

Pokus o formální popis českého slovosledu

Alexandr Rosen

Ústav teoretické a počítačnické lingvistiky

Universita Karlova v Praze

alexandr.rosen@ff.cuni.cz

<http://utkl.ff.cuni.cz/~rosen/public/THESIS>

26. března 2002

„A constraint-based approach to dependency syntax
applied to some issues of Czech word order“

*„Deklarativní formalizace teorie závislostní syntaxe
s aplikací na některé problémy českého slovosledu“*

Obsah

1	Proč?	4
2	Hypotézy	9
3	Teorie: FGP	10
4	Formalismus: RSRL	11
5	Reprezentace	17
6	Syntaktická „kostra“	25
7	Slovosled	27
8	Příklady	36
9	Výsledky a výhledy	43

1 Proč?

Předpoklad:

rozdíl mezi kompetencí (*langue*) a performancí (*parole*)

- popis (gramatika) kompetence:
 - které řetězce fonémů/grafémů patří a které nepatří do jazyka X
 - co tyto řetězce znamenají
(jaký mají vztah k reprezentaci významu)
- popis (gramatika) performance:
 - jak se tyto znalosti využívají při „jazykových aktivitách“

Může korpus nahradit gramatiku?

– Jak z něj „vycucnout“ implicitně obsažená pravidla gramatiky a jak je zobecnit?

Automaticky = statisticky?

- Neoznačovaný korpus:
jaké řetězce se vyskytly,
v jakém kontextu a situaci, jak často
- Označovaný korpus:
jaké kategorie, konstrukce, významy se vyskytly,
v jakém kontextu a situaci, jak často

Longum iter est per precepta, breve et efficax per exempla.

– Seneca

Es gibt nur die Beispiele

– Wittgenstein

Ale: statistické metody – jen simulace vědomého poznání

Gramatika:

- jak skládat delší výrazy z kratších
- jaký je vztah mezi povrchovým řetězcem a jeho reprezentací

Možnosti:

- Derivační (stratifikační) vs. nederivační (monostratální) přístup
- Složková vs. závislostní syntax
- Složková syntax a CF gramatika: derivace = reprezentace
- FGP: derivace \neq reprezentace, reprezentace abstraktnější
tradiční přístup: stratifikační a procedurální, jde to jinak?

Formalismus

- Standardní formalizace FGP je stratifikační s procedurálními prvky (generování tektogramatického zápisu, překladové složky)
- Problémy:
 - paralelní přístup k informacím z více rovin
 - interpretace neúplných výrazů
 - preference generování na úkor analýzy

Slovosled

- Slovosled odráží více faktorů z více rovin
(slovosledné faktory: aktuální členění, ‚povrchová‘ pravidla gramatiky)
- Slovosledné faktorů spolu ‚soutěží‘
(jejich role se liší v různých jazycích – V. Mathesius)

Řešení?

- ‚Constraint-based‘ formalismy umožňují paralelní přístup k více rovinám
- FGP nabízí adekvátní teoretickou koncepci a hloubkovou reprezentaci, včetně aktuálního členění

2 Hypotézy

- Teorie není nerozlučně svázána s formalismem
- FGP lze kombinovat s deklarativním formalismem
- Tato kombinace umožňuje lépe popsat působení slovosledných faktorů

„deklarativní formalismus“:

Relational Speciate Re-entrant Logic – RSRL

Paul King (1989), Frank Richter (2000)

- **Cíl práce:**
popsat podstatnou část slovosledných jevů v češtině pomocí:
 - FGP jako teoretického východiska
 - RSRL jako formálního jazyka

3 Teorie: FGP

- rozčlenění popisu jazyka na více rovin
- hranice mezi systémem jazyka a sémantickými + pragmatickými interpretacemi
- úloha výpovědní dynamičnosti a aktuálního členění v jazykovém významu
- reprezentace jazykového významu v podobě tektogramatického stromu: vztahy závislosti mezi autosémantickými slovy, hloubkový slovosled, kontextová zapojenost

4 Formalismus: RSRL

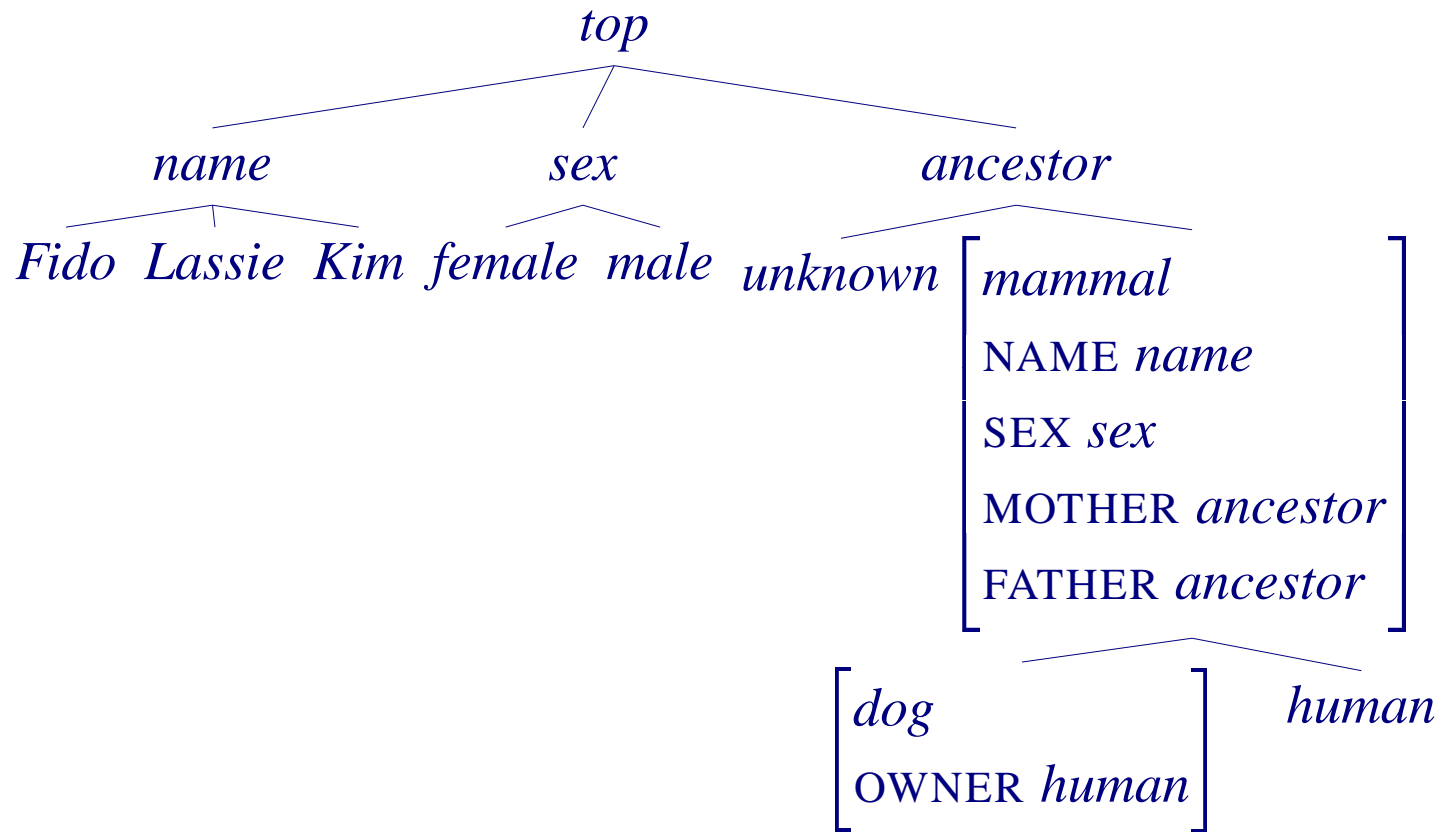
Gramatika popisuje jazyk nepřímo, pracuje s modelem jazyka. Skládá se ze „*signatury*“ a „*teorie*“.

