The background of the slide is an underwater scene with several fish swimming in clear blue water. The fish are of various sizes and are scattered throughout the frame, creating a sense of depth and movement.

# Dekomprese v popisu jazyka *aneb* hlubiny i mělčiny deklarativně

Alexandr Rosen  
Ústav teoretické a počítačnické lingvistiky  
Universita Karlova v Praze  
alexandr.rosen@ff.cuni.cz

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

4. 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“*

- Potíže při vynořování z hlubin – přechod mezi rovinami (z hloubkové na povrchovou)
- Složková syntax a CF gramatika: derivace = reprezentace
- Závislostní syntax: záleží na formalismu, navíc ve FGD je reprezentace abstraktnější (pomocná slova) → roviny popisu
- Vztah mezi povrchovým řetězcem a jeho reprezentací – záležitost teoretického popisu
- Tradiční přístup: stratifikační a procedurální, jde to jinak?

## Vývoj

„Lexikální informace ve strojovém překladu“



„Strojový překlad jako pokusný záhonek“



„Závislostní gramatika v deklarativní formalizaci, aplikovaná na kontrastivní popis dvou jazyků“



„Pokus o popis některých jevů českého slovosledu pomocí FGP v deklarativní formalizaci“

# Obsah

- 1 Motivace 5
- 2 Hypotézy 8
- 3 Teorie: FGP 9
- 4 Formalismus: RSRL 10
- 5 Rerezentace 16
- 6 Základy syntaxe 24
- 7 Slovosledná omezení 26
- 8 Výsledky a výhledy 32

# 1. Motivace

(jsou-li abstraktní entity pro poznání a popis jazyka nutné ...)

Jaký je vztah mezi povrchoým řetězcem fonémů/grafémů a jeho reprezentací?

- Stochastické metody – simulace vědomého poznání
- Teoretická lingvistika – různé hypotézy

