

# A learner corpus of Czech: Current state and future directions

Barbora Štindlová • Svatava Škodová  
Technical University of Liberec, Faculty of Education

Alexandr Rosen • Jirka Hana  
Charles University Prague, Faculty of Arts • Faculty of Mathematics and Physics

## Abstract

The paper describes *CzeSL*, a learner corpus of Czech, together with its design properties. We start with a brief introduction of the project within the context of *AKCES*, a program addressing Czech acquisition corpora; in connection with the programme we are also concerned with the groups of respondents, including differences due to their L1; further we comment on the choice of the sociocultural metadata recorded with each text and related both to the learner and the text production task. Next we describe the intended uses of *CzeSL*. The core of the paper deals with transcription and annotation. We explain issues involved in the transcription of handwritten texts and present the concept of a multi-level annotation scheme including a taxonomy of captured errors. We conclude by mentioning results from an evaluation of the error annotation and presenting plans for future research.

**Keywords:** learner corpus, Slavic languages, Czech, error annotation, error taxonomy, multi-level annotation.

## 1. Introduction – A learner corpus of Czech

The first learner corpus of *Czech as a Second Language (CzeSL)*,<sup>1</sup> with its size of 2 million words the only large learner corpus for a Slavic language at the moment,<sup>2</sup> is built as a joint project of Technical University Liberec and Charles University Prague. It is a part of the programme *Acquisition Corpora of Czech (AKCES)*, pursued at Charles University in Prague since 2005 (Šebesta 2010). In addition to *CzeSL*, *AKCES* includes the following subcorpora: (i) *SCHOLA 2010*<sup>3</sup> and *EDUCO* – recordings and transcripts capturing the language of Czech pupils attending primary school classes (about 800,000 words each, finished); (ii) *SKRIPT* – written texts produced by Czech

---

<sup>1</sup> The corpus is one of the tasks of the project Innovation of Education in the Field of Czech as a Second Language (project no. CZ.1.07/2.2.00/07.0259), a part of the operational programme Education for Competiveness, funded by the European Structural Funds (ESF) and the Czech government. The annotation tool was also partially funded by grant no. P406/10/P328 of the Grant Agency of the Czech Republic.

<sup>2</sup> To the best of our knowledge, there is only one learner corpus built for a Slavic language – *PiKUST* (Stritar 2009), a corpus of Slovene as a foreign language. However, it is of a modest size of 35,000 words, and its error annotation is adopted from the Norwegian project *ASK* (Tenfjord 2009).

<sup>3</sup> More details and a search interface are available at <http://ucnk.ff.cuni.cz/schola.php>.

students (about 600,000 words so far, in development); (iii) *ROMi* – texts and speech produced by young learners with Romani background (in development); and (iv) *IUVENT* – spoken corpus of language produced by young native Czechs (planned). A consistent methodology and set of tools used throughout the program represent a significant synergic effect, allowing for comparative analyses of native and non-native language in the corpora, built on identical principles.<sup>4</sup>

*CzeSL* is focused on four main groups of learners of Czech: (1) speakers of related Slavic languages, represented mainly by Russian, other Eastern Slavic languages and Polish, (2) speakers of other Indo-European languages, with a slight majority of German, (3) speakers of distant non-Indo-European languages, mainly Chinese, Vietnamese and Arabic, and (4) pupils in primary school age with Romani background.<sup>5</sup> About one half (1 million words) of the corpus consists of short essays written by non-native learners (1–3), while short essays written by Romani pupils (4) account for 25%. Theses for a university degree, written by non-native students, represent the remaining 25% of the corpus. Approximately 20% (300 thousand words) of the short essays are corrected and error-annotated.

In addition to its concern with several representative groups of speakers, *CzeSL* strives to cover as much ground as possible also in other aspects. This wide-scope design property, offering extensive data with rich annotation and metadata, is meant to serve a number of different users and to satisfy varied research requirements: (a) *CzeSL* consists of both spoken and written texts, produced during a range of situations throughout the language learning process, collected as manuscripts and transcribed into an electronic format. The transcription follows rules designed to preserve important features of handwritten texts such as self-corrections (see Štindlová 2011: 106). Apart from originally handwritten texts, *CzeSL* also includes Bachelors', Masters' and doctoral theses, written in Czech by non-native students and collected in an electronic format. (b) The data cover all language levels according to the Common European Framework of Reference for Languages (CEFR, 2001), from real beginners (level A1) to more advanced learners (level B2 and higher), with a balanced mix of levels as much as possible, although levels B1 and B2 prevail over the lower grades due to their easier availability. (c) The texts are elicited in various situations; they are not restricted to parts of written or oral examination, or to argumentative or reflective essays, as in many other learner corpora. (d) *CzeSL* also includes texts, collected at regular intervals from learners attending long-term language courses. All texts of a specific author can be retrieved using metadata (see below). This will support analyses