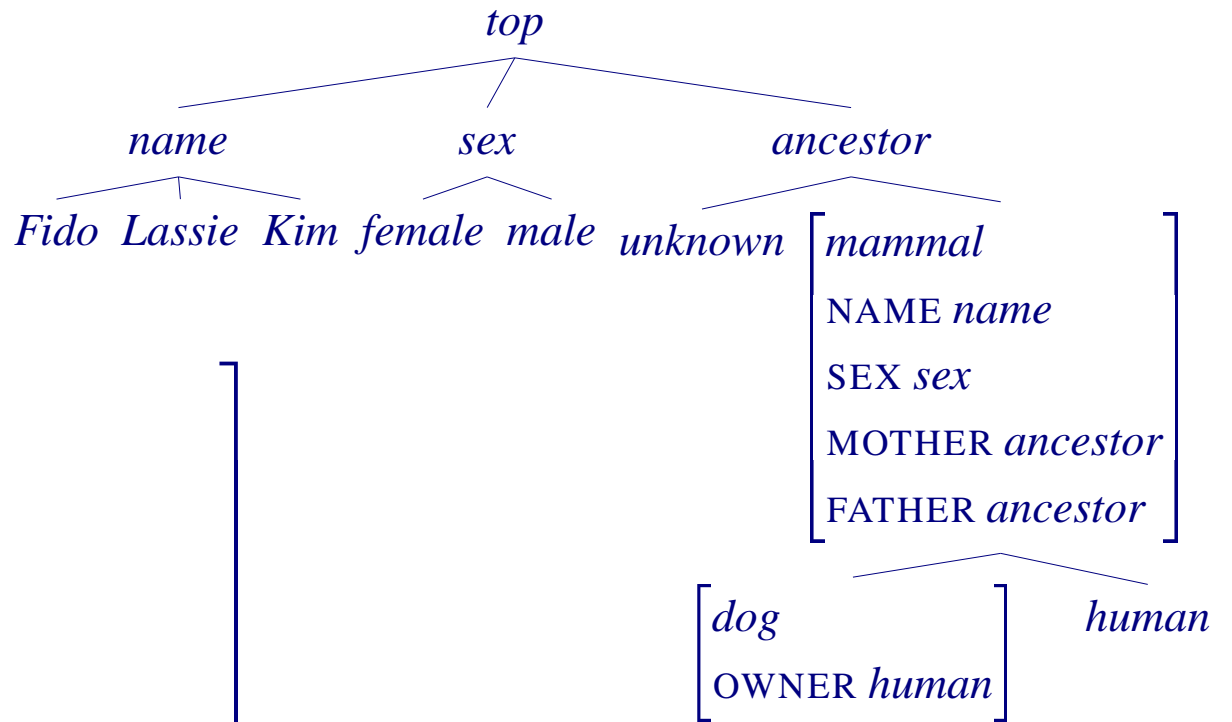
Signatura

definuje hierarchii typů – množin objektů modelu.

- Každému objektu modelu je přiřazen právě jeden maximálně specifický typ.
- Typy mají předepsané vlastnosti – atributy. Předepsané hodnoty atributů jsou opět typy.
- Objekty modelu musí mít všechny předepsané atributy, jejichž hodnoty musí být maximálně specifické typy.
- Typy stojící v hierarchii níže dědí všechny specifikace svých nadtypů.

Navíc: seznamy, množiny.





<i>dog</i>											
NAME	<i>Fido</i>										
SEX	<i>male</i>										
MOTHER	<i>unknown</i>										
FATHER	<i>unknown</i>										
OWNER	<table border="1"> <tr> <td><i>human</i></td> <td></td> </tr> <tr> <td>NAME</td> <td><i>Kim</i></td> </tr> <tr> <td>SEX</td> <td><i>female</i></td> </tr> <tr> <td>MOTHER</td> <td><i>unknown</i></td> </tr> <tr> <td>FATHER</td> <td><i>unknown</i></td> </tr> </table>	<i>human</i>		NAME	<i>Kim</i>	SEX	<i>female</i>	MOTHER	<i>unknown</i>	FATHER	<i>unknown</i>
<i>human</i>											
NAME	<i>Kim</i>										
SEX	<i>female</i>										
MOTHER	<i>unknown</i>										
FATHER	<i>unknown</i>										

Teorie

klade omezení na objekty modelu definované v signatuře.

- Teorie se skládá z formulí, každá formule musí být splněna všemi objekty modelu.
- Formule se skládají z typů a atributů (definovaných v signatuře) a logických symbolů.
- Logické symboly jsou obvyklé spojovací výrazy, proměnné, negace a kvantifikátory s dosahem přes komponenty daného objektu.
- V teorii lze definovat a používat relace.

	<i>dog</i>												
NAME	<i>Fido</i>												
SEX	<i>male</i>												
MOTHER	<table border="1"> <tr> <td><i>dog</i></td> <td></td> </tr> <tr> <td>NAME</td> <td><i>Lassie</i></td> </tr> <tr> <td>SEX</td> <td><i>male</i></td> </tr> <tr> <td>MOTHER</td> <td><i>unknown</i></td> </tr> <tr> <td>FATHER</td> <td><i>unknown</i></td> </tr> <tr> <td>OWNER</td> <td>①</td> </tr> </table>	<i>dog</i>		NAME	<i>Lassie</i>	SEX	<i>male</i>	MOTHER	<i>unknown</i>	FATHER	<i>unknown</i>	OWNER	①
<i>dog</i>													
NAME	<i>Lassie</i>												
SEX	<i>male</i>												
MOTHER	<i>unknown</i>												
FATHER	<i>unknown</i>												
OWNER	①												
FATHER	<i>unknown</i>												
OWNER	①	<table border="1"> <tr> <td><i>human</i></td> <td></td> </tr> <tr> <td>NAME</td> <td><i>Kim</i></td> </tr> <tr> <td>SEX</td> <td><i>female</i></td> </tr> <tr> <td>MOTHER</td> <td><i>unknown</i></td> </tr> <tr> <td>FATHER</td> <td><i>unknown</i></td> </tr> </table>	<i>human</i>		NAME	<i>Kim</i>	SEX	<i>female</i>	MOTHER	<i>unknown</i>	FATHER	<i>unknown</i>	
<i>human</i>													
NAME	<i>Kim</i>												
SEX	<i>female</i>												
MOTHER	<i>unknown</i>												
FATHER	<i>unknown</i>												

mammal → [MOTHER | SEX *female*
FATHER | SEX *male*]

