rangle \} \end{aligned}$$

$$\begin{aligned} \text{C20} = \{ & \langle -et, \text{CAS}=[145] \text{ NUM}=\text{S}, 0 \rangle, \\ & \langle -eti, \text{CAS}=[2367] \text{ NUM}=\text{S}, 0 \rangle, \\ & \langle -íti, \text{CAS}=[2367] \text{ NUM}=\text{P}, \text{d} \rangle \} \end{aligned}$$

Ostatní základní číslovky se buď skloňují podle substantivních vzorů (např. *milion*), nebo mají nepravidelné skloňování, které je vyřešeno výjimkou (v Praze rozepsáním slovních tvarů, v Brně zvláštním vzorem).

Tabulka 13.1 uvádí příklad pro číslovky *patnáct* a *dvacet*.

13 Vzory zájmen a číslovek

Kofix	Zakončení	Pád CAS	Číslo NUM	Flektivní mutace FMU	Slovní tvar
<i>patnáct</i>	ε	[145]	S	0	<i>patnáct</i>
	-e	[145]	S	D	<i>patnácte</i>
	-i	[2367]	P	0	<i>patnácti</i>
<i>dvacet</i>	-et	[145]	S	0	<i>dvacet</i>
	-eti	[2367]	S	D	<i>dvaceti</i>
	-íti	[2367]	P	0	<i>dvacíti</i>

Tabulka 13.1: Příklad tvarů číslovek *patnáct* a *dvacet*

Lemma	Kofix	Vzor	Poznámka
<i>první</i>	<i>prvn</i>	Ci	
	<i>prvn</i>	Cy23	<i>prvnější</i>
<i>druhý</i>	<i>dru</i>	Chy	<i>druhý, druzí</i>
<i>pátý</i>	<i>pát</i>	Cy	<i>pátý, pátí</i>
<i>dvojí</i>	<i>dvoj</i>	Ci	
<i>paterý</i>	<i>pater</i>	Cyy	* <i>pateří</i> v 1pl živ. se neuzívá

Tabulka 13.2: Příklady řadových a druhových číslovek se vzory. Číslovka *první* je zřejmě jediná, kterou lze stupňovat.

13.1.2 Číslovky řadové a druhové

Číslovky řadové a druhové se skloňují podle adjektivních vzorů. Příklady ukazuje tabulka 13.2.

13.1.3 Číslovky úhrnné a souborové

Mezi těmito číslovkami je poměrně nejasný rozdíl. Často se považují za jeden druh. Při jejich skloňování dochází k přechodům mezi oběma poddruhy, dokonce se může připlést ještě zakončení číslovek druhových, jak ukazuje tabulka 13.3 s příkladem skloňování *5 dveří* (záměrně nevypisujeme číslovku 5 slovy). Neurčíme zde ani hodnotu kategorie Číslo.

Pád	Alternativní tvary		
1	<i>paterý dveře</i>	<i>pateré dveře</i>	<i>patero dveří</i>
2	<i>patera dveří</i>	<i>paterých dveří</i>	<i>patero dveří</i>
3	<i>pateru dveří</i>	<i>paterým dveřím</i>	<i>patero dveřím</i>
4	<i>paterý dveře</i>	<i>pateré dveře</i>	<i>patero dveří</i>
6	<i>pateru dveří</i>	<i>paterých dveřích</i>	<i>patero dveřích</i>
7	<i>paterem dveří</i>	<i>paterými dveřmi</i>	<i>patero dveřmi</i>

Tabulka 13.3: Příklad neostré hranice mezi číslovkami úhrnnými, souborovými a druhovými

Číslovky druhové jsme vydělili zvlášť kvůli možnosti jejich čistě adjektivního skloňování. Číslovky úhrnné a souborové zde pro jejich obtížné vymezení rozdělávat nebudeme. Vzory by se mohly jmenovat opět podle kombinace kódů kategorií *Slovní druh* a *Poddruh*, tedy *Cu* pro číslovky úhrnné a *Cs* pro číslovky souborové. Je třeba pouze rozhodnout, které tvary do jednotlivých poddruhů patří.

13.1.4 Číslovky násobné, opakovací a výčtové

Číslovky násobné, opakovací a výčtové mají charakter příslovce a jako takové se neskloňují, nemají tedy ani vzor.

