

Annotating Foreign Learners' Czech

Alexandr Rosen¹ Svatava Škodová²
Barbora Štindlová² Jirka Hana¹

¹Charles University, Prague, Czech Republic

²Technical University, Liberec, Czech Republic

FDSL 8.5

Formal Description of Slavic Languages 8.5

Faculty of Arts, Masaryk University, Brno

26–27 November 2010

Outline of the talk

- 1 Introduction
- 2 Annotation scheme
- 3 Annotation process
- 4 Conclusion

Outline of the talk

- 1 Introduction
- 2 Annotation scheme
- 3 Annotation process
- 4 Conclusion

Learner Corpora

- Include texts produced by learners of a foreign language
- Early 1990s: used to compile learners' dictionaries (e.g., *Longman Learner Corpus*)
- Used by authors of textbooks and researchers in *2nd Language Acquisition*
- Deviant forms can be corrected and their error type identified
- There can be simultaneous deviations on multiple levels

Some currently available learner corpora

Size	L1	TL	TL proficiency	Error annotation
ICLE – <i>Internat'l Corpus of Learner English</i>				
3M	21	English	advanced	yes
CLC – <i>Cambridge Learner Corpus</i>				
30M	130	English	all levels	yes
USE – <i>Uppsala Student English Corpus</i>				
1.2M	Swedish	English	advanced	no
HKUST – <i>Hong Kong UST Corpus of Learner English</i>				
25M	Chinese	English	advanced	yes
CLEC – <i>Chinese Learner English Corpus</i>				
1M	Chinese	English	5 levels	yes
JEFL – <i>Japanese EFL Learner Corpus</i>				
0.7M	Japanese	English	advanced	yes
FALKO – <i>Fehlerannotiertes Lernerkorpus</i>				
1.2M	various	German	advanced	yes
FRIDA – <i>French Interlanguage Database</i>				
0.2M	various	French	intermediate	yes
CIC – <i>Chinese Interlanguage Corpus</i>				
2M	96	Chinese	intermediate	?
ASK – <i>Andersspråkskorpus</i>				
?	10	Norwegian	two levels	yes

A learner corpus of Czech

- The CzeSL Project: Czech as a Second Language
- Czech: rich inflection, derivation, complex agreement rules and information-structure-driven constituent order
- 2 million words to be transcribed, corrected and annotated within 3 years
- L1: Slavic (Russian, Ukrainian), Vietnamese, Romani, Chinese, ...
- Beginners to advanced learners
- Hand-written texts, elicited on various occasions in the class

Outline of the talk

- 1 Introduction
- 2 Annotation scheme**
- 3 Annotation process
- 4 Conclusion

Three-level format

- Level 0 for the original
- Successive corrections:
 - ▶ Level 1 – graphemics and morphology.
 - ▶ Level 2 – agreement, valency, complex verb forms, lexicon, word order and negative concord
- Able to capture errors in multi-word discontinuous expressions
- Errors due to missed agreement, valency and pronominal reference have links to the words responsible for the proper form
- Automatic assignment of error tags wherever possible, based on comparing faulty and corrected forms, sometimes using morphosyntactic tags, assigned by a tagger

A sample sentence

republicu a miluju. Tento ~~em~~ ^{em} ~~je~~ ^{že} potřebuji a a moje přítelkyně. Bojal ~~jsme~~ ^{že} se že ona se bude líbit prahu, proto to bylo velmi vadí pro mě. Česká republika je krásné místo.

Bojal jsme
*feared AUX-*PL
I was afraid

že ona se ne bude líbit slavnou prahu,
that she RFL not will *like *famous *prague,
that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mě.
therefore it was *very resent for me.
because I would be very unhappy about it.

A sample sentence

republiku a miluju. Tento ~~em~~ ^{em} ~~je~~ ^{že} potřebuju ja a moje přítelkyně. Bojá ~~jsme~~ ^{že} ona se bude líbit prahu, proto to bylo velmi vadí pro mně. Česká republika je krásné místo.

Bojál jsme
feared AUX-*PL
I was afraid

že ona se ne bude líbit slavnou prahu,
that she RFL not will *like *famous *prague,
that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mně.
therefore it was *very resent for me.
because I would be very unhappy about it.

A sample sentence

republicu a miluju. Tento ~~em~~ ^{em} ~~je~~ ^{je} ~~potřebuju~~ ^{potřebuju} ja a moje přítelkyně. Bojá ~~jsme~~ ^{jsme} se že ona se bude líbit prahu, proto to bylo velmi vadí pro mě. Česká republika je krásné místo.

Bojá ál jsmeem

feared AUX-SG

I was afraid

že ona se ne bude líbila slavnou prahu,
that she RFL not will *like *famous *prague,

that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mě.
therefore it was *very resent for me.

because I would be very unhappy about it.

A sample sentence

republicu a miluju. Tento ~~em~~ ^{em} ~~je~~ ^{je} ~~potřebuju~~ ^{potřebuju} ~~ja~~ ^{ja} ~~a~~ ^a ~~mám~~ ^{mám} ~~přítelkyni~~ ^{přítelkyni}. ~~Bojím~~ ^{Bojím} ~~se~~ ^{se} ~~že~~ ^{že} ~~ona~~ ^{ona} ~~se~~ ^{se} ~~ně~~ ^{ně} ~~bude~~ ^{bude} ~~líbit~~ ^{líbit} ~~prahu~~ ^{prahu}, ~~proto~~ ^{proto} ~~to~~ ^{to} ~~bylo~~ ^{bylo} ~~velmi~~ ^{velmi} ~~vadí~~ ^{vadí} ~~pro~~ ^{pro} ~~mně~~ ^{mně}. ~~Česká~~ ^{Česká} ~~republika~~ ^{republika} ~~je~~ ^{je} ~~krásné~~ ^{krásné} ~~místo~~ ^{místo}.

BojááI jsmeem se
feared AUX-SG RFL

I was afraid

že ona se ne bude líbila slavnou prahu,
that she RFL not will *like *famous *prague,
that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mě .
therefore it was *very resent for me.
because I would be very unhappy about it.

A sample sentence

republicu a miluju. Tento ~~em~~ ^{em} ~~je~~ ^{je} ~~se~~ ^{se} ~~že~~ ^{že} potřebuju a a moje přítelkyně. Bojal jsem se že ona se bude líbit prahu, proto to bylo velmi vadí pro mě. Česká republika je krásné místo.