---

<sup>4</sup> In this paper we use the term *second language* as denoting any language learned after first/native language or mother tongue (L1). Thus we use it as a hypernym of *foreign language*, a second language one learns outside of the environment the language is spoken. Some authors (e.g. Ellis 1994) use the terms second and foreign language in the same way as we do here, while others (e.g. Günther & Günther 2007) use them as complementary. *CzeSL* includes data from non-native residents of the Czech Republic, including those staying for a relatively short period (e.g. 1 year), and also from students of Czech abroad. This information is included in the metadata.

<sup>5</sup> Sometimes it is difficult to decide whether Czech is the first or second language of these children. Yet the sociocultural differences between the non-Roma and some Roma communities in the Czech Republic are such that the linguistic development of Roma children may show some traits of L2 acquisition. Because their linguistic integration represents a significant issue in the country's education system, this part of the *CzeSL* corpus will become a separate component of *AKCES*.

of temporal development of the author's interlanguage, providing the option of using some parts of *CzeSL* for longitudinal research.

Each text is equipped with detailed metadata records, for a total of 18 parameters. Some of them (12) relate to the respondent, while the remaining 6 specify the character of the text and circumstances of its production. All texts in the corpus produced by non-native speakers of Czech are assigned basic sociological data about the learners, such as age, gender, and language background (first language). Other obligatory variables describe (1) their proficiency level in Czech according to the CEFR, (2) conditions of the process of acquisition of Czech, including an indication of the institution, duration, or location – whether abroad or in the Czech Republic, (3) the textbooks used in learning Czech. More variables may be specified as an option,<sup>6</sup> such as the learner's knowledge of other (non-native) languages, her bilingual competence, length of stay in the country, or whether a family member has been a Czech speaker. In addition to standard metadata specifying temporal and size restrictions, we also register the availability of language reference tools and the extent and type of elicitation.

In Section 2, we sketch the intended use of the corpus, proceed to issues involved in the transcription of handwritten texts (Section 3) and argue for the specific design of our annotation scheme in Section 4. The concluding section, Section 5, deals with the error taxonomy.<sup>7</sup>

## 2. Intended use

Despite the fact that teaching Czech as a second language has acquired the status of a well-established field with a long tradition, a proper teaching methodology is not developed and available. Teachers often cope with this situation by adopting two possible strategies: (i) They have recourse to methods and techniques used for other languages, such as English or German. Since Czech is typologically different, mainly due to its rich morphology, this approach is grossly inadequate for training students in the use of rules of Czech grammar. (ii) They tend to transfer their detailed and fairly academic knowledge of Czech grammar to foreigners, in a way poorly structured for such a task, confusing the issues of presenting grammar to native and foreign students.

A specific problem is the issue of educating children with a native language other than Czech, whose presence at Czech primary schools is a recent phenomenon. Primary school teachers receive no training in teaching Czech as a foreign language, again resorting to an individual and intuitive approach. By its inclusion in *AKCES*, *CzeSL* will become a resource for research and design of teaching materials assisting teachers of young non-native speakers at different stages of the acquisition of Czech. At the same time, *CzeSL* should provide representative data that would help initiate and develop systematic and comprehensive research of Czech as a foreign language (so far, there are no monographs available dealing with this topic).

---

<sup>6</sup> The role of such variables has been emphasized by several authors (e.g. Granger 2003; 2008; Tono 3003).