13.1.5 Číslovky dílové

Číslovky dílové jsou jen tři, totiž *půl*, *čtvrt* a *třet*. *Půl* je nesklonné¹, druhé dvě, ač mají substantivní skloňování, pojmáme jako výjimky a nepřičítáme jim speciální vzor.

13.2 Zájmena

13.2.1 Zájmena substantivní

Zájmena substantivní určitá nemají společný vzor.

Vytvoříme vzory *kdo*, *co* pro zájmena substantivní tázací a vztažná. Podle nich se skloňují i ta zájmena substantivní neurčitá a záporná, která vznikla z těchto dvou zájmen pomocí předpony nebo přípony, tedy např. *kdekdo*, *nikdo*, *všelico*, *kdokoliv*, *cosi*. Předponu zde opět pojmáme šíře, než je obvyklé, jako počáteční řetězec.

Zájmena s předponou jsou bezproblémová. V podstatě bychom s nimi mohli nakládat stejně, jako s předponovým guessrem (viz poznámka pod čarou na str. 70). To však neděláme, neboť se jedná o slova poměrně častá.

Se zájmeny vzniklými pomocí přípony je třeba zacházet odlišně, neboť zde dochází k flexi „uvnitř“ slovního tvaru (*cosi*, *čehosi*, *čemusi*, atd.).

Máme v zásadě dvě možnosti. Buď každé takové zájmeno pojmeme jako výjimku a umístíme do slovníku všechny jeho slovní tvary s morfologickými značkami. Tak je to v současném pražském slovníku. Druhá možnost přenechává práci morfologickým nástrojům.

První přístup je obecnější, neboť je veškerá informace uložena ve slovníku. Nástroje, které se vyvíjejí pro takový slovník, jsou potom obvykle snadno přenositelné na jiné jazyky.

Druhý přístup je zase elegantnější a přehlednější.

Rozhodli jsme se zůstat u současné praxe a speciální vzory pro skloňování těchto zájmen (zatím) nezavádět.

¹Otázkou, zda slovní tvar *půli* náleží číslovce *půl*, nebo podstatnému jménu *půle*, se zde zabývat nebudeme.

13.2.2 Zájmena přivlastňovací

U některých zájmen přivlastňovacích můžeme využít adjektivního vzoru Pi. Jsou to tato zájmena:

- určité *její*,
- tázací a vztažné *čí*,
- záporné *ničí*,
- neurčitá *něčí*, *všeličí*, a další, viz tabulka 4.1.

Ostatní zájmena jsou výjimkami.

Neurčitá a záporná zájmena vzniklá z lemmatu *čí* pomocí přípon (*čísi*, *čípak* a další) se řeší stejně jako podobně tvořená zájmena substantivní, i když i zde by se dalo uvažovat o systémovém řešení (viz oddíl 13.2.1).

13.2.3 Zájmena ukazovací a vymežovací

Ta zájmena ukazovací a vymežovací, která mají adjektivní tvar, tedy *takový*, *onaký*, *taký*, *každý*, *samý*, *všeliký*, *veškerý*, se skloňují podle tvrdého adjektivního vzoru Pky, Pry a Py, podle zakončení svého lemmatu.

Pro ostatní se vzor nevytváří.

13.3 Ostatní zájmena

Zájmena zařazená do poddruhu ostatní (SUB=0) mají adjektivní skloňování podle vzorů Pky (např. *jaký*, *všelijaký*), Pry (např. *který*, *leckterý*) a Py (např. *žádný*).²

Problém se zájmeny odvozenými pomocí přípon (*jakýsi*, *kterýkoli*, ...) se řeší stejně jako u zájmen substantivních (viz oddíl 13.2.1).