Bojaál jsmeem se,
feared AUX-SG RFL

I was afraid

že ona se ne bude líbila slavnou prahu,
that she RFL not will *like *famous *prague,
that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mě.
therefore it was *very resent for me.
because I would be very unhappy about it.

A sample sentence

republicu jsem miluju. Tento...
že potřebuju a a moje přítelkyně. Bojal jsem
se že ona se bude líbit prahu, proto to bylo velmi
vadí pro mě. Česká republika je krásné místo.

Boja^al js^{me}em se,
feared AUX-SG RFL

I was afraid

že ona^í se ne bude líbila slavnou prahu,
that her RFL not will *like *famous *prague,

that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mě.
therefore it was *very resent for me.

because I would be very unhappy about it.

A sample sentence

republicu jsem miluji. Tento...
že potřebuji a a moje přítelkyně. Bojal jsem
se že ona se bude líbit prahu, proto to bylo velmi
vadí pro mě. Česká republika je krásné místo.

Boja^ál js^{me}em se,
feared AUX-SG RFL

I was afraid

že ~~ona~~^í se ne bude líbila slavnou prahu,
that her RFL not will *like *famous *prague,

that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mě .
therefore it was *very resent for me.

because I would be very unhappy about it.

A sample sentence

republicu a miluju. Tento ~~em~~ ^{em} ~~se~~ ^{že} potřebuju a a moje přítelkyně. Bojal jsem se že ona se bude líbit prahu, proto to bylo velmi vadí pro mě. Česká republika je krásné místo.

Bojaál js~~me~~em se,
feared AUX-SG RFL

I was afraid

že ~~ona~~í se nebude libila slavnou prahu,
that her RFL not will *like *famous *prague,
that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mě .
therefore it was *very resent for me.
because I would be very unhappy about it.

A sample sentence

republicu a miluju. Tento ~~em~~ ^{em} ~~je~~ ^{že} potřebuji a a moje přítelkyně. Bojál ~~jsa~~ ^{že} se že ona se bude líbit prahu, proto to bylo velmi vadí pro mě. Česká republika je krásné místo.

Bojaál js~~me~~em se,
feared AUX-SG RFL

I was afraid

že ~~ona~~í se nebude ~~líbit~~ slavnou prahu,
that her RFL not will like *famous *prague,

that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mě .
therefore it was very resent for me.

because I would be very unhappy about it.

A sample sentence

republicu a miluju. Tento ~~em~~ ^{em} ~~je~~ ^{je} ~~potřebuju~~ ^{potřebuju} ~~ja~~ ^{ja} ~~a~~ ^a ~~mám~~ ^{mám} ~~přítelkyni~~ ^{přítelkyni}. ~~Bojím~~ ^{Bojím} ~~se~~ ^{se} ~~že~~ ^{že} ~~ona~~ ^{ona} ~~se~~ ^{se} ~~bude~~ ^{bude} ~~líbit~~ ^{líbit} ~~prahu~~ ^{prahu}, ~~proto~~ ^{proto} ~~to~~ ^{to} ~~bylo~~ ^{bylo} ~~velmi~~ ^{velmi} ~~vadí~~ ^{vadí} ~~pro~~ ^{pro} ~~mně~~ ^{mně}. Česká republika je krásné místo.

Bojáál jsmeem se,
feared AUX-SG RFL

I was afraid

že ~~ona~~í se nebude ~~lí~~bilat slavnouá prahu,
that her RFL not will like famous *prague,

that she would not like the famous city of Prague,

proto to bylo velmi vadí pro mně .
therefore it was very resent for me.

because I would be very unhappy about it.

A sample sentence

republicu a miluju. Tento...
že potřebuju a a moje přítelkyně. Bojal jsem
se že ona se bude líbit prahu, proto to bylo velmi
vadí pro mě. Česká republika je krásné místo.

Bojaál jsmeem se,
feared AUX-SG RFL

I was afraid

že ~~ona~~jí se nebude ~~lí~~bit slavnouá ~~p~~Prahua,
that her RFL not will like famous prague,

that she would not like the famous city of Prague,

protože to bylo velmi vadí pro mě .
therefore it was very resent for me.

because I would be very unhappy about it.

A sample sentence

republicu a miluju. Tento ~~em~~ ^{em} ~~je~~ ^{je} ~~potřebuju~~ ^{potřebuju} ~~ja~~ ^{ja} a moje přítelkyně. Bojá ~~jsou~~ ^{jsou} ~~se~~ ^{se} ~~že~~ ^{že} ona se bude líbit prahu, proto to bylo velmi vadí pro mě. Česká republika je krásné místo.

Bojá ál jsmeem se,
feared AUX-SG RFL

I was afraid

že ~~ona~~ í se nebude ~~lí~~ bit slavnouá ~~ř~~ Prah u,
that her RFL not will like famous prague,

that she would not like the famous city of Prague,

proto že to by ~~to~~ velmi vadí pro mě .
therefore it was very resent me.

because I would be very unhappy about it.

A sample sentence

republicu a miluju. Tento člověk mi vadí
že potřebuju a a moje přítelkyně. Bojáť jsem
se že ona se bude líbit prahu, proto to bylo velmi
vadí pro mě. Česká republika je krásné místo.

Bojáť jsem se,
feared AUX-SG RFL

I was afraid

že ona se nebude líbit slavnou Prahu,
that her RFL not will like famous prague,

that she would not like the famous city of Prague,

protože to by velmi vadí pro mě.
therefore it would very resented me.

because I would be very unhappy about it.

A sample sentence

že potřebuju Já a moje přítelkyně. Budat Js^{em} se že ona se bude líbit prahu, proto to bylo velmi vadi pro mně. Česká republika je krásné místo.

Bojaál jsmeem se,
feared AUX-SG RFL

I was afraid

že ~~ona~~ji se nebude ~~li~~bí~~la~~t slavnouá ~~p~~Prahua,
that her RFL not will like famous prague,

that she would not like the famous city of Prague,

protože to by~~le~~ velmí vadílo pro mně .
therefore it would very resent me.

because I would be very unhappy about it.

A sample sentence

že potřebuju lásku a máš přítelkyni. Budeš se
se že ona se bude líbit práhu, proto to bylo velmi
vadí pro mě. Česká republika je krásné místo.