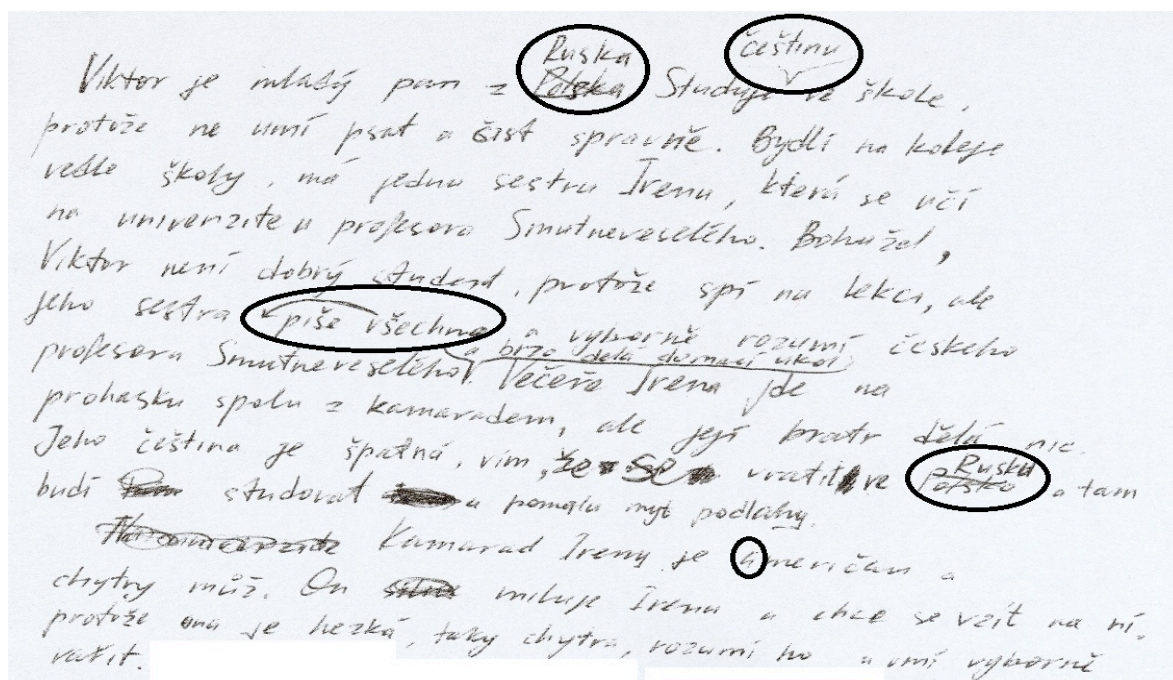
<sup>7</sup> For more details about some technical aspects of the compilation of the *CzeSL* corpus see Hana *et al.* (2012) and Jelínek *et al.* (2012).

The programme *Czech as a foreign language* has been available only recently as a three-year BA course at Technical University in Liberec (TU), and as a two-year MA course at Charles University in Prague (CU). Texts collected for *CzeSL* are already in use in the training of teachers both at TU and CU to give them an idea about the traits of the learner language in relation to the author's L1 and proficiency. This should help them to change perspective from viewing the language as an abstract system to approaching Czech as a sum of components acquired by learners at a specific stage of the development of their interlanguage.

In the following, the design and structure of the corpus is presented in more detail.

### 3. Transcription

Since most original texts are handwritten,<sup>8</sup> they are transcribed according to detailed rules using off-the-shelf tools (e.g., Open Office Writer or Microsoft Word). A set of codes is used to capture the author's corrections and other properties of the manuscript (e.g., for future research of handwriting of students with a different native writing system, for investigating the process of language acquisition, or to enable multiple interpretation – the same glyph may be interpreted as *i* in the handwriting of one student, *e* of another, and *a* of yet another). An additional reason for collecting handwritten texts is to avoid the use of a spell checker, because the result would not reflect the student's skills. In a highly inflectional language such as Czech, deviations in spelling very often do not only reflect wrong graphemics, but also indicate errors in morphology. In Figure 1 a sample text is presented with some of the author's self-corrections in circles.<sup>9</sup>



<sup>8</sup> Electronic texts (BA, MA and Ph.D. theses) represent a minority.

<sup>9</sup> For more details about the rules of transcription see Štindlová (2011).

*Figure 1: A sample handwritten text*

The handwritten text is transcribed as in Figure 2. Boldface highlights specific features of handwriting – self-corrections (deletion *Polska*, insertion {*češtinu*}<in>, text movement {*piše všechno -> všechno piše*}), and transcription codes in angle brackets. For example, the author replaced the form *Polska* ‘Poland’ by *Ruska* ‘Russia’, which is transcribed as strikeout text.

Viktor je mladý pan z **Polska** Ruska. Studuje {*češtinu*}<in> ve škole, protože ne umí psát a číst správně. Bydlí na koleje vedle školy, má jednu sestru Irenu, která se učí na univerzite u profesora Smutneveselého. Bohužel, Viktor není dobrý student, protože spí na lekci, ale jeho sestra {*piše všechno -> všechno piše*} a výborně rozumí českého profesora Smutneveselého {*a brzo delá domácí ukol*}<in>. Večeře Irena jde na prohasku spolu z kamaradem, ale její bratr dělá nic. Jeho čeština je špatná, vím, že se vrátit ve ~~Polsko~~ Rusko 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 ní, protože ona je hezká, taky chytra, rozumí ho a umí výborně vařit.

