

Snídaně  
s Českým akademickým korpusem

18. června 2007, 8:30

# Anotování Českého akademického korpusu

- PROČ? JAK?

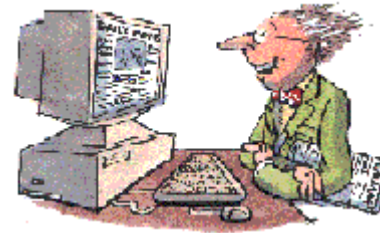
# Kdo

- **Hosté**
  - Veronika Čurdová
  - Lenka Žehrová
- **Domáci** ... viz dále



# Anotování korpusu. Proč?

- teoretická lingvistika
- počítačová lingvistika
  - strojové učení
  - trénovací data: anotované korpusu
  - aplikace: tagging, parsing, ...





*"The computer is claiming its intelligence is real, and ours is artificial."*



# Navštivte stránku

popularizačních článků a rozhovorů

[http://ufal.mff.cuni.cz/?a=popular&m=student\\_info](http://ufal.mff.cuni.cz/?a=popular&m=student_info)

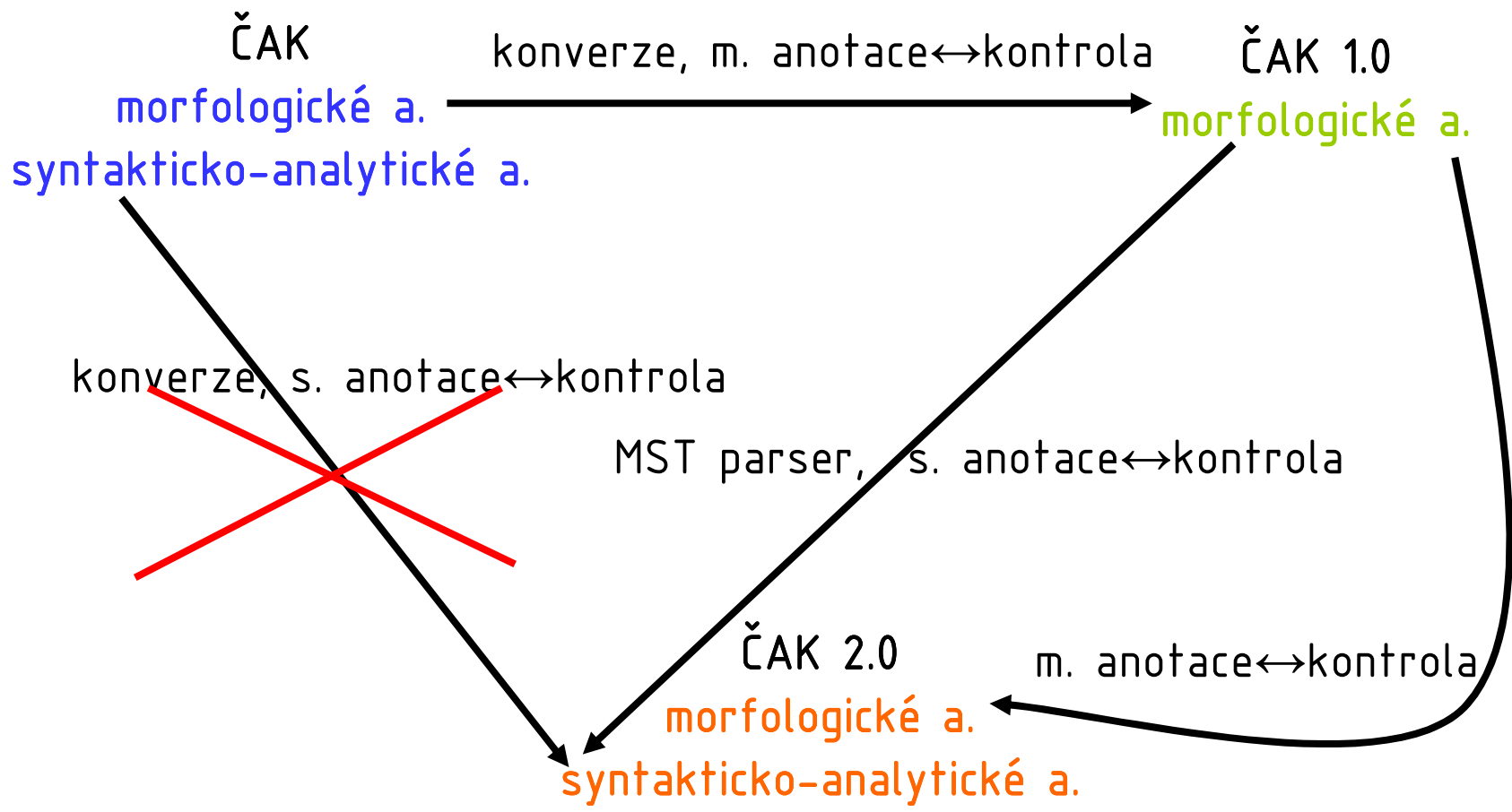
# projekt „Data a nástroje pro informační systémy“ (REST)

