

Building an Error-tagged Learner Corpus of Czech

Jirka Hana, Alexandr Rosen & Barbora Štindlová

Charles University, Prague & Technical University, Liberec

Institute of Formal and Applied Linguistics
Seminar of Formal Linguistics
April 22, 2013

Outline of the talk

- 1 Learner Corpora
- 2 CzeSL
- 3 Error Annotation of CzeSL
- 4 Evaluation
- 5 Postprocessing
- 6 Thanks

Outline of the talk

1 Learner Corpora

2 CzeSL

3 Error Annotation of CzeSL

4 Evaluation

5 Postprocessing

6 Thanks

Learner Corpora

- Include texts produced by learners of a foreign language
- Early 1990s: as a source of data for learners' dictionaries (e.g., *Longman Learner Corpus*)
- Used by authors of textbooks, methodologists and researchers in
 - Teaching X as a Foreign Language
 - Second Language Acquisition
- Deviant forms can be corrected and their error type identified
- There can be simultaneous deviations on multiple levels

Annotation of Learner Corpora

Learner corpora can be annotated in two independent ways:

Linguistic annotation

- Lemmatization, morphological tagging, syntactic structure, etc.
- On the original text or on the corrected text
- Usually automatic or semiautomatic

Error annotation

- Correcting and/or categorizing errors
- Diverse annotation systems
- Usually manual

Capturing errors

1. Implicit – errors are identified and corrected

- Pros:
 - faster training of annotators
 - faster process of annotation
- Cons:
 - results hard to search and analyze

2. Explicit – errors are identified and categorized

- Error categories (tags) reflect a specific theory

Outline of the talk

1 Learner Corpora

2 CzeSL

3 Error Annotation of CzeSL

4 Evaluation

5 Postprocessing

6 Thanks

CzeSL – A learner corpus of Czech

- The first large learner corpus of a Slavic language (*PiKUST* for Slovenian includes only 35KW)
- The first large learner corpus of Czech
- Part of the *AKCES* project – acquisition corpora of Czech (includes native speakers' spoken and written classroom language)

The parameters

- 2 MW (only a part annotated)
- 3 subcorpora by L1
 - Slavic language: Russian, Ukrainian, Polish
 - Other Indo-European language: German, English, French
 - Other language: Vietnamese, Arabic
- Written and spoken part
- All levels of proficiency according to CEFR
- Original documents are electronic and hand-written texts
- Elicited on various occasions in the class
- Metadata on learner and task (19 items)

Size

Size of various subcorpora (in thousand words, approximate)

	transcribed	annotated	doubly annotated
<i>CzeSL – Foreigners</i>			
spoken	11	0	
written	1,315	200	75
university	732	0	
<i>ROMi – Czech Roma</i>			
spoken	540 (270)	0	
written	450	170	110
<i>native Czech</i>			
spoken	1,046	?	
written	150	?	

Outline of the talk

- 1 Learner Corpora
- 2 CzeSL
- 3 Error Annotation of CzeSL
- 4 Evaluation
- 5 Postprocessing
- 6 Thanks

Workflow

- Acquisition
- Transcription
- Proofreading
- Conversion to PML
- Error annotation
- Revision
- *Adjudication*
- Postprocessing

Strategy

- Minimal correction
- Capture only grammatical and lexical characteristics of non-native language
- Relative to Literary Czech

Error Annotation of a Flective Language

Problems

- Inflection (nouns: 15 basic paradigms, subparadigms, subsub...)
- Derivation, agreement, word-order reflecting information structure, etc.

Solutions

- Multilevel annotation scheme
- Combining manual and automatic annotation

Error Annotation of a Flective Language

Problems

- Inflection (nouns: 15 basic paradigms, subparadigms, subsub...)
- Derivation, agreement, word-order reflecting information structure, etc.

Solution

- Multilevel annotation scheme
- Combining manual and automatic annotation

Ruska čeština

Víktor je intelektuální pam z ~~Ruského~~ Studeje ve škole, protože ne umí psat a číst správně. Bydlí na kolejce vedle školy, má jednu sestru Irenu, kterou se včetně univerzity u profesora Smutněveselého. Bohužel, Víktor není dobrý student, protože spí na lekcích, ale jeho sestra ~~píše všechno~~ a ^{a být všechno využitelné} rozumí českého profesora Smutněveselého. Včera Irena jde na prohádku spolu z kamarádem, ale její bratr dělá nic. Její čeština je špatná, vím, že ~~se~~ vratí ^{Rusko} do ~~Ruského~~ a tam budi ~~školovat~~ a pomoci myt podkladky.