*Figure 2: The sample transcribed*

#### 4. Annotation scheme

The language of a learner of Czech may deviate from the standard in a number of aspects at the same time: spelling, morphology, morphosyntax, semantics, pragmatics or style. To show some of the options, the transcribed example in Figure 2 is shown again in Figure 3 with the spelling, morphological and morphosyntactic problems identified and the handwriting-specific annotation resolved. Forms wrong in any context (due to an error in spelling or morphology) are set in boldface, forms wrong due to a morphosyntactic or lexical anomaly are underlined. Some forms may be faulty for both reasons; these are in bold and underlined.

Viktor je mladý pan z Ruska. Studuje češtinu ve škole, protože ne umí psát a číst správně. Bydlí na koleje vedle školy, má jednu sestru Irenu, která se učí na univerzite u profesora Smutneveselého. Bohužel, Viktor není dobrý student, protože spí na lekci, ale jeho sestra všechno **piše** a **výborně** rozumí českeho profesora Smutneveselého a brzo delá domácí ukol. Večeře Irena jde na **prohasku** spolu z kamaradem, ale její bratr dělá nic. Jeho čeština je špatná, vím, že se vrátit ve Rusku a tam budí studovat u pomalu myt podlahy. **Kamarad** Ireny je **Američan** a **chytry muž**. On miluje Irenu a chce se vzít na ní, protože ona je hezká, taky **chytra**, rozumí ho a umí **výborně** vařit.

*Figure 3: Spelling, morphological and morphosyntactic problems identified*

The highlighted and underlined parts are incorrect. The transcribed version (Figure 3) and a corrected version are shown in Table 1. Corrected characters in italics and

glossed translations are included for the benefit of the reader; they are not a part of the corpus.

Viktor je mladý <u>pan</u> z Ruska. [Viktor is a young <u>Mr.</u> from Russia.]	Viktor je mladý <u>pán</u> z Ruska. [Viktor is a young <u>man</u> from Russia.]
Studuje češtinu ve škole, protože <u>ne umí psát</u> a <u>číst správně</u> . [He studies Czech at school, because he <u>can not write</u> and <u>read correctly</u> .]	Studuje češtinu ve škole, protože <u>neumí psát</u> a <u>číst správně</u> . [He studies Czech at school, because he <u>cannot write</u> and <u>read correctly</u> .]
Bydlí na <u>koleje</u> vedle školy, má jednu sestru Irenu, která se učí na <u>univerzite</u> u profesora <u>Smutneveselého</u> . [He lives at <u>residence halls</u> <sub>GEN</sub> next to the school, has one sister Irena, who is a student of professor <u>Smutněveselý</u> at the <u>university</u> .]	Bydlí na <u>koleji</u> vedle školy, má jednu sestru Irenu, která se učí na <u>univerzitě</u> u profesora <u>Smutněveselého</u> . [He lives at <u>residence halls</u> <sub>LOC</sub> next to the school, has one sister Irena, who is a student of professor <u>Smutněveselý</u> at the <u>university</u> .]
Bohužel, Viktor není dobrý student, protože spí na lekci, ale jeho sestra všechno <u>píše</u> a <u>výborně</u> rozumí <u>českeho</u> profesora <u>Smutneveselého</u> a brzo <u>delá</u> domácí <u>ukol</u> . [Unfortunately, Viktor is not a <u>good</u> student, because he sleeps in the class, but his sister <u>writes</u> everything and <u>perfectly</u> understands the <u>Czech</u> professor Smutněveselý and <u>does</u> her <u>homework</u> soon.]	Bohužel, Viktor není <u>dobrý</u> student, protože spí na lekci, ale jeho sestra všechno <u>píše</u> a <u>výborně</u> rozumí <u>českému</u> profesorovi <u>Smutněveselému</u> a brzo <u>dělá</u> domácí <u>úkoly</u> . [Unfortunately, Viktor is not a <u>good</u> student, because he sleeps in the class, but his sister <u>writes</u> everything and <u>perfectly</u> understands the <u>Czech</u> professor Smutněveselý and <u>does</u> her <u>homework</u> soon.]
<u>Večeře</u> Irena jde na <u>prohasku</u> spolu <u>z kamaradem</u> , ale její bratr <u>dělá</u> nic. [ <u>Dinner</u> Irena goes for a <u>walk</u> with her <u>friend</u> , but her brother <u>does</u> nothing.]	<u>Večer</u> Irena jde na <u>procházku</u> spolu <u>s kamarádem</u> , ale její bratr <u>nedělá</u> nic. [In the <u>evening</u> Irena goes for a <u>walk</u> with her <u>friend</u> , but her brother <u>doesn't do</u> anything.]
Jeho čeština je špatná, vím, že se <u>vratit</u> ve <u>Rusku</u> a tam <u>budí</u> studovat <u>u</u> pomalu <u>myt</u> podlahy. [?] [His Czech is poor, I know that he will to <u>to return</u> to <u>Russia</u> and there he <u>wakes</u> study <u>at</u> slowly <u>wash</u> floors.]	Jeho čeština je špatná, vím, že se <u>vrátí do</u> <u>Ruska</u> a tam <u>bude</u> studovat <u>a</u> pomalu <u>myt</u> podlahy. [?] [His Czech is poor, I know that he will <u>return</u> to <u>Russia</u> and there he <u>will</u> study <u>and</u> slowly <u>wash</u> floors.]
<u>Kamarad</u> Ireny je Američan a <u>chytry muž</u> . [Irena's <u>boyfriend</u> is an American and a <u>smart</u> guy.]	<u>Kamarád</u> Ireny je Američan a <u>chytrý muž</u> . [Irena's <u>boyfriend</u> is an American and a <u>smart</u> guy.]
On miluje Irenu a chce <u>se vzít na ní</u> , protože ona je hezká, taky <u>chytra</u> , rozumí <u>ho</u> a umí <u>výborně</u> vařit. He loves Irena and wants to <u>marry on her</u> , because she is pretty, also <u>smart</u> , she understands <u>him</u> <sub>GEN</sub> and is an <u>excellent</u> cook.	On miluje Irenu a chce <u>si ji vzít, protože</u> ona je hezká, taky <u>chytrá</u> , rozumí <u>mu</u> a umí <u>výborně</u> vařit." He loves Irena and wants to <u>marry her</u> , because she is pretty, also <u>smart</u> , she understands <u>him</u> <sub>DAT</sub> and is an <u>excellent</u> cook.

