

# Hornjoserbski tekstowy korpus w nowej formje

Sonja Wölkowa

*Wot 15. měrca 2013 je so změnil přistup k hornjoserbskemu tekstowemu korpusej na stronach Serbskeho instituta. Při tutej skladnosći chcemy wo tutym projekće informować a někotre pokiwy za jeho wuživanje dać.*

Tekstowe korpusy su digitalne zběrki tekstow, kotrež su w džensnišim času zwjetša w interneče přistupne a kotrež služa přenjotnje jako zaklad za přeptytowanje rěče. W nich je zhromadženy awtentiski rěčny material, kotryž dowoli objektivniše posudžowanje rěčnych zjawow a je njeparujomny zaklad za moderne linguistiske slědženje a tohorunja za wudželjanje rěčnych přiručkow, słownikow abo wučnicow. Tola tež za literaturowědnikow, bibliotekarow abo stawiznarjow móže tajki korpus při wotpowědnym wuběrje tekstow wunošne žórlo być, a lektorojo abo redaktorojo móža so při rěčnych korekturach na njón zepěrać. Tekstowe korpusy eksistuju hižo za tójšto europskich rěčow, mjez nimi za nas wobdawacu němčinu<sup>1</sup> a za našej susodnej słowjanskej rěci pólščinu<sup>2</sup> a čěščinu<sup>3</sup>. A tež za delnjo- a hornjoserbščinu móžemy so na tajku informacisku resursu zepěrać. Wobaj korpusaj staj so wuwiłoj a hladatej a rozšerjatej so w zamołwitości Serbskeho instituta. Wo „DOlnoserbski Tekstowy KOrpus“ (skrótšeny jako DOTKO)<sup>4</sup> stara so wotrjad za delnjoserske slědženja w Choćebuzu, mjeztym zo je za HOrnjoserbski Tekstowy KOrpus (skrótka HOTKO)<sup>5</sup> załožity rěčespytny wotrjad Serbskeho instituta w Budýšinje w kooperacji ze Serbskej centralnej biblioteku. Tutón napisledk mjenowany je tema předležaceho přinoška.

## 1. Nastaće hornjoserbskeho tekstowego korpusa

Praktiske džělo na hornjoserbskim tekstowym korpusu je so započalo w lěće 1996, jeho zakłady je koncipowały a stworił tehdyši sobudželačer Serbskeho instituta Edward Wornar, nětko profesor na Lipščanskej uniwersiće. Po jeho wotchadže w lěće 2003 přewza zamołwitość za korpus wjednica rěčespytnego wotrjada Serbskeho instituta Sonja Wölkowa. Pjeć lět po zahajenju džěla na korpusu spřistupni so přenja wersija na internetowej stronje Serbskeho instituta. Tehdyši stav wuvića kompjutero-weho koděrowanja słowjanskich pismikow po wšelakorych zasadach bě při tym wosebite wužadanje. Zo by korpus zajimcam po cyłym swěće njewotwisnje wot wužiwaneho koděrowanskeho standarda přistupny był, wu-

wi jeho założer a wobdželer Edward Wornar za specifisce serbske pismiki z diakritiskimi znamješkami transkripciju, kotaž wuńdze ze 128 znamješkami tak mjenowanego ASCII-koda.<sup>6</sup> Pismiki z diakritiku rozpuščicu so na dwě znamješce – zakladny pismik z předchadžacym symbolom za hóčku ( $\check{c}$  =  $^{\wedge}c$ ), smužku ( $\acute{c}$  =  $/c$ ) resp. nakósnu smužku pola  $\dot{t}/\ddot{z}$  ( $\dot{t}$  =  $_l$ ). Tak wupadachu drje serbske teksty chětro njezwučene, za to pak na wšěch kompjuterach jenak.<sup>7</sup>