~~Hodinopisec~~ Kamarád Ireny je Američan a chytrý muž. On ~~stále~~ miluje Irenu a chce se vrátit na ni návštěvu.

Viktor je mladý pan z Polska Ruska. Studuje češtinu</in> ve škole, protože ne umí psat a čist spravně. Bydlí na koleje vedle školy, má jednu sestru Irenu, která se učí na univerzite u profesora Smutneveselého. Bohužel, Viktor není dobrým studentem, protože spí na lekci, ale jeho sestra {piše všechno -> všechno piše} a vyborně rozumí českého profesora Smutneveselého {a brzo delá domácí úkol}</in>. Večeře Irena jde na prohaska spolu z kamaradem, ale její bratr dělá nic. Jeho čeština je špatná, vím, že se vrátit ve Polsku Rusku a tam budí studovat u pomalu myt podlahy.

Kamarad Ireny je {A|a}meričan a chytry muž. On miluje Irenu a chce se vzít na ni. protože ona je hezká, taky chytra, rozumí ho a umí vyborný vařit.

Viktor je mladý pan z Polska Ruska. Studuje {češtinu}<in> ve škole, protože ne umí psat a čist spravně. Bydlí na koleje vedle školy, má jednu sestru Irenu, která se učí na univerzite u profesora Smutneveselého. Bohužel, Viktor není dobrým studentem, protože spí na lekci, ale jeho sestra {piše všechno -> všechno piše} a vyborně rozumí českého profesora Smutneveselého {a brzo delá domácí úkol}<in>. Večeře Irena jde na prohaska spolu z kamaradem, ale její bratr dělá nic. Jeho čeština je špatná, vím, že se vratit ve Polsko Rusku a tam budí studovat u pomalu myt podlahy.

Kamarad Ireny je {A|a}meričan a chytry můž. On miluje Irenu a chce se vzít na ní. protože ona je hezká, taky chytra, rozumí ho a umí vyborný vařit.

Viktor je mladý **pan** z **Polska** Ruska. Studuje {češtinu}<in> ve škole, protože ne umí psat a čist spravně. Bydlí na kolej vedle školy, má jednu sestru Irenu, která se učí na univerzite u profesora Smutneveselého. Bohužel, Viktor není **dobrym** student, protože spí na lekci, ale jeho sestra {piše všechno -> všechno piše} a vyborně rozumí **českeho** profesora **Smutneveselého** {a brzo delá domácí ukol}<in>. Večeře Irena jde na **prohaska** spolu z kamaradem, ale její bratr **dělá** nic. Jeho čeština je špatná, **vím**, že se **vratit** ve **Polsko Rusku** a tam **budí** studovat u pomalu **myt** podlahy.

Kamarad Ireny je {A|a}meričan a chytry můž. On miluje Irenu a chce se vzít na ní. protože ona je hezká, taky chytra, rozumí ho a umí **vyborný** vařit.

Multilevel Annotation Scheme

Level 0

- Original text (transcribed, self-corrections inlined)

Level 1

- Corrections disregarding word context
- Spelling, form of stems and endings
- Result: sequence of existing Czech forms

Level 2

- Remaining errors: syntactic, lexical, word-order, style, referential, negation, ...
- Result: grammatically correct sentence

**Bojal jsme se že ona se ne bude libila slavnou prahu,
proto to bylo velmí vadí pro mně.**

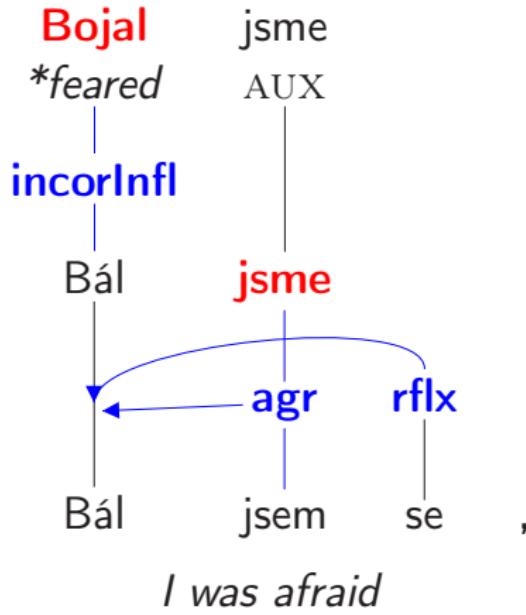
Bál jsem se, že se jí nebude líbit slavná Praha,
protože to by mi velmi vadilo.