Table 1: The transcribed and the corrected versions with translations

To cope with the multi-level options of erring in Czech and to satisfy the goals of the project, our annotation scheme answers the following requirements:

1. Preservation of the original text alongside with the emendations
2. Successive emendations
3. Ability to capture errors in single forms as well as in multi-word discontinuous expressions
4. Syntactic relations as supplementary information for some error types: agreement, valency, pronominal reference
5. Automatic assignment of errors when possible, based on comparing faulty and corrected forms, using morphosyntactic tags, assigned by a tagger

To meet these requirements, we use a multilevel annotation scheme, supporting successive emendations. As a compromise between several theoretically motivated levels and practical concerns about the process of annotation, the scheme offers two annotation levels. This enables the annotators to register anomalies in isolated forms separately from the annotation of context-based phenomena but saves them from difficult theoretical dilemmas.

Level 0 is the level of transcribed input, where the words represent the original strings of graphemes, with some properties of the handwritten original preserved in the mark-up. Level 1 gives orthographical and morphological emendation of isolated forms as a text consisting of existing Czech forms; the sentence as a whole can still be incorrect. A formally correct form *weak* in a sentence such as *I'll see you in a weak* would be corrected since the author clearly misspelled the form she intended to use, creating an unintended homograph. On the other hand, the form *week* in *I'll see you in two week* is an error in morphosyntax and will be corrected at Level 2, where all other types of deviations are treated, resulting in a grammatically correct sentence. This includes errors such as those in syntax (agreement, government), lexicon, word order, usage, style, reference, or negation.

Levels of annotation are represented as a graph consisting of a set of interlinked parallel paths, where a path is a sequence of word forms corresponding to a sentence at a given level. Each word in the input text is represented at every level, unless it is split, joined, deleted or added by the annotator. Whenever a word form is emended, the type of error can label the link connecting the incorrect form with its emended version (such as *incorInfl* or *incorBase* for morphological errors in inflectional endings and stems). The sample text is shown in Figure 4 as displayed by the tool used by the annotators.<sup>10</sup>

The whole annotation process proceeds as follows:

1. The transcript is converted into the annotation format, where Level 0 roughly corresponds to the tokenized transcript and Level 1 is set as equal to Level 0 by

---

<sup>10</sup> The tool *feat* (*Flexible Error Annotation Tool*) is an environment for layered error annotation of learner corpora, see Hana *et al.* (2010). It is freely available from <http://purl.org/net/feat>.



default. Both are encoded as PML (an XML-based format for structural linguistic annotation, see Pajas & Štěpánek 2006).

