

## Slovenský národný korpus (2002 – 2012): východiská, ciele a výsledky pre výskum a prax

Mária Šimková, Radovan Garabík

Jazykovedný ústav Ľ. Štúra SAV, Bratislava

korpus@korpus.sk

Rozvoj jazykových technológií a zdrojov na Slovensku a vo svete  
Bratislava, 7. 6. 2012 – 8. 6. 2012

# Východiská



- **Vonkajšie faktory**
  - technologizácia a informatizácia
  - rozvoj korpusov a korpusovej lingvistiky
  - prístupové konania do Európskej únie
- **Vnútorne faktory**
  - tvorba nového výkladového slovníka
  - koncepcia starostlivosti o štátny jazyk

# Východiská



- **interný korpus textov slovenského jazyka v JÚLŠ SAV**
  - základné počítačové vybavenie
  - voľne dostupné programy
  - texty bez licencie – cca 25 mil. jednotiek
  - 3 prac. s čiastkovou kapacitou
- **snahy o podporu v grantovej agentúre VEGA**
- **zapájanie sa do medzinárodných projektov**

# Východiská



- **2000 – 2001 vypracovanie a podanie projektov:**
  - *Budovanie Národného korpusu slovenského jazyka*
  - *Elektronizácia jazykovedného výskumu*
- **MK SR, MŠ SR, SAV**
- **uznesenie vlády SR č. 137 z 13. 2. 2002**

# Ciele



- **čo najviac čo najkvalitnejších dát**
- **potreby používateľov a možnosti oddelenia**
- **vývoj korpusov a korpusovej lingvistiky**
  
- **primárne vedecko-výskumné a učebné využitie (lingvistické/lexikografické)**

# Ciele



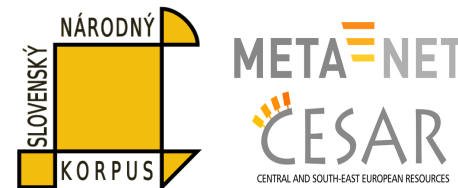
- **všeobecný jednojazyčný korpus**
- **databáza lexikografických diel**
- **terminologická databáza**
- **korpus diachrónnych textov**
- **korpus nárečových textov**
- **paralelné korpusy**
- **korpus hovorenej slovenčiny**

# Jazykové zdroje



- **prim-5.0 (720 mil.) + podkorporusy + skweb**
- **Slovenská terminologická databáza**
- **databáza lingvistických zdrojov**
- **<http://slovniky.korpus.sk/>**
- **<http://www.juls.savba.sk/ediela/>**

# Využitie



- **národné výskumné projekty**
- **500 registrovaných používateľov ročne, 40 000 dopytov denne**
- **slovníky (výkladový, kolokačný, prekladové), monografie, štúdie, postupové práce**



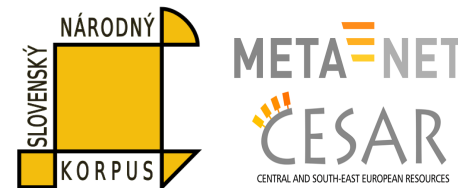
# Medzinárodné projekty



- **Mondilex (FP7)**
  - Conceptual Modelling of Networking of Centres for High-Quality Research in Slavic Lexicography and Their Digital Resources
- **Slovak Online (Lifelong Learning)**
- **EuroMatrixPlus (FP7)**
  - Bringing Machine Translation for European Languages to the User
- **CESAR (ICT PSP)**
  - Central and South-east european Resources
- **MAD: SAV + BAH**
  - Elektronický korpus – konfrontačná štúdia so zameraním na návrh bulharsko-slovenských elektronických jazykových zdrojov

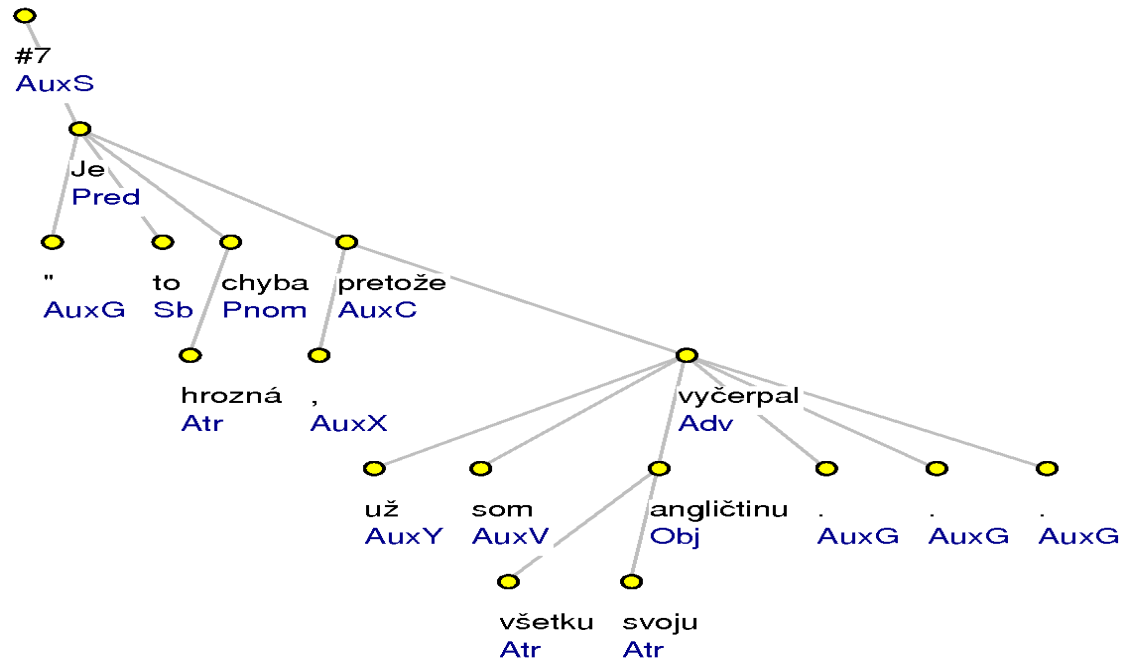
**Pôvodné východiská a ciele sa vzhľadom na súvislosti s týmto mierne posunuli (vzhľadom na podobné zahraničné projekty, zdroje a nástroje)**

