

DOI: 10.11649/sfps.2014.008

Sonja Wölkowa
(Serbski institut, Budyšin)

Tekstowy korpus a dalše informaciske srědki wo hornjoserbskej rěči w interneće

W dobje globalizacije a přiběraceje omniprezency interneta steji wopisanje hornjoserbšćiny před nowymi wužadanjemi – Franc Šěn je situaciju 2012 na konferency w Gdańsku jara markantnje charakterizował: „štož njeje w interneće, njeeksistuje” (Šěn, 2013, s. 131). A tak nastawa a rosće potřeba za digitalnje a přez internet přistupnymi informaciskimi źórłami wo hornjoserbšćinje kaž za praksu tak tež za wědomostne slědženje. Na aktualnosć tuteje problematiki pokazuje w léće 2013 w nadawku Załožby za serbski lud załožena džěłowa skupina „Serbšćina w nowych mediach”. Wona je jako tuchwilu najwažniši nadawk identifikowała a wobzamknyła wudžěłanje noweho němsko-hornjoserbskeho słownika w interneće, a tutomu nadawkej mamy so jako rěčespytnicy stajić. Jako palaca začuwa so tuta potřeba wosebje tež hladajo na to, zo je so za delnjoserbšćinu w poslednim lětdžesatku w delnjoserbskim wotrjedže Serbskeho instituta zrealizował a w syći spřistupnił cyły system informaciskich srědkow, ke kotremuž słuša nimo němsko-delnjoserbskeho słownika tež digitalny přistup k najwažnišim delnjoserbsko-němskim słownikam Zwahra, Muki, Šwjele a Starosty kaž tež delnjoserbski tekstowy korpus (Dolnoserbski tekstowy korpus, n.d.; www.dolnoserbski.de).

W swojim přednošku chcu tu předstajić hornjoserbski tekstowy korpus HoTKo w konteksće dalšich digitalnych projektow za hornjoserbšćinu,

This is an Open Access article distributed under the terms of the Creative Commons Attribution 3.0 PL License (creativecommons.org/licenses/by/3.0/pl/), which permits redistribution, commercial and non-commercial, provided that the article is properly cited. © The Author(s) 2014.

Publisher: Institute of Slavic Studies, PAS & The Slavic Foundation
[Wydawca: Instytut Sławistyki PAN & Fundacja Sławistyczna]

realizowanych, přihotowanych a planowanych w Serbskim instituče. Wobšěrniši přehlad wo prezency serbskich wobsahow w interneće powšitkownje je podať w hižo mjenowanym přednošku w Gdańsku Franc Šěn, při tym wšak njebě městna na nadrobníše předstajenje jednotliwych poskitkow (Šěn, 2013). Přenje plany stworjenja tajkich internetowych pomocnych srědkow za hornjoserbsčinu sahaja do přenjeje połojcy 90tych lět 20. lětstotka. Tehdy naćisny so w Serbskim instituče projekt jednorěčneho hornjoserbskeho słownika, kotryž mjeješe primarnje w digitalnej formje nastać (Šoćina & Wornar, 1996). Bohužel njeje so tutón projekt personalnych pričín dla realizować móhł. Tola w přihotowanskej fazy w zwisku z tutymi słownikarskimi planami započea so w léće 1996 džělo na digitalnym hornjoserbskim tekstowym korpusu, kotryž eksistuje hač do džzensnišeho a so dale wuwija (přir. Wölkowa, 2013). Jeho zakłady (runje tak kaž mjenowany słownikarski projekt) je koncipowať a stworiť Edward Wornar, tehdy sobudžělačer Serbskeho instituta, nětko profesor Instituta za sorabistiku na Lipsčanskej uniwersiće. Po jeho wotchadže na profesuru w Lipsku w léće 2003 je zamołwitosć za korpus přešla na awtorku předležaceho přinoška w kooperaciji z nawodu Serbskeje centralneje biblioteki Francom Šěnom.

Upper Sorbian Text Corpus -- Hornjoserbski tekstowy korpus

Tutón interface služi k wuhodnoćenju tekstoweho korpusa přez WWW. Móžeće zapodać, kotre teksty maja so přeptać, hač ma so za jednornymi słowami abo za pravidlownymi wurazami pytać, kelko linkow konteksta so podawa, abo hač ma so statistika formow napisać.

This interface is for exploiting the text corpus via WWW. You can specify which texts you want to search, whether to search for strings or regexps, how many lines of context you want or whether you want statistics of forms.

Questions/Prašnja? sonja@serbski-institut.de

Wo koderowanju hlej [na tutej stronje](#).

Information on the encoding of the texts you'll find [here](#).

Pytanje/Search

Pytaj/Search:

Typ pytaneho wuraza/Type of expression:

- jedhory/fixed string pravidlowny wuraz/regular expression

Pravidlowny wuraz rěka tu *rozšěrjeny* pravidlowny wuraz kaž `grep -E`. Informacija wo `grep` je [tu](#) (jenož jendźelsce)

Regular expression means here *extended* regular expression as understood by `grep -E`. Information can be found [here](#).