²Ač jsme na několika místech tvrdili, že se nechceme zabývat výjimkami, upozorňujeme zde na neobvyklý slovní tvar *žáden*. Vyskytuje se jen v mužském rodě (u tvaru *žádna*, který se též vyskytuje, jde pravděpodobně vždy o překlep, stejně jako *žádno*, *žádni*, *žádnou* i *žádný*). Od tvaru *žádný* ho odlišujeme hodnotou FMU=e (vložené e)

14 Závěr

Předložená práce se zabývá systémem morfologického popisu češtiny přesto, že na toto téma bylo napsáno i řečeno již mnoho. Po více než deseti letech intenzivního užívání elektronických morfologicky anotovaných korpusů češtiny se totiž ukazuje, že leckteré detaily popisu potřebují revizi, doplnění, případně i zcela odlišný přístup. Všechny návrhy vznikly z připomínek, stížností i nápadů ze strany uživatelů korpusů řady SYN, i na základě vlastní práce s těmito korpusy a s pražským morfologickým slovníkem.

Cílem práce bylo navrhnout systém, který jednotně, konzistentně a co nejúplněji popíše všechny morfologické jevy, které jsou potřebné při práci s českým jazykovým korpusem.

V prvních kapitolách se zabýváme systémem kategorií, pomocí kterých se české slovní tvary popisují. Patří sem i pojednání o lemmatizaci. V druhé části potom navrhujeme nový způsob zápisu flektivních a derivačních vzorů.

V čem spočívá přínos předkládané práce:

Definujeme přesné vymezení jednotlivých morfologických kategorií a jejich hodnot. Přitom se snažíme být v souladu s tradičními lingvistickými popisy. V některých případech však navrhujeme vlastní, netradiční řešení, protože tradiční popis nevyhovuje požadavkům na použití při automatickém zpracování jazyka. Vždy však dbáme na to, aby veškerá lingvistická informace zůstala zachována.

Stanovujeme základní princip pro budování morfologického slovníku, a to Zlaté pravidlo morfologie. S tím souvisí i důsledné zpracování slovních variant. Předkládáme jejich formální popis, který není závislý na jejich neobjektivním stylovém hodnocení. Zavádíme termín mutace, který pojetí varianty rozšiřuje. Mutace potom dělíme na flektivní a globální, což nám usnadní zachytit jejich variabilitu a snadnou kombinovatelnost.

Pro globální mutace zavádíme tzv. vícenásobné lemma, které umožní zahrnout tyto mutace pod společné lemma, ale přesto je popisuje tak, aby zůstaly rozlišeny a neporušily tak Zlaté pravidlo morfologie.

Při popisu jednotlivých morfologických kategorií jsme narazili na několik případů, které dosavadní systémy pro automatické zpracování češtiny buď zcela ignorují, nebo popisují nekonzistentně, nepřehledně nebo ne příliš šikovně vzhledem k dalšímu zpracování morfologicky označených textů. Tyto problematické případy se snažíme řešit lépe. Toho dosahujeme většinou zavedením nových kategorií pro popis některých jevů, nebo zavedením nových hodnot kategorií tradičních.

Mezi nově zavedené kategorie patří kategorie Duál, která umožní snadnější práci při analýze češtiny na syntaktické rovině a rovinách vyšších.

Dalšími kategoriemi jsou Flektivní mutace a Globální mutace, které rozšiřují známé termíny varianta a dubleta. Pomocí těchto kategorií popisujeme ty slovní tvary, které se v hodnotách ostatních kategorií neliší.

Nová je i kategorie Slovesný tvar, která slučuje tradiční slovesné kategorie,

neboť jejich hodnoty se většinou nedají vzájemně kombinovat. Některé hodnoty této kategorie jsou relevantní pro jiný slovní druh než pro sloveso. Konkrétně jde o pasivum, které anotujeme jako jmenný tvar přídavného jména, a kondicionál, který je relevantní pouze pro tři lemmata: částici *by* a spojky *aby*, *kdyby*.

Nové hodnoty jsme přidali kategorii **Slovní druh**. Důležitá je hodnota **cizí slovo**, která umožní popsat cizojazyčná slova, především vlastní jména, u nichž nemá smysl se snažit o zařazení do českého morfologického systému. Cizí jména přinášejí v dalších rovinách zpracování mnoho problémů. Jejich vyjmutím z množiny tradičních slovních druhů s nimi můžeme zacházet podle potřeby odlišně.

Podáváme rozbor tzv. složenin, které se nedají jednoznačně zařadit do systému slovních druhů, neboť ve svém tvaru jich sdružují vícero. V popisu složenin použijeme též nově zavedený koncept vícenásobného lemmatu.