Bojaál jsmeem se,
feared AUX-SG RFL

I was afraid

že ~~ona~~ji se nebude ~~li~~bi~~la~~t slavnouá ~~p~~Prahua,
that her RFL not will like famous prague,

that she would not like the famous city of Prague,

protože to by~~le~~ velmí vadílo pro mně.
therefore it would very resent me.

because I would be very unhappy about it.

A sample sentence

republicu se miluju. Tento...
že potřebuju ja a moje přítelkyně. Bojáť ~~jsme~~^{em}
se že ona se bude líbit prahu, proto to bylo velmi
vadí pro mne. Česká republika je krásné místo.

Bojáť | js~~me~~em | se,
feared AUX-SG RFL

I was afraid

že ~~ona~~í se nebude ~~líbit~~ slavnou~~á~~ ~~p~~Praha,
that her RFL not will like famous prague,

that she would not like the famous city of Prague,

protože to by~~to~~ ~~velm~~i vadí~~lo~~ mi.
therefore it would very resent me.

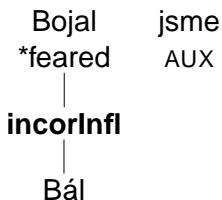
because I would be very unhappy about it.

Annotation of a sample sentence, part I

Bojal	jsme
*feared	AUX

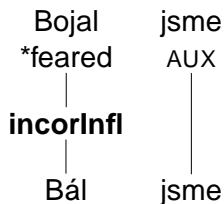
I was afraid

Annotation of a sample sentence, part I



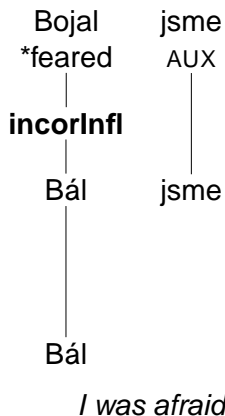
I was afraid

Annotation of a sample sentence, part I

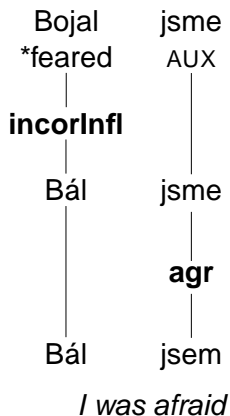


I was afraid

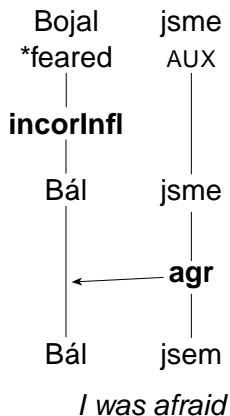
Annotation of a sample sentence, part I



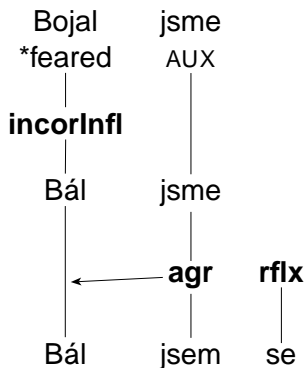
Annotation of a sample sentence, part I



Annotation of a sample sentence, part I

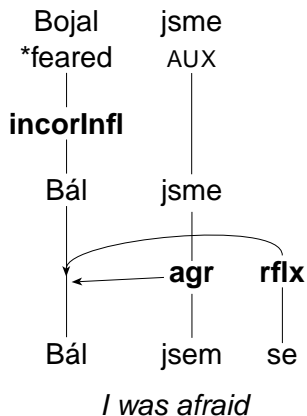


Annotation of a sample sentence, part I

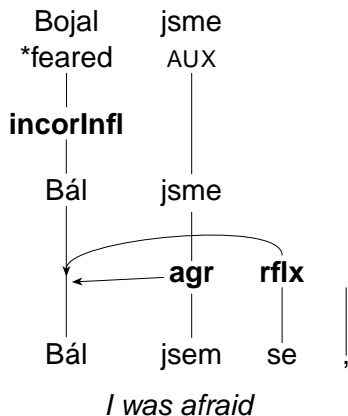


I was afraid

Annotation of a sample sentence, part I



Annotation of a sample sentence, part I



Annotation of a sample sentence, part II

že ona se ne bude libila slavnou prahu ,
that she RFL not will *like famous *prague ,

that she would not like the famous city of Prague,

Annotation of a sample sentence, part II

že	ona	se	ne	bude	libila	slavnou	prahu	,
that	she	RFL	not	will	*like	famous	*prague	,

že

that she would not like the famous city of Prague,

Annotation of a sample sentence, part II

že	ona	se	ne	bude	libila	slavnou	prahu	,
that	she	RFL	not	will	*like	famous	*prague	,
že	ona							

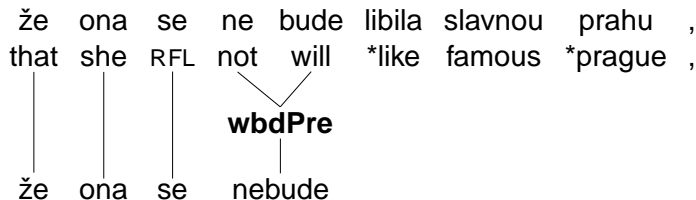
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II

že	ona	se	ne	bude	libila	slavnou	prahu	,
that	she	RFL	not	will	*like	famous	*prague	,
že	ona	se						

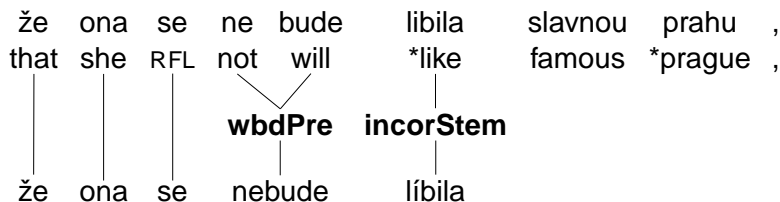
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



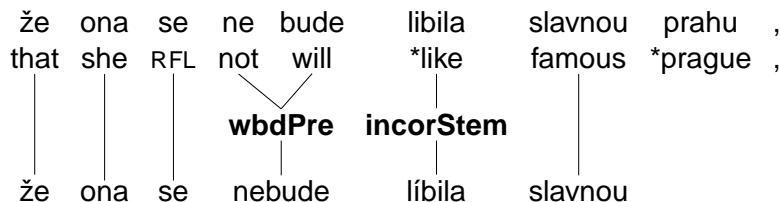
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