## 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

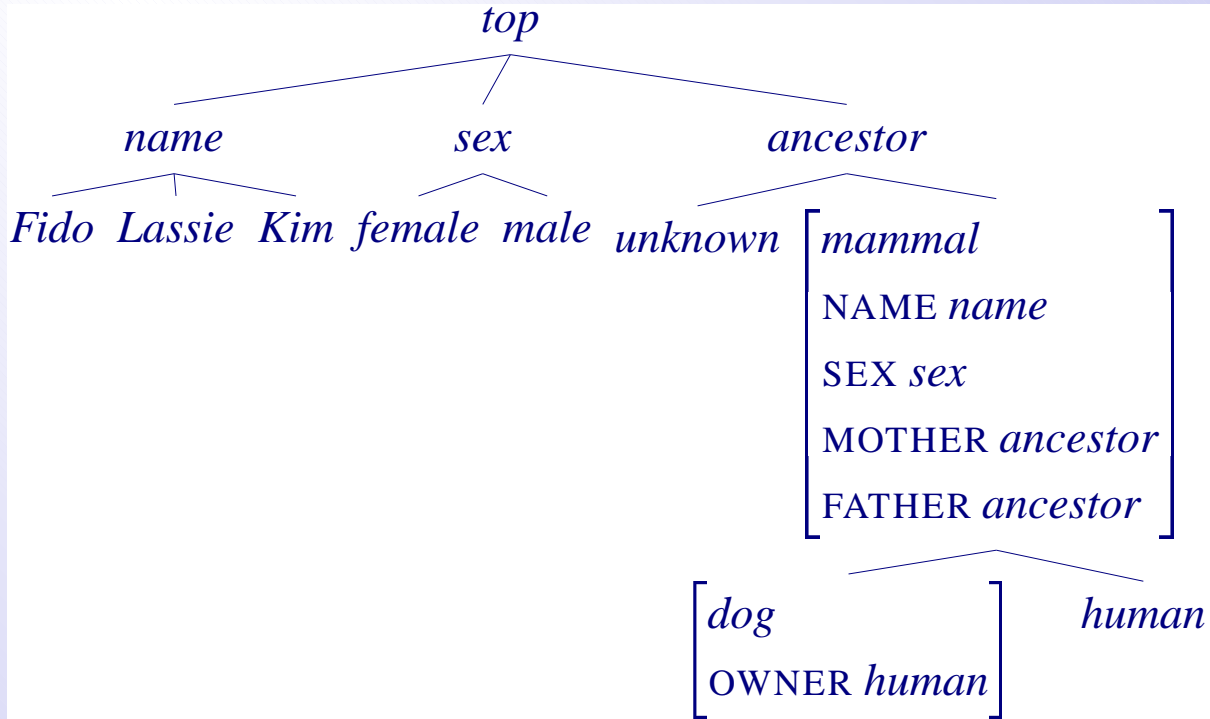
## 4. Formalismus: RSRL

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.



<i>dog</i>																	
NAME	<i>Fido</i>																
SEX	<i>male</i>																
MOTHER	<i>unknown</i>																
FATHER	<i>unknown</i>																
	<table> <tr> <td colspan="3"><i>human</i></td> </tr> <tr> <td>NAME</td> <td colspan="2"><i>Kim</i></td> </tr> <tr> <td>SEX</td> <td colspan="2"><i>female</i></td> </tr> <tr> <td>MOTHER</td> <td colspan="2"><i>unknown</i></td> </tr> <tr> <td>FATHER</td> <td colspan="2"><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>																
OWNER																	