[MOTHER | SEX ①] → ① *female*

[FATHER | SEX ①] → ① *male*

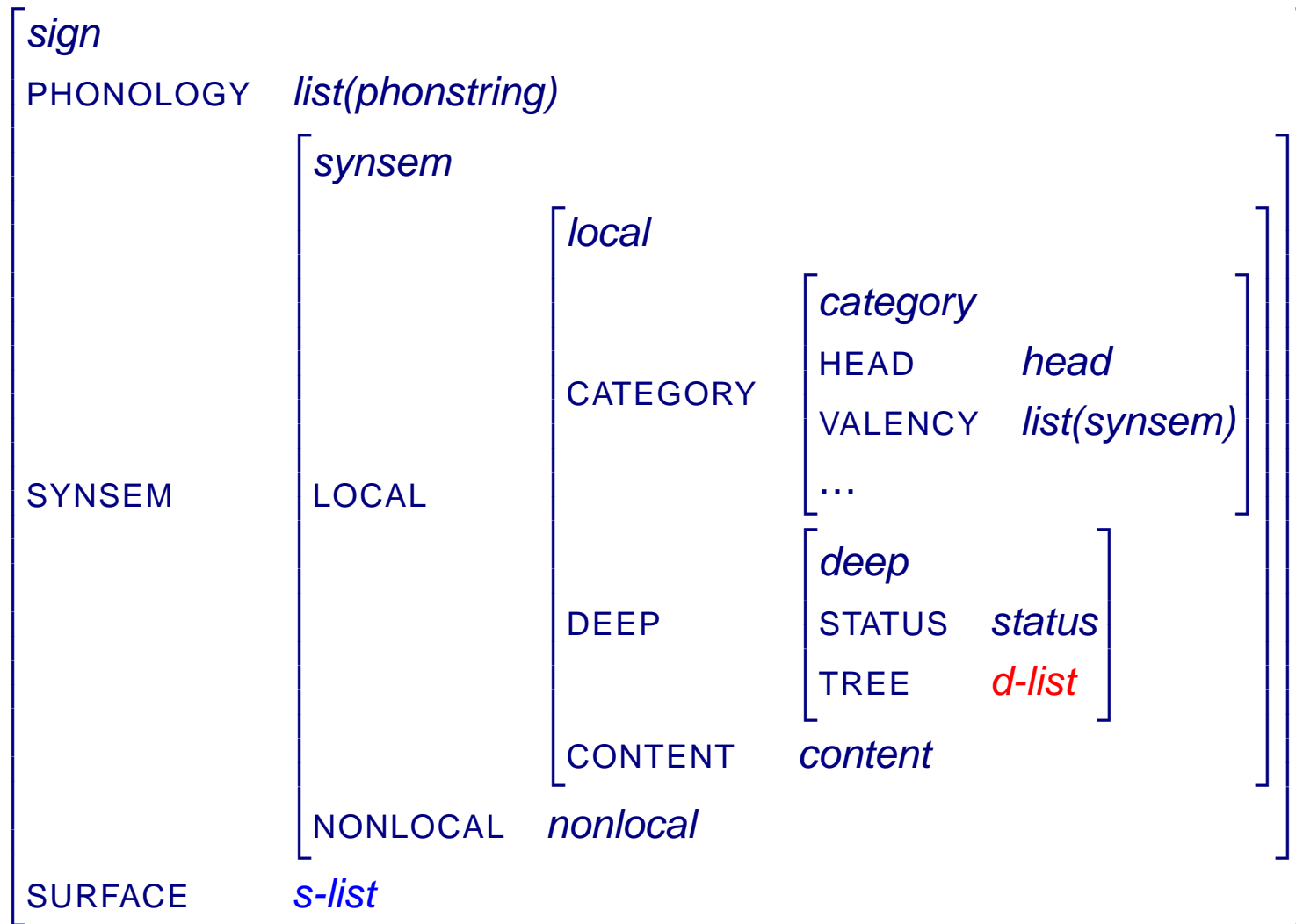
human → [MOTHER *human*
FATHER *human*]