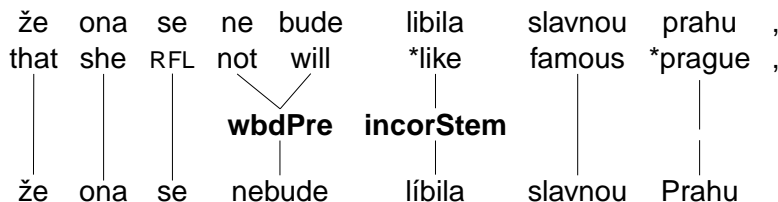
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



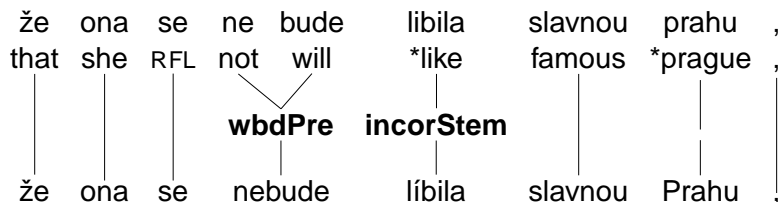
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



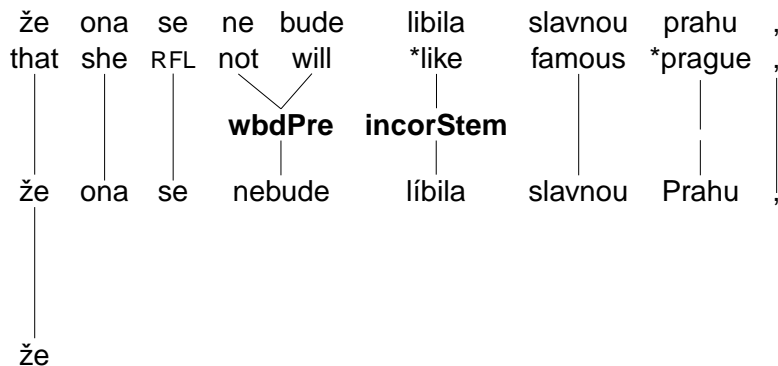
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



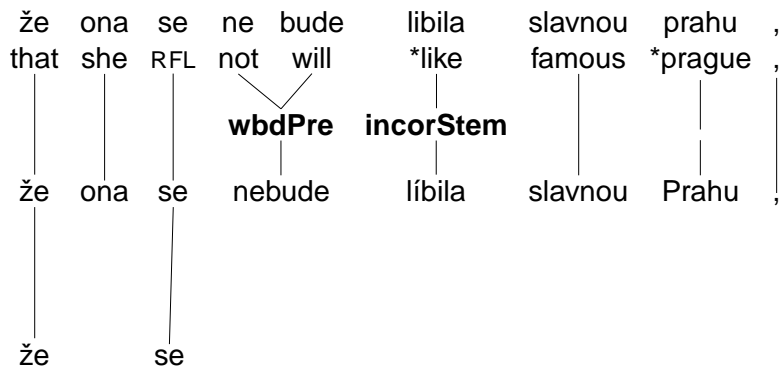
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



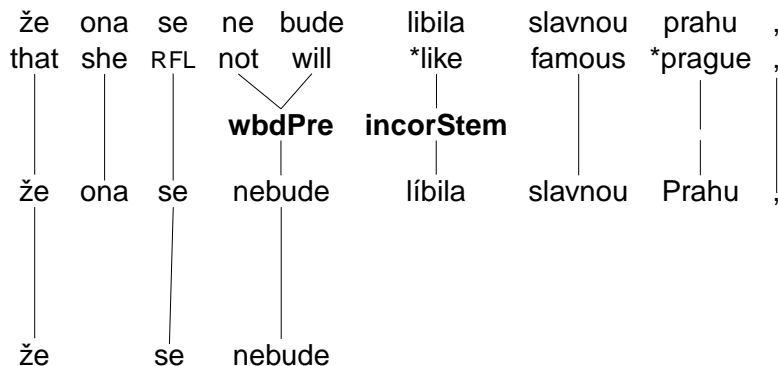
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



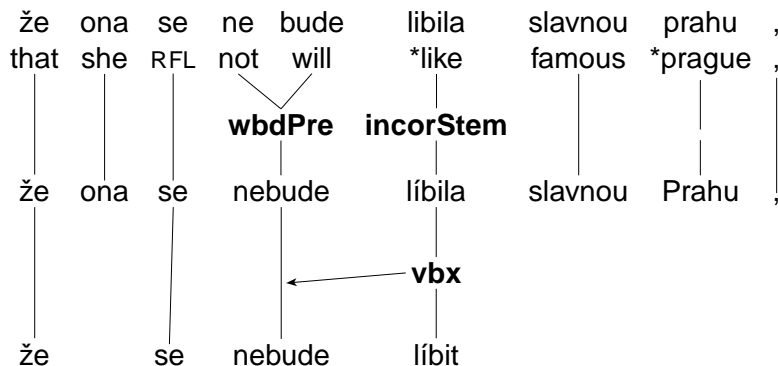
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



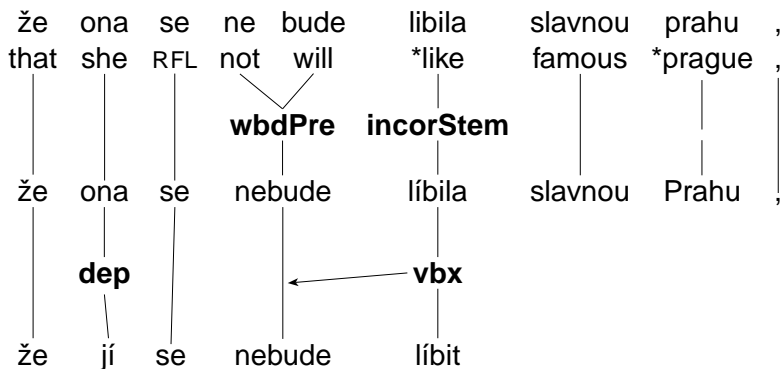
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



