

O jednej skratke

Radovan Garabík

JÚLŠ SAV

813 64 Bratislava, Slovakia

korpus@korpus.juls.savba.sk, <http://korpus.juls.savba.sk>

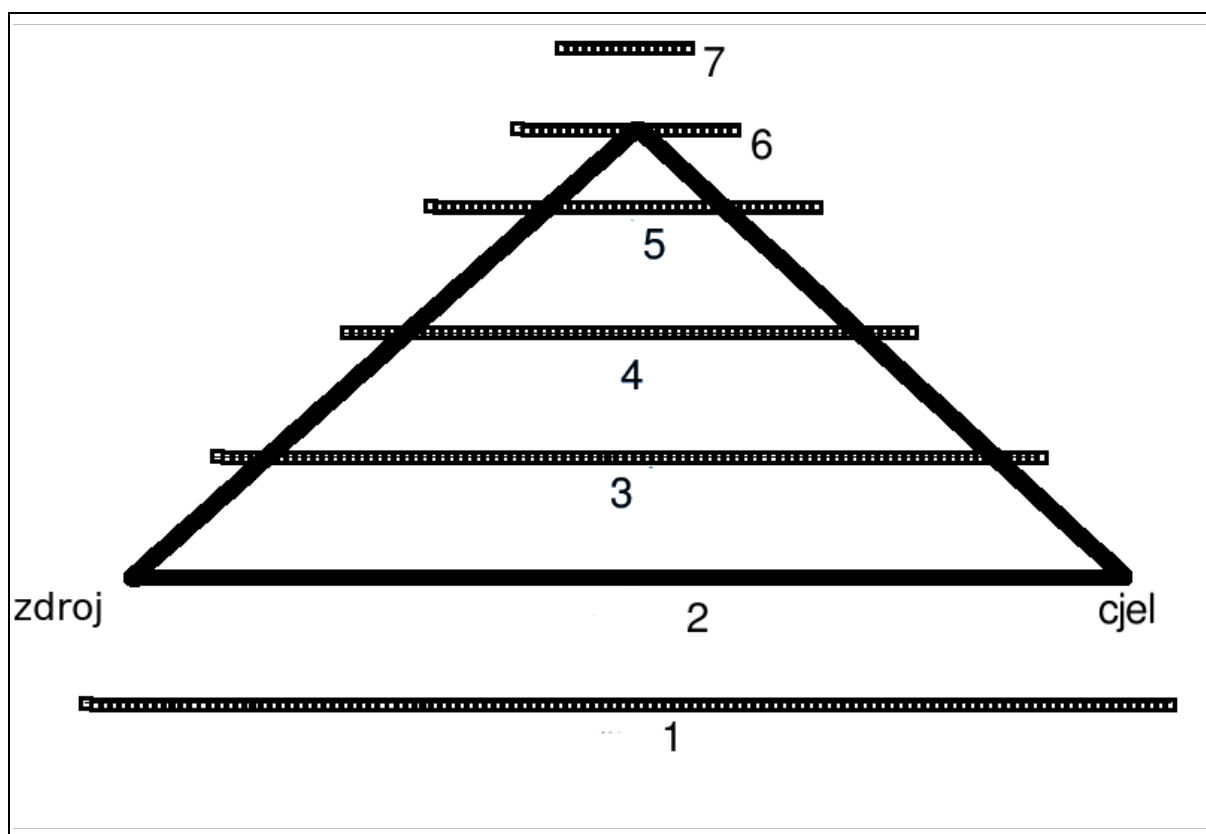
XVI. kolokvium mladých jazykovedcov, Častá-Papiernička 8. – 10. 11. 2006

Abstract. Machine translation systems tend to be rather complicated and the results are often disappointing. However, the difficulties involved in a successful translation diminish when dealing with a pair of very close languages, and the translation can be ameliorated by strategic use of common morphological, grammar and lexical features of the languages involved. Presented system can be used for translation on the orthographic and lexical level between very close languages and was successfully applied to translation from standard Slovak into the L. Štúr's Slovak language.