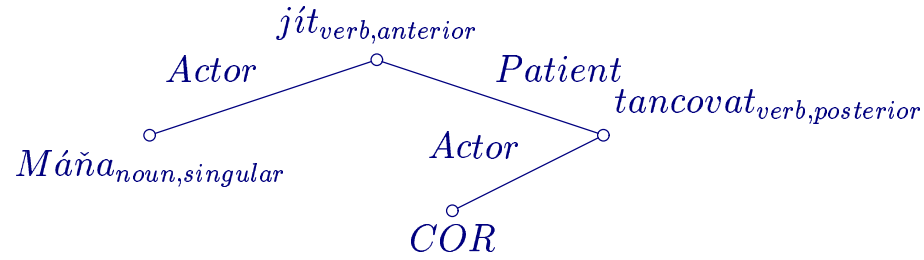
[*human*
MOTHER ① *mammal*] → ① *human*

[*human*
FATHER ① *mammal*] → ① *human*

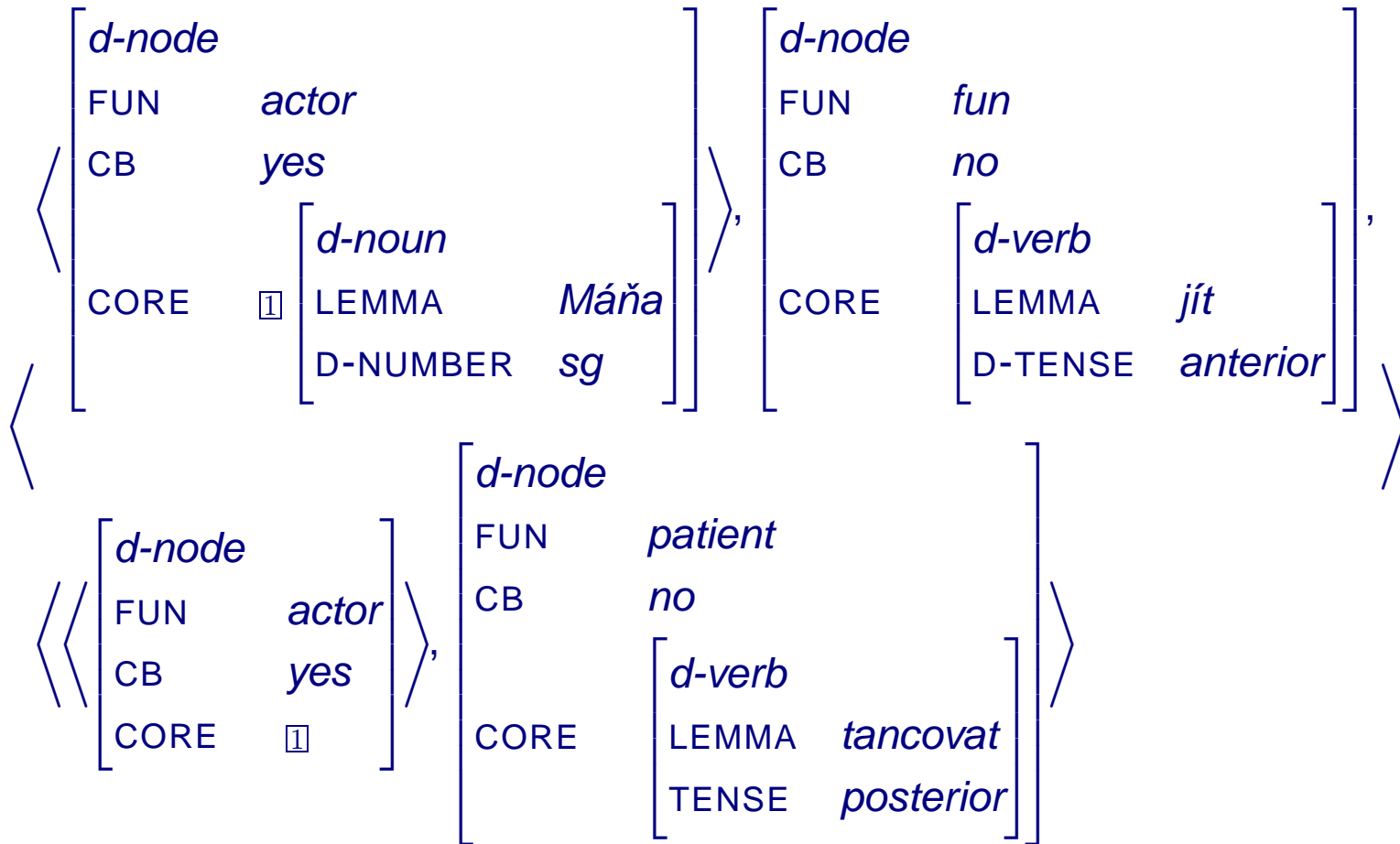
[OWNER ①] → [MOTHER | OWNER ①]

5 Rezentance





⟨ ⟨ Máňa ⟩, jít, ⟨ ⟨ COR ⟩, tancovat ⟩ ⟩

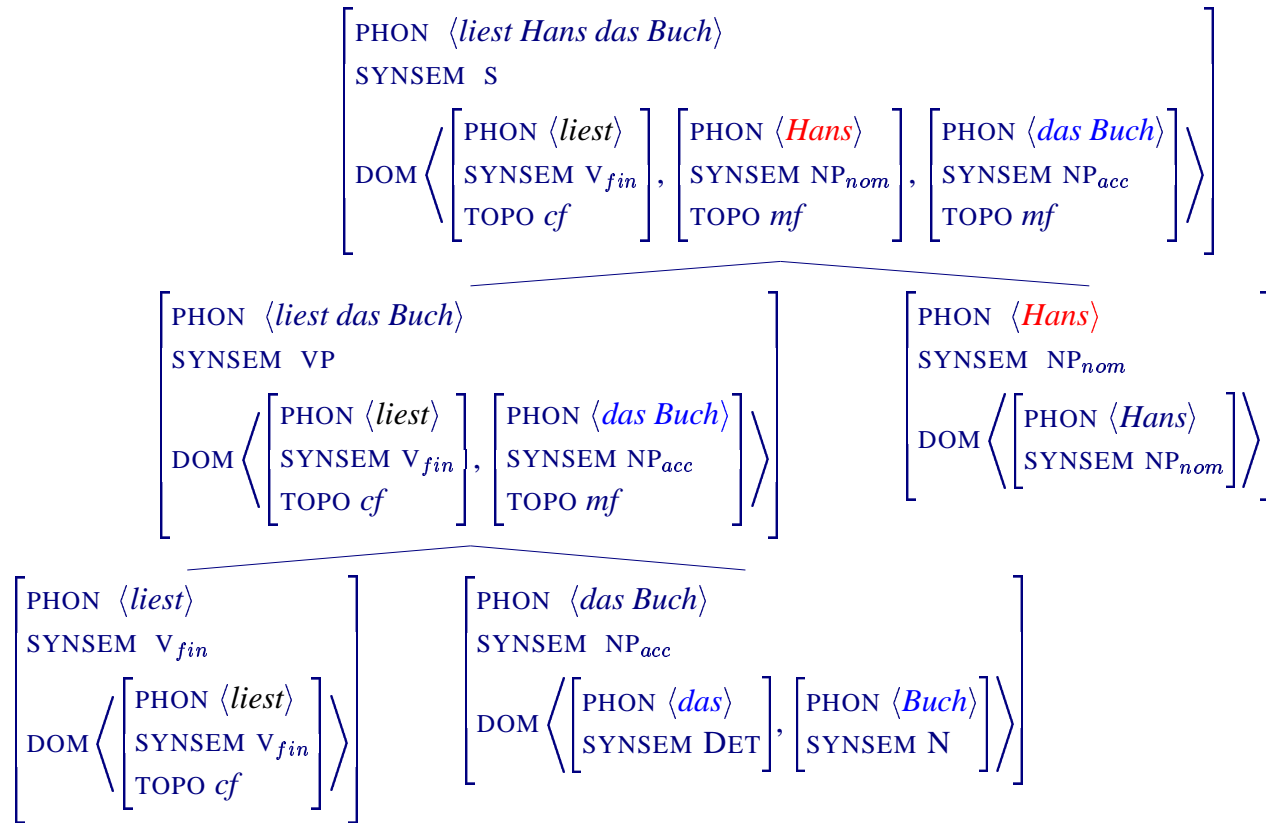


Povrchový řetězec

Inspirace č. 1: Mike Reape, Andreas Kathol

- Horizontální pořadí terminálů v derivačním stromě se nerovná povrchovému slovosledu
- Slovosledné domény, objekty domén
- Relace *shuffle* a *compaction*
- Topologická pole

$$sign \rightarrow \left[\begin{array}{c} \text{PHONOLOGY } \boxed{1} \oplus \dots \oplus \boxed{n} \\ \text{DOMAIN } \left\langle \left[\begin{array}{c} \textit{domain-object} \\ \text{PHONOLOGY } \boxed{1} \end{array} \right], \dots, \left[\begin{array}{c} \textit{domain-object} \\ \text{PHONOLOGY } \boxed{n} \end{array} \right] \right\rangle \end{array} \right]$$



$$\left(\begin{array}{l} \textit{phrase} \\ \text{PHON } \boxed{2} \oplus \boxed{3} \oplus \boxed{1} \oplus \boxed{4} \\ \text{DOM } \boxed{5} \circ \left\langle \begin{array}{l} \textit{dom-obj} \\ \text{PHON } \boxed{2} \oplus \boxed{3} \\ \text{SS } \boxed{7} \end{array} \right\rangle, \boxed{8} \right) \\ \text{SS } \boxed{L} \mid \boxed{C} \mid \boxed{H} \boxed{11} \\ \wedge \text{p-compactio}n(\boxed{9}, \langle \boxed{6} \rangle, \langle \boxed{8} \rangle) \end{array} \right)$$