Mezi hodnoty kategorie **Osoba** zahrnujeme též zdvořilostní formu druhé osoby singuláru, tedy vykání, které sice do paradigmatu zájmen a sloves patří, ale nebylo do něj dosud formálně zařazeno.

Podobný je případ stupně **s** kategorie **Stupeň**, tedy tvarů typu *sebekrásnější*, *sebekrásněji*, které se také dosud neanalyzovaly jako součást paradigmatu příslušného pozitivu přídavného jména nebo příslovce, k němuž přirozeně náleží.

Dále se zabýváme tzv. „stupňováním sloves“, které je sice spíše okrajovým jevem, ale velmi pravidelným, takže jeho zařazení do paradigmatiky sloves je také přirozené.

Nově navrhujeme systém flektivních vzorů, a to tak, aby byly parametrizovatelné. Pomocí vhodného nastavení parametru jednotlivým vzorům dosáhneme lepšího pokrytí systémových slovních tvarů. Flektivní vzory jsou, podobně jako v dosavadním pražském systému vzorů, doplněny o vzory umožňující pravidelné tvoření slov odvozených. Na rozdíl od pražského systému však i zde zavádíme parametrizaci, která umožní volbu, které odvozeniny tvořit a které ne.

Je velmi pravděpodobné, že se najdou slova, která náš návrh úplně nepokryje. Jsme však přesvědčeni, že návrh je dostatečně obecný na to, aby umožnil i řešení složitých výjimek.

Návrh, který jsme zpracovali, je v současné době již částečně implementován:

- Vytvořili jsme nástroje pro převod vzorů ze současného pražského morfologického slovníku do nového systému.
- Nově zpracovaný nástroj pro morfologickou analýzu pracuje s prefixovým guessrem navrženým na základě zmiňovaného výzkumu předpon.
- Postfixový guesser jsme použili při morfologické anotaci jedné z verzí korpusu SYN2000.
- Začali jsme pracovat na návrhu konkrétního schématu pro uchování morfologického slovníku v systému PML (Prague Markup Language) pro ukládání jazykových dat. V novém formátu budou především snadno zaznamenatelné derivační vztahy mezi jednotlivými lemmaty.

Chtěli bychom implementaci co nejdříve dokončit, aby mohly morfologické nástroje začít pracovat již s novými kategoriemi a jejich hodnotami. S tím sou-

visí i prapůvodní impuls k započetí této práce, totiž sjednocení pražského a brněnského pohledu na morfologické anotace, zejména vytvoření jednoznačného převodu mezi morfologickými značkami pražskými a brněnskými. Výrazně se tak zjednoduší vzájemná spolupráce.

Literatura

- Akademická mluvnice. *Mluvnice češtiny 2*. Praha, ACADEMIA, 1986.
- BRABCOVÁ, R. Kolísání rodu substantiv. In *Korpus jako zdroj dat o češtině*, s. 47–50. Brno: Masarykova univerzita, 2004.
- CRUSE, D. A. *Lexical Semantics*. Cambridge, UK, Cambridge University Press, 1986.
- ČERMÁK, F. Povaha a úzus interjekcí: případ češtiny. In *Computer Treatment of Slavic and East European Languages*, s. 299–307. Slovak Academy of Sciences, 2007.
- HAIJČ, J. *Disambiguation of Rich Inflexion*. Praha, Karolinum, 2004.
- HAIJČ, J. *Unification Morphology Grammar*. PhD thesis, Matematicko-fyzikální fakulta Univerzity Karlovy v Praze, 1994.
- HANKS, P. – PUSTEJOVSKY, J. Common Sense About Word Meaning: Sense in Context. In *Lecture Notes in Artificial Intelligence, Proceedings of the 7th International Conference, TSD 2004*, Berlin Heidelberg, 2004. Springer-Verlag.
- HAVRÁNEK, B. – JEDLIČKA, A. *Česká mluvnice*. Praha, Státní pedagogické nakladatelství, 1981.
- HLAVÁČOVÁ, J. Pravopisné varianty a morfologická anotace korpusů. In *Grammar & Corpora / Gramatika a korpus*, s. 161–168. Praha: Academia, 2008.
- HLAVÁČOVÁ, J. – HRUŠECKÝ, M. Affisix — Tool for Prefix Recognition. In *Lecture Notes in Artificial Intelligence, Proceedings of the 11th International Conference, TSD 2008*, s. 85–92, Berlin Heidelberg, 2008. Springer-Verlag.
- HLAVÁČOVÁ, J. – KOLOVRATNÍK, D. Morfologie češtiny znovu a lépe. In *Informačné Technológie – Aplikácie a Teória. Zborník príspevkov, ITAT 2008*, s. 43–47, 2008.
- HLAVÁČOVÁ, J. – LOPATKOVÁ, M. Variants and Homographs: Eternal Problem of Dictionary Makers. In *Proceedings of the 11th International Conference, TSD 2008*, s. 93–100, Berlin Heidelberg, 2008. Springer-Verlag.
- HLAVÁČOVÁ, J. Morphological Guesser of Czech Words. In *Proceedings of TSD 2001*, s. 70–75. Springer-Verlag Berlin Heidelberg, 2001.
- JIRANOVÁ, P. Morfologická a syntaktická charakteristika českých číslovek vyjadřujících počet entit, jejich souborů a druhů. Diplomová práce, Filosofická fakulta UK v Praze, 2008.