Úvod

Sistemi automatickeho prekladu patria k najkomplikovanejším aplikáciám v oblasti počítačového spracovania prirodzeného jazyka. Toto vypláva z potreby urobiť hĺbkovú analýzu zdrojového nárečia a transformovať zmysel puvodného textu do celového nárečia. Schematicky muožeme proces prekladu znázorniť diagramom podobným tomu na obr. 1. Plocha trojuholníka vijadruje oblasť, v ktorej pracujú tipické systémy automatickeho prekladu. Každá vodorovná čiara zodpovedá abstraktnej úrovni transferu medzi zdrojovým a celovým nárečím. Čím vyššia úroveň, tým abstraktnejší transfer sa uskutočňuje, a výsledok je tým bližší prirodzenému nárečju. Úroveň 1vá zodpovedá fonetike a na obrázku je uvedená len kvuoli úplnosti, pretože vo väčšine systémou automatickeho prekladu (ako aj v našom článku) ide o písaní text. Úroveň 2há zodpovedá ortografii, transfer na tejto úrovni znamená len zmenu ortografického systému (takíto transfer je použiteľní napríklad pri zmeňe pravopisu jedného nárečia, alebo preklad medzi nárečjami, ktoré sa líšja iba ortografiou). Úroveň 3tja zodpovedá morfológii a je použiteľná pre preklad medzi nárečjami, ktoré sa odlišujú maximálne morfológiou (s istými obmedzenjami muožže ísť o dve veľmi blízke príbuzné nárečia). Pri odlišnejších nárečjach dostaňeme na výstupe syntakticki a semanticki nežmyselní text. Úroveň 4tá zodpovedá sintaxi, na výstupe dostaňeme text syntakticki správni, aj keď možno s nežmyselním významom (alebo s významom nežzodpovedajúcim originálnemu textu). Modernje špičkovje systémy automatickeho prekladu sa k tejto úrovni iba približujú. Úroveň 5ta, semantika, zodpovedá pochopenju významu slov a slovních spojení originálneho textu a ich preklad na slovnje spojeňa s rovnakím významom – na tejto úrovni pracujú prekladateľá-ludja. Úroveň 6ta, na diagrame znázorňená vrcholom trojuholníka zodpovedá užitju interlingvi (medzireči), pri ktorom preklad prebehou už po stranách trojuholníka a

transfer sa zredukovau na identicku operáciu, pretože všetky črti puvodnjeho aj preloženjeho textu sú obsahnutje v medziprodukt'e. Do diagramu sme ešte doplnili sjedmu úroveň, ležjacu nad vrcholom trojuholníka. Táto úroveň bi sa dala opísať ako „pochopeňja toho, čo chceu autor povedať“ a jej znázorňeňja je vjacmeňej iba akademickje, pretože k dosjahnúťju tejto úrovne dochádza veľmi zriedka.



Obrázok 1 ví: Schematicki znázorňeni trojuholník prekladu

Preklad medzi veľmi blízkimi nářečjami

Blízkje (geneticki aj štrukturňe) nářečja majú vela podobních čft. Pri vzd'alovaní nářečí rozd'jeli medzi nimi celkom dobre sledujú úrovňe v uved'enom trojuholníku – najprú sa zjavja rozd'jeli v fonetike (aj v rámci jednjeho nářečja či dokonca rozličnorečja), potom v ortografii (pri kodifikácii alebo odšt'jepení nářečja, často s politickou motiváciou). Pri morfológickích rozd'jeloch sme už oprávn'eni hovoriť o ruoznich nářečjach. Sintax často zostáva kompatibilná aj pri nářečjach od seba značne vzd'jaleních, a v prípade dramatických rozd'jelou v lexike už ňemožeme hovoriť o blízkích nářečjach v našom poňímaní. Z automatických prekladovích

sistemou medzi blízkimi nárečjami muožeme spomenúť preklad medzi Češtinou a Slovenčinou[1] a preklad medzi Turečtinou a krimskou Tatárčinou[2].