Statistika -- Příkladny -- Kontekst / Statistics -- Examples -- Context

Informacija/Information:

- příklady/examples statistika/statistics linkow konteksta/Lines of Context

Teksty wubrać/Text Selection

Wobr. 1. Přistupowy formular za hornjoserbski tekstowy korpus 2001

Pjeć lět po zahajenju džěla na korpusu, potajkim w léce 2001, spřístupni so přerňa wersija zjawnje na internetowej stronje Serbskeho instituta.

Přístupowy formular bě dwurěčny serbsko-jendźelski a njebě tehdy ani estetisce wosebje naročny ani jara komfortabelny za wužiwanje. Tehdyši staw wuwĩa kompjuteroweho koděrowanja słowjanskich pismikow po wšelakorych zasadach a kodowych stronach bě při tym wosebite wužadanje. Zo by so korpus wot zajimcow po cyłym swěće njewotwisnje wot wužiwaneho koděrowanskeho standarda wužiwač hodźał, wuwi jeho załožer a wobdźělar Edward Warnar za specifisce serbske pismiki z diakritiskimi znamješkami transkripciju, kotraž wuńdže ze 128 znamješkami koda ASCII. Pismiki z diakritiku rozpušćichu so na dvě znamješce – zakładny pismik z předchadźacym symbolom za hóčku (č = ^c), smužku (ć = /c) resp. nakósnu smužku pola l/ḷ (ł = _l). Tak wupadachu drje serbske teksty chětro njezwučene, za to pak na wšěch kompjuterach jenak. K znazornjenju a wujasnjenju poda so w interneće přikład za transkripciju.

Example for Text Encoding -- Příklad koděrowanja

<p>Rjana _Lu^zica, sprawna, p^re/celna, mojich serbskich w/otcow kraj, mojich zb/o^zných sonow raj. Šwjate su mi twoje hona.</p>	<p>Rjana Łužica, sprawna, přećelna, mojich serbskich wótcow kraj, mojich zbóžnych sonow raj. Šwjate su mi twoje hona.</p>
--	---

(_ for lslash/Lslash, ^ for caron, / for acute)

Kedźbu/Attention!

Hdyž jako "typ wuraza" wubjerjeće "prawidłowny wuraz", ma třěška ^ funkciju wosebiteho znamješka. Tutu funkciju hasnje backslash před třěšku (\^):
 česki = \^c\^eski.

If you select "regular expression" from "type of expression", caron ^ gets a special meaning. To turn off this meaning, type a backslash before caron (\^):
 česki = \^c\^eski.

Wobr. 2. Pokazka koděrowanja słowjanskich pismikow z diakritiskimi znamješkami

Kompjuterowa technika pak je so dale wuwijała, w koděrowanju je so mjezynarodnje přesadził standard Unicode, kotryž wobsahuje tež wosebite pismiki słowjanskich rěčow. Tohodla dotalny přistup k hornjoserbskemu tekstowemu korpusej nowym móžnosćam po lětach wjace njewotpowědowaše. K tomu přińdže fakt, zo je wotrjad za delnjoserbske slědženja Serbskeho instituta delnjoserbski

digitalny korpus kónc lěta 2010 spřístupnił w spodobnej a za wužiwarja přijomnej formje na stronje www.dolnoserbski.de a zdobom na zakładze kooperacije z Institutom za Čěski narodny korpus (ÚČNK) na tamnišej stronje www.korpus.cz z přidatnymi móžnosćemi za přepytowanje a rešeršowanje. Zo móžemy nětko tohorunja hornjoserbski tekstowy korpus na tutej stronje a ze samsnymi móžnosćemi wužiwać, za to mamy so kooperaciskim počaham a sobudžěfu wjacorych partnerow džakować. Kolegaj z Choćebuskeho wotrjada Serbskeho instituta Fabian Kaulfürst a Marcin Szczepański posrědkowaštaj kontakt a zhromadne džěło z Praskim ÚČNK a pomhaštaj formu podawanja tekstow a trěbnych informacijow wo nich zjednotnić. Michal Křen z Praskeho instituta postara so mjez druhim wo segmentowanje běžnych tekstow na sady a wo jich přihotowanje za wšelake přepytowanske móžnosće, na př. za analyzowanje słownych skupin (kolokacijow). Zjednorjeni přistup za njelinguistow, kaž poskića so wón pod www.dolnoserbski.de za delnjoserbski korpus, tuchwilu za hornjoserbski hišće móžny njeje, za to dyrbjja so hišće financielne srědki a programowar namakać.