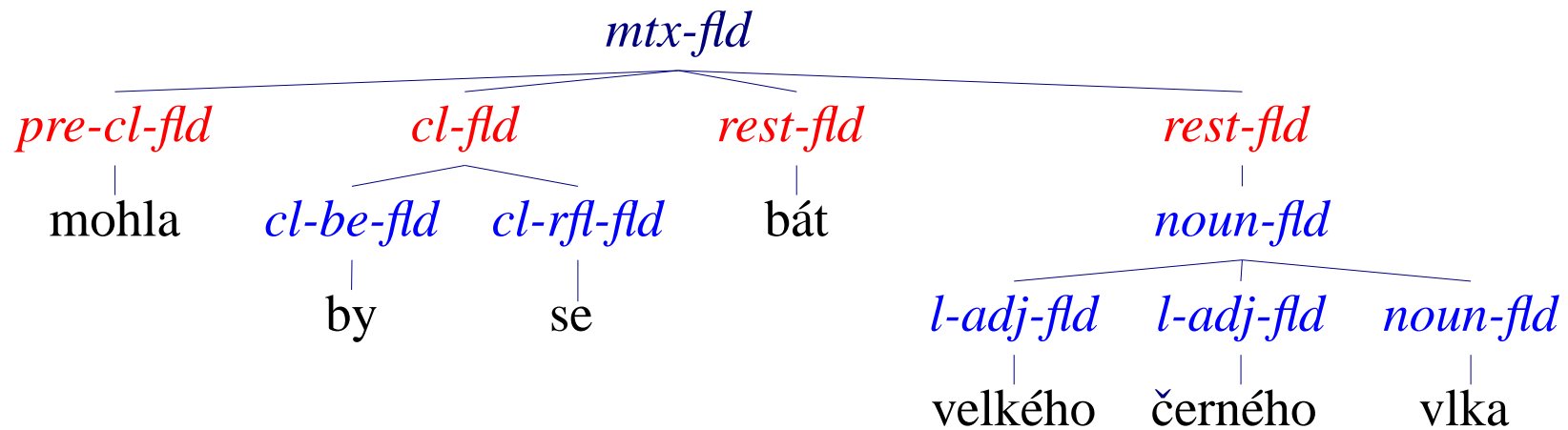
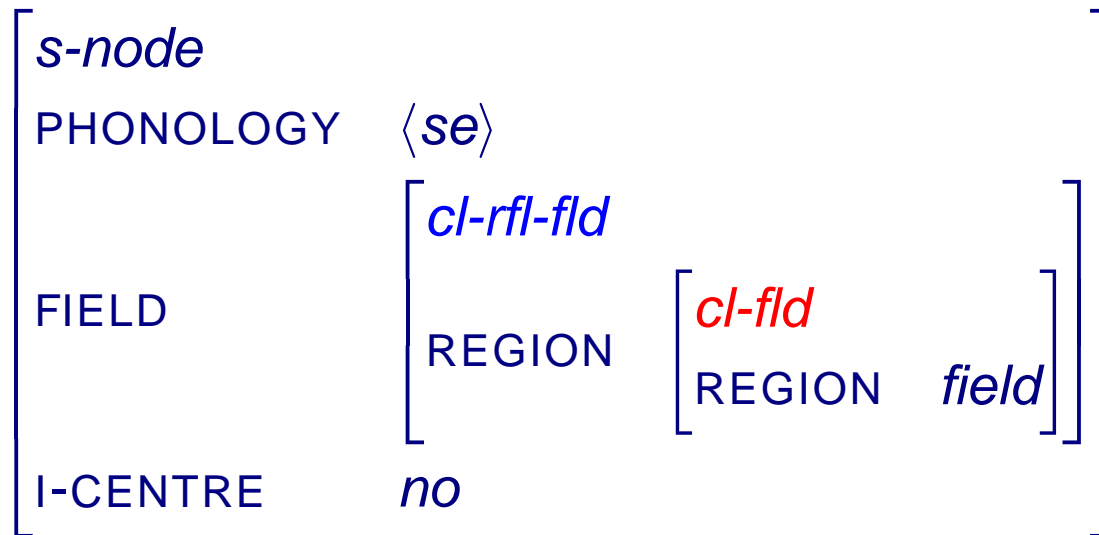
$$\left[\begin{array}{l} \textit{word} \\ \text{PHON } \boxed{1} \langle \textit{mieszka} \rangle \\ \text{DOM } \boxed{5} \left\langle \begin{array}{l} \textit{dom-obj} \\ \text{PHON } \boxed{1} \\ \text{SS } \boxed{10} \end{array} \right\rangle \\ \text{SS } \boxed{10} \left[\boxed{L} \mid \boxed{C} \mid \boxed{H} \boxed{11} \textit{verb} \right] \end{array} \right]$$