Kompjuterowa technika pak je so dale wuwiwała, w koděrowaniu je so mjezynarodne přesadžíł standard Unicode, kotryž tež wosebite pismiki słowjanskich rěčow wobsahuje. Tohodla dotalny přistup k hornjoserbskemu tekstowemu korpusej nowym móžnosćam poněčim wjace njewotpowědowaše. K tomu příndže fakt, zo je wotrjad za delnjoserske slědženja kónc lěta 2010 spřistupnił delnjoserski tekstowy korpus w spodobnej a za wužiwarja přijomnej formje na stronje [www.dolno-serbski.de](http://www.dolno-serbski.de)<sup>8</sup> a zdobom na zakladże kooperacije z Institutem za Česki narodny korpus (ÚČNK) na tamnej stronje [www.korpus.cz](http://www.korpus.cz) z přidatnymi móžnosćemi přeptytowanja a rešeršowanja. Zo móžemy nětko tohorunja hornjoserbski tekstowy korpus na tutej stronje a ze samsnymi móžnosćemi wužiwać<sup>9</sup>, za to mamy so kooperaciskim počaham a sobudželu wjacorych partnerow džakować. Kolegaj z Choćebuskeho wotrjada Serbskeho instituta Fabian Kaulfürst a Marcin Szczepański sposřdkowaštaj kontakt a zhromadne džělo z Praskim ÚČNK a pomhaštaj formu podawanja tekstow a trěbnych informacijow w nich zjednotnić. Michal Křen z Praskeho instituta postara so mjez druhim wo segmentowanje běžnych tekstow na sady a wo jich přihotowanje za wšelake přeptytowske móžnosće, na př. za analyzowanje słownych skupin (koločkow).

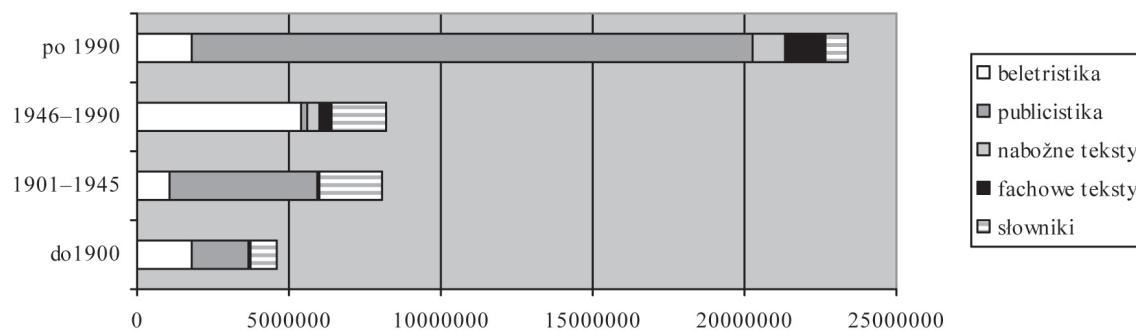
## 2. Tuchwilny staw Hornjoserbskeho rěčnego korpusa

W běhu lět je Hornjoserbski tekstowy korpus dosć nahladne rozrostł: Tuchwilu je w nim 384 datajow wšelakeho razu a wšelakeje wulkosće, dohromady je to wo koło 44 milionow tak mjenowanych „tokenow“ – z tym su měnjene jednotliwe separatne słowne formy.<sup>10</sup> Jako tokeny liča tež w tekscie wustupowace ličby a nimo toho kóžde interpunkciske znamješko, tak zo je w korpusu trochu mjenje słowow hač tokenow (ca. 36 milionow). Teksty hornjoserbskeho korpusa słušeja do wšelakich kategorijow (hlej wobr. 1): Najwjetši je podzél publicistiky (57%) a beletristiky (23%), mjenje je nabožnych (4%) a wědomostnych (4%) tekstow, k tomu příndu

někotre słowniki a rjad terminologijow za jednotliwe šulske předmjety (12%).<sup>11</sup> Najstarše teksty sahaja hač do přeňeje połojcy 19. lětstotka, su to spisy Handrija Zejlerja a basnje Rudolfa Mjenja w nowowudaćach z 20. lětstotka, najmłodše su Serbske Nowiny a Katolski Posoł z lěta 2010. Podzél jednotliwych časowych dobow je wšelaki: Najwjace mamy modernych tekstow z časa po přewróće 1989/90 (54%). Z 19. lětstotka pochadža 10 %, z přeňeje połojcy 20. lětstotka hač do 1945 je jich 18% a z lět mjez 1945 a 1990 19% – tak reprezentuje hornjoserbski tekstowy korpus z nimale třomi štvrćimi nami tekstow z časa po Druhej swětowej wójne předewšem rěč přitomnosće.<sup>12</sup>