- KARLÍK, P. – HLADKÁ, Z. Kam s ním? (Problém stupňování adjektiv). In *Život s morfémy. Sborník studií na počest Zdenky Rusínové*, s. 73–93. Brno: Masarykova univerzita v Brně, 2004.
- KIRSCHNER, Z. MOSAIC — A Method of Automatic Extraction of Significant Terms from Texts. Technical report, Faculty of Mathematics and Physics, Charles University, Prague, 1983.
- KOPEČNÝ, F. *Základy české skladby*. Praha, Státní pedagogické nakladatelství, 1962a.
- KOPEČNÝ, F. *Slovesný vid v češtině*. Praha, Nakladatelství ČSAV, 1962b.
- KOSKENNIEMI, K. Two-level morphology: a general computational model for word-form recognition and production. Technical Report Publication No. 11, Helsinki: University of Helsinki Department of General Linguistics, 1983.
- KŘÍSTEK, M. Způsoby vymezení stylové příznakovosti v lexiku (na materiálu současné češtiny). In *Varia IX: zborník materiálů z IX. kolokvia mladých jazykovedcov*, s. 102–112. Slovenská jazykovedná spoločnosť pri SAV, 2002.
- KUČERA, K. Hyperlemma: A Concept Emerging from Lemmatizing Diachronic Corpora. In *Computer Treatment of Slavic and East European Languages*, s. 121–125. Slovak Academy of Sciences, 2007.
- LOPATKOVÁ, M. – ŽABOKRTSKÝ, Z. – BENEŠOVÁ, V. Valency Lexicon of Czech Verbs VALLEX 2.0. Technical Report 34, UFAL MFF UK, 2006.
- OSOLSOBĚ, K. O rozdílech mezi pražským a brněnským značkováním. Nепublikováno.
- OSOLSOBĚ, K. *Algoritmický popis české formální morfologie a strojový slovník češtiny*. PhD thesis, Filosofická fakulta Masarykovy univerzity v Brně, 1996.
- PANEVOVÁ, J. Honorifika v češtině (České vykání - teorie a korpusová data). In *Vybrané kapitoly z české gramatiky*. Praha: Academia, 2008. v tisku.
- Pravidla. *Pravidla českého pravopisu*. Ústav pro jazyk český, Praha, Pansofia, 1993.
- PRZEPIÓRKOWSKI, A. *The IPI PAN Corpus: Preliminary version*. Warszawa, IPI PAN, 2004.
- ROMPORTL, S. *Struktura gramatické složky slovesných tvarů určitých v češtině*. Praha, ACADEMIA, 1970.
- ROSEN, A. – SALONI, Z. Kategorie honorativu v českých konjugačních paradigmatech. *Slovo a slovesnost*. 2006, , 1, s. 36–45.
- SEDLÁČEK, R. Morfologický analyzátor češtiny. Diplomová práce, Fakulta informatiky Masarykovy univerzity v Brně, 1999. <http://nlp.fi.muni.cz/projekty/ajka>.
- SGALL, P. – HRONEK, J. *Čeština bez příkras*. Praha, H&H, 1992.