## 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> <td></td> </tr> <tr> <td>NAME</td> <td><i>Lassie</i></td> <td></td> </tr> <tr> <td>SEX</td> <td><i>male</i></td> <td></td> </tr> <tr> <td>MOTHER</td> <td><i>unknown</i></td> <td></td> </tr> <tr> <td>FATHER</td> <td><i>unknown</i></td> <td></td> </tr> <tr> <td>OWNER</td> <td>1</td> <td></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> <td></td> </tr> <tr> <td>NAME</td> <td><i>Kim</i></td> <td></td> </tr> <tr> <td>SEX</td> <td><i>female</i></td> <td></td> </tr> <tr> <td>MOTHER</td> <td><i>unknown</i></td> <td></td> </tr> <tr> <td>FATHER</td> <td><i>unknown</i></td> <td></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 *mammal*] → [MOTHER | SEX *female*]

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

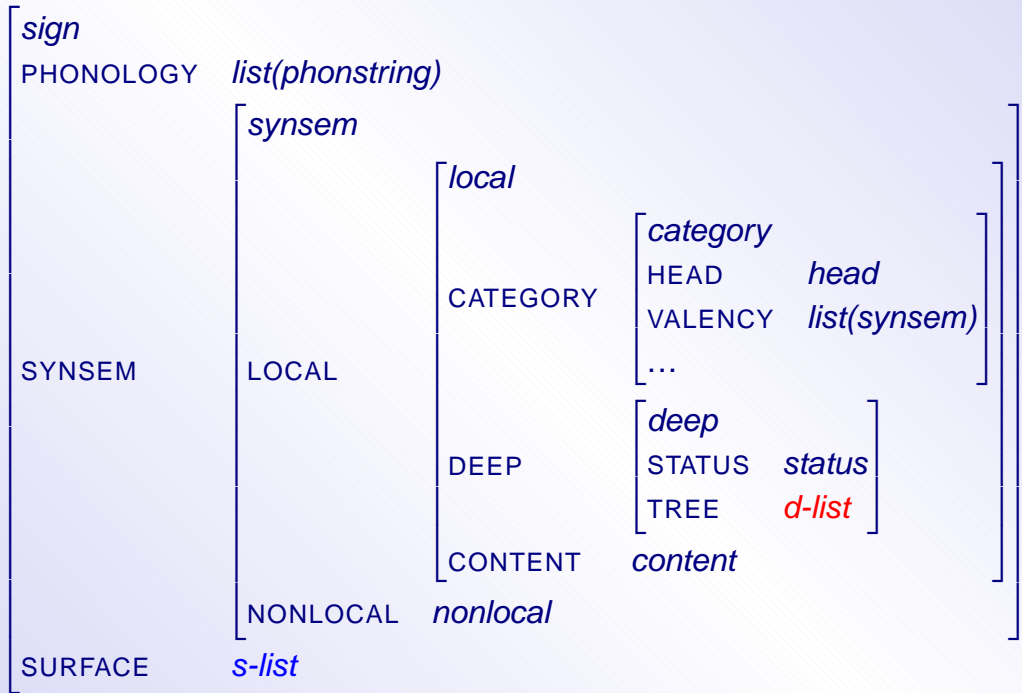
*human* → [MOTHER *human*  
FATHER *human*]

