



Словацкий национальный корпус

Radovan Garabík

Jazykovedný ústav Ľ. Štúra

Panská 26

813 64 Bratislava 1

Slovakia

e-mail: korpus@korpus.juls.savba.sk

www: <http://korpus.juls.savba.sk>

Принципы и цели корпуса

- корпус текстов современного словацкого языка (1955 – 2005)
- лингвистическая разметка (лемматизация, морфология)
- представительный корпус
- библиографическая аннотация
- доступен

- ◄ до 1989 – ничего
- ◄ 1989 – 1991 – первые разговоры о необходимости корпуса
- ◄ 1991 – Текстовый корпус словацкого языка
 - без лингвистической обработки
 - лексикография
- ◄ 2002 – проект Словацкого национального корпуса
- ◄ половина 2003 – 30 млн слов
- ◄ конец 2003 – 110 млн слов
- ◄ начало 2004 – лемматизация, морфосинтаксическая разметка
- **сегодня – 192 млн слов**
- *конец 2004 – сбалансированный корпус 30 млн слов*
- *2005 – параллельные корпуса*
- *2006 – представительный корпус 200 млн слов*

Структура данных в корпусе

- архив
 - документы в оригинальном виде
- банк
 - общий формат текстов (почти XML)
- корпусоид
 - стандарт XCES (TEI), лингвистическая разметка
- дата
 - цифровой формат для корпусного менеджера