Žórła korpusa a wosebitosće tekstow

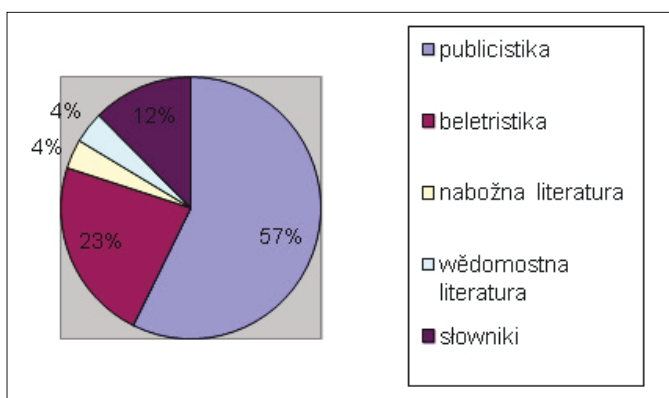
Džakowano zrěčenjam z Ludowym nakładnistwom Domowina w Budyšinje, w kotrymž wuchadža wjetšina publikacijow w hornjoserbskej rěči, a z Rěčnym centrumom WITAJ – to je wotrjad Domowiny, kotryž mjez druhim wudawa šulske wučbnicy a druhe wučbne srědki – smy móhli za rešeršowanje spřístupnić wobšěrnju zběrku aktualnych tekstow – nowinow, časopisow a knihow. Při tym ma so wězo wobkedźbowanje awtorskich prawow zaručić. Tohodla je maksimalna wulkosć pokazowaneho konteksta pytaneho słowa wobmjezowana na 100 tokenow. Najwjetši džěl tekstow za korpus pak je so z pomocu wědomostno-techniskeho personala Serbskeho instituta zaskenował a wobdžěłał z pismo spóznawacymi OCR-programami. Dokelž njejsu so wuslědki tuteje procedury njedosahaceho personala dla w dosahacej měrje korigowali, njemóža so bohužel w tekstach korpusa zmylki wuzamknýć, za korektne citowanje je trjeba sahać po original. Tu čaka nas w přichodže hišće wjele korekturneho džěła, tola tuchwilu ma hišće prioritu kwantitatiwny přirost korpusa.

Wěsty problem za lingwistiske korpusowe slědženje je fakt, zo njejsu citowane pasaže w druhich rěčach, mjez druhim tež w delnjoserbšćinje a němčinje, wosebje markěrowane. Za serbskeho wužiwarja abo dobreho znajerja serbšćiny pak je poprawom lochko spóznac, hač jedna so wo hornjoserbšćinu abo hinašu rěč.

Mjeńši džěl tekstow, předewšëm z 19. lětstotka, předleži w historiskim prawopisu (na př. pječ lětnikow časopisa „Serbski hospodar”) – džensnišemu š wotpowěduje tu na př. *sch* (*mysch* město *mys*), za *ć/č* pisa so *cž/cž* (*cžakacž* město *čakać*), za *s/z* pisa so *β/s* (*βasy* město *sazy*). Pisanje *kh* za džensniše *h* na spočatku morfemow a mjehke *ř* pak namakatej so tež hišće w tekstach nastatych w 20. lětstotku do lěta 1945. Na internetnej stronje Serbskeho instituta smy wo tym zaměstnili nadrobnu informaciju, kotraž móže wužiwarjej pomhać wšitke přikłady za swoje naprašowanje namakać, hačrunjež jewja so prawopisne warianty.

W běhu lět je Hornjoserbski tekstowy korpus dosć nahladnje rozrostł: Tuchwilu je w nim 384 datajow wšelakeho razu a wšelakeje wulkosće, dohromady je to wokoło 44 milionow „tokenow”, potajkim jednotliwych separatnych słownych formow¹. Dokelž liča jako tokeny tež w teksće wustupowace ličby a nimo toho kóždy interpunkciske znamješko, je w korpusu trochu mjenje slowow hač tokenow (ca. 36 milionow).

Teksty hornjoserbskeho korpusa slušeja do wšelakich kategorijow: Najwjetši je podžěl publicistiki (57 %) a beletristiki (23 %), mjenje je nabožnych (4 %) a wědomostnych (4 %) tekstow, k tomu přińdu někotre słowniki a rjad terminologijow za jednotliwe šulske předmjety (12 %)².



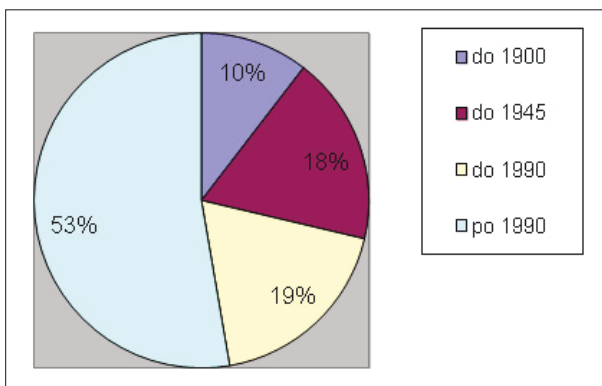
Wobr. 3. Podžěl tekstowych družin na korpusu

Najstarše teksty sahaja hač do přenjeje połojcy 19. lětstotka, su to spisy Handrija Zejlerja a basnje Rudolfa Mjenja w nowowudačach z 20. lětstotka a lětniki „Časopisa Maćicy Serbskeje” wot lěta 1848, najmłódše su „Serbske Nowiny”

¹ Zestajane časowe formy kaž na př. perfekt *sym spěwał* liča jako dvě slowje.