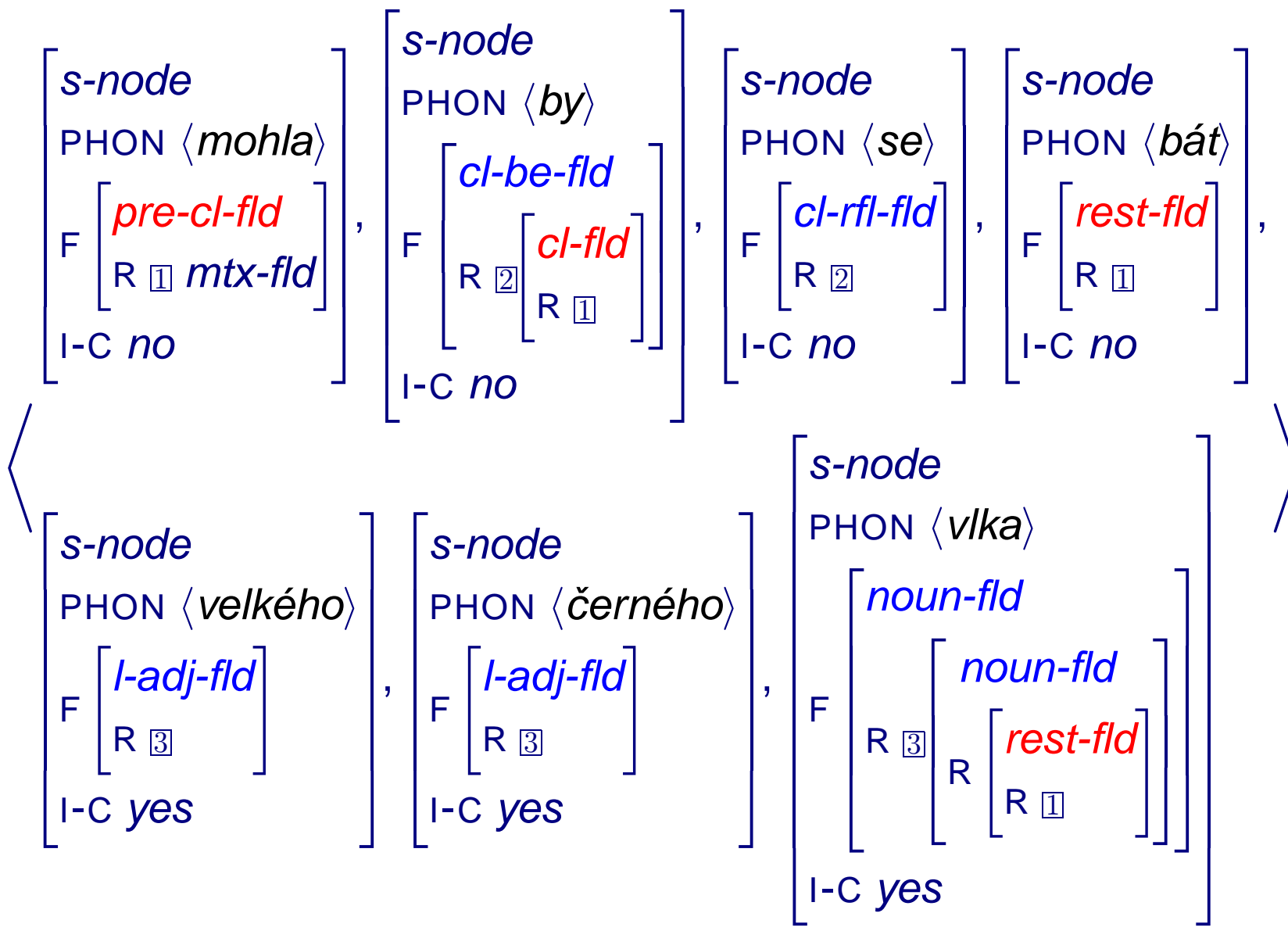
$$\left[\begin{array}{l} \textit{phrase} \\ \text{PHON } \boxed{2} \langle \textit{w} \rangle \oplus \boxed{3} \langle \textit{duzym} \rangle \oplus \boxed{4} \langle \textit{domu} \rangle \\ \begin{array}{l} \textit{dom-obj} \\ \text{PHON } \boxed{2} \\ \text{SS } \boxed{L} \mid \boxed{C} \mid \boxed{H} \boxed{12} \textit{prep} \end{array} \right], \\ \boxed{9} \text{DOM} \left\langle \begin{array}{l} \textit{dom-obj} \\ \text{PHON } \boxed{3} \\ \text{SS } \boxed{L} \mid \boxed{C} \mid \boxed{H} \textit{adj} \end{array} \right\rangle, \right. \\ \left. \begin{array}{l} \textit{dom-obj} \\ \boxed{8} \text{PHON } \boxed{4} \\ \text{SS } \boxed{L} \mid \boxed{C} \mid \boxed{H} \textit{noun} \end{array} \right] \\ \text{SS } \boxed{7} \left[\boxed{L} \mid \boxed{C} \mid \boxed{H} \boxed{12} \right] \end{array} \right]$$

Inspirace č. 2: Gerald Penn

- Objekty domén odpovídají jednotlivým slovům
- Příslušnost ke slovosledné oblasti je vyjádřena sdílením hodnot
- To umožňuje hierarchickou strukturu topologických polí a nezávislost na syntaktických objektech

<i>s-node</i>					
PHONOLOGY	<i>list(phonstring)</i>				
FIELD	<table><tr><td><i>field</i></td><td></td></tr><tr><td>REGION</td><td><i>field</i></td></tr></table>	<i>field</i>		REGION	<i>field</i>
<i>field</i>					
REGION	<i>field</i>				
I-CENTRE	<i>boolean</i>				





6 Syntaktická „kostra“

Vlastně pravidla derivační struktury, tedy skladby nelexikálních objektů typu *sign*.

- Deep List Composition Principle
- Surface List Composition Principle
- Valency Principle
- Head Principle
- Phonology Principle

Derivační struktura je plochá (všechny závislé uzly, včetně funkčních slov, jsou sestry řídicího uzlu). Výjimku funkční slova modifikovaná jiným funkčním slovem.

Deep List Composition Principle

V každém objektu typu *non-lexical* je *d-list* syntaktické matky roven seznamu *d-list* řídící dcery, do něhož jsou vloženy seznamy *d-list* ostatních dcer.

non-lexical →

$$\left(\begin{array}{l} \left[\text{SYNSEM} \mid \text{LOCAL} \mid \text{DEEP} \mid \text{TREE} \text{ [5]} \right. \\ \left. \text{HEAD-DAUGHTER} \mid \text{SYNSEM} \mid \text{LOCAL} \mid \text{DEEP} \mid \text{TREE} \text{ [1]} \right. \\ \left. \text{NONHEAD-DAUGHTERS} \text{ [2]} \right. \\ \wedge \text{COLLECT_DLISTS}(\text{[2]}, \text{[3]}) \\ \wedge \text{APPEND}(\text{[1]}, \text{[3]}, \text{[4]}) \\ \wedge \text{PERMUTE}(\text{[4]}, \text{[5]}) \end{array} \right)$$

7 Slovosled

- Omezení hloubkového slovosledu
- Omezení povrchového slovosledu
- Omezení vztahu mezi hloubkovým a povrchovým slovosledem

Omezení hloubkového slovosledu

- V každém nezanořeném stromě musí být aspoň jeden NB uzel.
- V každém stromě musí řídicí uzel v horizontálním pořadí předcházet všechny NB uzly.
- V každém stromě musí NB uzly v horizontálním pořadí následovat po CB uzlech podle systémového uspořádání.

Omezení vztahu mezi hloubkovým a povrchoým slovosledem

Uplatní se pro každou dvojici tektogramatických uzlů, není-li povrchová pozice žádného z nich určena povrchoými slovoslednými pravidly (dáno topologickým polem).

