

# O jednej skratke

Radovan Garabík

JÚLŠ SAV

813 64 Bratislava, Slovakia

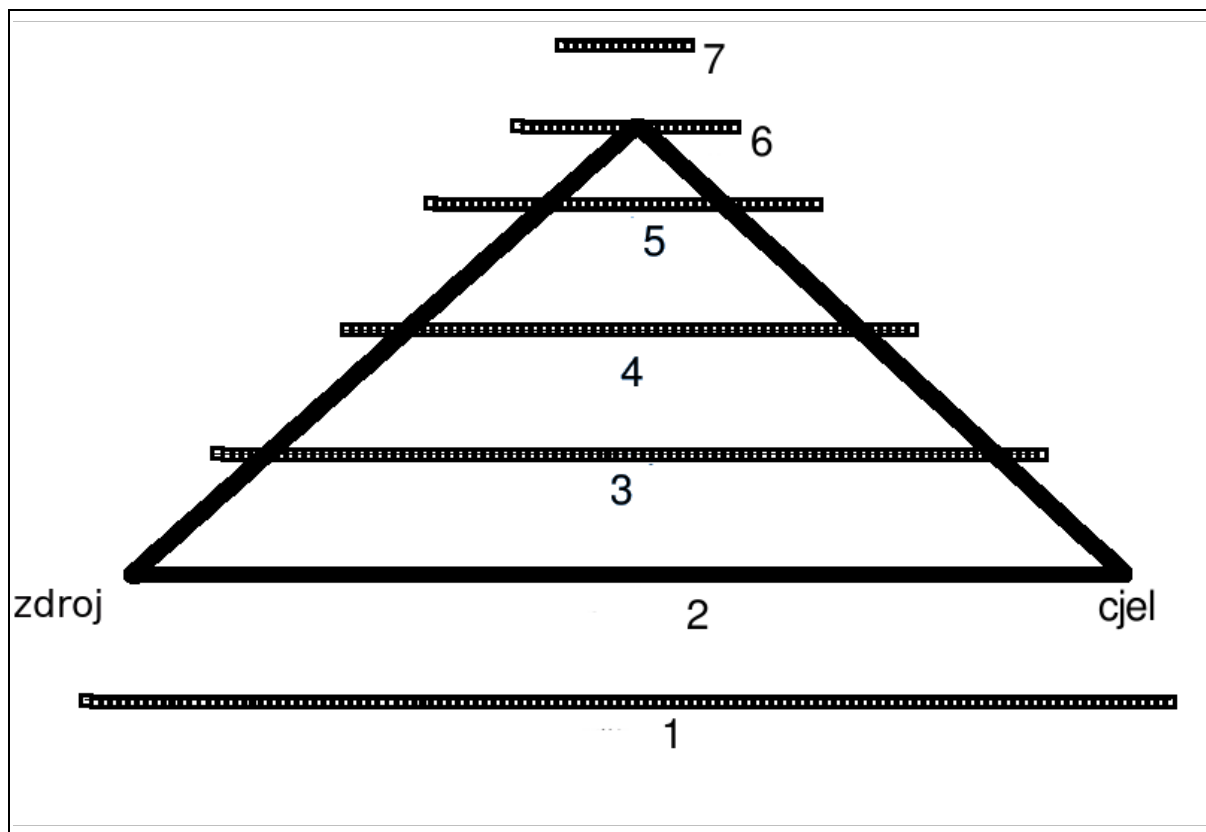
[korpus@korpus.juls.savba.sk](mailto:korpus@korpus.juls.savba.sk), <http://korpus.juls.savba.sk>