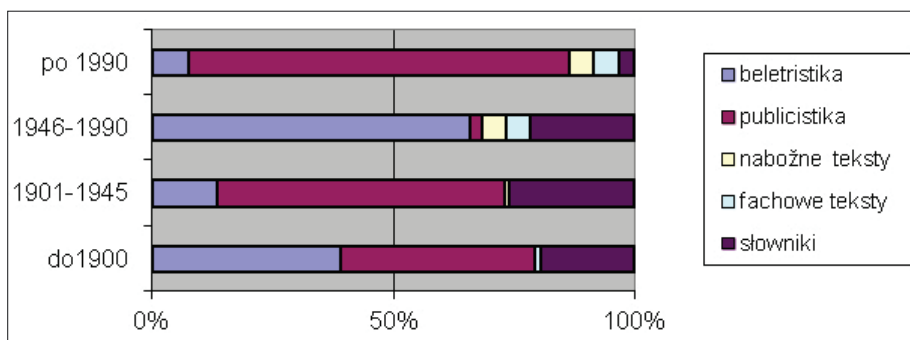
² Procentualne podžěle su so wuličili na zakładze ličby tokenow (hlej horjeka).

a časopis „Katolski Posoł” z lěta 2010. Podźěl jednotlivych časowych dobow je wšelaki: Najwjace mamy modernych tekstow z časa po přewróće 1989/90 (54 %). Z 19. lětstotka pochadza 10 %, z přenjeje 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 štwórcinami tekstow z časa po Druhej swětowej wójnje předewšěm rěč přitomnosće, tola wón dowola we wěstej měrje tež diachroniski přistup.



Wobr. 4. Podźěl tekstow ze wšelakich dobow

Hdyž wobhladamy sej rozdźelenje tekstowych družin na jednotlive časowe wotrězki, pytnjemy, zo njeje jich podźěl přeco jenaki. Při dalšim rozšěrjenju korpusa budže so wosebje na wurunanje disproporcijow mjez wšelakimi tekstowymi typami a časowymi doбами džiwać dyrbjeć, najtrěbniše je wudospołnjenje korpusa za druhu połojcy 20. lětstotka ze zběrku publicistiki.




Wobr. 5. Podźěl tekstowych družin po časowych dobach

Dokelž so hornjoserbski tekstowy korpus běžnje wudospołnja, dyrbi so rjec, zo njejedna so wo tak mjenowany referenčny korpus, to rěka wo statiski korpus,

z kotrehož móža so konstantne statistiske daty dobywać. Wuhódnoćenja wotbłyščuja jenož aktualny staw w tekstach korpusa.

Přístup k hornjoserbskemu tekstowemu korpusej

Hornjoserbski Tekstowy KOrpus, skrótka HOTKO, je zaměstnjeny na stronje Serbskeho instituta pod menijowym dypkom *online-publikacije*. Za serbskich wužiwarjow smy přihotowali nadrobnu dokumentaciju a pomocne informacije w hornjoserbskej rěči, kotrež namakaće, hdyž nakliknjeće jednotliwe dypki pod krótkim zawodnym tekstom.



The screenshot shows the website interface for the Hornjoserbski tekstowy korpus HOTKO. At the top, there is a blue header with the logo of the Serbski Institut (Sorbisches Institut) on the left and handwritten text in Serbo-Croatian on the right. Below the header, a navigation menu on the left lists various sections, with 'online-publikacije' highlighted. The main content area features the title 'Hornjoserbski tekstowy korpus HOTKO' and a short introductory text. On the right, there is a search bar and a section titled 'Serbska lětnja šula' with a photo of a group of people. Below that, another section is titled 'Nowa wustajeŕca z wobstatkow archiwa a'.

Wobr. 6. Přístup k HoTKo na stronje <http://www.serbski-institut.de/cms/os/48/hornjoserbski>

Dokelž zmóžnja so přepytowanje hornjoserbskeho tekstoweho korpusa přez kooperaciju z Institutom za Čěski narodny korpus na Filozofiskej fakulće Karloweje uniwersity w Praze (www.korpus.cz), dyrbyja sej wužiwarjo na jeho stronach wosobinske wužiwarke konto založić. Za tutu proceduru namakaja so na mjenowanej internetowej stronje Serbskeho instituta wotpowědne serbskorěčne wujasnjenja (menijowy dypk: přístup ke korpusej).

Štóž je sej na te wašnje přístup zarjadowaŕ, móže z pomocu wotprašowankeho programa Čěskeho narodneho korpusa NoSketchEngine za wustupowanjom wěstych słowow w zapřijatyh tekstach pytać – při tym namakaja so awtomatisce konkordancy za słowa, jich formy abo kombinacije. Koděrowane su tež metainformacije wo žórłach: Podaća wo awtorje, titlu, městnje a času wuchadženja móža so jako skrótsenka na lěwym boku abo tež dokładnišo

wosebitym woknješku pokazować. Tež šěrši kontekst hač do 100 slowow hodži so pokazać.