Džakowano zrěčenjam z Ludowym nakładnistwom Domowina a Rěčnym centrumom WITAJ smy móhli wobšernu zběrku aktualnych tekstow za rešeršowanje spřistupnić, při čimž ma so wězo wobkedžbowanie autorskich prawow zaručić.<sup>13</sup> Najwjetsi džel tekstow pak je so zaskenował a potom z programami za spóznawaće teksta (t. mj. OCR-programy) wobdželał. Njedosahaceho personala dla njejsu so začitane teksty bohužel w dosahcej měrje skorigować móhli, tak zo njemóžemy wuzamknyc, zo jewja so zmylki. Za serbskeho wužiwarja mjenje problematiski drje je fakt, zo njejsu citowane pasaže w druhich rěčach, mjez druhim w delnjoserbščinje, wosebje markěrowane – wšak je za njeho lochko spóznać, hač jedna so wo hornjoserbščinu abo hinašu rěč.

wobraz 1



wobraz 2

Mjeřísi džel tekstow, předewšem z 19. lětstotka, předleži w historiskim prawopisu (na př. pjeć lětnikow časopisa „Serbski hospodar“), pisanje kh za džensniše ch na spočatku morfemow a mjechke í pak namakatej so tež hišće w tekstach, nastatych do lěta 1945. Na internetnej stronje Serbskeho instituta smy wo tym zaměstnili nadrobnu informaci (hlej wobr. 2, dypk: prawopisne warianty), kotaž móže wužiwarjej pomhać wšitke příklady za swoje naprašowanje namakać, hačrunjež jewja so prawopisne warianty.

### 3. Wužiwanje hornjoserbskeho tekstowego korpusa

Hornjoserbski tekstowy korpus namakaće na stronje Serbskeho instituta pod menijowym dypkom *online-publikacije* (<http://www.serbski-institut.de/cms/os/48/hornjoserbski>, hlej wobr. 2). Za wužiwarjow smy přihotowali nadrobnu dokumentaciju a pomocne informacie, kotrež namakaće, hdźiž nakliknjeće jednotliwe dypki pod krótkim zawodnym tekstem.

Dokelž zmónjna so přeptytowanje hornjoserbskeho tekstowego korpusa přez kooperaciju z Institutom za Češki narodny korpus ([www.korpus.cz](http://www.korpus.cz)), dyrbja sej wužiwarjo na jeho stronach wosobinske wužiwerske konto załožić. Za tutu proceduru smy na horjeka mjenowanej stronje spřihotowali wotpowědne wujasnjenja (wobr. 2, dypk: přistup ke korpusej).

Štóż je sej na te wašnje přistup zarjadował, móže z pomocu wotpråšowanskich programow<sup>14</sup> za wustupowanjom wěstych słowow w zapřijatych tekstuach pytać – při tym namakaja so awtomatisce wobšeńe lisčiny příkladow (tak mjenowane konkordancy) za słowa, jich formy abo kombinacije, hromadže z dokumentaciju, z kotrych žórlów pochadžeja. W tutejch lisčinach hodži so zaso přepytować, kak husto wěste słowa wustupuja a w kajkim susodstwje so wosebje jewja, móže so tež zwěścic, hač su za wěstych awtorow typiske, na wěste držiny tekstu abo snano na wěstu časowu dobu wobmjezowane.