# Jazykové zdroje



- **Morfologická databáza META-NET SHARE**
  - pokrýva základnú slovnú zásobu jazyka (96000 lem,  $1.15 \cdot 10^6$  tvarov,  $3.25 \cdot 10^6$  záznamov)
- **Slovenský morfológický (morfosyntaktický) tagset**
- **MULTEXT East (Mondilex)**
  - MULTEXT East tagset
  - 1984

- Syntakticky anotovaný korpus (EuroMatrixPlus)
  - Kompatibilný s PDT1.0 (analytická rovina)
  - 50 tisíc viet, 2× nezávisle anotované



# Jazykové zdroje



- Paralelný slovensko-český korpus (EuroMatrix+) META-NET SHARE
- Paralelný slovensko-anglický korpus (EuroMatrix+) META-NET SHARE
- Paralelný slovensko-ruský korpus
- Paralelný slovensko-bulharský korpus (MAD: JÚLŠ+ИМИ БАН)
- Paralelný slovensko-francúzsky korpus
- Paralelný slovensko-latinský korpus

# Jazykové zdroje



<a href="#">4457</a>	Obraz v zrkadle sa strácal . „ <b>Somár</b> ! Čo si urobil s tou topánkou ? “ skríkol otec , ktorý sa strácal v dyme , stúpajúcim z podrážky .	Rozplýval se i obraz v zrcadle . „ <b>Blbe</b> ! Cos udělal s tou botou ? “ křikl tatínek , který zmizel v kouři , stoupajícím z podrážky .
<a href="#">1682</a>	Nikto tomu nerozumel , ani tučný básnik ; čítal moje básne prasačimi očkami a kričal , ty sviňa prekliata , kde si toto nabral ? <b>A potom sa šiel ožrat</b> na oslavu poézie a plakal : pozrite na toho ťulpasa , to je <b>básnik</b> ! Taký potmehúd , a čo vie napísať !	Nikdo tomu nerozuměl , ani tlustý básník ne ; četl mé básně prasečíma očkama a křičel , ty svině zatracená , kdes tohle vzal ? <b>A pak se šel ožrat</b> na oslavu poezie a plakal : koukejte se na toho <b>blba</b> , to je <b>básník</b> ! Takový tichošlápek , a co dovede napsat !
<a href="#">8193</a>	Fakt je jeden , ak tu chce niekto ukázať , že je schopný niečo urobiť , šanci je neúrekom . Bohužial' , nemožno sa vyhovárať , že šéf je hlúpy a šéfuje len preto , lebo je v partaji . Hádám preto je všade ako - tak a doma najlepšie . “	Fakt je ten , pokud tu chce někdo dokázat , že je schopen něčeho dosáhnout , šanci je nad hlavu . Bohužel , nelze se vymlouvat , že šéf je <b>blb</b> a šéfuje jen proto , že je v partaji . Myslím , že proto je všude dobře , a doma nejlíp . “
<a href="#">9189</a>	ZASIAHNUTIE DOBRÉHO VOJÁKA ŠVEJKA DO SVETOVEJ VOJNY „ Zabili nám teda Ferdinanda , “ povedala posluhovačka pánu Švejkovi , ktorý pred rokmi , keď ho vojenská lekárska komisia definitívne vyhlásila za hlupáka , opustil vojenskú službu a živil	ZASÁHNUTÍ DOBRÉHO VOJÁKA ŠVEJKA DO SVĚTOVÉ VÁLKY “ Tak nám zabili Ferdinanda , “ řekla posluhovačka panu Švejkovi , který opustiv před léty vojenskou službu , když byl definitivně prohlášen vojenskou lékařskou komisí za <b>blba</b> , živil se

# Jazykové zdroje



- WordNet (Slovak Online)
- Web korpus
- Korpus právnych textov
- Slovenský hovorený korpus

META-SHARE

# Perspektívy



- rozširovanie a skvalitňovanie zdrojov
- nárečový korpus
- historický korpus
- špecializované slovníky

Ďakujem(e) za pozornosť