The screenshot shows the NoSketch Engine search results for the term "rěčespytnicy". The interface includes a search bar at the top right, a navigation menu on the left, and a main results area. The search results are displayed as a list of entries, each with a document ID and a snippet of text. A detailed view of a specific entry is shown below the list, displaying fields such as document ID, full title, title, author, publisher, place, year, year of origin, and original font.

Document ID	Snippet
Černý1	bě to , štož prjedy wo nas wědžachu ; lědma rěčespytnicy něšto wo našej stowjanskej prastarej rěči spomnichu . Ale što
Groß11	, kotraž je drje pomjenowana po černjach abo kaž družu rěčespytnicy mjenja po čornej pjeršci , na kotrejš sej holanscy burja
SSP1	! Přeni , kiž na nju dźiwać počachu , njeběchu rěčespytnicy , kaž bychmy wočakować móhli , ale přirodowědnicy . Jako
Lp95_1	liča k wobdźělnikam sorabizća-slawišća w najšěršim zmysle , historikarjo , rěčespytnicy , ludowědnicy , literarni a kulturni slědźerjo , studenca a
Lp95_1	aktualne za dzensniši rěčespyt . Metodogramatikarjo běchu nimale wuwučajće němscy rěčespytnicy . Najznačšie mjena mjez Indoeuropeistami resp . germanistami su :
R1990	, wobšerna a popularna biografija pobrachowaše , to njezačuwaču jenož rěčespytnicy , ale tež kulturni stawiznarjo a stawiznarjo sojca - bě
R1990	třo daši studenca , kiž su pozdźišo skutkowali jako sorabisticy rěčespytnicy ; Pawel Wirtl , Hans Holm Bielefeldt a Arnulf Schröder
R1991	Moskowskich fachowcow . Na přeslědženju serbskich dialektow su so wobdźěllili rěčespytnicy z Pólskeje , Čěskeje a Sowjetskeho zwjazka . Tutón bjazposrědni
R1992	žonow njeje to hač dotaj hišće pytnyl , ani naši rěčespytnicy . Hewak bychu sej wězo hižo dawno blowu lamali a
R1993	stron institucionalizowaneho rěčespyta . Poměr wšak je dwustronski . Tež rěčespytnicy Serbskeho Instituta trjebaja w swojim powołanskim dźěle radu a pomoc
R1993	Zoležow a pod . Što mjenja družu ? Redaktorjo a rěčespytnicy ? P . S . : Mjez Ukrajinanami su tež
R1994	wobste- jacy styki wožiwć a nowe zwiski nawjazac . Chorwatscy rěčespytnicy su wosebje na wobstajnej wsměnje literatury na polu filologije z
R1994	dyskeho Serbskeho Instituta měli někotři serbscy a pólscy wědomostnicy (rěčespytnicy a rěčni sociologjo) wuwić serb- skeje rěče připytaować
R1994	w serbšćinje po l . 1991 dale što . Serbscy rěčespytnicy so z tuzym wěsćim rozestajaja , za zjawnosć buchuwu hłownje
R1995	a w kažubskej rěči . Poslednju hódnotu pólscy a kažubscy rěčespytnicy jako nastawcu spisownu rěč . Přinohki na konferency bachu wěnowane
R1995	a HSRK su (haž na jednoho) wěšty serbscy rěčespytnicy , někotři lektorjo , redaktorjo , wěboarjo a daši multiplikatorjo
R1996	to tuchwilu w DSRK praktikuja . Wuznawaja da so sojca rěčespytnicy - fachowcy sami dosć we wěšch tych komplikowanych , mjez sobu
R1997	Stajimy drje před nowej tendencu w serbšćinje , kotruž nam rěčespytnicy zawěšće bōrzu wopodstatnja . Mozak Lěto 1996 w SLA z
R1997	zjawnosći a jeho rōlu při zwěšćenju zakadaw narodnjeje eksistencny su rěčespytnicy a stawiznarjo w dalekej měrje přeslědžili a pokazali . Mjenje
R1998	hu , notj danakim kněžstwom do wěsnelne rōle zwěšćarje . Přini rěčespytnicy z chroamtichu w XIX . stóleću taja žurkni w zanjesenyh wěšch

R	document.id	Černý1	
R	document.full_title	Adolf Černý, Serbske wobrazki / přeložił a zawod napisał Ota Wičaz. Wušće w Praze 1890, Wozjewjene w „Dom a swět“ 1923, znowa w SN 2(25.2.-5.5.1992)	
S	document.title	Serbske wobrazki,	
SI	document.author	Černý, Adolf	
SI	document.publisher		
SI	document.place		
SI	document.year	1992	
SI	document.year_orig	1923	
SI	document.orig_font	latin	

Wobr. 7. Konkordanca zdobyta w HoTKo z rozšěrjenym podaćom žórła

W tutych konkordancach płaća pytanse a statistiske móžnosće programa NoSketchEngine: Hodži so na př. přepytować, kak husto wěšte słowa abo słowne formy wustupuja a w kajkim susodstweje so wosebje jewja, móže so tež zwěšćić, hač su za wěštych awtorow typiske, na wěšte družiny tekstow abo snano na wěštu časowu dobu wobmjewowane. Program móžnja tež pytanje slowow w definowanych kontekstach a dowoli individuělny wuběr, w kotrych tekstach ma so pytać, z pomocu definowanja wosobinskih subkorpusow.