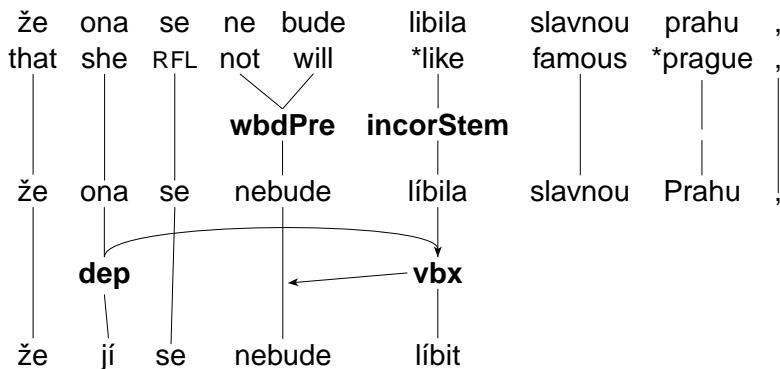
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



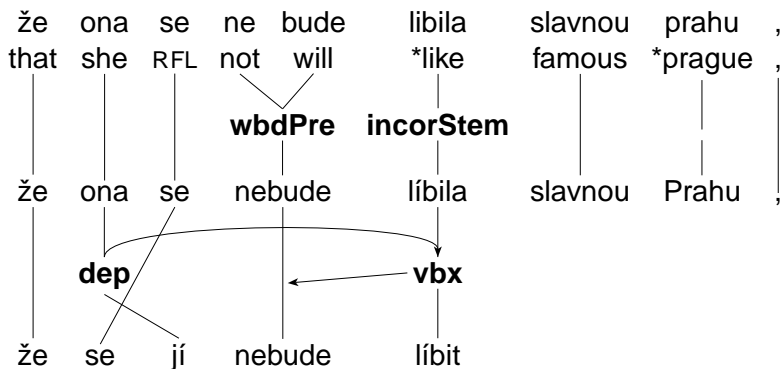
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



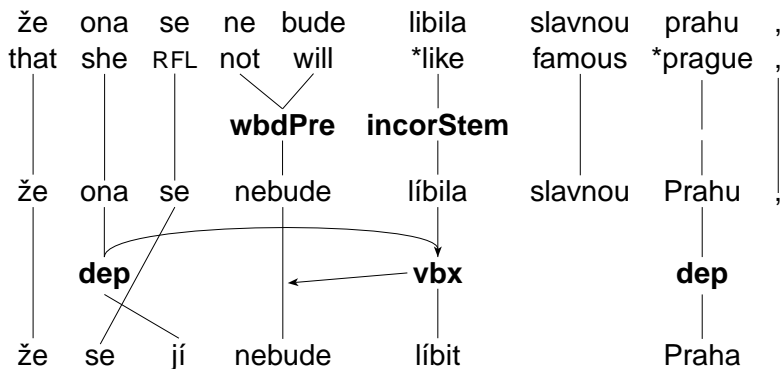
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



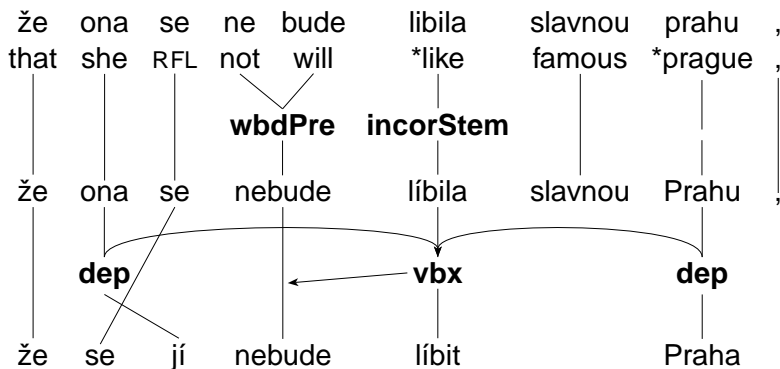
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



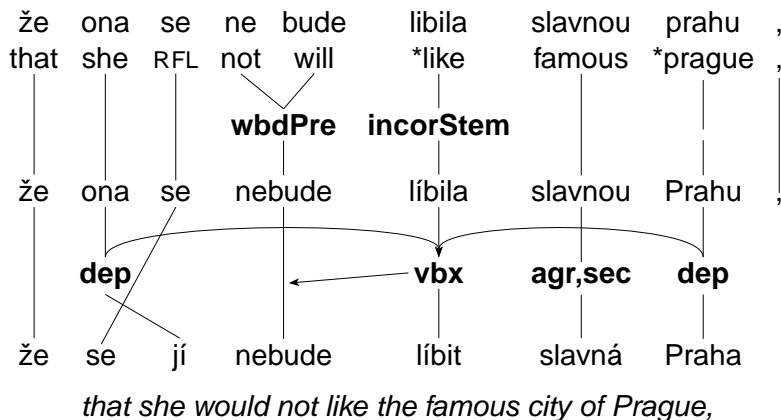
that she would not like the famous city of Prague,

Annotation of a sample sentence, part II

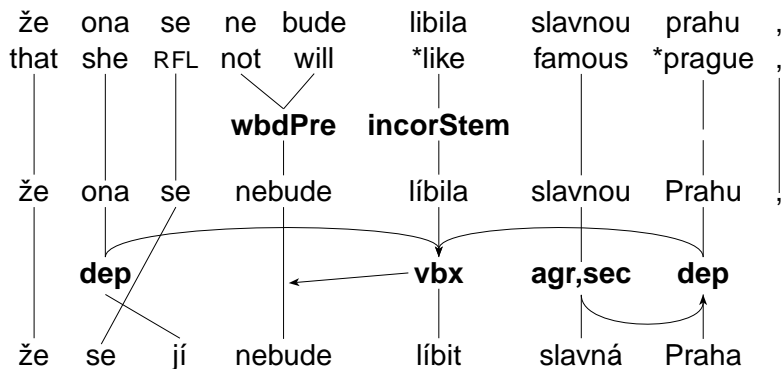


that she would not like the famous city of Prague,

Annotation of a sample sentence, part II

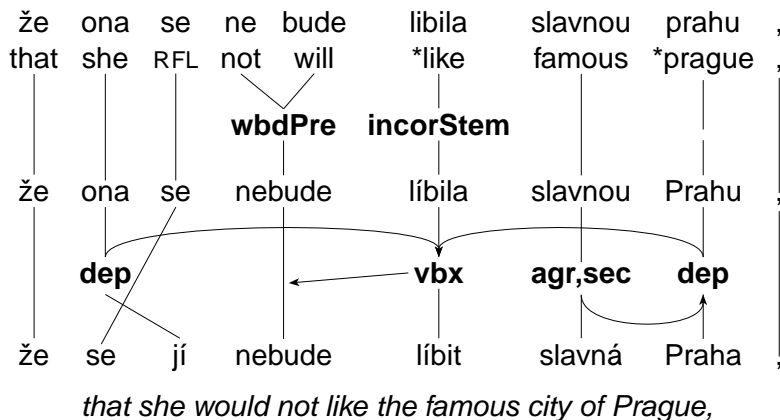


Annotation of a sample sentence, part II



that she would not like the famous city of Prague,

Annotation of a sample sentence, part II



Annotation of a sample sentence, part II

proto	to	bylo	velmí	vadí	pro	mně	.
therefore	it	was	*very	resent	for	me	.