Štúrovská Slovenčina

Spisovnuo Slovenskuo nárečje, tak ako ho definoval Ludevít Štúr v [3] sa od modernej Slovenčini [4] líši na prví pohľad prevažne ortografiou, pričom rozďjeli sú lahko algoritmicki popísaťelňje. Hlavňje ortografickje rozďjeli spočívajú v absencii grafemi „y“, v inej realizácii dvojhlasok a v explicitnom povinnom značení mekkosťi spoluhlások d, t, n.

Lexikálnje rozďjeli sú subtílňejšje, na prví pohľad badaťelňje len v ňjektorích najčast'ejších slovách, ale v skutočnosťi mjerne posúvajúce semantickí význam celích trjed slov.

Technická realizácia

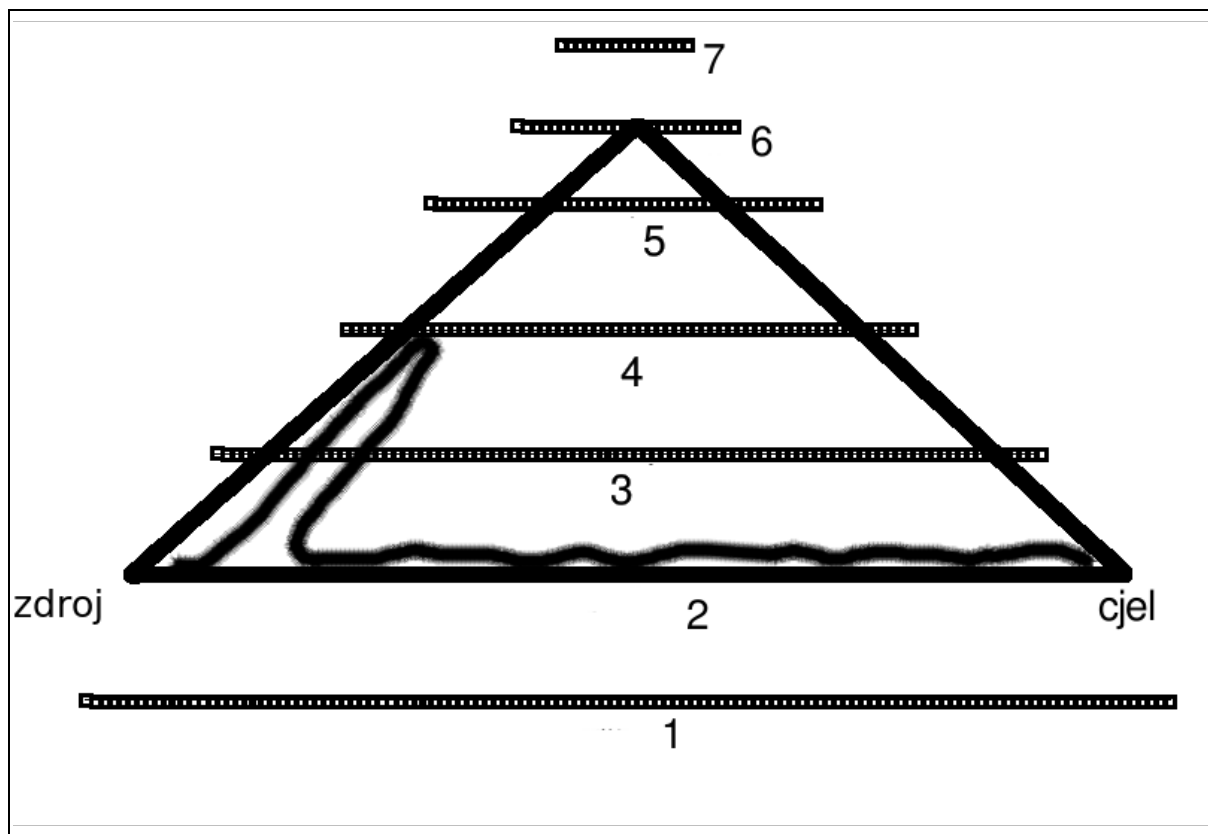
Pred prekladom je text najprú skonvertovaní zo vstupnjeho kódovaňja do Unicode, potom normalizovaní na NFKC normalizáciu Unicode a všetki ďalšje operácie prebjeajú duosledňje v Unicode. Text je tokenizovaní na základňje jednotki – tokeni (slová), ku každjemu tokenu je priradená informácia o prípadných bjelich znakoch (whitespace; Leerraum) pred slovom, abi sa po preklade mohlo zrekonštruovať vernuo rozložejja textu. Po preklade je veľkosť písmen preloženjeho slova upravená tak, abi kopírovala veľkosť písmen puovodnjeho slova – ak je prelozenuo slovo dlhšje ako puovodnuo, veľkosť „nadbitočnich“ písmen kopíruje veľkosť poslednej písmeni puovodnjeho slova. Toto zabezpečí verní preklad vlastních mjen a prípadných slov písaních kapitálkami. Ako víňimka sú koreňe slov „Sloven(čina, skí)“ a „Vlád(a)“ vždi v preklade písane so začjatočným veľkým písmenom, podľa úzu uživanjeho L. Štúrom.