2. The annotator manually corrects the document and provides some information about errors using our annotation tool *feat*.
3. Automatic post-processing provides additional information about lemma, part-of-speech and morphological categories for emended forms.
4. Error information that can be inferred automatically is added by comparing original and corrected strings: type of spelling alternation, missing/redundant expression, and inappropriate word order.

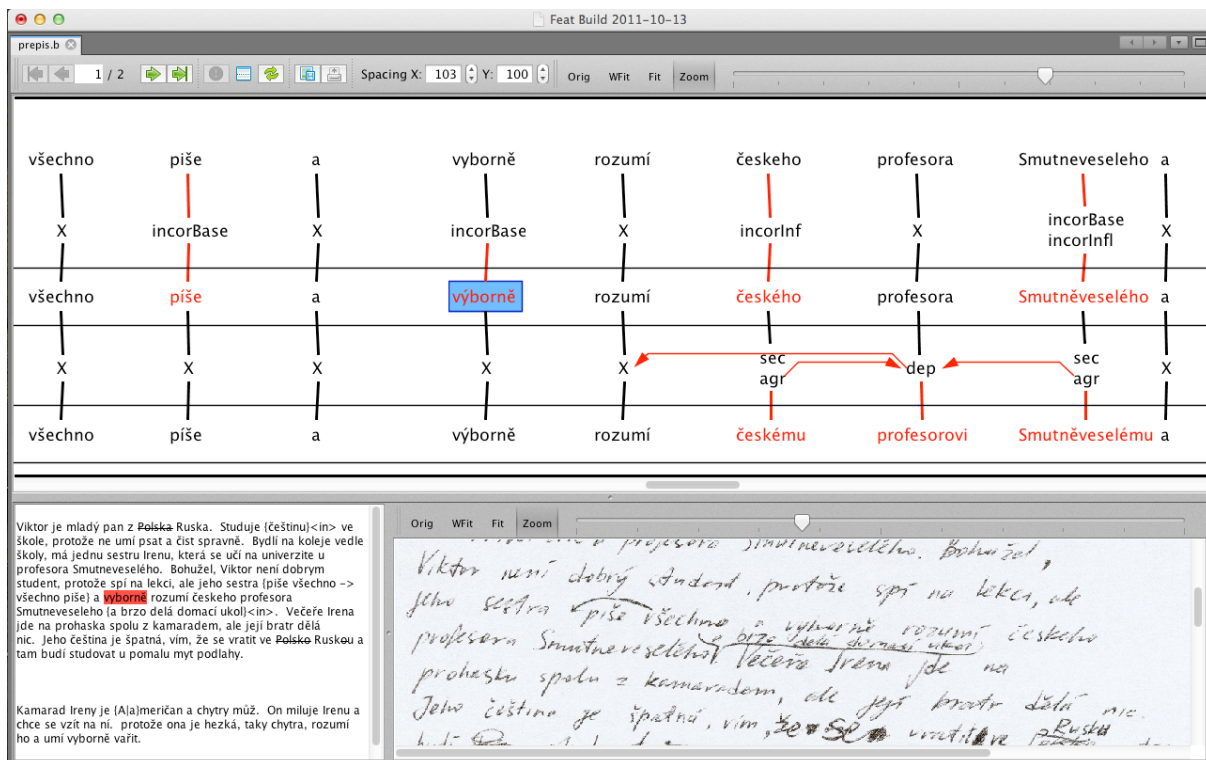


Figure 4: The sample in the annotation tool *feat*

## 5. Error taxonomy

The taxonomy of errors is based on previous research of frequent error types and reflects implicitly stated research hypotheses about the acquisition of an inflectional language.<sup>11</sup> To a large extent the taxonomy uses standard linguistic categories, complemented by a classification of superficial alternations of the source text, such as missing, redundant, faulty or incorrectly ordered element. The tagset consists of 22 error tags, 8 for Level 1, 11 for Level 2, and 3 that can be used at both levels. They are

<sup>11</sup> For some taxonomies used in previous projects see, e.g., Díaz-Negrillo *et al.* (2006), Nicholls (2003), Izumi *et al.* (2005), or Granger (2003a).



supplemented by tags generated automatically by comparing original with emended forms and manually assigned tags.

