

# Štruktúra dát v Slovenskom národnom korpuse a ich vonkajšia anotácia

**Radovan Garabík:** Slovenský národný korpus, oddelenie Jazykovedného ústavu Ľ. Štúra SAV

Príspevok odznel na konferencii Slovenčina na začiatku 21. storočia, Prešov 7.3.2003.

## Abstract

The paper "Data structure in Slovak National Corpus and its external annotation" presents proposed system for data storage and annotation in Slovak National Corpus. Taking into account the desirability to ensure readability of the data in future times, the system is targeted towards platform independent, easily decipherable annotation format. The data in corpus are organised into hierarchical structure divided into four different levels. Archive contains original documents without any other processing, Bank contains documents converted into common text format, Corpusoid adds linguistic annotation to the data, and finally the texts are converted into data suitable for corpus manager. Further we discuss the proposed common text format for documents in the Bank, as well as the structure of style-genre annotation of documents.

Výhľadovo si treba uvedomiť, že existuje pomerne veľká pravdepodobnosť, že v budúcnosti (perspektívne aj v značne vzdialenej) bude potrebné sa k dátam uloženým v korpuse vrátiť a ďalej s nimi pracovať. Preto sa snažíme zvoliť taký formát ukladania dát, ktorý je čo najmenej závislý na použitej počítačovej architektúre a dostupných programových prostriedkoch, s perspektívou udržania čitateľnosti podľa možnosti v priebehu nasledujúcich desaťročí. Preto bol pri návrhu štruktúry korpusu kladený obzvláštny dôraz na formálny návrh ukladania údajov.

Korpus je rozdelený na viac úrovní, v každej úrovni sa uchovávajú dáta v inom stupni spracovania.

Hierarchia sa skladá z nasledujúcich stupňov:

- archív
- banka
- corpusoid
- dáta pre korpusový manažér

## Archív

Archív zdrojových textov je tvorený súbormi, ktoré sú priamo získané z externých zdrojov. Každý takýto text, získaný z jedného zdroja v jednom časovom okamihu, tvorí jednu základnú nedeliteľnú jednotku archívu, nazývanú pre naše účely „dokument“. Fyzicky môže byť dokument tvorený viacerými súbormi, niekedy aj zložitou adresárovou štruktúrou. Tieto súbory sú v zdrojovom tvare bez akejkoľvek konverzie alebo úpravy a sú takto aj uchovávané. Napríklad dokumenty získané skenovaním sa ukladajú ako obrázky, ak to dovoľia technické možnosti.

Archivované texty sú nedotknuteľné, pôvodné tlačové a pravopisné chyby sa neopravujú.

## Banka

Banka obsahuje texty v spoločnom formáte inšpirovanom jazykom XML. V textoch v banke sú už odstránené všetky nadbytočné informácie, napríklad typografické údaje neobsahujúce informáciu o logickej štruktúre textu, vzorce,

tabuľky, obrázky a podobne. V texte sú vyznačené časti písané iným fontom a typom písma, ak sú použité na zvýraznenie alebo iné upozornenie na odlišnú funkciu v texte. Banka neobsahuje žiadne dodatočné (lingvistické) informácie, tie budú pridané až v ďalšej úrovni. Nad bankou je naplánované implementovať fulltextové vyhľadávanie, z licenčných dôvodov zrejme len pre vnútornú potrebu Slovenského národného korpusu.

Základnú nedeliteľnú jednotku banky nazývame „text“, pričom existuje jednoznačná väzba medzi textom v banke a dokumentom v archíve. V prípade, že dokument archívu sa skladá z viacerých bibliograficky identifikovateľných jednotiek, snažíme sa tento dokument rozdeliť na viacej textov, každý z nich jednoznačne identifikovateľný. Avšak aj v prípade že dokument je tvorený jedinou ucelenou bibliograficky určenou množinou dát, môže byť rozdelený na viacero textov, ak sú napríklad od rôznych autorov alebo rôzneho žánru, ako typický príklad uveďme knihu a doslov ku knihe.

Texty v banke sa občas môžu meniť, najmä v prípade opravovania chýb vzniknutých konverziou z dokumentov archívu.

## **Corpusoid**

Corpusoid kopíruje svojou štruktúrou banku. Texty v corpusoide sú vnútorne anotované, väčšinou automaticky použitím externých programových nástrojov. Ručná vnútorná anotácia je osobitne označená. Predpokladá sa, že tieto texty budú často regenerované zo zdrojových textov banky.