Hornjoserbski korpus njeje hišće lematizowany a gramatisce anotěrowany – to rěka, zo njejsu wšelake słowne formy jednotnemu hesłu přiřadowane kaž w słowniku (na př. *ruku*, *ruce*, *rukow*, *rukomaj* k hesłu *ruka*) a njejsu po słownych družinach abo gramatiskich formach analyzowane. To wobmjewuje pytanse móžnosće wosebje za rěčespytne prašenja z wobłuka syntaksy a morfologije. Za namakanje po móžnosći wšěch formow někajkeho pytaneho słowa móža sej wužiwarjo pomhać z tak mjenowanymi regularnymi wurazami (jendź. regular expressions), w kotrychž wužiwarja so znamješka, kotrež funguja jako naměstniki za wšelake warianty – za

to wobsteja wšelake w nadrobnosćach so wotchilace standardy. Za HoTKo płaćace zasady, kotrež bazuja wězo na zasadach Čěskeho narodneho korpusa, su w lisčinje na stronje Serbskeho instituta pod dypkom „regularne wurazy” wopisane.

Wězo njesmědža so při wužiwanju ženje zabyć wěste wobmjezowanja za wuslědki 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 zaleži na tym, zo móža teksty korpusa nimo prostych OCR-zmylkow tež rozšěrjene rěčne zmylki dokumentować – hranica mjez zmylkom a noweju, so šěrjaceju formu tež za rěčespytnika druhdy cyle wótra njeje.

Dalše informaciske srědki wo hornjoserbskej řeči w interneće

Přestajeny hornjoserbski tekstowy korpus dyrbi so widzieć w konteksće dalšich projektow Budyskeho rěčespytneho wotrjada Serbskeho instituta. Čežišćo tworja tu prašenja słowoskłada a leksikografije. Centralny zaměr předstaja projekt wobšěrneho noweho němsko-hornjoserbskeho słownika, kotryž tuchwilu w našim wotrjedže koncipujemy – z tym reagujemy na potreby serbskeje rěčneje praksy, kotraž žada sej wjac hač dwaj lětdžesatkaj po wudaću poslednjeho tajkeho kompendija (Jenč, Michałk & Šěrakowa, 1989–1991) nowy a metodologisce dalewuwity pomocny grat. Planowany słownik je primarnje za digitalne wužiwanje w interneće mysleny. Njerěka pak to, zo njebudže so z njeho sekundarnje tež čišćany słownik generěrować dać. Technologisce budže so wón zepěrać na nazhonjenja wotrjada za delnjoserbske slědženje, a tež koncepcionelnje matej wobaj słownikaj mjezsobu kompatibelnej być. Tole wšak je – hladajo na hornjoserbšćinu – tuchwilu hišće hudźba přichoda. Zrealizowane su porno tomu hižo wěste předdžěła.

Projekt *Němsko-hornjoserbskeho słownika noweje leksiki* wobdžěłany wot Helmuta Jenča, Anje Pohončoweje a Jany Šołćineje (Jenč, Pohončowa & Šołćina, 2006), wobsahowacy někak 12000 hesłow, drje njeje ani za internet koncipowany ani w nim přistupny, jako zběrka, registrowaca nowu leksiku po wudaću němsko-hornjoserbskeho słownika 1989–1991, pak ma so wón jako krok po puću k planowanemu internetowemu słownikeju widzieć.

Prěni z digitalnych, online přistupnych poskitkow SI za hornjoserbšćinu, kiž chcu tu mjenować, je *Hornjoserbski frazeologiski słownik* w interneće. Wón je nastal na zakładze hornjoserbsko-němsko-ruskeho frazeologiskeho słownika,

kotryž smój hromadže z Anatolijom Iwčenku w léće 2004 wudałoj (Ivčenko & Wölke, 2004). Internetna wersija, kotraž je wot 2005 přistupna (<http://www.serbski-institut.de/cms/de/50/Obersorbisches-phraseologisches-Woerterbuch>), wobmjězuje so na hornjoserbšćinu a němčinu, dowoli pak za to wobšěrnije pytanske móžnosće wuchadžejo z wobeju rěčow. Nimo toho so tuta digitalna wersija běžnje z přidatnymi frazeologizmami wudospołnja.

Specialnu leksikografisku zběrku předstaja datowa banka geografiskich mjenow w hornjoserbšćinje a němčinje (Geografiske mjena hornjoserbsce, n.d.; <http://www.serbski-institut.de/cms/os/396/Geografiske-mjena-hornjoserbsce>). Wona je nastala w zhromadnym džěle z Rěčnym centrumom WITAJ, hdžež nastawaja serbske wučbnicy, mjez nimi tež tajke za geografiju. Tuchwilu su w datowej bance předewšěm zezběrane tak mjenowane eksonymy (něhdže 2000), potajkim geografiske mjena za objekty zwonka Serbow, tola planujemy ju rozšěrić tež wo hornjoserbske geografiske mjena z Łužicy. Wotprašowanje datow je móžne wuchadžejo z němčiny, ale tež wuchadžejo z hornjoserbšćiny. Systematiske wobdžělanje tutoho wosebiteho džěla słowoskłada w cyłku je zakład za jeho po móžnosći dospołne a zdobom konsekwentne a we sebi konsistentne předstajenje w planowanym němsko-hornjoserbskim słowniku a słuži na druhim boku stabilizaciji normy spisowneje rěče.