**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 spracúvania prirodzeného jazyka. Toto vyplíva z potreby urobiť hĺbkovú analýzu zdrojového nárečia a transformovať zmysel puvodneho textu do cjelovjeho nárečia. Schematicki muožeme proces prekladu znázorniť diagramom podobným tomu na obr. 1. Plocha trojuholníka vijadruje oblasť, v ktorej pracujú tipickje sistemi automatickeho prekladu. Každá vodorovná čjara zodpovedá abstraktnej úrovni transferu medzi zdrojovým a cjelovým nárečím. Čím viššja úroveň, tím abstraktnejší transfer sa uskutočňuje, a výsledok je tím bližší prirodzenjemu nárečju. Úroveň 1vá zodpovedá fonetike a na obrázku je uvedená len kvuoli úplnosťi, pretože vo väčšine sistemou automatickeho prekladu (ako aj v našom článku) ide o písaní text. Úroveň 2há zodpovedá ortografii, transfer na tejto úrovni znamená len zmenu ortografickjeho systému (takíto transfer je použitelní napríklad pri zmeňe pravopisu jedného nárečia, alebo preklad medzi nárečjami, ktorje sa líšja iba ortografiou). Úroveň 3ťja zodpovedá morfológii a je použitelná pre preklad medzi nárečjami, ktorje sa odlišujú maximálne morfológiou (s istými obmedzeňjami muože ísť o dve veľmi blízke príbuzňje nárečia). Pri odlišnejších nárečjach dostaňeme na výstupe syntakticki a semanticki ňezmyselní text. Úroveň 4tá zodpovedá sintaxi, na výstupe dostaňeme text syntakticki správni, aj keď možno s ňezmyselním významom (alebo s významom ňezodpovedajúcim originálnemu textu). Modernje špičkovje sistemi automatickeho prekladu sa k tejto úrovni iba približujú. Úroveň 5ta, semantika, zodpovedá pochopeňju významu slov a slovních spojení originálnjeho textu a ich preklad na slovňje spojeňja s rovnakým významom – na tejto úrovni pracujú prekladateľá-luďja. Úroveň 6ta, na diagrame znázornená vrcholom trojuholníka zodpovedá užítju interlingvi (medzireči), pri ktorom preklad prebehou už po stranách trojuholníka a transfer sa zredukovau na identicku operáciu, pretože všetki črti puvodneho aj preloženjeho textu sú obsiahnutje 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 1v1: Schematicki znázorňeňja trojuholník prekladu

### **Preklad medzi veľmi blízkimi nářečjami**

Blízke (geneticki aj štrukturňe) nářečja majú vela podobných črt. Pri vzdalovaní nářečí rozďjeli medzi nimi celkom dobre sledujú úrovňe v uvedenom trojuholníku – najprú sa zjavja rozďjeli v fonetike (aj v rámci jedneho nářečja či dokonca rozličnorečja), potom v ortografii (pri kodifikácii alebo odšťjepení nářečja, často s politickou motiváciou). Pri morfologických rozďjeloch sme už oprávněni hovoriť o ruoznich nářečjach. Syntax často zostáva kompatibilná aj pri nářečjach od seba značne vzdjalených, a v prípade dramatických rozďjelou v lexike už ňemuožeme hovoriť o blízkich nářečjach v našom poňimaní. Z automatických prekladových sistemou medzi blízkimi nářečjami muožeme spomenúť preklad medzi Češťinou a Slovenčinou[1] a preklad medzi Turečťinou 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 rozdiely sú ľahko algoritmicky popísateľnejšie. Hlavné ortografické rozdiely spočívajú v absencii grafemu „y“, v inej realizácii dvojhlasok a v explicitnom povinnom značení mekkosti spoluhlások d, t, n.

Lexikálne rozdiely sú subtilnejšie, na prví pohľad badaťelne len v niektorých najčastejších slovách, ale v skutočnosti mierne posúvajúce semantický význam celých trjed slov.