## **Dáta pre korpusový manažér**

Dáta, vhodné na spracovanie korpusovým manažérom automaticky generujeme z corpusoidu. Pri tejto konverzii môže prichádzať k strate informácií, pretože korpusový manažér nemusí podporovať všetky typy informácií nachádzajúcich sa v corpusoide. Aké informácie budú prístupné pre používateľov, záleží na použitom korpusovom manažéri.

V Slovenskom národnom korpuse predpokladáme v prvej etape nasadenie korpusového manažéra Manatee (<http://www.textforge.cz>).

## **Externá anotácia**

Anotácia je tvorená štruktúrou kľúč-príznak-hodnota. Hodnota je vždy zapísaná ako reťazec znakov ukončený koncom riadku, čo vylučuje viacriadkové názvy. Hodnota môže byť buď voľná (napr. meno autora), alebo vybraná z fixnej množiny povolených hodnôt (napr. žáner).

Tieto hodnoty majú špeciálny význam:

- prázdne miesto: neurčená hodnota, túto hodnotu sme neurčovali, prípadne ju môžeme určiť neskôr, ale apriori o nej nič nevieme povedať.
- *XXX*: neznáma hodnota. Nevieme určiť (napr. meno autora v novinovom článku).
- *YYY*: neurčiteľná hodnota. Z princípu sa nedá určiť alebo nemá význam (napr. pohlavie autora, ak ide o kolektívne dielo, pohlavie prekladateľa, ak nejde

- o preložené dielo).
- *MIX*: v prípade hodnoty vyberanej z fixnej množiny je hodnota zmes iných.
  - *MSC*: v prípade hodnoty vyberanej z fixnej množiny je hodnota, ktorá nie je v danej množine určená.
  - *TTT*: neznáma hodnota ktorú treba určiť – Značka, že k danej anotácii sa nevyhnutne treba neskôr vrátiť a hodnotu doplniť.

Hodnoty môžu byť bližšie označené týmito príznakmi:

- *fuzzy*: označuje nejasnosť hodnoty. Spolu s hodnotou označuje, že je určená len približne (a na rozdiel od TTT nie je potrebné alebo neplánujeme sa k záznamu vrátiť a hodnotu spresniť).
- *auto*: označuje, že daná hodnota bola automaticky vygenerovaná na rozdiel od hodnôt, ktoré zadal človek.

## Implementácia archívu

Každý dokument je umiestnený v samostatnom adresári. S ohľadom na obmedzenia filesystému a na lepšiu orientáciu sa používa štvorstupňová hierarchia adresárov podľa dátumu *Acqdate*. V adresári zodpovedajúcom danému dokumentu sa nachádza zdrojový dokument. Okrem toho sa tam nachádza súbor *.annotation*, obsahujúci externú anotáciu. Zdrojový dokument môže byť skomprimovaný štandardným komprimačným programom *gzip(1)*, v takom prípade bude mať k pôvodnému menu pripojenú príponu *.gz*. Iné komprimačné programy (napr *bzip2*) nepoužívame. V prípade, že je dokument tvorený zložitou adresárovou štruktúrou, môže byť skomprimovaný vo formáte *.tar.gz*.

## Implementácia banky

Pri implementácii externej bibliografickej a štýlovo-žánrovej anotácie sme vychádzali zo štruktúry anotácie v Českom národnom korpuse, s niektorými zmenami reflektujúcimi skúsenosti kolegov z ČNK a špecifické podmienky SNK.

Všetky texty v banke a corpusoide sú bez výnimky v kódovaní UTF-8, normalizácia C. Konce riadkov používajú unixovskú konvenciu, t. j. ASCII znak 0x0A, '\n'. Snažíme sa zbytočne nezneužívať kontrolné znaky v unicode, pretože tieto nebývajú dostatočne dobre zobraziteľné v bežných textových editoroch. Namiesto toho zavedieme XML tagy. Text v banke by mal odrážať logickú štruktúru dokumentu, nie typografické formátovanie. Avšak v mnohých prípadoch je logická štruktúra tesne zviazaná s typografickou a sformátovať súbor tak, aby odrážal logickú štruktúru, často vyžaduje neprímerane veľa manuálnej práce. Preto text v banke zachytáva podľa možnosti pôvodnú typografickú štruktúru dokumentu.

V banke uchováваме informáciu o dvoch typoch koncov riadku, mäkkom konci bez významnej logickej hodnoty, a o tvrdom konci s významnou hodnotou (napr. konce veršov v básňach).