Zwonka rěčespytneho wotrjada je so w nadawku Serbskeho instituta zdžěłała zběrka w interneće přistupnych hornjoserbskich dwurěčnych słownikow SERDIS – Serbske digitalne słowniki. SERDIS nasta na iniciatiwu a z koordinaciju přez Franca Šěna w zhromadnym džěle Serbskeho instituta w Heidelbergu skutkowacym neurobiologom Jurjom Brankačkom a ukrainskim programowarjom Vladimirom Kukušku. Online přistupne, programowane w Java®, namakaja so tam digitalne wersije hornjoserbsko-němskeju słownikow Křescana Bohuwěra Pfula (Pful, 1866) a Jurja Krala (Kral, 1927), hornjoserbsko-ruskeho słownika Konstantina Trofimowiča (Trofimowič, 1974) a dvě dalšej mjeńšej leksikografiskej zběrce. W SERDIS móže so principielnje w dwěmaj směromaj pytać – tola pytanje wuchadžejo z ciloweje rěče je skerje wobmjězowane dla přesnadneje hłubokosće a nadrobnosće analyzy słownikowych artiklow. Za dalewuwiće tutoho online-poskitka by derje było, so orientować na metodach, kiž su so nałożowali za wotpowědny delnjoserbski poskitk.

Hornjoserbski tekstowy korpus twori wažny materialowy zakład za wšě leksikografiske džěla rěčespytneho wotrjada Serbskeho instituta. Intensiwnje je so wón hižo wužiwał za wšitke tři runje mjenowane hižo zrealizowane pro-

jekty: słownik noweje leksiki (Jenč et al., 2006), frazeologiski słownik (Ivčenko & Wölke, 2004) a datowu banku geografiskich mjenow, nimo toho tež zwonka SI za jendźelsko-hornjoserbski słownik Měrcina Straucha a Edwarda Wornarja (Wornar & Strauch, 2007). Porno tradicionalnym metodam leksikografije před dobu kompjutera – jako džěłaše so ze statiskimi papjerjanyimi kartotekami – mamy z tym wulke lěpšiny. Hladajo na wobjim (tuchwilu 36 mio slowow porno 1 mio kartkow za dwuzwjazkowy němsko-hornjoserbski słownik Jenča et al. 1989–1991) je digitalny korpus wo wjele spuščomniše wuchadžišćo. Za digitalny korpus steja k dispoziciji wjele diferencowaniše móžnosće analyzy, na př. po frekwency a konteksće. Nimo toho njeje digitalny korpus lokalnje wjazany na městno kartoteki, ale je přistupny zjawnje po cyłym swěće, a to nic jenož za leksikografiske, ale tohorunja za druhe rěčespytne předewzaća, a z přiběracym mnóstwom zapřijatych tekstow tež přeco lěpje za rešerše literaturowědnego abo stawizniskeho razu.

Bibliography

- Dolnoserbski. (n.d.). Retrieved December 19, 2013, from www.dolnoserbski.de
- Geografiske mjena hornjoserbsce. (n.d.). Retrieved from <http://www.serbski-institut.de/cms/os/396/Geografiske-mjena-hornjoserbsce>
- Ivčenko, A., & Wölke, S. (2004). *Hornjoserbski frazeologiski słownik* [Oborsorbisches phraseologisches Wörterbuch; *Верхнелужицкий фразеологический словарь*]. Budyšin: Ludowe nakładnistwo Domowina. Retrieved from <http://www.serbski-institut.de/cms/de/50/Oborsorbisches-phraseologisches-Woerterbuch>
- Jenč, H., Michałk, F., & Šěrakowa, I. (1989–1991). *Němsko-hornjoserbski słownik* (Vol. 1–2). Budyšin: Ludowe nakładnistwo Domowina.
- Jenč, H., Pohončowa, A., & Šołćina, J. (2006). *Němsko-hornjoserbski słownik noweje leksiki*. Budyšin: Ludowe nakładnistwo Domowina.
- Kral, J. (1927). *Serbsko-němski Słownik hornjołužiskeje rěče* [Wendisch-deutsches Wörterbuch der oberlausitzer Sprache]. Budyšin: Maćica Serbska.
- Pful, K. B. (1866). *Łužiski serbski słownik: Lausitzisch Wendisches Wörterbuch*. Budyšin: Maćica Serbska.
- Šěn, F. (2013). Sorabistika a nowe medije. In M. Milewska-Stawiany, & S. Wölkowa (Eds.), *Leksikologiske přinoški II: IV. seminar serbskeje słowotwórby a leksiki* [IV Seminarium Słowotwórstwa i Słownictwa Łużyckiego] : *Uniwersytet Gdański. Serbski institut 31.5. – 1.6.2012* (pp. 131–141). Budyšin: Serbski institut.