Existují 3 možnosti:

- Relativní pořadí uzlů je stejné (uplatní se princip aktuálního členění).
- Je-li první z uzlů NB a na prvním místě v seznamu *d-list*, může být v *pre-cl-flt* vyšší klauze (narušení principu členské sounáležitosti).
- Je-li druhý z uzlů na posledním místě v seznamu *d-list*, může být v seznamu *s-list* na jiném než posledním místě a označen jako součást intonačního centra (uplatní se princip důraznosti).

Obecná omezení povrchového slovosledu

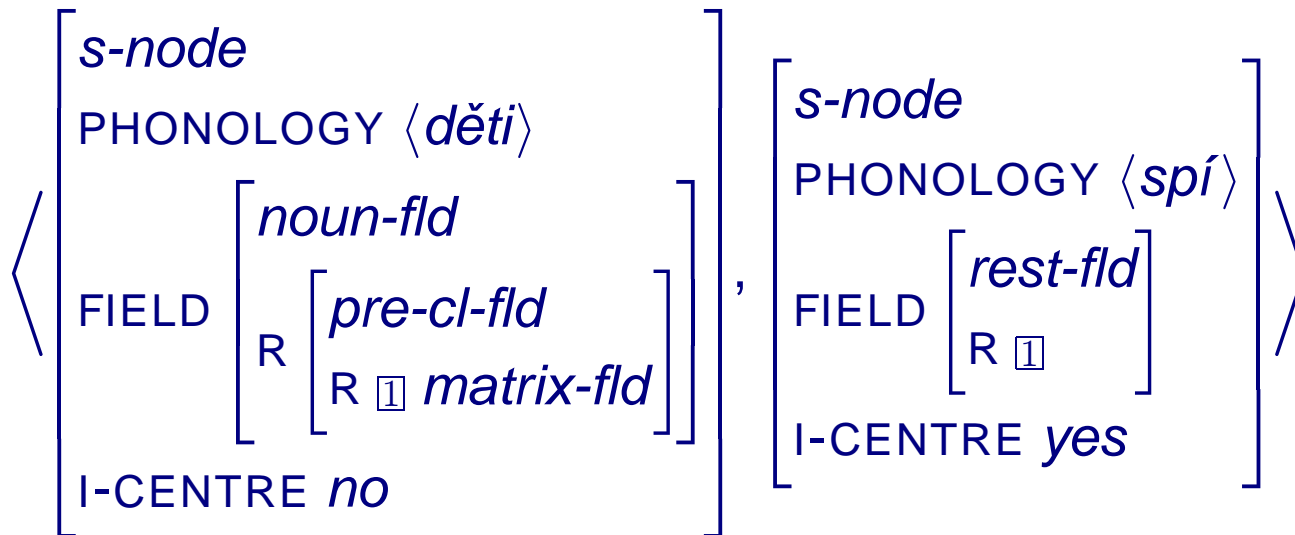
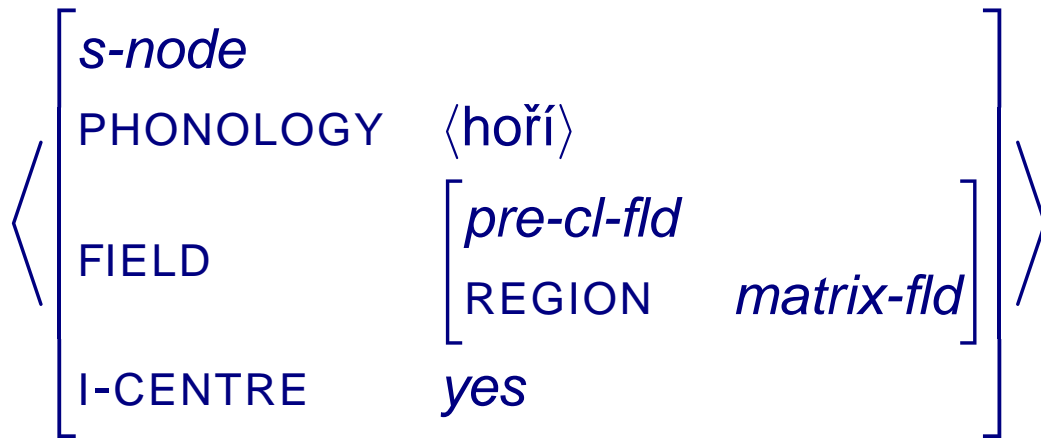
Matrix Compaction Principle – všechny struktury *s-node* ve větě jsou částí jednoho pole *matrix fld*

Planarity Principle – spojitost polí („projektivita“)

Topological Order, Field Existence, Field Uniqueness Principles – definice skladby topologických oblastí:

- pořadí, existence (≥ 1) a jedinečnost ($=1$) polí v oblasti
- konkrétní hodnoty definovány relacemi specifikovanými pro přehlednost v podobě tabulek

Oblast	Pole	Pořadí	Obsazení
<i>matrix-flđ</i>	<i>pre-cl-flđ</i>	1	1
	<i>cl-flđ</i>	2	≤1
	<i>rest-flđ</i>	3	neurčeno
	<i>fin-flđ</i>	4	≤1



field

matrix-flid

embedded-flid REGION *field*

pre-cl-flid

cl-flid

rest-flid

fin-flid

dep-flid

noun-flid

pp-flid

adj-flid

adv-flid

inf-flid

aux-flid

emb-clis-flid

sconj-clis-flid

wh-clis-flid

Oblast	Pole	Pořadí	Obsazení
<i>pre-cl-flid</i>	<i>dep-flid</i>	1	1
<i>rest-flid</i>	<i>dep-flid</i>	any	any

Oblast	Pole	Pořadí	Obsazení
<i>noun-flđ</i>	<i>l-adj-flđ</i>	1	any
	<i>noun-flđ</i>	2	≥ 1
	<i>emb-clš-flđ</i>	3	≤ 1
	<i>h-adj-flđ</i>	≥ 4	any
	<i>pp-flđ</i>	≥ 4	any
<i>pp-flđ</i>	<i>prep-flđ</i>	1	1
	<i>noun-flđ</i>	2	1
<i>l-adj-flđ</i>	<i>adv-flđ</i>	1	any
	<i>adj-flđ</i>	2	1
<i>h-adj-flđ</i>	<i>adv-flđ</i>	1	any
	<i>adj-flđ</i>	2	1
	<i>pp-flđ</i>	≥ 3	any
	<i>emb-clš-flđ</i>	≥ 3	any
	<i>noun-flđ</i>	≥ 3	any
<i>adv-flđ</i>	<i>adv-flđ</i>	1	≤ 1
	<i>adv-flđ</i>	2	1