Z ďalších typografických údajov zachytávame konce odsekov (značka **<p>**), nejasností v texte (značka **<noise>**, nadpisy (**<h1>**Nadpis**</h1>**), zvýraznený text (**<em>**Zvýraznený text**</em>**).

Vzorcie (matematické, chemické a iné) nahradíme značkou **<equation>**, jednoduché vzorce, zápisy fyzikálnych veličín a chemické zlúčeniny a reakcie (ak nesú podstatnú lingvistickú informáciu, napríklad H<sub>2</sub>O je dostatočne vžitá vo

všeobecnom povedomí slovensky hovoriacej komunity a je chápané ako slovo) zapíšeme použitím príslušných unicodových symbolov.

Tabuľky nahradíme značkou <table>>), obrázky značkou <picture>, malé obrázky v texte značkou <icon>.

Rozlišujeme nasledujúce typy pomlčiek:

- U+002D HYPHEN-MINUS (-) používame v banke na zápis nejednoznačnej pomlčky, keď nevieme pri konverzii textu automaticky rozhodnúť o aký typ ide – napríklad ak bol v zdrojovom texte použitý iba jeden typ znaku aj pre spojovník, aj pre pomlčku.
- U+2010 HYPHEN (-) používame v banke na zápis spojovníka. Príklad: „Rakúsko-Uhorsko“.
- U+2014 EM DASH (–) používame na zápis pomlčky. Príklad: „Peniaze – radosť“.
- U+2212 MINUS SIGN (−) používame na zápis unárneho alebo binárneho matematického operátora mínus, ak je v zdrojovom texte tento operátor odlišný. Ak nie, používame U+002D HYPHEN-MINUS (-).
- U+00AD SOFT HYPHEN je dosť nejednoznačne definovaný. Interne ho používame na zápis hyphenation point na mieste, kde prišlo naozaj k rozdeleniu slova (inými slovami, výlučne iba na konci riadku). Tento znak nebude mať dôvod sa objaviť v banke.
- U+2011 NON-BREAKING HYPHEN považujeme za ekvivalent U+2010 HYPHEN a v banke ho nepoužívame.

## Príloha

### Anotácia archívu

Anotácia archívu obsahuje tieto kľúče:

- *Name*: meno dokumentu.
- *Acqdate*: dátum získania dokumentu pre korpus.
- *Id*: identifikačný reťazec dokumentu. Dokument nemusí byť fyzicky umiestnený na serveri, môže sa nachádzať na externom médiu s označením identifikačného reťazca.
- *Source*: zdroj, odkiaľ bol získaný dokument, zvyčajne názov vydavateľstva alebo redakcie.
- *License*: licencia dokumentu – rozlišujeme tri úrovne licencie, len pre vnútorné použitie pracovníkmi JÚLŠ SAV, licenciu umožňujúcu zaradenie dokumentu na CD nosič, a nakoniec licenciu umožňujúcu sprístupnenie prehľadávania dokumentu na Internete.
- *Comment*: ľubovoľný komentár určený pre človeka (nie pre počítačové spracovanie).

### Anotácia banky

Anotácia banky obsahuje tieto kľúče:

- *Name*: názov textu.
- *Origname*: pôvodný názov textu, ak ide o preklad. Zapíname *fuzzy* príznak, ak je názov transliterovaný.
- *Author*: meno autora. Tak, ako je uvedené v zdrojoch k dokumentu, podľa normy pre bibliografické zápisy.
- *Origauthor*: originálne meno autora v pôvodnom jazyku. Z technických dôvodov zrejme zostane väčšinou prázdne, vyplňať sa bude ak bolo meno autora uvedené v preklade chybné.
- *Translator*: meno prekladateľa. *YYY*, ak nejde o preložené dielo.
- *Translation*: určuje či bol text preložený.

<i>trn</i>	preklad
<i>org</i>	originálny text
<i>ftt</i>	voľne preložený, prerozprávaný text

- *ISBN*: isbn číslo (ak existuje).
- *ISSN*: issn číslo (ak existuje).
- *SourceId*: Id dokumentu archívu, z ktorého dokument banky pochádza.
- *Rhyme*: rýmovanie dokumentu.

<i>nrh</i>	nerýmovaný
<i>rhy</i>	rýmovaný
<i>MIX</i>	kombinovaný

- *Type*: typ textu, presnejšiu špecifikáciu typu textu určuje kľúč *Subtype*.