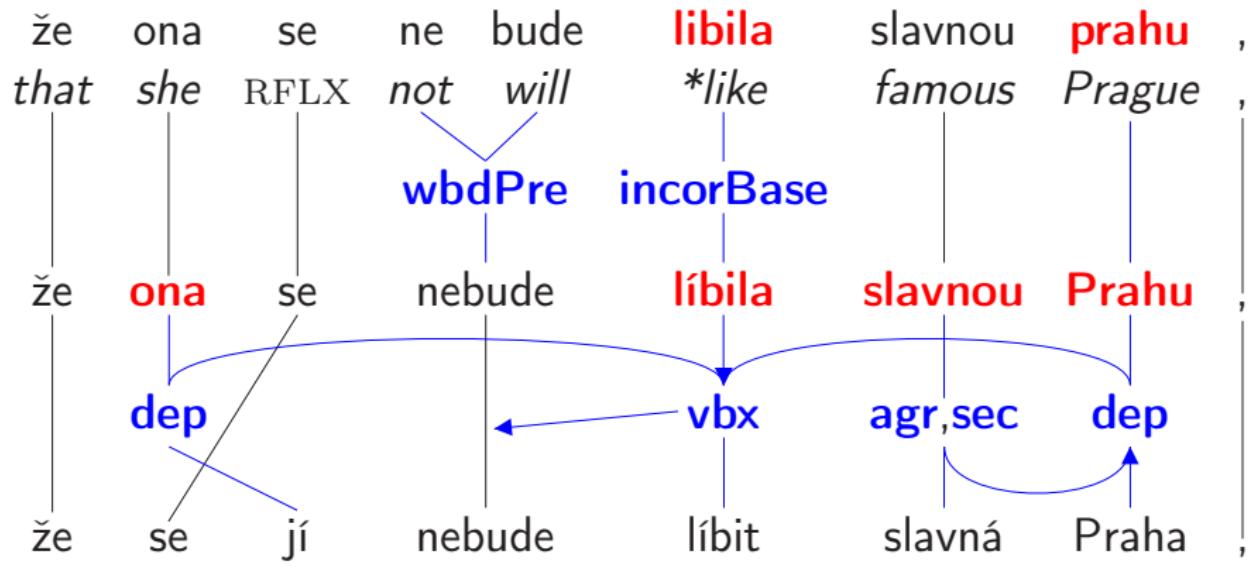
'I was afraid that she would not like the famous city of Prague,
because I would be very unhappy about it.'