- <http://ufal.mff.cuni.cz/rest>
- datová komponenta – **Český akademický korpus**
- komponenta nástrojů



# Český akademický korpus

- 1971–1985 – ÚJČ AV
- 540 000 morfologicky a syntakticky anotovaných slov
- psané a mluvené texty
- publicistika, administrativa, odborné statě
- **anomálie**
  - vymazaná interpunkce
  - vymazané ciferné výrazy



# Na cestě k syntakticko-analytickým a. ČAK 2.0

- MST parser

VSTUP: ČAK 1.0 (w, m soubory)

VÝSTUP: a' soubory

ÚSPĚŠNOST: 84,6% (na hranách)

- s. anotace ↔ kontrola

VSTUP: a' soubory

VÝSTUP: a soubory

## s. anotace ↔ kontrola

1. dvojitá anotace v editoru TrEd (anotátoři)
2. technická podpora v TrEd pro řešení mezianotátorských odlišností; (ne)shody v číslech (Kiril Ribarov)
3. vyřešení odlišností v TrEd (Alla Bémová, Zdeňka Urešová)
4. automatické kontrolní skripty à kontrola PDT 2.0 (Jiří Mírovský)

paralelně m. anotace ↔ kontrola (Jiří Mírovský)

# s. anotace ↔ kontrola: TrEd

## Zvláštní pozornost:

- Správné použití tzv. kontextů v TrEdu. Nově upravený kontext pro snadnou anotaci PML\_CAC\_A\_Edit
- `guessed_form` je součástí `m-souboru`. Slouží k upřesnění, dle chápání anotátora, obsahu prázdných uzlů (původně chybějících slov, či někdy špatně vložených slov).  
*Pozor: Nutně ukládat i m-soubor -> viz „Save“*
- Při anotaci koordinace a apozice, nastavit u všech relevantních uzlů podstromu hodnotu `is_member` na 1.
- Při anotaci parentéze, nastavit hodnotu u všech relevantních uzlů podstromu hodnotu `is_parenthesis_root` na 1.

## s. anotace ↔ kontrola: od anotátorů

- m-soubor (původní název)
- a-soubor (původní název doplněný iniciálami)
- poznámky

# Termín konce ...

- nesdělují. Stačí, když jsem nervózní já, příp. pane ředitel 😊
- Při dodržování dílčích termínů vše stihneme v termínu!

# Odměna

- 5 Kč à věta + odměny
- DP[ČP]



# Co dál?

- demonstrace TrEd
- debata nad odlišnostmi v anotacích

TTree Editor Default(1/1): C:\Documents and Settings\baradmin\Plocha\porovna...

File View Node Session Bookmarks User-defined Help Context: PML\_A\_View

Především to je problém političnosti školy 17/148

```

graph TD
    Root((a-s03w-s17)) --- AuxS[AuxS]
    AuxS --- Node1(( ))
    Node1 --- je[je  
VB-S-3P-AA  
Pred]
    Node1 --- Node2(( ))
    je --- Node3(( ))
    Node3 --- Predvem[Především  
Db  
AuxY  
-a(AuxZ)]
    Node3 --- to[to  
PDNS1  
Sb]
    Node2 --- problem[problém  
NNIS1-A  
Pnom]
    problem --- Node4(( ))
    Node4 --- politicnosti[političnosti  
NNFS2-A  
Atr]
    Node4 --- skoly[školy  
NNFS2-A  
Atr]
  
```

id: a-s03w-s17