Samotní preklad prebjea v dvoch fázach: najprú sa aplikuje lexikálna transformácia, pri ktorej sa nahrádzajú slová, ktorje sú v štúrovej Slovenčine inak reprezentovaňje. Víhodňje sa dá viužiť prevažná ekvivalencja morfem medzi súčasnuo a štúrovskou Slovenčinou a v prekladovej tabulke stačí povečšijne uvjest' iba prekladi koreňovích morfem, iba ňjekedi je potrebnuo uvjest' preklad celích tvarov slov.

Druhú fázju prebjea na ortografickej úrovni. Prekladi v oboch fázach sú realizovaňje jednoduchím nahrádzaním originálnich reťazcou prekladovimí. Začjatki a konce slov sú označejne špeciálnimi znakmi (^ začjatok, \$ koňjec), čo umožňuje efektívňje spracovať transformácie v príveskách slov a zabraňuje možným ňesprávnim nahradeňjam. Vzhľadom na

duoslednje značeŋja palatalizovanih spoluhlások v štúrovskej Slovenčiŋe je potrebno duokladne rozlišovat' tvrdje a mekkje „i“ podla vislovnosti (čo vede k prekladom politi→polity, diplo→dyplo, poezia→poesya) a takt'jež bolo potrebno zavjest' tvrdje „e“ na označeŋja ňepalatalizujúceho „e“ (toto písmeno sme arbitrárne označili znakom „ě“ – U+00EB LATIN SMALL LETTER E WITH DIAERESIS, príkladi prekladou: internet→intěrnět). V druhej úrovni budú tjeťo slová transformovaŋje na štúrovskí pravopis (polity→politi, dyplo→diplo, poesya→poesia, intěrnět→internet).

Naša skratka v prekladovom trojuholŋíku potom sleduje transfer na ortografickej úrovni, s krátkim vibočeŋím do oblasti semantiki (vlastŋe iba zámena ňjektorích lexikálnich jednot'jek).



Obrázok 2hí: Trojuholŋík prekladu so znázorŋeŋím našej skratki

lexikální preklad	ortografický preklad
u'grék' : u'rék',	u'ov\$' : u'ou\$',
u'gréc' : u'réc',	u'né\$' : u'ňje\$',
u'gréč' : u'réc',	u'é\$' : u'je\$',
u'maďar\$' : u'uher\$',	u'ého\$' : u'jeho\$',
u'maďar' : u'uhr',	u'ému\$' : u'jemu\$',
u'maďara' : u'uhra',	u'é' : u'e',
u'talian' : u'talyan',	u'ý' : u'í',
u'^ludovít' : u'^ludévít',	u'y' : u'i',
u'slávneho' : u'slávňého',	u'ô' : u'uo',
	u'l' : u'l',
# pieseň -> peseň	u'ä' : u'e',
u'^pies' : u'^pes',	u'ë' : u'e',
u'vidiet\$' : u'videt\$',	
u'vedietet\$' : u'vedet\$',	u'ia' : u'ja',
u'vedie' : u'vede',	u'dia' : u'dja',
u'erie\$' : u'ere\$',	u'diakon' : u'diakon',
u'erieš\$' : u'ereš\$',	u'tia' : u'tja',
u'^zmenši' : u'^umenši',	u'nia' : u'ňja',