### Technická realizácia

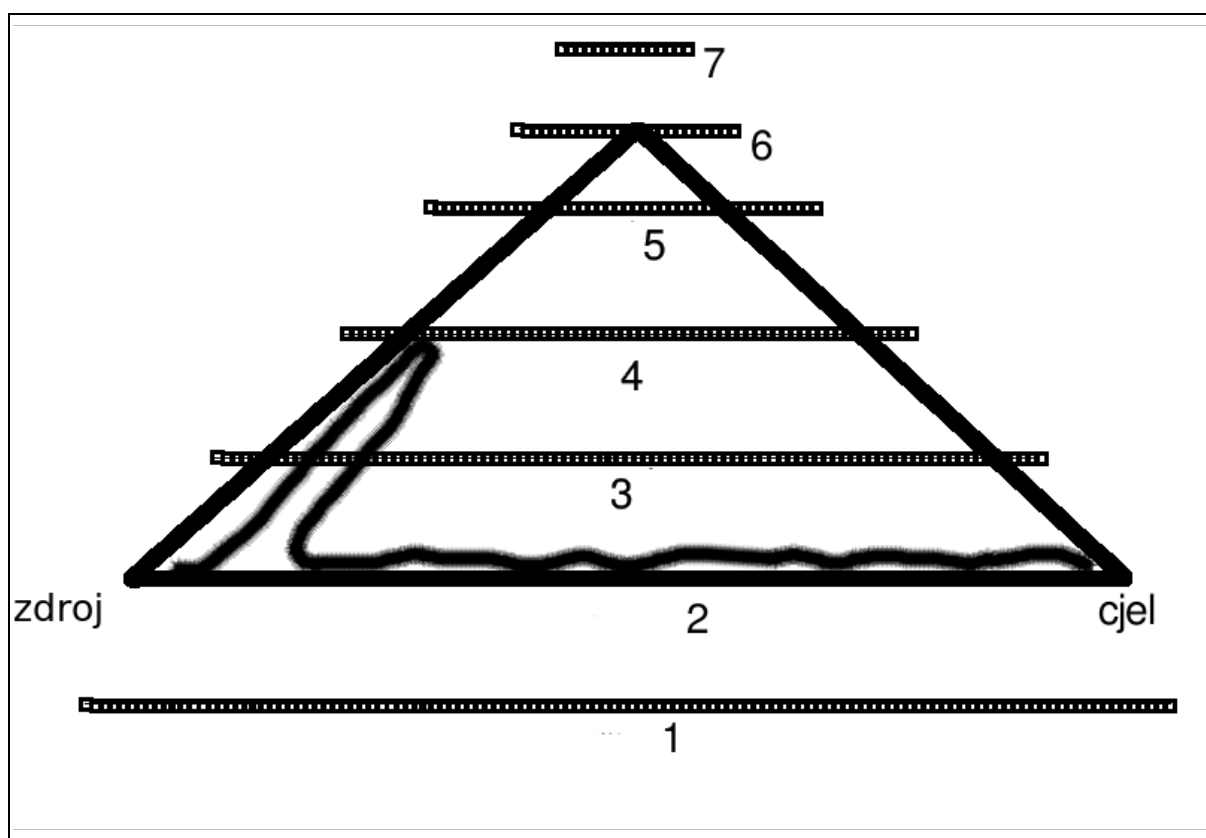
Pred prekladom je text najprú skonvertovaný zo vstupného kódovania do Unicode, potom normalizovaný na NFKC normalizáciu Unicode a všetky ďalšie operácie prebiehajú duosledne v Unicode. Text je tokenizovaný na základné jednotky – tokeny (slová), ku každému tokenu je priradená informácia o prípadných bielich znakoch (whitespace; Leerraum) pred slovom, aby sa po preklade mohlo zrekonštruovať vernuo rozloženie textu. Po preklade je veľkosť písmen preloženého slova upravená tak, aby kopírovala veľkosť písmen puovodného slova – ak je prelozeno slovo dlhšie ako puovodnuo, veľkosť „nadbitočných“ písmen kopíruje veľkosť poslednej písmeni puovodného slova. Toto zabezpečí verní preklad vlastných mjen a prípadných slov písaných kapitálkami. Ako výnimka sú korene slov „Sloven(čina, skí)“ a „Vlád(a)“ vždy v preklade písané so začiatčím veľkým písmenom, podľa úzu užívaného L. Štúrom.

