

Pokus o formální popis českého slovosledu

Alexandr Rosen
Ústav teoretické a počítačové 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 Repräsentace	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 ≠ 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

Slovsled

- Slovosled odráží více faktorů z více rovin (slovsledné faktory: aktuální členění, ‚povrchová‘ pravidla gramatiky)
- Slovsledné 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 slovsledných faktorů

„deklarativní formalismus“:
Relational Speciate Re-entrant Logic – RSRL
Paul King (1989), Frank Richter (2000)

- **Cíl práce:**
popsat podstatnou část slovsledných 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 „signature“ a „teorie“.

Signatura

definuje hierarchii typů – množin objektů modelu.

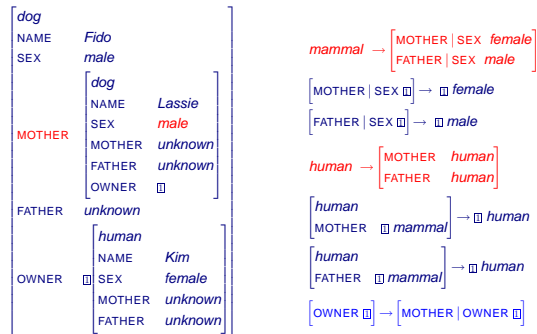
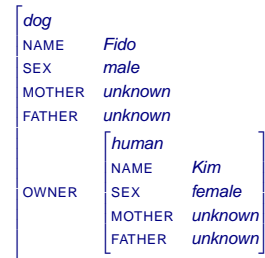
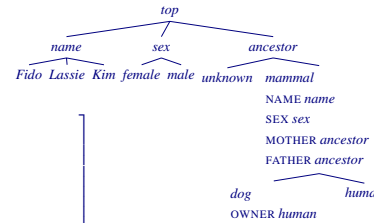
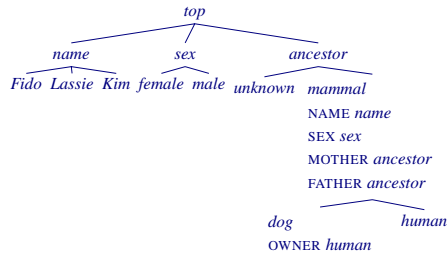
- Každému objektu modelu je přiřazen právě jeden maximálně specifický typ.
- Typy mají předepsané vlastnosti – atributy. Předepsané hodnoty atributů jsou opět typy.
- Objekty modelu musí mít všechny předepsané atributy, jejichž hodnoty musí být maximálně specifické typy.
- Typy stojící v hierarchii níže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

Navíc: seznamy, množiny.

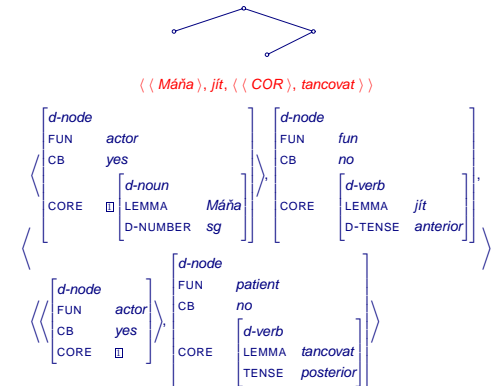
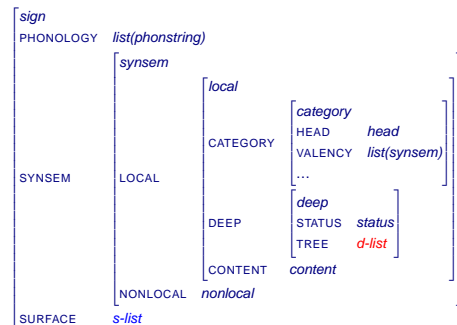
Teorie

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.



5 Repräsentace

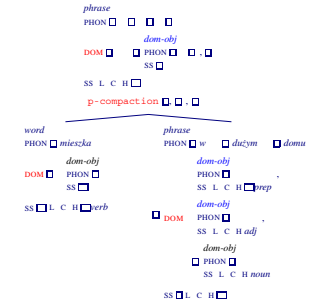
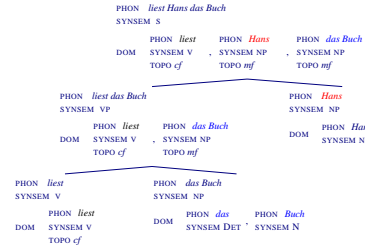


Povrchový řetězec

Inspirace č. 1: Mike Reape, Andreas Kathol

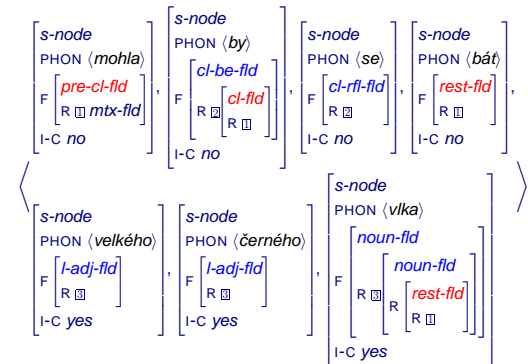
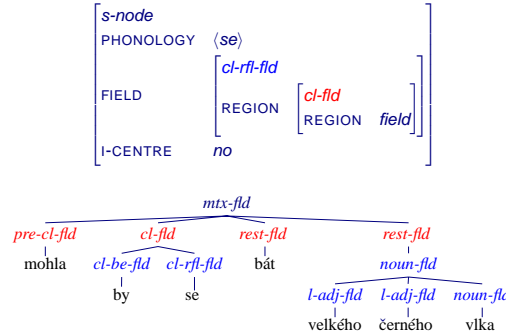
- Horizontální pořadí terminálů v derivačním stromě se nerovná povrchovému slovosledu
- Slovosledné domény, objekty domén
- Relace *shuffle* a *compaction*
- Topologická pole