	Type		Subtype
<i>po</i>	<b>poézia</b>		
		<i>ver</i>	báseň
		<i>son</i>	pieseň
<i>dra</i>	<b>dráma</b>		
		<i>scr</i>	scenár, prepis rozhlasového vysielania
<i>pro</i>	<b>próza</b>		
		<i>nov</i>	román
		<i>col</i>	súbor poviedok, jednotlivá poviedka
		<i>fac</i>	literatúra faktu
<i>inf</i>	<b>informatívny text</b>		
		<i>pub</i>	publicistika
		<i>adv</i>	reklamné materiály, reklama, inzeráty
<i>prf</i>	<b>odborný text</b>		
		<i>sci</i>	vedecká literatúra, články, odborné časopisy, učebnice pre VŠ
		<i>pop</i>	populárnovedecká literatúra, záujmové časopisy
		<i>txb</i>	učebnice pre stredné a základné školy
		<i>enc</i>	encyklopédie a podobné abecedne usporiadané diela
		<i>adm</i>	administratíva
		<i>man</i>	manuály, návody na obsluhu
<i>liv</i>	<b>živá komunikácia</b>		
		<i>spk</i>	hovorený
		<i>wri</i>	písomný

- *Genre*: žáner/odborná oblasť, *Subgenre* určuje presnejšiu špecifikáciu

	<b>Genre</b>		<b>Subgenre</b>
<i>img</i>	<b>imaginatívne</b>		
		<i>bel</i>	beletria
		<i>crm</i>	detektívky, špionážne romány
		<i>scf</i>	sci-fi, fantasy
		<i>jun</i>	pre deti a mládež, rozprávky
<i>int</i>	<b>interstage</b>		
		<i>trv</i>	cestopisy
		<i>mem</i>	(auto)biografie, memoáre
		<i>chr</i>	kroniky
		<i>let</i>	listy
		<i>ess</i>	eseje
<i>pol</i>	<b>politika</b>		
<i>ars</i>	<b>umelecké vedy</b>		
		<i>mus</i>	hudba, opera, opereta, balet
		<i>cin</i>	film
		<i>arc</i>	architektúra
		<i>art</i>	výtvarné umenie, fotografia, sochárstvo
		<i>the</i>	divadlo, divadelná veda a kritika
		<i>lit</i>	literatúra, literárna veda a kritika
<i>hum</i>	<b>humanitné vedy</b>		
		<i>his</i>	história, archeológia
		<i>psy</i>	psychológia
		<i>edu</i>	pedagogika
		<i>soc</i>	sociológia, komunikácia, médiá
		<i>phi</i>	filozofia
		<i>inf</i>	knihovníctvo a informačné zdroje
		<i>pol</i>	politológia
		<i>lin</i>	lingvistika
		<i>eth</i>	etnológia, etnografia
<i>law</i>	<b>právo</b>		
		<i>jur</i>	právo
<i>nat</i>	<b>prírodné vedy</b>		
		<i>agr</i>	poľnohospodárstvo
		<i>med</i>	medicína
		<i>zoo</i>	zoológia
		<i>bot</i>	botanika
		<i>bio</i>	biológia
		<i>che</i>	chémia
		<i>mat</i>	matematika
		<i>ggr</i>	geografia
		<i>phy</i>	fyzika
		<i>met</i>	meteorológia
		<i>geo</i>	geológia
		<i>env</i>	environmentalistika, ekológia
<i>tec</i>	<b>technika</b>		
		<i>tra</i>	doprava, spoje, telekomunikácie
		<i>ene</i>	energetika
		<i>ind</i>	priemysel
		<i>com</i>	informatika, počítače

	Genre		Subgenre
		<i>bui</i>	stavebníctvo
		<i>sta</i>	normalizácia, štandardizácia
<i>ecn</i>	<b>ekonómia, riadenie</b>		
		<i>eco</i>	ekonómia, bankovníctvo, obchod
		<i>mng</i>	manažment, riadenie
		<i>mer</i>	tovaroznalectvo, spotrebiteľská oblasť
<i>blf</i>	<b>náboženstvo, viera</b>		
		<i>rel</i>	náboženstvo, viera, sekty
		<i>teo</i>	teológia
		<i>exc</i>	nadprirodzeno, okultné vedy, mágia
<i>lif</i>	<b>life style</b>		
		<i>hou</i>	domácnosť (byt, záhrada, ručné práce)
		<i>fsh</i>	odievanie, móda
		<i>spo</i>	šport
		<i>sct</i>	spoločenský život
		<i>amu</i>	zábava, hra, koníčky, voľný čas
		<i>min</i>	etnické minority
		<i>reg</i>	región
<i>ins</i>	<b>interdisciplinárne vedy</b>		