Samotný preklad prebieha v dvoch fázach: najprú sa aplikuje lexikálna transformácia, pri ktorej sa nahrádzajú slová, ktoré sú v štúrovej Slovenčine inak reprezentované. Výhodne sa dá využiť prevažná ekvivalencia morfev medzi súčasnou a štúrovskou Slovenčinou a v prekladovej tabulke stačí povečšine uvjesť iba preklady koreňových morfev, iba niekedy je potrebnuo uvjesť preklad celých tvarov slov.

Druhá fáza prebieha na ortografickej úrovni. Preklady v oboch fázach sú realizované jednoduchým nahrádzaním originálnich reťazcov prekladovými. Začiatky a konce slov sú označené špeciálnymi znakmi (^ začiatok, \$ koniec), čo umožňuje efektívne spracovať transformácie v príveskách slov a zabránuje možným nesprávnym nahradeniam. Vzhľadom na duosledné značenie palatalizovaných spoluhlások v štúrovej Slovenčine je potrebnuo duokladne rozlišovať tvrdje a mekkje „i“ podľa výslovnosti (čo vedie k prekladom

politi→polity, diplo→dyplo, poézia→poesya) a takt'jež bolo potrebnuo zavjest' tvrdje „e“ na označeňja ňepalatalizujúceho „e“ (toto písmeno sme arbitrárne označili znakom „ë“ – U+00EB LATIN SMALL LETTER E WITH DIAERESIS, príkladi prekladou: internet→intërnët). V druhej úrovni budú tjeťto slová transformovaňje na štúrovskí pravopis (polity→politi, dyplo→diplo, poesya→poesia, intërnët→internet).

Naša skratka v prekladovom trojuholníku potom sleduje transfer na ortografickej úrovni, s krátkim vibočeňím do oblast'i semantiki (vlastne iba zámëna ňjektorích lexikálnich jednot'jek).



Obrázok 2hí: Trojuholník prekladu so znázorňëním našej skratki

lexikální preklad	ortografický preklad
u'grék' : u'rék',	u'ovš' : u'ouš',
u'gréc' : u'réc',	u'něš' : u'ňješ',
u'grěč' : u'rěč',	u'éš' : u'ješ',
u'maďarš' : u'uherš',	u'éhoš' : u'jehoš',
u'maďar' : u'uhr',	u'émuš' : u'jemuš',
u'maďara' : u'uhra',	u'é' : u'e',
u'talian' : u'talyan',	u'ý' : u'í',
u'^ludovít' : u'^luděvít',	u'y' : u'i',
u'slávneho' : u'slávňeho',	u'ô' : u'uo',
	u'l' : u'l',
# pieseň -> peseň	u'ä' : u'e',
u'^pies' : u'^pes',	u'ë' : u'e',
u'vidietš' : u'videtš',	
u'vedjetš' : u'vedetš',	u'ia' : u'ja',
u'vedie' : u'vede',	u'dia' : u'dja',
u'erieš' : u'ereš',	u'diakon' : u'diakon',
u'eriešš' : u'erešš',	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 jaziku Python. Hlavnou zamerajúcou programom 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 – písmeni d' a t' budú v NFKD normalizácii, ostatne v NFKC
- e ENCODING alebo --encoding ENCODING – miesto štandardného kodujúca utf-8, predpokladá vstup a výstup v kodujúci ENCODING, ktorú môže byť hociká kodujúca podporovaná pythonom, ale pravdepodobne význam má len jedno z utf-8, iso8859\_2, cp1250, cp852 alebo mac\_latin2. Kodujúca iná než utf-8 nie je kompatibilná s voľbami -D a -d.