because I would be very unhappy about it.

Annotation of a sample sentence, part II

proto to bylo velmi vadí pro mně .
therefore it was *very resent for me .

proto

because I would be very unhappy about it.

Annotation of a sample sentence, part II

proto	to	bylo	velmí	vadí	pro	mně	.
therefore	it	was	*very	resent	for	me	.
proto	to						

because I would be very unhappy about it.

Annotation of a sample sentence, part II

proto	to	bylo	velmí	vadí	pro	mně	.
therefore	it	was	*very	resent	for	me	.
proto	to	bylo					

because I would be very unhappy about it.

Annotation of a sample sentence, part II

proto	to	bylo	velmí	vadí	pro	mně .
therefore	it	was	*very	resent	for	me .
proto	to	bylo	incorStem			
proto	to	bylo	velmi			

because I would be very unhappy about it.

Annotation of a sample sentence, part II

proto	to	bylo	velmí	vadí	pro	mně .
therefore	it	was	*very	resent	for	me .
proto	to	bylo	incorStem	vadí		
proto	to	bylo	velmi	vadí		

because I would be very unhappy about it.

Annotation of a sample sentence, part II

proto	to	bylo	velmí	vadí	pro	mně .
therefore	it	was	*very	resent	for	me .
proto	to	bylo	incorStem	vadí	pro	
proto	to	bylo	velmi	vadí	pro	

because I would be very unhappy about it.

Annotation of a sample sentence, part II

proto	to	bylo	velmí	vadí	pro	mně .
therefore	it	was	*very	resent	for	me .
			incorStem			
proto	to	bylo	velmi	vadí	pro	mně

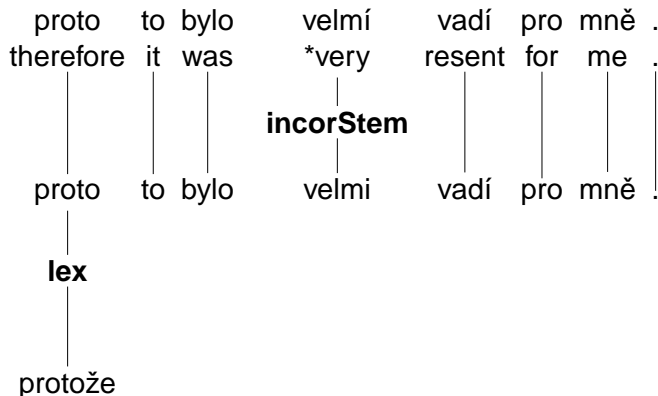
because I would be very unhappy about it.

Annotation of a sample sentence, part II

proto	to	bylo	velmí	vadí	pro	mně .
therefore	it	was	*very	resent	for	me .
			incorStem			
proto	to	bylo	velmi	vadí	pro	mně !

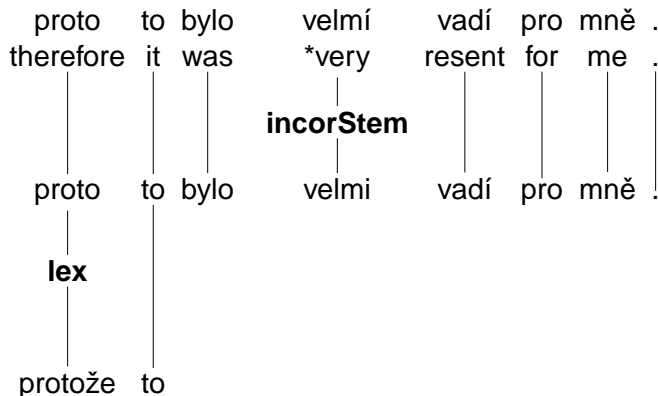
because I would be very unhappy about it.

Annotation of a sample sentence, part II



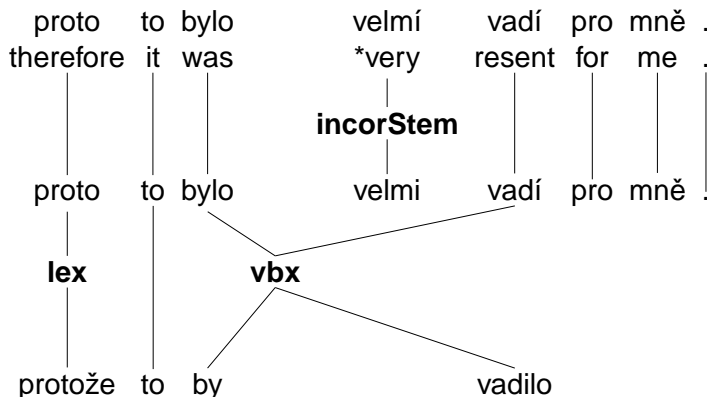
because I would be very unhappy about it.

Annotation of a sample sentence, part II

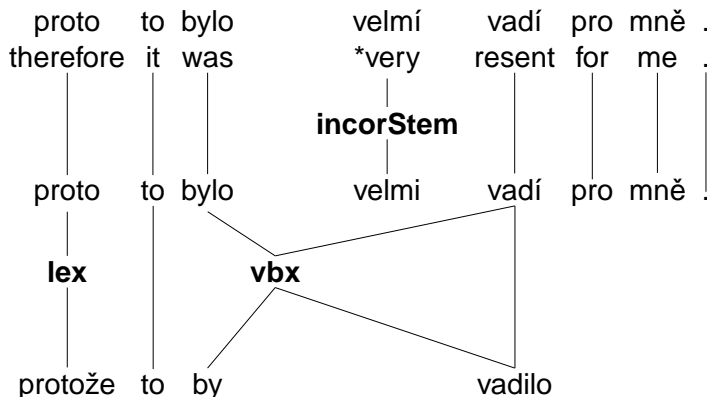


because I would be very unhappy about it.