Tabulka 1vá – časť prekladovej lexikálnej a ortografickej tabulky

Popis funkcií programu

Program (nazvaný ludevít) je napísaný v programovacom jazyku Python. Hlavnou zameranou programu je pre unixovú systém, ašak, sú napísané len s užitím štandardných pythonovských knižníc, funguje na veľmi širokej množine systémov a platform. Program funguje ako filter, čítajúc štandardný vstup a zapisujúc preložený text na štandardný výstup.

Výstup je možno modifikovať ďalšími argumentami k programu:

- o súbor alebo --output-file súbor – výstup zapíše do súboru miesto na štandardný výstup
- D alebo --nfkd – výstup buď v NFKD normalizácii
- d alebo --nfkd-hack – písmená d' a t' budú v NFKD normalizácii, ostatné v NFKC
- e ENCODING alebo --encoding ENCODING – miesto štandardného kodovania utf-8, predpokladá vstup a výstup v kodovaní ENCODING, ktorú môže byť hociká kodovanie podporované pythonom, ale pravdepodobne význam má len jedno z utf-8, iso8859_2, cp1250, cp852 alebo mac_latin2. Kodovanie iné než utf-8 nie je kompatibilné s voľbami -D a -d.

Kde sa zvuky mekko vyslovujú takto sa zmekčujúcou čarkou viznačujú, ale písmená „d“ a „t“ ju v dobe modernej inakšie označujú, značka táto skoro ako dlhá čarka má podobu. Aby sa historická vernosť zachovala, tieto dve písmená je možno normalizovať na unicodovské

„NFKD“ spôsob (parameter -d, prípadne normalizovať všetci písmeni parametrom -D), to značí že zmekčujúce čjarki sú ako samostatne kombinujúce písmeni (combining characters, kombinierende diakritische Zeichen) reprezentované, keď sa s predchádzajúcou písmenou vjažu, v renderovacích systémoch zriedka býva úplná podpora kombinujúcich písmen, a tak sa často písmena nad predchádzajúcou nežmeňená zobrazia, čo vizerá temer ako puovodní historicki správni spôsob písania. Žjal, mnohokrát sú tjetto čjarki zle zobrazenje, alebo naopak tak ako majú byť skombinované dobre a správne (na moderní spôsob) zobrazenje, a teda tento spôsob ňje vždi dobrje výsledki dáva.

Nedostatki

Ňedostatki uvedĕnjeho systému prekladu (či prepisu) muožeme rozdeliť na dve skupini. Prvú skupinu tvorja ňedostatki teoreticki ňepodstatňje, ktorje je aspon teoreticki možnuo lahko odstrániť aplikovaním dostatočnjeho množstva ľudskej práce. Sem patrí hlavňe malí rozsah prekladovjeho slovníka (na lexikálnej úrovni) a chibi v slovníkoch na lexikálnej aj ortografickej úrovni. Chibi je možnuo odstrániť duokladným skontrolovaním slovníkou, a malí rozsah slovnej zásobi samozrejme doplnením – okrem potrebi ňevihnutej ľudskej práce tu ňje sú žjadnje problemi, ktorje tomuto princípijálne bráňja.