Literatura

- SKOUMALOVÁ, H. Czech lexicon by two-level morphology. In *Proceedings of the Second European Seminar of TELRI – Language Applications for a Multilingual Europe*, s. 123–145. IDS/VDU, Mannheim/Kaunas, 1997.
- Slovenská morfológie. *Morfológia slovenského jazyka*. Bratislava, Vydavateľstvo SAV, 1966.
- ŠEVČÍKOVÁ, M. Pronouns Introducing Content Clauses. In *Grammar & Corpora / Gramatika a korpus*, s. 277–284. Praha: Academia, 2008.
- ŠIMANDL, J. Kvantifikátory v korpusech ÚČNK a možnosti jejich značkování. Nepublikováno, 2007.
- ŠIMKOVÁ, M. O lexikálnom význame častíc. *Slovenská reč*. 2001, , 66, s. 37–51.
- ŠMILAUER, V. *Novočeské tvoření slov*. Praha, Státní pedagogické nakladatelství, 1971.
- TUŠKOVÁ, J. M. *Variantní a dubletní tvary v současné deklinaci apelativních feminin*. Spisy Pedagogické fakulty MU, sv. č. 98. Brno, Masarykova univerzita, 2006.
- URREA, A. M. – HLAVÁČOVÁ, J. Automatic Recognition of Czech Derivational Prefixes. In *Proc. CICLING 2005*, s. 189–197. Springer-Verlag Berlin Heidelberg, 2005.

A Přehled kategorií a jejich hodnot

Uvádíme souhrnný přehled kategorií a jejich hodnot z kapitoly 4. Číslování odpovídá návrhu morfologické značky z oddílu 4.3, souhrn tedy neobsahuje hodnoty kategorií **Flektivní mutace** a **Globální mutace**.

U hodnot kategorie **Slovní druh** uvádíme rovnou všechny poddruhy, které jsou pro daný slovní druh relevantní. Souhrnný přehled všech poddruhů je uveden dále pod číslem 2, kde je v závorce místo příkladu kód slovních druhů, pro které je daná hodnota relevantní.

1. Slovní druh (POS)

- N: podstatné jméno
 - S: deverbativní typu *věznění, pokrytí,...*
 - 0: ostatní
- A: přídavné jméno
 - U: přivlastňovací (*matčín, otcův,...*)
 - G: od přechodníku přít. (*mající, sedící, beroucí,...*)
 - M: od přechodníku min. (*ušedší, nakupovavší,...*)
 - S: ostatní deverbativní (*namazaný, zemřelý, nakousnutý, namažán, nakousnut,...*)
 - 0: ostatní (*jarní, starý,...*)
- P: zájmeno
 - Z: substantivní (*já, kdo, nikdo, oni,...*)
 - U: přivlastňovací (*můj, čí,...*)
 - D: ukazovací (*ten, takový,...*)
 - V: vymešovací (*každý, všechen, týž, sám*)
 - 0: ostatní
- C: číslovka
 - 1: základní (*jedna, sto,...*)
 - r: řadové (*druhý, pátý,...*)
 - u: úhrnné (*dvé, patero,...*)
 - s: souborové (*dvoje, paterý,...*)
 - d: druhové (*dvojí, paterý,...*)
 - n: násobné (*dvakrát, pětkrát,...*)
 - o: opakovací (*podruhé, popáté,...*)
 - v: výčtové (*zadruhé, zapáté,...*, ale i *druhé* z dvojice *za druhé,...*)
 - p: dílové (*půl, čtvrt, třet*)
- V: sloveso
 - m: modální (*moci/moct, mít, mívat, muset, musívat, smět, chtít, hodlat, dát se, dávat se, dovést, umět*)
 - f: fázová (*začít, začínat, přestat, přestávat, zahájit, skončit,...*)
 - b: pomocná (*být, bývat, mít, dostat*)
 - 0: ostatní (*navštívit, koupat se,...*)