Hornjoserbski korpus njeje hišće lematizowany a gramatisce anoterowany – to rěka, zo njejsu wšelake słowne formy jednotnemu heslu přirjadowane kaž w słowniku (na př. *ruku*, *ruce*, *rukow*, *rukomaj* k heslu *ruka*) a njejsu po słownych družinach abo gramatiskich formach analyzowane. To wobmjezuje pytanske móžnosće wosebje za rěčespytné prašenja z wobluka syntaksy a morfolođije. Za namakanje po móžnosći wšich formow někajkeho pytaneho słowa pak móža sej wužiwarjo pomhać z tak mjenowanymi regularnymi wurazami (jendź. regular expressions). To su znamješka, kotrež funguju jako naměstniki za wšelake warianty, takrjec kaž jokery w

kartowej hrě. Tak steji kombinacija .\* za 0 abo wjace pismikow abo cyfrow. Hdyž zapodamy potajkim naprasowanje *měsac.\**, namaka pytanski program tež příklady za wotchilne formy tutoho słowa, na př. *měsaca*, *měsacej*, *měsacy*, *měsacow*, *měsacami*. Naprasowanje .\*spěwam přinjese příklady za *spěwam*, *N/njespěwam*, *Z/zaspěwam*, *W/wuspěwam*, *D/dospěwam* atd. Nadrobne serbske wujasnjenja k tutej znamješkam namakaće na stronje Serbskeho instituta (wobr. 2, dypk: regularne wurazy).

Nimo ryzy wědomostnych zaměrow hodži so tekstowy korpus tohorunja za serbskorěčnu praksu zvužitkować. Awtorojo wučbnych materialijow namakaja z jeho pomocu zlochka mnóstwo awtentiskich příkladow za słowa abo słowne skupiny. Tute hodža so na př. wužiwać za zestajanje lisčinow sadow z džerami k wupjelnjenju za leksikaliske zvučowanja. Tež za rěčnymi wobrotami hodži so derje pytać, hdyž kombinuje so naprasowanje we hłownym woknješku z kontekstowym filtrom, kaž sčehowacy příklad pokazuje (wobr. 3).

Hdyž zapodamy na př. we hłownym woknješku sekwencu *[Kk]joza*<sup>15</sup> a w kontekstowym woknješku (w menuju wubjerje so *context*) zapisamy .\*liz[nl].\* a wuzamnjemy pod menijowym dypkom *subcorpora* jako žórló

wobraz 3

## wobraz 4

The screenshot shows a computer screen displaying a web-based concordance tool. The title bar says 'Concordance - Opera'. The main content area is titled 'NoSketch Engine' with a logo featuring a stylized 'E' made of Chinese characters. Below this, it displays 'User: hotko Corpus: hotko Description: Hornolužický textový korpus, verze 1 z 6. 3. 2013 Size: 44,367,372 positions? Hits: 72'. On the left, a sidebar has a tree view with nodes like 'Concordance', 'Word List', and 'Help'. The main pane lists 'Hits: 72 ( 1.62 i.p.m.; related to the whole corpus) | ARF: 4'. It shows a list of words from the corpus followed by their context in a sentence. For example, 'Rad15' is followed by 'takemu wosobnemu a žadnemu hošcej , popřeja . Ale koza mje lizny ! Na blidě stejachu w mlóce jahfy a za štyrače dnjow za chribjetom . Ale to bět mje koza lizna . Tak jednoho dnja w kofejowni nowinki studuju a wubrác , žony kapitalistow a tak . Z tym jich koza lizny . Za to su jim tych čerwjenych dali , » Ja so hišće netkłe wjesel , zo je SA-kow koza lizda « , rjekny wón zaskle . Wona delnju hubu chcyše z luteje nadutosc partout našich wózhríwcow spěwać wučić . Koza je če liznya ! A ordny twojego přichodného nana dí a so nadžješ , zo žadn liscík wjace nedoštanje . Koza jeho lizny , wón dosta město w njewobsadzenym poslednim rynku zo bytžk z lochkošu sama dočni . » Tebje je koza liznya . Sy wopak zastal « , rjekny . » wušparanje nós , » nadžjomne , ménju , če zaso koza njelizuje . « Jeho hudlenje bě přezerzawé , ani hižo z deňčkom na blido . » Ně , hospoza , koza je waz liznya ! Namrét sym ! = wotmolwic . po poważu zaso dele puščit . A Lulu je hakle koza liznya . Je so při wobjedze ta běrnjacej políkwi natykať budu ja w pincy rejwać . Teho je tež prawje koza liznya . 35 Tej dawnwo póčki šćerkotaja . Tej kóžda budu ja w pincy rejwać . Teho je tež prawje koza liznya . 35 Tej dawnwo póčki šćerkotaja . Tej kóžda » ( w Rab . ) Prajenje : Nas je koza lizla my smy so prawje zjebali . Prajenje : Tón mětk , tež přewodzjerjej džéč . A mje bětē potajkim koza lizna ? Běte to snaž kwitwownaka tufoho mřodženca , přewodzjerja .