Ňedostatki teoreticki podstatňje sú horšje, pretože viplívajú buď prjamo z návrhu prekladacjeho sistjemu, alebo z vrodĕných vlastností oboch verzií Slovenčini a vzťahu medzi ňimi. Tjeto ňedostatki ňje je možnuo jednoducho odstrániť. Najzávažnejšje z ňich sú:

- Absencja kontextovjeho prekladu. Sistem sa vždi pozerá iba na jedno konkrétne slovo. Toto zabraňuje možnosťi lexikálneho prekladu slovami s iním pohlavím, pretože ňje je možnuo súčasne preložiť prípadňje adjektíva a slovesá (v menosloví) tak, aby bola zachovaná zhoda pohlaví. V programe sme urobili jedĕnú výnimku, slovo „Bratislava“ prekladáme slovom „Prešporok“ (spolu s patričnými tvarmi skloňeňja), pretože idĕ o slovo dost’ známe a často užívaňje, a občas sa viskitnuvšju chibu v ňesúlade pohlavja prípadnjeho mena prídavnjeho sme považovali za menšje zlo ako poňechať tvar „Brat’islava“ (s ňejasním užívaním v štúrovskích dobách)
- Ňerozlíšiteľná homonímja v štandardnej Slovenčine. Najčast’ejším príkladom sú prídavnje mená ňijakjeho pohlavja v nominative a akusative v jednotnom počťe (príveska -uo) a prídavnje mená ženskjeho a ňijakjeho pohlavja v množnom počťe (príveska -je), ktorje v štandardnej Slovenčine majú rovnakú prívesku (-je, alebo -e ak bola predchádzajúca silaba dlhá) a ňje je ňijakí spôsob, ako bez duokladnej semantickej analisi určiť správni preklad (taketo slovňje spojeňja bývajú často

nerozhodnuteľnejšie aj skúsením čitateľom-človekom).

- Absencia úpravy syntaxe. Štúrovská Slovenčina sa od štandardnej odlišuje aj mjerne inou skladbou veti, náš sistem neobsahuje nijakje prostriedki aňi na syntaktickú analisu originálu aňi úpravu prekladu.

Zhrnut'ja

Uved'ení sistem umožňuje základní preklad zo štandardnej Slovenčini do štúrovskej na ortografickej a čjastočne lexikálnej úrovni. Preklad ňje je celkom dokonalí a od originálnej štúrovskej Slovenčini sa odlišuje hlavne v syntactickej skladbe vjet a v lexike, ale je dostatočne dobrí na občasňje užít'je, na demonstrácju štúrovskej Slovenčini a ako pomuocka pre prekladaťelou do štúrovskej Slovenčini. Po jednoduchej úprave (náhrada slovníka) je možnuo program užít' pre preklad medzi podobne odlišnými jazikovými sistemami (napríklad v prípad'e závažnejších zmjen v Slovenskom pravopise pri preklad'e do novej normi).

Program je dostupní pod licenciou GNU GPL v. 2.0, a jeho demoversiu prístupnú cez WWW rozhraňja prostredníctvom jednoduchjeho CGI skriptu je možnuo si prezreť na stránke Jazikovednjeho ústavu Ludevíta Štúra SAV[5]. O potrebe a užitočnosťi takjeho prekladu svedčí aj neočakávaná popularita, ktorej sa istú dobu uved'enuo WWW rozhraňja t'ešilo[6].

Literatura

1. Hajič, Jan – Hric, Ján – Kuboň, Vladislav: Machine translation of very close languages. In: Proceedings of the sixth conference on Applied natural language processing, Seattle, Washington, s. 7 – 12. Morgan Kaufmann Publishers Inc., San Francisco, 2000.
2. Altıntaş, Kemal: Turkish To Crimean Tatar Machine Translation System, MSc Thesis, Bilkent University Computer Engineering Department, July 2001
3. Štúr, Ludevít: Nauka reči Slovenskej. Prešporok, Tatrín, 1846.
4. Morfológia slovenského jazyka. Red. J. Ružička. Bratislava, Vydavateľstvo Slovenskej akademie vied 1966. 896 s.
5. <http://vvv.juls.savba.sk/ludevít/>
6. Štúrovská slovenčina. In: SME, 20. 12. 2006, s. 29