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

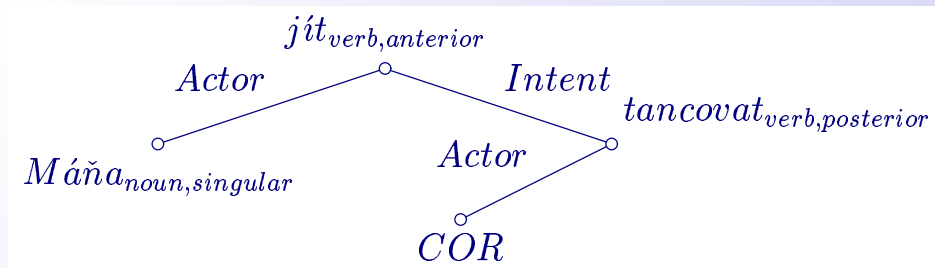
[*human*  
FATHER *mammal*] → [FATHER *human*]

*dog* → [OWNER 1  
MOTHER | OWNER 1]

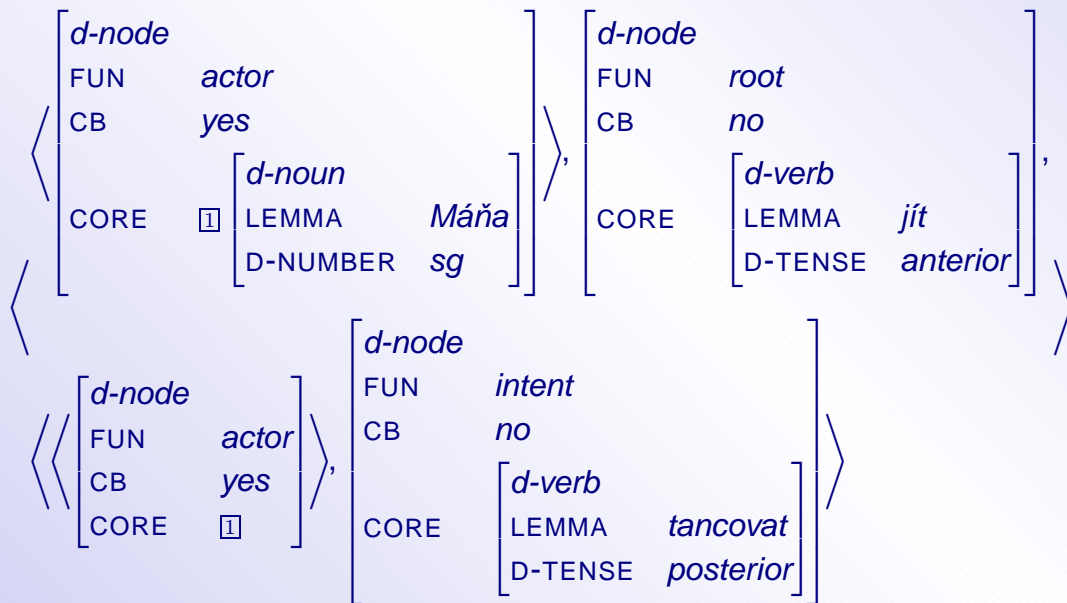
# 5. Rerezentace







*<< Máňa >>, jít, << COR >>, tancovat >>*

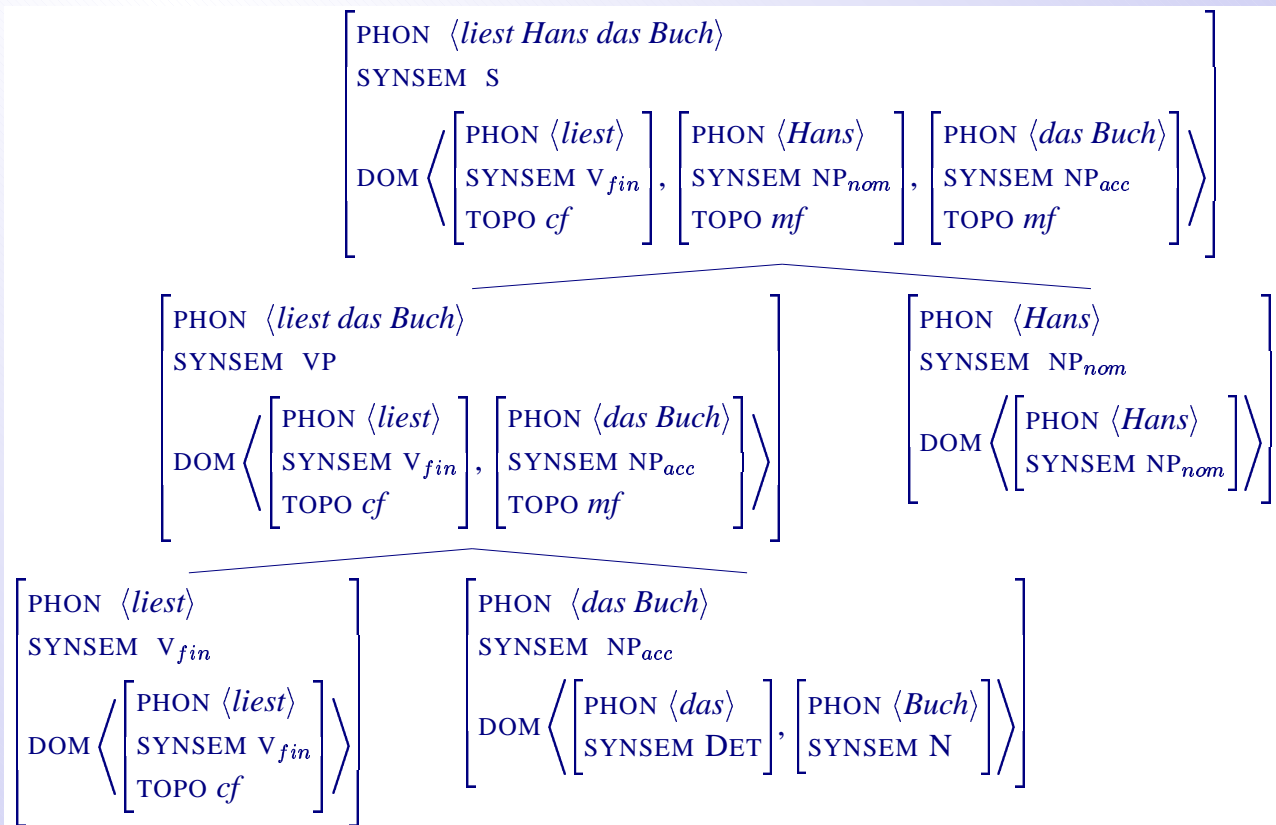


## Povrchový řetězec

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

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