Annotation of a sample sentence, part II

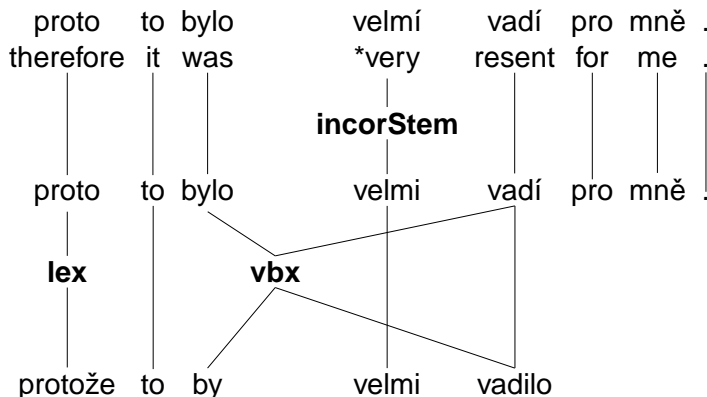


Annotation of a sample sentence, part II



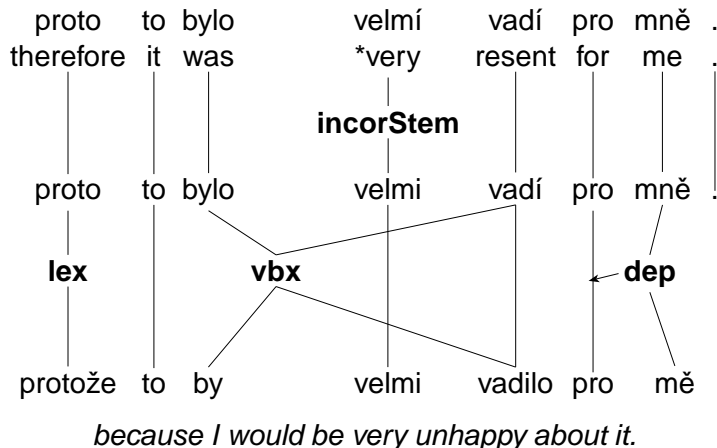
because I would be very unhappy about it.

Annotation of a sample sentence, part II

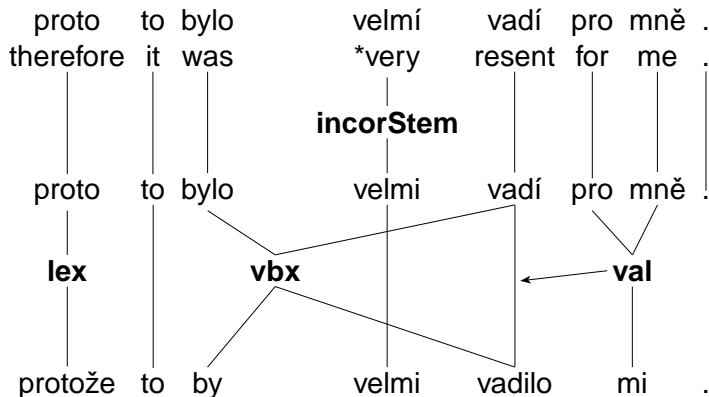


because I would be very unhappy about it.

Annotation of a sample sentence, part II

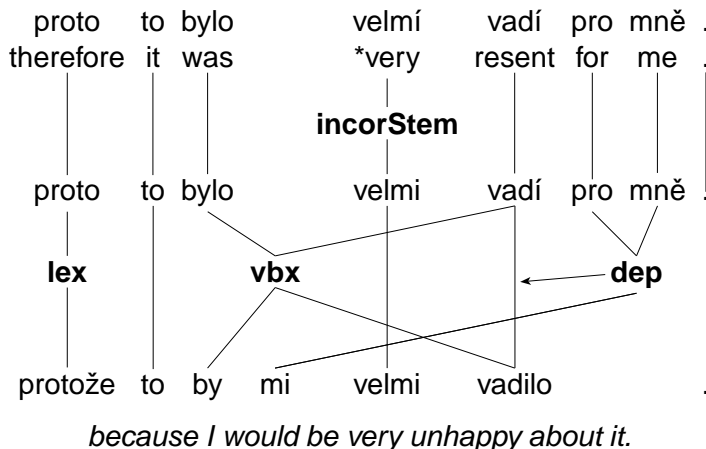


Annotation of a sample sentence, part II

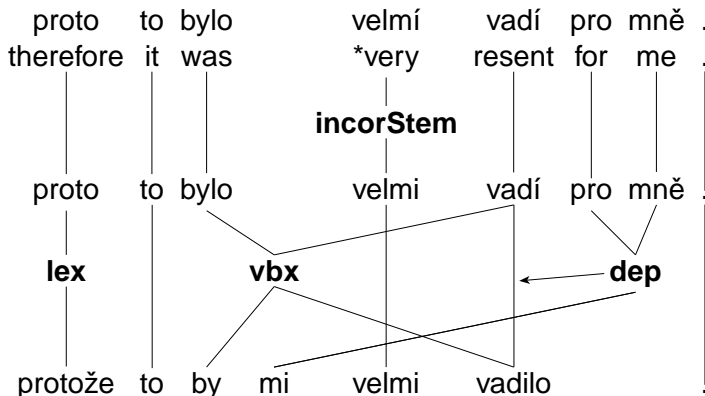


because I would be very unhappy about it.

Annotation of a sample sentence, part II



Annotation of a sample sentence, part II



because I would be very unhappy about it.