**Bojal jsme se že ona se ne bude libila slavnou prahu,
proto to bylo velmí vadí pro mně.**

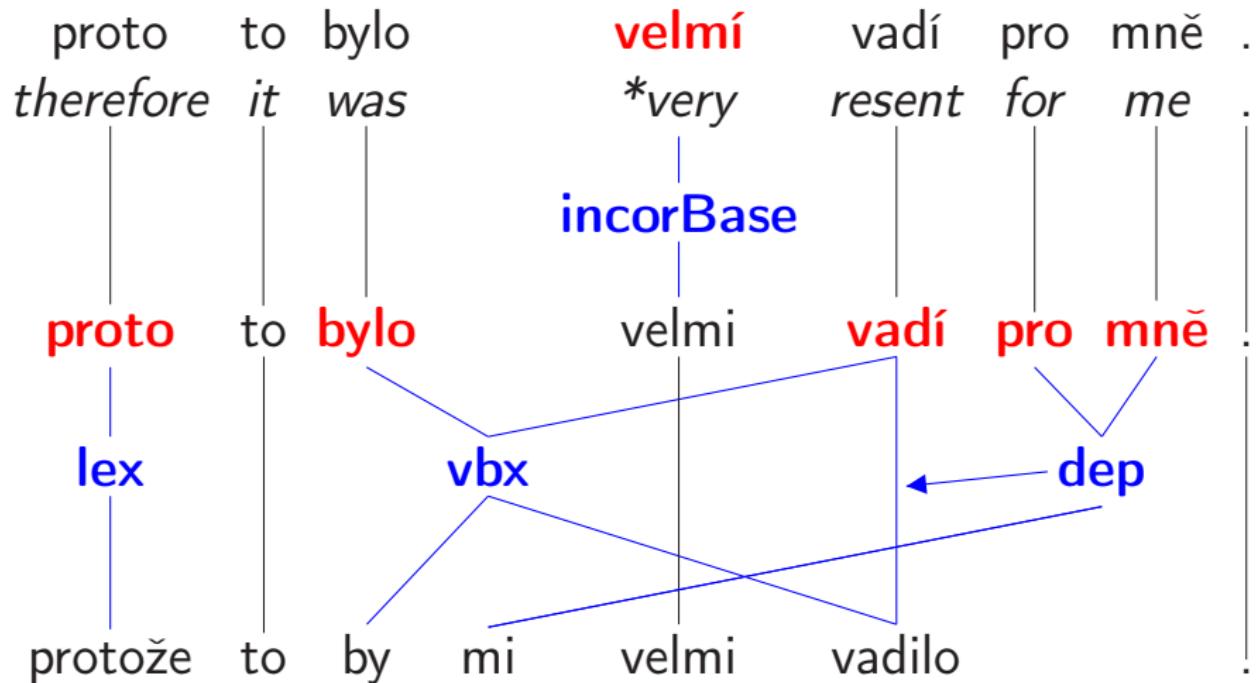
Bál jsem se, že se jí nebude líbit slavná Praha,
protože to by mi velmi vadilo.

'I was afraid that she would not like the famous city of Prague,
because I would be very unhappy about it.'

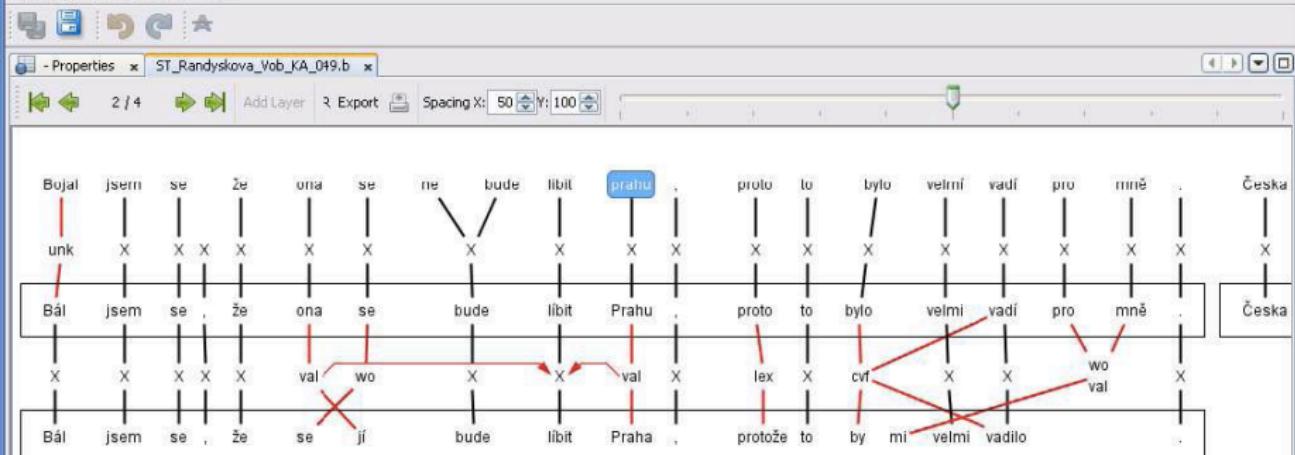




that she would not like the famous city of Prague,



because I would be very unhappy about it.



Proč mám/nemám rád (Č|č) eskou republiku?