- *Medium*: médium

<i>lib</i>	kniha
<i>nws</i>	noviny
<i>jou</i>	časopis
<i>ste</i>	skriptá
<i>net</i>	internet a iné (pre-internetové) siete
<i>for</i>	formulár
<i>occ</i>	príležitostné
<i>npu</i>	nepublikované texty, rukopisy, školské práce
<i>tvf</i>	televízia
<i>rad</i>	rozhlas

- *Authsex*: pohlavie autora

<i>msc</i>	muž
<i>fem</i>	žena

- *Lang*: jazyk diela, trojpísmenová skratka podľa ISO 639-2 (slovenčina *svk*). Prípadné inojazyčné súvislé časti textu budú označené v corpusoide. Kompletne inojazyčný text nebude mať dôvod vyskytnúť sa v SNK, kým sa nezačne budovať paralelný korpus.
- *Varieta*: variant jazyka, v ktorom je dokument. Väčšinou bude spisovná slovenčina.

<i>std</i>	spisovná slovenčina
<i>nst</i>	neštandardná slovenčina

- *Transsex*: pohlavie prekladateľa.

<i>msc</i>	muž
<i>fem</i>	žena

- *Origlang*: pôvodný jazyk diela.
- *Date*: dátum vydania.
- *Dateorig*: dátum prvého vydania (vzniku diela).
- *Comment*: komentár.



## Príklad záznamu jedného dokumentu v archíve

```
Comment: Pravda
Source: Perex, a. s.
Name: Pravda 1991-04
Licence(auto): 1
Acqdate(auto): 2003-02-19
Id(auto): 2003-02-019-A
```

## Prezeranie banky:

FormEdit 0.01  
Editácia existujúceho súboru:  
/tmp/annotation-bank

Author	Klára Buzássyová	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Authsex	fem	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Comment	Kultúra slova 2000/1	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Date	2000-01-01	<input checked="" type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Dateorig	2000-01-01	<input checked="" type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Genre	Dátum prvého vydania	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
ISBN		<input type="checkbox"/> fuzzy	<input checked="" type="checkbox"/> auto
Lang	svk	<input type="checkbox"/> fuzzy	<input checked="" type="checkbox"/> auto
Medium	jou	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Name	Na osemdesiatku profesora Jána Horváta	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Origauthor		<input type="checkbox"/> fuzzy	<input checked="" type="checkbox"/> auto
Origlang		<input type="checkbox"/> fuzzy	<input checked="" type="checkbox"/> auto
Origname		<input type="checkbox"/> fuzzy	<input checked="" type="checkbox"/> auto
Rhyme	nrh	<input type="checkbox"/> fuzzy	<input checked="" type="checkbox"/> auto
Sourceid	2000-01-01-A	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Subgenre		<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Subtype	pub	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Translation	org	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Translator	YYY	<input type="checkbox"/> fuzzy	<input checked="" type="checkbox"/> auto
Transsex	YYY	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Type	inf	<input type="checkbox"/> fuzzy	<input type="checkbox"/> auto
Varieta	std	<input type="checkbox"/> fuzzy	<input checked="" type="checkbox"/> auto

Abort

Save

## Príklad záznamu jedného dokumentu v banke

```
Name: Na osemdesiatku profesora Jána Horeckého
SourceId: 2000-01-01-A
Date(fuzzy): 2000-01-01
Dateorig(fuzzy): 2000-01-01
ISBN(auto):
Author: Klára Buzássyová
Authsex: fem
Origauthor(auto):
Origlang(auto):
Origname(auto):
Translator(auto): YYY
Transsex: YYY
Translation: org
Lang(auto): svk
Type: inf
Subtype: pub
Genre: int
Subgenre:
Medium: jou
Varieta(auto): std
Rhyme(auto): nrh
```

## **Bibliografia**

**Leech, Geoffrey** (2000): *Anotační systémy pro značkování korpusu*. In: *Acta Universitatis Carolinae – Philologica 3–4 – Studie z korpusové lingvistiky*. Praha: Univerzita Karlova v Praze, Nakladatelství Karolinum, s. 185–197.

**Kol.** (2001): *Slovenčina a čeština v počítačovom spracovaní*. Bratislava: VEDA Vydavateľstvo Slovenskej akadémie vied.