słowniki, namaka program 72 przykładowych sadow z frazeologizmom *někoho koza lizne* we wšelakich tempusowych formach, kaž pokazuje sčěhowacy wurězk (wobr. 4).

Wězo njesmědža so při wužiwanju ženje zabyć wěste wobmjezowanja za wuslědkи rešeršow. Prěnje tajke wobmjezowanje rezultuje z toho, zo nimamy korpus wšěch hornjoserbskich tekstow. Hdyž potajkim někajku formu njenamakamy, rěka to jenož, zo njeje wona w korpusu. Dopokaz, zo wona njeeksistuje, to hišće njeje.

Druhe wobmjezowanje leži w tym, zo móža teksty korpusa rozšerjene rěčne zmylki dokumentować – hranica mjez zmylkom a nowej, so šerjacej formu tež za rěčespytnika druhdy cyle wótra njeje. A třeće wobmjezowanje zaleži na wužiwarju samym – kwalita wuslědkow korpusowego pytanja wotwisuje wot kwality prašenjow, kotrež so jemu stajeja. A tu móže so kóždy najlepje přez swójske nazhonjenja wukmanić. Nadžijamy so, zo namaka hornjoserbski tekstowy korpus tež w serbskorěčnej praksy swojich zajimcow a wužiwarjow.

- 1 <http://www.ids-mannheim.de/kl/projekte/korpora/> kaž tež <http://wortschatz.uni-leipzig.de/>
- 2 <http://nkjp.pl/>
- 3 [www.korpus.cz](http://www.korpus.cz)
- 4 Přístup přez <http://www.dolnoserbski.de/korpus/> abo přez <http://www.korpus.cz/dotko.php>.
- 5 Přístup pod <http://www.serbski-institut.de/cms/os/48/hornjoserbski> abo direktnje přez <http://www.korpus.cz/hotko.php>
- 6 Skróšenka za: American Standard Code for Information Interchange.
- 7 K znazornjenju a wujasnenju poda so w interneće příklad za transkripciju: Rjana \_Lu^zica, sprawna, p^re/celna, mojich serbskich w/otcow kraj, mojich zb/o^znych sonow raj. Swjate su mi twoje hona.
- 8 Hlej Nowy Casnik 28.12.2010, str. 2.
- 9 Zjednorjeny přístup kaž pod [www.dolnoserbski.de](http://www.dolnoserbski.de) za delnjoserbski korpus za hornjoserbski tuchwilu hišće móžny njeje, za to dyrbja so hišće pjenjezy a programowar namakać.
- 10 Zestajane časowe formy kaž na př. perfekt *sym spěwał* liča jako dwě slowje.
- 11 Procentualne podžele su so wuličili na zakladže ličby tokenow (hlej horjeka).
- 12 Při dalšim rozšerjenju korpusa budže so wosebje na wurunanje disproporcijow mjez wšelakimi tekstowymi typami a časowymi dobami džiwać dyrbjeć.
- 13 Tohodla je maksimalna wulkosć pokazowanego konteksta pytaneho słowa wobmjezowanego na 100 tokenow.
- 14 Za wužiwanje wšěch pytanskich a přepytowanskich opcijow je trěbna znajomość jendželščiny.
- 15 Róžkate spinki signalizują naprašowanskemu programej, zo maja so w nich stejace znamješka jako warianty rozumić.