Už se nacházíme v české republice až půl roku. toho mě musilo by stačit, abych rozuměl, mám rád to země nebo ne rád. teďko mužů učítříčku, že českou republiku já miluju. tento zámcí má všechna že potřebuju ja a moje přítelkyně. Bojal jsem se že ona se ne bude líbit **práh**, proto to bylo velmi vadi pro mně. Česká republika je krásné místo, tady je hodně hezké památek například pražský hrad a vyšehrad. libím se moc pražský hrad, protože tam je zamky, který velmi krásne a hezke. take v čechach je dobrá příroda a když jsme se procházeli na divoké řárce byly šokovani ečť z těch krásnych pohledů. Je to nekrásneší místo ve všem bilém světě. take rád že Česi je dobrí

Fit WFit Orig Zoom
LiLuJu. tento země na všechna
ja a moje přítelkyně. Boží Ja ~~sem~~
že libit počátku, proto to bylo velmi
Česká republika je krásné místo,
hezké památky, například pražský
libím se moc pražský hrad, proto

Error tags

- 22 error tags added manually
- 7 error tags added automatically

	manual	automatic	total
level 1	8	1	9
level 2	11	6	17
level 1 & 2	3	0	3
total	22	7	29

- Additional formal tags are added automatically by comparing original and corrected forms

incor			incorrect form	
	incorInfl	M	inflection error	L1
	incorBase	M	stem error	L1
	incorOther	A	other	L1
fw			foreign word, neologism, unidentifiable	
	fwFab	M	newly created “Czech” word	L1
	fwNc	M	foreign word	L1
	flex	M	inflection of fw	L1
wbd			word-boundary error	
	wbdPre	M	separate prefix, attached preposition	L1
	wbdComp	M	incorrectly separated/joined composites	L1
	wbdOther	M	other word-boundary errors	L1
styl			colloquial, bookish, regional expression	
	stylColl	M	colloquial expression	L1,L2
	stylOther	M	bookish, regional, slang expression	L1,L2
	stylMark	M	filler	L2
problem		M	problem	L1,L2

agr	M	agreement error	L2	
dep	M	structural error	L2	
ref	M	pronominal reference error	L2	
vbx	M	complex verb error	L2	
	cvf	A	analytical verb form error	L2
	mod	A	modal verb error	L2
	vnp	A	copula	L2
rflx	M	reflexive form error	L2	
neg	M	negation error	L2	
odd	A	extra word	L2	
miss	A	missing word	L2	
wo	A	word-order error	L2	
lex	M	lexical and idiomatic error	L2	
use	M	incorrect use of a category	L2	
sec	M	secondary error	L2	
disr	M	word salad	L2	

Outline of the talk

1 Learner Corpora

2 CzeSL

3 Error Annotation of CzeSL

4 Evaluation

5 Postprocessing

6 Thanks

Evaluation – Sample

- 67 texts, about 150 words each
- 9373 tokens
- CEFRL level A2–B1
- Various L1s
- 14 annotators, each text by two

A measure of IAA: Kappa

$$\bullet \quad \kappa = \frac{\Pr(a) - \Pr(e)}{1 - \Pr(e)}$$

$\Pr(a)$ – observed agreement

$\Pr(e)$ – chance agreement

- $\kappa = 1$ – perfect agreement
- $\kappa = 0$ – random agreement
- $\kappa > 0.4$ – reasonable

Evaluation

Higher IAA on formally well defined tags

error tag	κ
<i>incorBase</i>	0.75
<i>agr</i>	0.54
<i>styl</i>	0.38

Inter-annotator agreement

9848 words

tag	only A1	only A2	A1 & A2	κ
incor	168	130	894	0.84
incorBase	167	165	559	0.75
incorInfl	173	130	250	0.61
wbd	14	21	45	0.72
fw	25	17	18	0.46
agr	82	99	110	0.54
dep	99	118	87	0.43
neg	11	9	9	0.47
styl	19	14	10	0.38
lex	107	131	74	0.37
use	60	74	19	0.21
sec	45	18	4	0.11

Examples of high IAA

Agreement error

$\kappa = 0.54$

- (1) Viděl malého Petra
- (2) Viděl *malou Petra

Why not higher?

Different intonations

L0: Věci budou *težki

A1 – L1: těžký, L2: těžké + AGR

A2 – L1: těžké, L2: těžké

Examples of high IAA

Agreement error

$\kappa = 0.54$

- (1) Viděl malého Petra
- (2) Viděl *malou Petra

Why not higher?