- Šołćina, J., & Warnar, E. (1996). Čehodla trjebamy jednorěčny hornjoserbski słownik? In E. E. Siatkowska, & J. Molas (Eds.), *Sprawy lużyckie w ich słowiańskich kontekstach* (pp. 15–22). Warszawa: Instytut Filologii Słowiańskiej UW.
- Trofimowič, K. (1974). *Hornjoserbsko-ruski słownik [Верхнелужицко-русский словарь]*. Budyšin: Ludowe nakładnistwo Domowina.
- Wölkowa, S. (2013). Hornjoserbski tekstowy korpus w nowej formje. *Serbska Šula*, 66(2), 44–47.
- Warnar, E., & Strauch, M. (2007). *Jendźelsko-hornjoserbski šulski słownik [English-Upper Sorbian Learner's Dictionary]*. Budyšin: Ludowe nakładnistwo Domowina.

Bibliography (transliteration)

- Dolnoserbski. (n.d.). Retrieved December 19, 2013, from www.dolnoserbski.de
- Geografiske mjena hornjoserbsce. (n.d.). Retrieved from <http://www.serbski-institut.de/cms/os/396/Geografiske-mjena-hornjoserbsce>
- Ivčenko, A., & Wölke, S. (2004). *Hornjoserbski frazeologiski słownik [Oborsorbisches phraseologisches Wörterbuch; Verkhneluzhitskii frazeologičeskii slovar']*. Budyšin: Ludowe nakładnistwo Domowina. Retrieved from <http://www.serbski-institut.de/cms/de/50/Oborsorbisches-phraseologisches-Woerterbuch>.
- Jenč, H., Michałk, F., & Šerakowa, I. (1989–1991). *Němsko-hornjoserbski słownik* (Vol. 1–2). Budyšin: Ludowe nakładnistwo Domowina.
- Jenč, H., Pohončowa, A., & Šołćina, J. (2006). *Němsko-hornjoserbski słownik noweje leksiki*. Budyšin: Ludowe nakładnistwo Domowina.
- Kral, J. (1927). *Serbsko-němski Słownik hornjołužiškeje rěče [Wendisch-deutsches Wörterbuch der oberlausitzer Sprache]*. Budyšin: Maćica Serbska.
- Pful, K. B. (1866). *Lužiski serbski słownik: Lausitzisch Wendisches Wörterbuch*. Budyšin: Maćica Serbska.
- Šěn, F. (2013). Sorabistika a nowe medije. In M. Milewska-Stawiany & S. Wölkowa (Eds.), *Leksikologiske přinoški II: IV. seminar serbskeje słowotwórby a leksiki [IV Seminarium Słowotwórstwa i Słownictwa Łużyckiego]: Uniwersytet Gdański. Serbski institut 31.5. – 1.6.2012* (pp. 131–141). Budyšin: Serbski institut.
- Šołćina, J., & Warnar, E. (1996). Čehodla trjebamy jednorěčny hornjoserbski słownik? In E. E. Siatkowska & J. Molas (Eds.), *Sprawy lużyckie w ich słowiańskich kontekstach* (pp. 15–22). Warszawa: Instytut Filologii Słowiańskiej UW.
- Trofimowič, K. (1974). *Hornjoserbsko-ruski słownik [Verkhneluzhitsko-russkii slovar']*. Budyšin: Ludowe nakładnistwo Domowina.
- Wölkowa, S. (2013). Hornjoserbski tekstowy korpus w nowej formje. *Serbska Šula*, 66(2), 44–47.
- Warnar, E., & Strauch, M. (2007). *Jendźelsko-hornjoserbski šulski słownik [English-Upper Sorbian Learner's Dictionary]*. Budyšin: Ludowe nakładnistwo Domowina.

The Upper Sorbian text corpus and further sources of information with regard to Upper Sorbian in the Internet

Summary

In the present era of globalisation and the omnipresence of the Internet, Sorbian linguistics faces new challenges along the lines “What is not in the Internet, does not exist” The demand for digital sources of information with regard to Upper and Lower Sorbian and those accessible online as working tools and reference points for language practice and as a source for academic research increases. As a result of this ongoing development, the *Foundation for the Sorbian People* established a workgroup called “Sorbian in the new media” at the end of 2012, which has pointed out the creation of an online German-Upper Sorbian dictionary as the major task in this field of activities. The focus of this article, however, is the Upper-Sorbian text corpus HoTKo, which has been created by the *Sorbian Institute* and which has been made available in co-operation with the Institute of the Czech National Corpus at the Charles University in Prague. The article presents the history and development of the corpus, its extent and shape as well as its link to or incorporation into further planned digital projects of the *Sorbian Institute* with regard to the Upper Sorbian language.

Keywords: computational lexicography; corpus linguistics; digitalization; text corpus; Upper Sorbian language

Słowa kluczowe: digitalizacja; język górnołużycki; korpus językowy; leksykografia komputerowa; lingwistyka korpusowa