Errors in individual word forms, treated at Level 1 (see Table 2), include misspellings (also diacritics and capitalization), misplaced word boundaries but also errors in inflectional and derivational morphology and unknown stems – fabricated or foreign words. Except for misspellings, all these errors are annotated manually.

<b>Error type</b>	<b>Description</b>	<b>Example</b>
<i>incorInfl</i>	incorrect inflection	<i>pracovají</i> v továrně; bydlím s <i>matkoj</i>
<i>incorBase</i>	incorrect word base	lidé jsou moc <i>mérný</i> ; musíš to <i>posvětlit</i>
<i>fwFab</i>	non-emendable, “fabricated” word	pokud nechceš slyšet <i>smášky</i>
<i>fwNC</i>	foreign word	váza je na <i>Tisch</i> ; jsem v <i>truong</i>
<i>flex</i>	with <i>fwFab</i> and <i>fwNC</i> : inflected	jdu do <i>shopa</i>
<i>wbdPre</i>	word boundary: prefix or preposition	musím to <i>při</i> <i>pravít</i> ; <i>veškole</i>
<i>wbdComp</i>	word boundary: compound	<i>český</i> <i>anglický</i> slovník
<i>wbdOther</i>	other word boundary error	<i>mocdobře</i> ; <i>atak</i> ; <i>kdykoli</i>
<i>stylColl</i>	colloquial form	<i>dobrej</i> film
<i>stylOther</i>	bookish, dialectal, hypercorrect	holka s <i>hnědými</i> <i>očimi</i>
<i>problem</i>	problematic cases	

Table 2: Errors at Level 1

Emendations at Level 2 concern errors in agreement, valency, analytical forms, pronominal reference, negative concord, the choice of aspect, tense, lexical item or idiom, and also in word order. For the agreement, valency, analytical forms, pronominal reference and negative concord cases, there is usually a correct form, which determines some properties (morphological categories) of the faulty form. Table 3 gives a list of error types manually annotated at Level 2. The automatically identified errors include word order errors and subtypes of the error in analytical verb forms (*vbX*).

<b>Error type</b>	<b>Description</b>	<b>Example</b>
<i>agr</i>	violated agreement rules	to jsou <i>hezké</i> chlapci; Jana <i>čtu</i>
<i>dep</i>	error in valency	bojí se <i>pes</i> ; otázka <i>čas</i>
<i>ref</i>	error in pronominal reference	dal jsem to jemu i <i>jejího</i> bratrovi
<i>vbX</i>	error in analytical or compound verb form	musíš <i>přijdeš</i> ; kluci <i>jsou</i> běhali
<i>rflx</i>	error in reflexive expression	dívá na televizi; Pavel <i>si</i> raduje
<i>neg</i>	error in negation	žádný to <i>ví</i> ; <i>půjdu</i> <i>ne</i> do školy
<i>lex</i>	error in lexicon or phraseology	jsem <i>ruská</i> ; dopadlo to <i>přírodně</i>

<i>use</i>	error in the use of a grammar category	pošta je nejvíc <i>blízko</i>
<i>sec</i>	secondary error	stará se o <i>našich holčičkách</i>
<i>stylColl</i>	colloquial expression	viděli jsme <i>hezky</i> holky
<i>stylOther</i>	bookish, dialectal, hypercorrect	rozbil se mi <i>hadr</i>
<i>stylMark</i>	redundant discourse marker	<i>no; teda; jo</i>
<i>disr</i>	disrupted construction	<i>kratka jakost vyborné ženy</i>
<i>problem</i>	problematic cases	

Table 3: Errors at Level 2

Rather than aiming at perfect Czech, we emend the input conservatively, modifying incorrect and inappropriate forms and expressions to arrive at a coherent and well-formed result, without any ambition to produce a stylistically optimal solution. Word order, for example, is corrected only when the input is ungrammatical.

Overall, we are convinced that annotation guided by formal criteria is useful at least as a base for comparison with native speakers' language, automatic (error) annotation, and for annotating communicative adequacy, style, etc. in the future. As a further step towards a common ground for the comparison and guidance for the annotators, grammatical and lexical aspects of the learner language are emended and tagged to conform to the rules of Standard Czech.