Different corrections

L0: *Věci budou *težki*

A1 – L1: *těžký*, L2: *těžké + AGR*

A2 – L1: *těžké*, L2: *těžké*

Examples of low IAA

Lexical error

$\kappa = 0.37$

Due to semantic proximity of lexemes annotators disagree about the need for correction:

- (3) *když se dívám na *?druhý/jiný kultury*
‘when I look at other cultures’

On the other hand, some lexemes are distant enough and annotators agree about the need for correction:

- (4) **housenky/housky kupuju v pekařství*
‘I buy caterpillars in the baker’s shop’

Some reasons for low IAA

- Errors of type **lex** involve a high degree of subjective judgement, thus cannot aim at high IAA.
- Errors of type **sec** – highly formal specific, due to primary errors.

Evaluation – Conclusion

- Morphosyntactic errors are easy to formalize and lead to a high κ – incor, agr, dep
- Semantic errors depend on subjective judgement, should standard measures be applied?

Outline of the talk

1 Learner Corpora

2 CzeSL

3 Error Annotation of CzeSL

4 Evaluation

5 Postprocessing

6 Thanks

Postprocessing

- Formal error tags
- Morphology

Formal error tags

Error type	Error description	Example
Cap0	capitalization: incor. lower case	<i>evropě/Evropě; štědrý/Štědrý</i>
Cap1	capitalization: incor. upper case	<i>Staré/staré; Rodině/rodině</i>
Vcd0	voicing assimilation: incor. voiced	<i>stratíme/ztratíme; nabítku/nabídku</i>
Vcd1	voicing assimilation: incor. voiceless	<i>zbalit/sbalit; nigdo/nikdo</i>
VcdFin0	word-final voicing: incor. voiceless	<i>kdyš/když; vztach/vztah</i>
VcdFin1	word-final voicing: incor. voiced	<i>přez/přes; pag/pak</i>
Vcd	voicing: other errors	<i>protošel/protože; hodili/chodili</i>
Palat0	missing palatalization (<i>k,g,h,ch</i>)	<i>amerikě/Americē; matkě/matce</i>
Je0	<i>je/ě:</i> incorrect ě	<i>ubjehlo/uběhlo; Nejvjetší/Největší</i>
Je1	<i>je/ě:</i> incorrect je	<i>vjeděl/věděl; vjeci/věci</i>
Mne0	<i>mě/mně:</i> incorrect mě	<i>zapoměla/zapomněla</i>
Mne1	<i>mě/mně:</i> incor. mně, mňe, mňě	<i>mněla/měla; rozumněli/rozuměli</i>
ProtJ0	protethic <i>j</i> : missing <i>j</i>	<i>sem/jsem; menoval/jmenoval</i>
ProtJ1	protethic <i>j</i> : extra <i>j</i>	<i>jse/se; jmé/mé</i>
ProtV1	protethic <i>v</i> : extra <i>v</i>	<i>vosm/osm; vopravdu opravdu</i>
EpentE0	<i>e</i> epenthesis: missing <i>e</i>	<i>domček/domeček</i>
EpentE1	<i>e</i> epenthesis: extra <i>e</i>	<i>rozeběhl/rozběhl; účety/účty</i>

Morphology

- L2 (correct Czech): disambiguated lemma + tag
- L1: morphological analysis
- If possible, disambiguated lemma+tag is transferred from L2 to L1
- If forms are different, but L2 lemma is among L1 lemmas:
use L2 lemma for L1 $L1=má$ 'has' or 'my'
 $L2=mou$ 'my': use the lemma $můj$ 'my' on L1

Outline of the talk

1 Learner Corpora

2 CzeSL

3 Error Annotation of CzeSL

4 Evaluation

5 Postprocessing

6 Thanks

Thanks to...

... other members of the team, esp.:

Barbora Štindlová, Tomáš Jelínek, Svatava Škodová, Karel Šebesta,
Vladimír Petkevič, Hana Skoumalová, Milena Hnátková, Jan Štěpánek,
Zuzanna Bedřichová, Kateřina Šormová

... the sponsors:

- The European Social Fund and the Czech government:
Education for Competitiveness – Innovation in Education in the Field of Czech as a Second Language (CZ.1.07/2.2.00/07.0259)
- Large Research, Development and Innovation Infrastructures:
The Czech National Corpus (LM2011023)

... and you!