A Přehled kategorií a jejich hodnot

- D: příslovce
 - P: místní (*kudy, tudy, odkud, nikudy, nikam; daleko, nedaleko,...*)
 - T: časová (*kdy, nikdy; včera, odpoledne,...*)
 - D: způsobová (*jak, všelijak; krásně, velmi, široce,...*)
 - R: predikativní (*jasno, možno, teplo, volno,...*)
- R: předložka
- J: spojka
 - ^ (stříška): souřadící (*a, ale, nebo,...*)
 - , (čárka): podřadící (*protože, když, až, -li,...*)
 - * (hvězdička): matematické operace (*plus, minus/mínus, krát, děleno* — neplést s *děleno* jako jmenný tvar přídavného jména *dělený*)
- I: citoslovce
- T: částice
 - 7: zvrtné (*se, si*)
 - c: kondicionálová (*pouze by*)
 - 0: ostatní (*ba, ano, bože, nechť, ať,...*)
- F: cizí slovo (K)
- G: prefixový segment (K)
- S: složenina (K)
- X: neznámé slovo

2. Slovní poddruh (SUB)

- \wedge (stříška): spojka souřadící (J)
- , (čárka): spojka podřadící (J)
- * (hvězdička): matematické operace (J)
- b: pomocné sloveso (V)
- c: kondicionálová částice (T)
- d: druhová číslovka (C)
- D: ukazovací zájmeno (P)
- f: fázové sloveso (V)
- G: adjektivum od přechodníku přít. (A)
- J: způsobové příslovce (D)
- m: modální sloveso (V)
- M: adjektivum od přechodníku min. (A)
- n: násobná číslovka (C)
- o: opakovací číslovka (C)
- p: dílová číslovka (C)
- P: místní příslovce (D)
- R: predikativní příslovce (D)
- r: řadová číslovka (C)
- S: deverbativní adjektivum/substantivum (AN)
- s: souborová číslovka (C)
- T: časové příslovce (D)
- U: přivlastňovací adjektivum/zájmeno (AP)
- u: úhrnná číslovka (C)
- v: výčtová číslovka (C)
- V: vymežovací zájmeno (P)
- Z: substantivní zájmeno (P)
- 1: základní číslovka (C)
- 7: zvrtná částice (T)
- 0: ostatní (vše)

3. Funkce (FCE)

- U: určitá (všechna osobní zájmena, určité číslovky, *tady, teď,...*)
- N: neurčitá (*někdo, čísi, několik, někdy,...*)
- Z: záporná (*nikdo, ničí, nijak,...*)
- T: tázací (*kdo, čí, kolik, kde,...*)
- V: vztažná (*kdo, čí, jenž, kdy,...*)
- S: zvrtná (*se, si, sobě, sebe, sebou*)

4. Slovesný vid (ASP)

- D: dokonavý (*koupit, napsat, doručit, narodit se,...*)
- N: nedokonavý (*kupovat, psát, doručovat, chodívat,...*)
- O: obouvidý (*referovat, absolvovat, izolovat,...*)

5. Zkratka (ABR)

- + : ano

6. Rod (GEN)

- M: mužský životný
- I: mužský neživotný
- F: ženský
- N: střední
- X: sdružená hodnota

7. Číslo (NUM)

- S: jednotné
- P: množné
- X: sdružená hodnota

8. Duál (DUA)

- +

9. Pád (CAS)

- 1 až 7
- X: sdružená hodnota

10. Osoba (PER)

- 1 až 3
- v: vykání

11. Stupeň (DEG)

- 1: pozitiv
- 2: komparativ
- 3: superlativ
- s: typ sebe + komparativ

12. Negace (NEG)

- N: pro záporné slovní tvary, které začínají záporkou *ne-*
- A: pro ostatní slovní tvary

13. Slovesný tvar (VRB)

- P: indikativ přítomného času (*kolíbá*)
- B: budoucí čas (*ponese, bude*)
- F: infinitiv (*otevřít*)
- I: imperativ (*peč*)
- L: přídavné sloveso (*strouhal*)
- T: přídavné sloveso (*zavřen*)
- K: kondicionál (*aby, kdyby, by*)
- p: přechodník přítomný (*starajíc*)
- m: přechodník minulý (*vtoupiv*)

14. Jmenný tvar přídavných jmen (NOM)

- J: jmenný tvar
- 0: ostatní přídavná jména

15. Stupeň intenzity slovesného děje (INT)

- r: pro předponu *roz-*
- p: pro předponu *po-*
- z: pro předponu *za-*
- n: pro předponu *na-*
- v: pro předponu *vy-*
- u: pro předponu *u-*