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

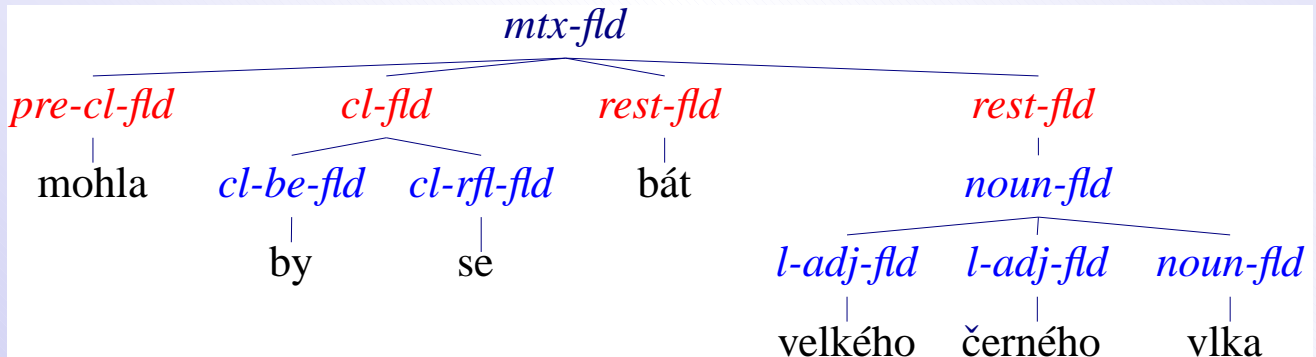
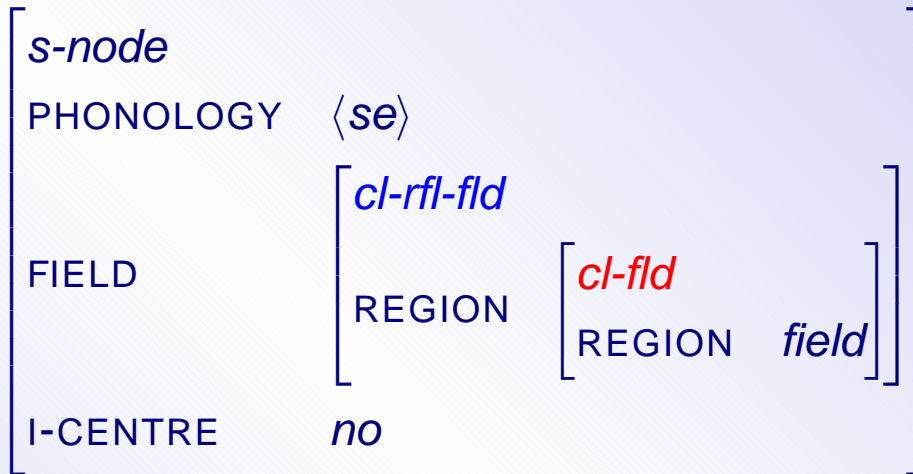
$$\left[ \begin{array}{l} \textit{word} \\ \text{PHON } \boxed{1} \langle \textit{mieszka} \rangle \\ \\ \text{DOM } \boxed{5} \left\langle \begin{array}{l} \textit{dom-obj} \\ \text{PHON } \boxed{1} \\ \text{SS } \boxed{10} \end{array} \right\rangle \\ \text{SS } \boxed{10} \left[ \boxed{L} \mid \boxed{C}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