Types of errors at Level 1

Error type	Tag	Links	Assignment
Word boundary	bnd	m:n	Auto
Punctuation	p	0:1, 1:0	Auto
Capitalisation	cap	1:1	Auto
Diacritics	dia	1:1	Auto
Character(s)	char	1:1	Auto
Inflection	infl	1:1	Auto
Unknown lexeme	unk	1:1	Manual
Foreign word	fw	1:1	Manual

Types of errors at Level 2

Error type	Tag	Links	Ref	Assignment
Agreement	agr	1:1	1	Manual
Valency	val	1:1	1	Manual
Pronominal reference	ref	1:1	1	Manual
Complex verb forms	cvf	m:n	0,1	Manual
Negation	neg	m:n	0,1	Manual
Missing constituent	miss	0:1	0	Manual
Odd constituent	odd	1:0	0	Manual
Modality	mod	1:1	0	Manual
Word order	wo	m:n	0	Manual
Lexis & phraseology	lex	m:n	0,1	Manual

Annotation policy

Minimal intervention: corrected text need not be perfect, grammatical is enough

To do

We still need to provide annotators with guidelines on how to:

- handle uncertainty about the author's intended meaning,
- identify false-friends errors,
- handle colloquial language.

Data format

- Prague Markup Language
(PML, used in *Prague Dependency Treebank*)
- Generic, XML-based, for rich layered annotation
- A higher level contains information about words on that level, about errors and about relations to tokens on lower levels
- Portion of Level 1 of the sample sentence encoded in the PML data format – see next slide

```

<?xml version="1.0" encoding="UTF-8"?>
<adata xmlns="http://utkl.cuni.cz/czesl/">
  <head>
    <schema href="adata_schema.xml" />
    <references>
      <reffile id="w" name="wdata" href="r049.w.xml" />
    </references>
  </head>
  <doc id="a-r049-d1" lowerdoc.rf="w#w-r049-d1">
    ...
    <para id="a-r049-d1p2" lowerpara.rf="w#w-r049-d1p2">
      ...
      <s id="a-r049-d1p2s5">
        <w id="a-r049-d1p2w50">
          <token>Bál</token>
        </w>
        <w id="a-r049-d1p2w51">
          <token>jsem</token>
        </w>
        ...
      </s>
      ...
      <edge id="a-r049-d1p2e54">
        <from>w#w-r049-d1p2w46</from>
        <to>a-r049-d1p2w50</to>
        <error> <tag>unk</tag> </error>
      </edge>
      <edge id="a-r049-d1p2e55">
        <from>w#w-r049-d1p2w47</from>
        <to>a-r049-d1p2w51</to>
      </edge>
      ...
    </para>
    ...
  </doc>
</adata>

```

Outline of the talk

- 1 Introduction
- 2 Annotation scheme
- 3 Annotation process**
- 4 Conclusion

The annotation workflow

- 1 A handwritten document is transcribed into HTML using off-the-shelf tools.
- 2 The information in the html document is used to generate Level 0 and a default Level 1 encoded in the PML format.
- 3 An annotator manually corrects the document and provides some information about errors using our annotation tool.
- 4 Error information that can be inferred automatically is added.
- 5 See next slide for a sample sentence in the annotation tool.

feat 201006101454

File View Tools Window Help

- Properties x ST_Randyskova_Vob_KA_049.b x

2 / 4 Add Layer R Export Spacing X: 50 Y: 100

Bojal	jsem	se	že	ona	se	ne	bude	líbit	prahu	,	proto	to	bylo	velmi	vadí	pro	mně	Česka
unk	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X
Bál	jsem	se	že	ona	se	bude	líbit	Prahu	,	proto	to	bylo	velmi	vadí	pro	mně	Česka	
X	X	X	X	X	val	wo	X	val	X	lex	X	cvř	X	X	wo	val	X	
Bál	jsem	se	že	se	ji	bude	líbit	Praha	,	protože	to	by	mi	velmi	vadilo			

Proč mám/nemám rád (Č)ěskou republiku?

Už se nacházím v česce republice až půl roku. toho mě musilo by stačit, abych rozuměl, mám rád to země nebo ne rád. tedko mužů určitě řeknu, že českou republiku ja miluju. tento země ma všechna že potřebuju ja a moje přítelkyně. Bojal jsem se že ona se ne bude líbit **prahu**, proto to bylo velmi vadí pro mně. Česka republika je krásne místo. tady je hodne hezké pamatek. například pražský hrad a vysehrad. líbim se moc pražský hrad, protože tam je zámky, který velmi krásne a hezke. take v čechach je dobra příroda a když jsme se procházili na divoke šarce byli šokováni o4 z tech krásnych pohledů. Je to nekrasneší místo ve všem bílém světě. take rad že Češi ie dobri

Fit WFit Orig Zoom

miluju. tento země ma všechna
ja a moje přítelkyně. Bojal jsem se
že líbit prahu, proto to bylo velmi
Česka republika je krásne místo.
hezke pamatek, například pražský
líbim se moc pražský hrad, proto

	0	1	2	3	4	5	6	7	8
L0	proto	to	bylo	velmí	vadí	promně	.		
gloss	<i>therefore</i>	<i>it</i>	<i>was</i>	<i>*very</i>	<i>resent</i>	<i>for</i>	<i>me</i>	.	
errid				dia					
L1	proto	to	bylo	velmi	vadí	promně	.		
errid	lex		cvf		2	val 4			
L2	protože	to	by	velmi	vadilo	mi	.		
errid				wo					
L3	protože	to	by	mi velmi vadilo				.	

Done.

Postprocessing

Manual annotation is followed by automatic post-processing, providing the corpus with additional information:

- ➊ Level 1: lemma and morphosyntactic tags (not disambiguated)
- ➋ Level 2: lemma and morphosyntactic tags (disambiguated)
- ➌ Level 1: type of error (by comparing the original and corrected strings) (e.g. **libit – líbit* ‘like’ – error in diacritics)
- ➍ Level 2: type of morphosyntactic errors caused by agreement or subcategorisation error (by comparing morphosyntactic tags at Level 1 and 2)
- ➎ Formal error description: missing/extra expression, wrong order
- ➏ In the future, we plan to automatically tag errors in verb prefixes, inflectional endings, spelling, palatalisation, metathesis, etc.

Outline of the talk

- 1 Introduction
- 2 Annotation scheme
- 3 Annotation process
- 4 Conclusion**

Conclusion

- Error annotation is a very resource-intensive task,
- But an error-tagged corpus is an invaluable tool:
 - ▶ to obtain a reliable picture of the learners' interlanguage and
 - ▶ to adapt teaching methods and learning materials.

Acknowledgments

Thanks to

other members of the project team, namely Karel Šebesta, Milena Hnátková, Tomáš Jelínek, Vladimír Petkevič, and Hana Skoumalová