Oblast	Pole	Pořadí	Obsazení
<i>cl-fld</i>	<i>cl-lis-fld</i>	1	≤ 2
	<i>cl-be-fld</i>	2	≤ 1
	<i>cl-rfl-fld</i>	3∨4	≤ 1
	<i>cl-ethdat-fld</i>	3∨4	≤ 1
	<i>cl-freedat-fld</i>	5	neurčeno
	<i>cl-dat-fld</i>	6	neurčeno
	<i>cl-acc-fld</i>	7	neurčeno
	<i>cl-gen-fld</i>	8	neurčeno
	<i>cl-ins-fld</i>	9	neurčeno
	<i>cl-nom-fld</i>	neurčeno	≤ 1
	<i>cl-uz-fld</i>	neurčeno	≤ 1
	<i>cl-pry-fld</i>	neurčeno	≤ 1
	<i>cl-vsak-fld</i>	neurčeno	≤ 1

Specifická omezení povrchového slovosledu

- komparativní konstrukce
(*menší vesnice než Lhota*)
- některé případy vzdálených závislostí
(*Koho jsi říkal, že Marie myslela, že Pavel pozve?*)
- rozdělené předložkové fráze
(*O jakou se jedná soutěž?*)
- postavení příklonek, haplologie
(*Miloš se jim nakonec rozhodl omluvit.*)

8 Příklady

Příklad 1: šplhání příklonek

1. A clitic can climb to a higher 2P unless it is governed by a finite verb (1), a deverbative (gerund) (2), an adjectival participle (3), or an adverbial participle (4).

- (1) a. Šéf *ho* nařídil zbavit všech výsad.
b. Šéf nařídil, aby *ho* zbavili všech výsad.
c. *Šéf *ho* nařídil, aby zbavili všech výsad.
- (2) a. Dědeček nemá rád dětské ušklíbání se nad polévkou.
b. *Dědeček se nemá rád dětské ušklíbání nad polévkou.
- (3) a. Uvíтали bychom více takových kajících se hříšníků.
b. *Uvíтали bychom se více takových kajících hříšníků.
- (4) a. Ředitel vzhlédl od dopisu, tváře se ustaraně.
b. *Ředitel se vzhlédl od dopisu, tváře ustaraně.

2. A clitic may only climb through a domain governed by an infinitive.
3. A more deeply embedded clitic cannot climb over a less deeply embedded clitic (5).

- (5)
- a. Pavel se snažil *mu* pomoci *ho* najít.
 - b. Pavel se *mu* snažil *ho* pomoci najít.
 - c. Pavel se *mu ho* snažil pomoci najít.
 - d.*?Pavel se *ho* snažil *mu* pomoci najít.

4. Two phonologically identical clitics with different governors either do not co-occur in a single clitic cluster or haplologize (6).

- (6)
- a. Kamila mi slíbila to vrátit MNĚ.
 - b. *Kamila mi mi to slíbila vrátit.
 - c. Kamila mi to slíbila vrátit.
 - d. Kamila mi slíbila mi to vrátit.

5. Two reflexives *si* and *se* can haplogogize yielding *si*, if the reflexive *si* originates in a more embedded domain (7).

- (7)
- a. Jan *si* bál vzít kravatu. [KO]
 - b. Snažím *si* to představit.
 - c. *Snažím *se* to představit.
 - d. Styděla *si* sednout do první řady.
 - e. *Styděla *se* sednout do první řady.
 - f. *Troufla *si* usadit v první řadě.
 - g. *Troufla *se* usadit v první řadě.

6. The order of pronominal clitics with the same morphological case co-occurring in a single cluster due to clitic climbing corresponds to the level of embedding of their governors.

Jeden z příkladů, který se zatím nepodařilo formálně popsat:
pořadí nepřízvučných zájmen v dativu

- (8) a. [Podařilo se *mi*] [*mu* amputovat pravou zadní nohu].
b. [Podařilo se *mu*] [*mi* udělat velkou radost].
- (9) a. Včera se *ti jí* to konečně povedlo vysvětlit.
b. Včera se *jí ti* to konečně povedlo vysvětlit.
- (10) a. *Vysvětlit* se *mi jí* to konečně povedlo až včera.
b. *Vysvětlit* se *jí mi* to konečně povedlo až včera.

Příklad z korpusu:

- (11) Naštěstí se mu jim podařilo chlapce v bezvědomí vyrvat a utéci s ním za plot. (LN 96)

Příklad 2: modifikátory substantiv v první pozici

- (12) a. Jaká se *vám* vybaví představa?
b. Vážně nevím, kterou *si* vybral nevěstu.
- (13) a. PĚKNÁ se nám vybaví představa.
b. BOHATOU *si* vybral nevěstu.
c. TAKOVÝ se *mi* líbí básník.

Ale proč ne taky tohle:

- (14) a. *TENHLE *mi* člověk slíbil peníze.
b. *VYSOKÝ *mi* slíbil peníze člověk.

Příklad 3: rozdělené předložkové fráze

- (15) O jakou se jedná soutěž?
- (16) a. O jak dotovanou soutěž se jedná?
b. O jak dotovanou se jedná soutěž?
c. *O jak se jedná dotovanou soutěž?
d. ?O velmi dobře dotovanou se jedná soutěž.
- (17) O jakou jsi myslel, že se jedná soutěž?

(18) PP COMPACTION:

$$\left(\begin{array}{l} \left[\begin{array}{l} \text{SYNSEM} \mid \text{LOCAL} \mid \text{CATEGORY} \mid \text{HEAD } \textit{noun} \\ \text{NONHEAD-DAUGHTERS } \boxed{1} \\ \text{SURFACE } \boxed{2} \langle [\text{FIELD } \textit{prep-flt}] \rangle \bigcirc \boxed{3} \bigcirc \boxed{4} \end{array} \right] \\ \wedge \text{MEMBER}([\text{SURFACE } \boxed{2}], \boxed{1}) \\ \wedge \text{MEMBER}([\text{SURFACE } \boxed{3}], \boxed{1}) \end{array} \right)$$

$$\rightarrow \exists \boxed{5} \left(\begin{array}{l} \left(\begin{array}{l} \text{APPEND}(\boxed{2}, \boxed{3}, \boxed{5}) \\ \wedge \text{REGION}(\boxed{5}, \textit{pre-cl-flt}) \\ \wedge \text{REGION_SETUP}(\boxed{4}, \textit{noun-flt}) \end{array} \right) \\ \vee \left(\begin{array}{l} \text{SHUFFLE}(\boxed{3}, \boxed{4}, \boxed{5}) \\ \wedge \text{REGION_SETUP}(\boxed{5}, \textit{noun-flt}) \end{array} \right) \end{array} \right)$$

9 Výsledky a výhledy

- formalizace FGP pomocí RSRL
- popis vztahu hloubkového a povrchového slovosledu
- popis interakce hloubkového slovosledu s povrchovými pravidly
- kompatibilita s Mathesiovými slovoslednými principy
- formalizace několika zákonitostí povrchového slovosledu s interakcí syntaktických, diskursních a prozodických faktorů

Co dál:

- implementace
- více jevů
- více jazyků
- modifikace?

A to je konec ...