<tok>

<orth>meč</orth>

<disamb>

<base>meč</base>

<ctag>SSis1</ctag>

</disamb>

<lex>

<base>meč</base>

<ctag>SSis1</ctag>

</lex>

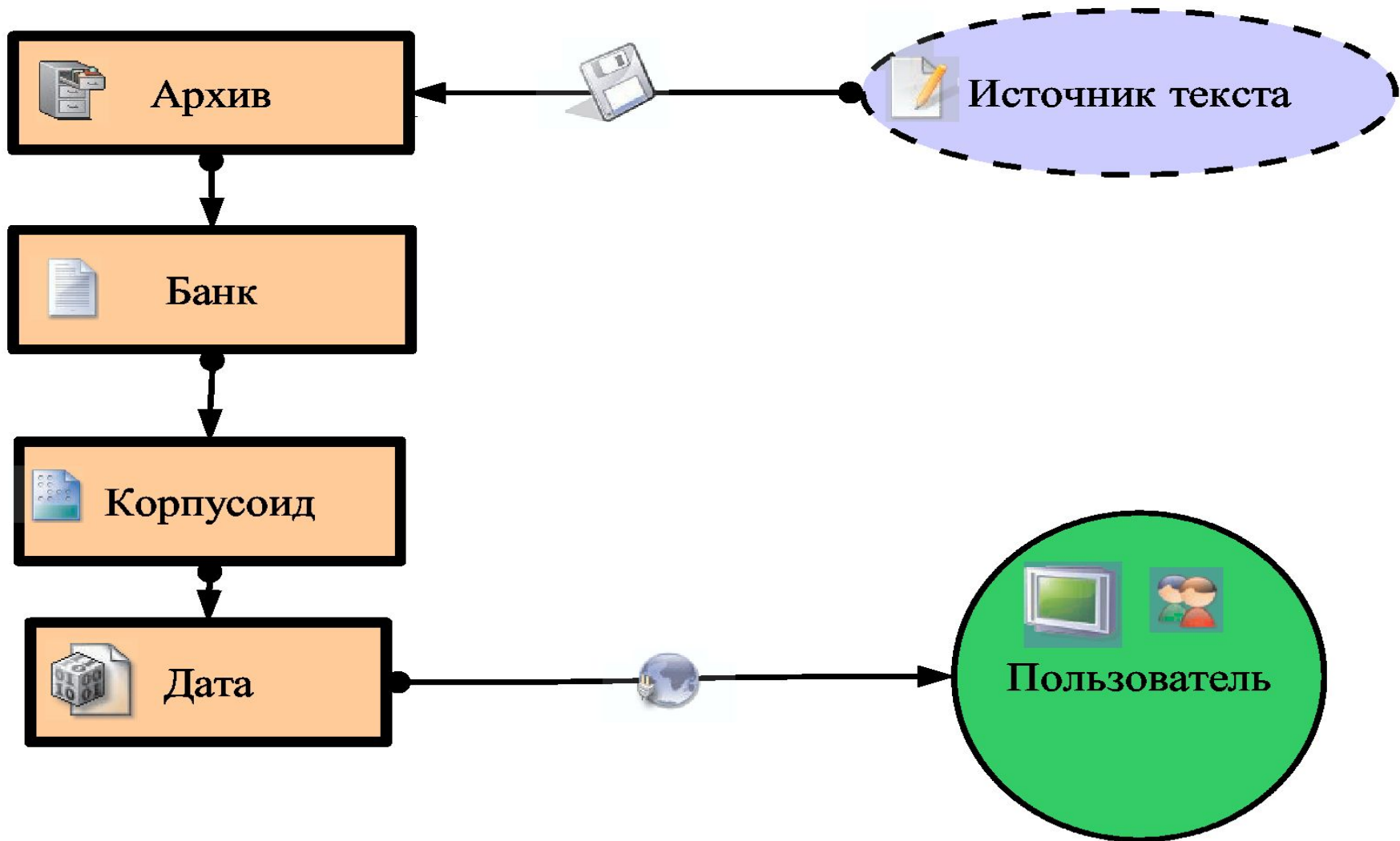
<lex>

<base>mečat'</base>

<ctag>VMesb+</ctag>

</lex>

</tok>



Поиск в корпусе

Система Manatee (<http://www.textforge.cz>)

Поиск слова, фразы

Регулярные выражения

Поиск по лемме, тэгу

Частотный анализ

Колокации

MI-score, T-score

Элементы структурной разметки

Библиографическая информация & стиль, жанр текста

Лингвистическая разметка

лемма

морфология

синтаксис?

морфосинтаксический тэг

1.автоматически

- существующие инструменты
 - Hajič & Hric (Прага)
 - Ažka (Брно)
- ~95% аккуратность

2.вручную

- небольшое количество текстов
- цель: несколько сто тысяч слов
- «ядро» корпуса

Система морфосинтаксических ТЭГОВ

строго морфологический принцип (неужели?)

один знак \Leftrightarrow одна грамматическая категория

знаки по определённой очереди

знаки (почти) не повторяются

„čas“ SSis1 – *S-существительное имя, S-существительный парадигм, i-неодушевлённый мужской род, s-единственное число, 1-именительный падеж*

1.		2.		3.		4.		5.		6.		7.	
S	Substantives	S A F U	Paradigm substantive adjective mixed incomplete	m i f n	Gender masc. animate masc. inanimate feminine neutrum	s p o	Number singular plural unknown	1 2 3 4 5 6 7 o	Case nominative genitive dative accusative vocative locative instrumental unspecified				
A	Adjectives	A F U	Paradigm adjective mixed incomplete	m i f n o	Gender masc. animate masc. inanimate feminine neutrum unspecified	s p o	Number singular plural unknown	1 2 3 4 5 6 7 o	Case nominative genitive dative accusative vocative locative instrumental unspecified	x y z	Grade positive/irrelevant comparative superlative		
P	Pronouns	S A P F U D	Paradigm substantive adjective pronoun mixed incomplete adverbial	m i f n o h	Gender masc. animate masc. inanimate feminine neutrum unspecified general	s p o	Number singular plural unknown	1 2 3 4 5 6 7 o	Case nominative genitive dative accusative vocative locative instrumental unspecified	g	Agglutinated agglutinated		
N	Numerals	S A N F U D X	Paradigm substantive adjective numeral mixed incomplete adverbial solitaire use	m i f n o	Gender masc. animate masc. inanimate feminine neutrum unspecified	s p o	Number singular plural unknown	1 2 3 4 5 6 7 o	Case nominative genitive dative accusative vocative locative instrumental unspecified				

1.		2.		3.		4.		5.		6.		7.	
V	Verbs	I K M H L B	Form infinitive indicative imperative transgressive <i>l</i> -participle futurum	d e j	Aspect perfective imperfective ambivalent	s p	Number singular plural	a b c	Person first second third	m i f n o h	Gender masc. animate masc. inanimate feminine neutrum unspecified general	+ -	Negation affirmative negative
G	Participles	k t	Type active passive	m i f n o	Gender masc. animate masc. inanimate feminine neutrum unspecified	s p o	Number singular plural unknown	1 2 3 4 5 6 7 o	Case nominative genitive dative accusative vocative locative instrumental unspecified	x y z	Grade positive/irrelevant comparative superlative		
D	Adverbs	x y z	Grade positive/irrelevant comparative superlative										
E	Prepositions	v u	Form vocalised non-vocalised	2 3 4 6 7 o	Binds with genitive dative accusative locative instrumental unspecified								
O	Conjunctions	Y	contains conditional morpheme <i>by</i>										
T	Particles	Y	contains conditional morpheme <i>by</i>										

J	Interjection	#	Not a word
R	Reflexive particle/pronoun <i>sa, si</i>	%	Citation element (e.g. foreign language word)
Y	Morpheme <i>by</i>	0	Digits
Z	Punctuation	:r	Proper noun
W	Abbreviation	:q	Incorrect spelling
Q	Unknown POS type		