16. Typ složeniny (CMP)

- n: zájmenný (*proň, zaň*)
- c: zájmenný (*oč, zač, začpak*)
- t: zájmenně-slovesný (*toť*)
- Z: zkratkový (*atd., apod.*)
- A: slovesný, 1. složka je A (*krásnýs*)
- N: slovesný, 1. složka je N (*latínys*)
- P: slovesný, 1. složka je P (*jemus*)
- C: slovesný, 1. složka je C (*kolíks*)
- V: slovesný, 1. složka je V (*zavřelas*)
- D: slovesný, 1. složka je D (*včeras*)
- T: slovesný, 1. složka je T (*sis, ses*)
- J: slovesný, 1. složka je J (*protožes*)
- S: slovesný, 1. složka je S (*načs*)

B Kopie účastnického slibu z Konkláve

Přisáhám na svou čest,

že během celého konkláve budu jednat podle svého nejlepšího
morfologického svědomí,

zejména že

zapomenu na veškeré aplikace, na kterých jsem dosud pracoval(a),

a obrátím svou mysl výhradně do budoucnosti,
aby výsledkem našeho jednání byla skutečně čistá a průzračná
morfologie českého jazyka,
nezatížená jakýmikoliv partikulárními zájmy.

Karel Pala



Vladimír Petkevič



Karel Oliva



Jan Hajič



Klára Osolsobě



Jaroslava Hlaváčová



Cikháj, 21. října 2005

Rejstřík

- číslo, 38
- ABR, 34
- ASP, 34
- budoucí čas, 46
- CAS, 41
- cizí slovo, 23
- DEG, 42
- derivační odkaz, 73, 80
- derivační vzor, 77, 80
- duál, 38
- duálová slova, 39
- DUA, 38
- dubleta, 15
- flektivní morfologická kategorie, 21
- flektivní mutace, 17, 47
- flektivní vzor, 77
- FMU, 47
- funkce, 30
- GEN, 38
- generování, 1
- globální morfologická kategorie, 21
- globální mutace, 17
- guesser, 69
- hyperlemma, 7
- INT, 47
- jmenný tvar přídavných jmen, 46
- kofix, 77
- kompaktní slovník, 74
- kompaktní systém značek, 5
- Konkláve, 2
- kritická kombinace, 82
- lemma, 3
- lemma vícenásobné, 7
- lemma variantní, 7
- lemmatizace, 5
- morfologická analýza, 1, 6
- morfologická kategorie, 4
- morfologická syntéza, 1
- morfologická značka, 4, 5, 49
- morfologická značka hodnotová, 50
- morfologická značka kompaktní, 50
- morfologická značka poziční, 50
- morfologické konkláve, 2
- morfologický slovník, 68
- mutace, 15
- mutace flektivní, 17
- mutace globální, 17
- NEG, 43
- negace, 43
- NOM, 46
- nulové mutace, 49
- NUM, 38
- oddělovače, 2
- osoba, 41
- pád, 41
- příčestí trpné, 45
- příbuzná lemmata, 73
- přechodníky, 46
- paradigma, 4
- paradigma rozšířené, 8
- parametr negace, 79
- PER, 41
- PML, 74
- poddruh, 24
- POS — slovní druh, 22
- poziční systém značek, 5
- prefix negace, 80
- relevantní morfologická kategorie, 4, 51
- rod, 38

rozšířené paradigma, 8

S-formy, 22

složenina, 56

složka složeniny, 56

slovesný tvar, 43

slovní druh, 22

slovní tvar, 3

slovo, 2

stupňování sloves, 9

stupeň, 42

stupeň intenzity slovesného děje, 47

stylový příznak, 19

typ složeniny, 58

undef, 4, 21

význam slova, 3

vícenásobné lemma, 8

vícenásobné lemma složeniny, 8

varianta, 15

variantní lemma, 7

vid, 34

VRB, 43

vykání, 41

zájmenná příslovce, 29

záznam morfologického slovníku, 71

zakončení, 77

zakončení slovního tvaru, 77

zkratka, 34

Zlaté pravidlo morfologie, 6

zvratná slovesa, 8