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



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



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* řídicí dcery, do něhož jsou vloženy seznamy *d-list* ostatních dcer.

$$\text{non-lexical} \rightarrow \left(\begin{array}{c} \text{SYNSEM | LOCAL | DEEP | TREE } \square \\ \text{HEAD-DAUGHTER | SYNSEM | LOCAL | DEEP | TREE } \square \\ \text{NONHEAD-DAUGHTERS } \square \\ \wedge \text{ COLLECT_DLISTS } (\square, \square) \\ \wedge \text{ APPEND } (\square, \square, \square) \\ \wedge \text{ PERMUTE } (\square, \square) \end{array} \right)$$

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 řídící 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-flid* 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-flid*

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-flid</i>	<i>pre-cl-flid</i>	1	1
	<i>cl-flid</i>	2	≤ 1
	<i>rest-flid</i>	3	neurčeno
	<i>fin-flid</i>	4	≤ 1

s-node
PHONOLOGY (hoří)
FIELD [*pre-cl-flid*]
REGION [*matrix-flid*]
I-CENTRE yes

s-node
PHONOLOGY (děti)
FIELD [*noun-flid*]
R [*pre-cl-flid*]
R [*matrix-flid*]
I-CENTRE no

s-node
PHONOLOGY (spí)
FIELD [*rest-flid*]
R []
I-CENTRE yes

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
scorj-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-flid</i>	<i>l-adj-flid</i>	1	any
	<i>noun-flid</i>	2	≥ 1
	<i>emb-clis-flid</i>	3	≤ 1
	<i>h-adj-flid</i>	≥ 4	any
	<i>pp-flid</i>	≥ 4	any
<i>pp-flid</i>	<i>prep-flid</i>	1	1
	<i>noun-flid</i>	2	1
<i>l-adj-flid</i>	<i>adv-flid</i>	1	any
	<i>adj-flid</i>	2	1
<i>h-adj-flid</i>	<i>adv-flid</i>	1	any
	<i>adj-flid</i>	2	1
	<i>pp-flid</i>	≥ 3	any
	<i>emb-clis-flid</i>	≥ 3	any
	<i>noun-flid</i>	≥ 3	any
<i>adv-flid</i>	<i>adv-flid</i>	1	≤ 1
	<i>adv-flid</i>	2	1

Oblast	Pole	Pořadí	Obsazení
<i>cl-flid</i>	<i>cl-lis-flid</i>	1	≤ 2
	<i>cl-be-flid</i>	2	≤ 1
	<i>cl-rtl-flid</i>	3∨4	≤ 1
	<i>cl-ethdat-flid</i>	3∨4	≤ 1
	<i>cl-freedat-flid</i>	5	neurčeno
	<i>cl-dat-flid</i>	6	neurčeno
	<i>cl-acc-flid</i>	7	neurčeno
	<i>cl-gen-flid</i>	8	neurčeno
	<i>cl-ins-flid</i>	9	neurčeno
	<i>cl-nom-flid</i>	neurčeno	≤ 1
	<i>cl-uz-flid</i>	neurčeno	≤ 1
	<i>cl-pry-flid</i>	neurčeno	≤ 1
	<i>cl-vsak-flid</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, haplogogie (*Miloš se jim nakonec rozhodl omluvit.*)

8 Příklad

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 ušklibání se nad polévkou.
b. *Dědeček se nemá rád dětské ušklibání nad polévkou.
 - (3) a. Uvítali bychom více takových kajících se hříšníků.
b. *Uvítali 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 haplogogize (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 nepřizvučných zájmen v dativu

- (8) a. [Podařilo se *mí*] [*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 *ji tí* 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:

$$\left(\begin{array}{l} \left[\begin{array}{l} \text{SYNSEM} \mid \text{LOCAL} \mid \text{CATEGORY} \mid \text{HEAD } \textit{noun} \\ \text{NONHEAD-DAUGHTERS } \square \\ \text{SURFACE } \square([\text{FIELD } \textit{prep-fld}] \circ \square \circ \square) \\ \wedge \text{ MEMBER}(\{\text{SURFACE } \square, \square\}) \\ \wedge \text{ MEMBER}(\{\text{SURFACE } \square, \square\}) \end{array} \right] \\ \rightarrow \exists \square \left(\begin{array}{l} \left(\begin{array}{l} \text{APPEND}(\square, \square, \square) \\ \wedge \text{ REGION}(\square, \textit{pre-cl-fld}) \\ \wedge \text{ REGION_SETUP}(\square, \textit{noun-fld}) \end{array} \right) \\ \vee \left(\begin{array}{l} \text{SHUFFLE}(\square, \square, \square) \\ \wedge \text{ REGION_SETUP}(\square, \textit{noun-fld}) \end{array} \right) \end{array} \right)$$

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 ...