Kde sa zvuky mekko vislovujú takto sa zmečujúcou čarkou viznačujú, ale písmeni „d“ a „t“ ju v dobe modernej inakšie označujú, značka táto skoro ako dlhá čarka má podobu. Aby sa historická vernosť zachovala, tieto dve písmeni 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 tjeo čjarki zle zobrazenje, alebo naopak tak ako majú byť skombinované dobre a správne (na moderní spôsob) zobrazenje, a teda tento spôsob nje vždi dobrje výsledki dáva.

## Nedostatki

Ňedostatki uved'enjeho systému prekladu (či prepisu) muožeme rozďelit' na dve skupini. Prvú skupinu tvoria Ňedostatki teoreticki Ňepodstatnje, ktorje je aspon teoreticki možnuo lahko odstrániť aplikovaním dostatočnjeho množstva ľudskej práce. Sem patrí hlavne 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 Ňevihnutnej ľudskej práce tu nje sú žjadnje problemi, ktorje tomuto princípijálne bráňja.

Ňedostatki teoreticki podstatnje sú horšje, pretože viplívajú buď prjamo z návrhu prekladacjeho systému, alebo z vroďeních vlastností oboch verzií Slovenčini a vzt'ahu medzi nimi. Tjeo Ňedostatki nje je možnuo jednoducho odstrániť. Najzávažnejšje z nich sú:

- Absencja kontextovjeho prekladu. Sistem sa vždi pozerá iba na jedno konkrétne slovo. Toto zabraňuje možnosťi lexikálneho prekladu slovami s iním pohlavím, pretože nje je možnuo súčasne preložiť prípadnje adjektíva a slovesá (v menosloví) tak, aby bola zachovaná zhoda pohlaví. V programe sme urobili jed'ínú výnimku, slovo „Bratislava“ prekladáme slovom „Prešporok“ (spolu s patričnými tvarmi skloňenja), pretože id'e o slovo dost' známe a často užívané, a občas sa viskitnuvšju chibu v Ňesúlade pohlavja prípadnjeho mena prídavnjeho sme považovali za menšje zlo ako poňchať tvar „Brat'islava“ (s Ňejasným užívaním v štúrovských dobách)
- Ňerozlíšiteľná homonímja v štandardnej Slovenčine. Najčastejším príkladom sú prídavnje mená Ňijakjeho pohlavja v nominative a akusative v jednotnom počte (príveska -uo) a prídavnje mená ženskjeho a Ňijakjeho pohlavja v množnom počte (príveska -je), ktorje v štandardnej Slovenčine majú rovnakú prívesku (-je, alebo -e ak bola predchádzajúca silaba dlhá) a nje je Ňijakí spôsob, ako bez duokladnej semantickej analísi určiť správni preklad (taketo slovnje 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 občasne užít'je, na demonštráciu štúrovskej Slovenčini a ako pomuocka pre prekladateľou do štúrovskej Slovenčini. Po jednoduchej úprave (náhrada slovníka) je možnuo program užít' pre preklad medzi podobne odlišnými jazikovými sistemami (napríklad v prípade závažnejších zmjen v Slovenskom pravopise pri preklade do novej normi).

Program je dostupní pod licenciou GNU GPL v. 2.0, a jeho demoversiu prístupnú cez WWW rozhraňa prostredníctvom jednoduchjeho CGI skriptu je možnuo si prezreť na stránke Jazikovednjeho ústavu Ludevíta Štúra SAV[5]. O potrebe a užitočnosti takjeho prekladu svedčí aj neočakávaná popularita, ktorej sa istú dobu uved'enuo WWW rozhraňa tešilo[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, Tatrin, 1846.
4. Morfológia slovenského jazyka. Red. J. Ružička. Bratislava, Vydavateľstvo Slovenskej akademie vied 1966. 896 s.
5. <http://vvv.juls.savba.sk/ludevit/>
6. Štúrovská slovenčina. In: SME, 20. 12. 2006, s. 29