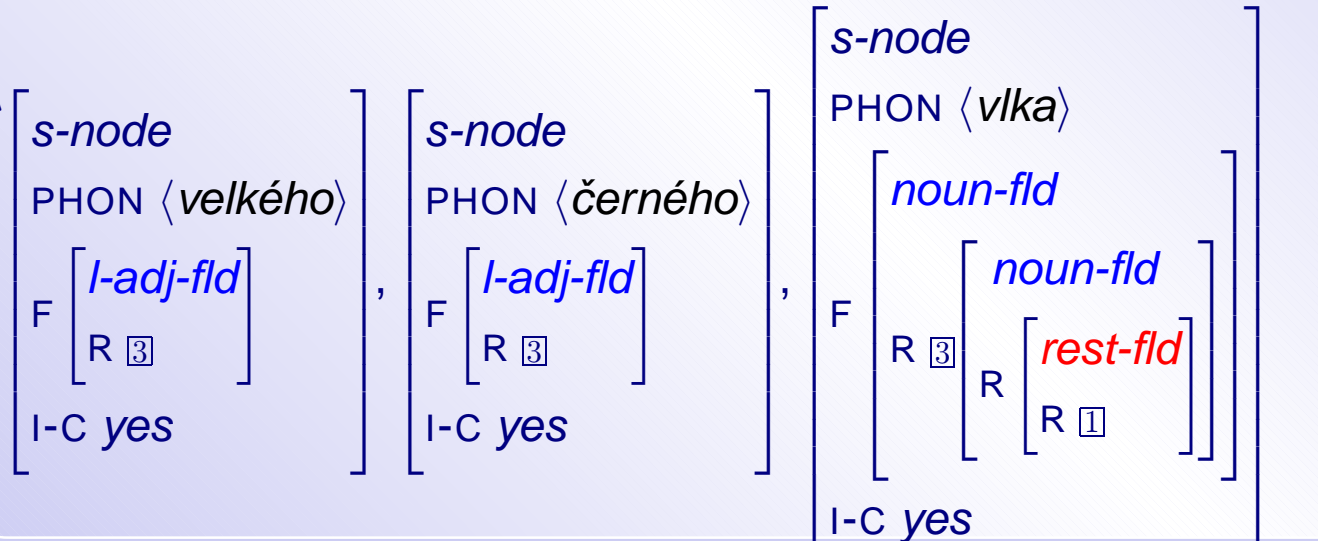
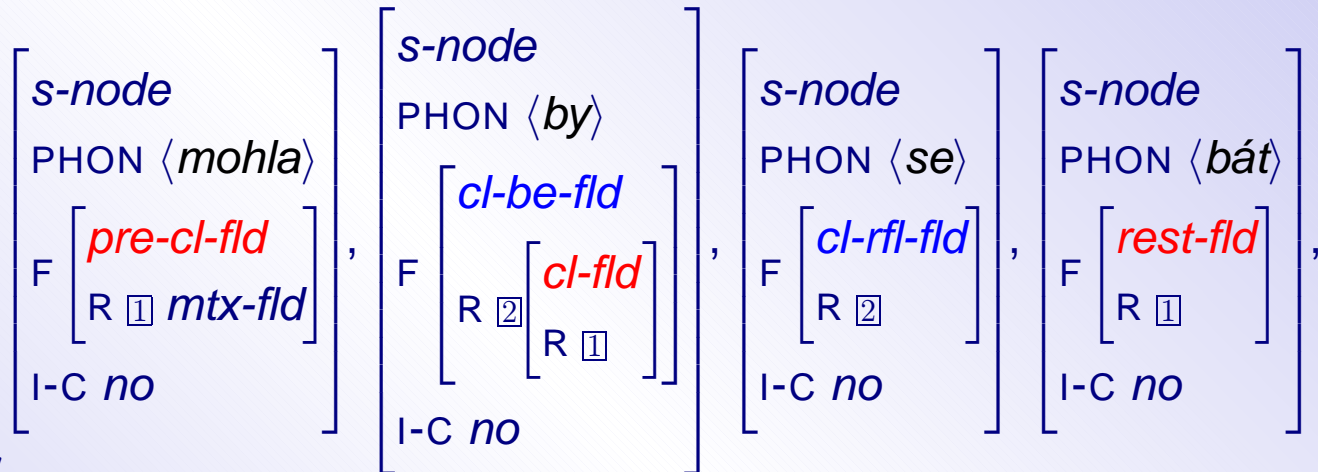
$$\left[ \begin{array}{l} \textit{phrase} \\ \text{PHON } \boxed{2} \langle \textit{w} \rangle \oplus \boxed{3} \langle \textit{dużym} \rangle \oplus \boxed{4} \langle \textit{domu} \rangle \\ \\ \begin{array}{l} \textit{dom-obj} \\ \text{PHON } \boxed{2} \\ \text{SS } \boxed{L} \mid \boxed{C} \mid \boxed{H}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

## 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. Základy syntaxe

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* →

$$\left( \begin{array}{l} \text{SYNSEM | LOCAL | DEEP | TREE } \boxed{5} \\ \text{HEAD-DAUGHTER | SYNSEM | LOCAL | DEEP | TREE } \boxed{1} \\ \text{NONHEAD-DAUGHTERS } \boxed{2} \\ \wedge \text{ COLLECT\_DLISTS } (\boxed{2}, \boxed{3}) \\ \wedge \text{ APPEND } (\boxed{1}, \boxed{3}, \boxed{4}) \\ \wedge \text{ PERMUTE } (\boxed{4}, \boxed{5}) \end{array} \right)$$

## 7. Slovosledná omezení

- 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-fld* 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 ( $\geq 1$ ) a jedinečnost ( $=1$ ) polí v oblasti
- konkrétní hodnoty definovány relacemi specifikovanými pro přehlednost v podobě tabulek

viz table.pdf  
viz table1.pdf

## 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. 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 ...