A doubly annotated sample (10,000 word forms) was evaluated for inter-annotator agreement to verify that the annotation scheme and taxonomy are sufficiently robust to be used in the corpus. Higher agreement was found for formally well-defined error categories, with satisfactory results even for categories requiring subjective judgment. For more details see Štindlová *et al.* (2012).

## 6. Conclusion and outlook

Experience from teaching Czech as a foreign language clearly indicates the need for a rich source of data on the language of learners, one which would help to design optimal presentation of the Czech language for non-native speakers. A learner corpus is the answer also because the typological properties of Czech as a highly inflectional language make the use of experience from other, better positioned languages at least questionable. In this sense, Czech may serve as a testbed for the development of methods and tools targeting inflectional languages.

In order to support varied types of use and to maintain consistency of annotation, we opted for an annotation scheme and error taxonomy based on grammatical deviations from the standard, without one specific focus. This strategy fits well with the typological properties of Czech and allows for extensions now or in the future – both

into new domains of annotated phenomena, and into more efficient annotation processes, such as automatic assignment of more detailed error categories (implemented), automatic morphological analysis (implemented) and syntactic analysis or semi-automatic emendation and error tagging using a spell and grammar checker, integrated with the annotation tool (in preparation).

## References

- CEFR (2001). *Common European Framework of Reference for Languages: Learning, Teaching, Assessment*. Applied Linguistics Non Series. Cambridge University Press.
- Díaz-Negrillo, A. and Fernández-Domínguez, J. (2006). Error tagging systems for learner corpora. *Resla*, 19:83–102.
- Ellis, R. (1994). *The Study of Second Language Acquisition*. OUP, Oxford, 12.
- Granger, S. (2003). The International Corpus of Learner English: a new resource for foreign language learning and teaching and second language acquisition research. *TESOL Quarterly*. 37(3):538–545.
- Granger, S. (2003a). Error-tagged learner corpora and CALL: A promising synergy. *CALICO journal*, 20:465–480.
- Granger, S. (2008). Learner corpora. In Lüdeling, A. and Kytö, M., editors, *Corpus Linguistics. An International Handbook*. Mouton De Gruyter, Berlin/New York, 259–274.
- Günther, B. and Günther P. (2007). *Erstsprache, Zweitsprache, Fremdsprache: Eine Einführung*. Beltz Pädagogik, Belz.
- Hana, J., Rosen, A., Škodová, S. and Štindlová, B. (2010). Error-tagged learner corpus of Czech. In *Proceedings of the Fourth Linguistic Annotation Workshop*. Uppsala, Sweden: Association for Computational Linguistics, 11–19.
- Hana, J., Rosen, A., Štindlová, B., and Jäger, P. (2012). Building a learner corpus. In Calzolari, N. et al., editors, *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, Istanbul, Turkey. European Language Resources Association (ELRA).
- Günther, B., Günther P. (2004). *Erstsprache und Zweitsprache: Einführung aus pädagogischer Sicht*. Weinheim, Basel, Beltz.
- Izumi, E., Uchimoto, K., and Isahara, H. (2005). Error annotation for corpus of Japanese learner English. In *Proceedings of the Sixth International Workshop on Linguistically Interpreted Corpora (LINC 2005)*, Korea, 71–80.
- Jelínek, T., Štindlová, B., Rosen, A., and Hana, J. (2012). Combining manual and automatic annotation of a learner corpus. In *Proceedings of the 15th International Conference TSD 2012*. To appear.

- Nicholls, D. (2003). The Cambridge Learner Corpus: error coding and analysis for lexicography and ELT. *Proceedings of the Corpus Linguistics 2003 Conference, 28–31 March*. Lancaster, 572–581.
- Pajas, P. and Štěpánek, J. (2006). XML-Based Representation of Multi-Layered Annotation in the PDT 2.0. In R. E. Hinrichs, N. Ide, M. Palmer, J. Pustejovsky (eds.) *Proceedings of LREC 2006 Workshop on Merging and Layering Linguistic Information, Genoa, Italy: ELRA*, 40–47.
- Stritar, M. (2009). Slovene as a foreign language: The pilot learner corpus perspective. *Slovenski jezik / Slovene Linguistic Studies* 7, 135–152.
- Šebesta, K. (2010). Korpusy češtiny a osvojování jazyka [Corpora of Czech and language acquisition]. *Studie z aplikované lingvistiky/Studies in Applied Linguistics* 1, 11–34.
- Štindlová, B. (2011). *Evaluaace chybové anotace v žákovském korpusu češtiny [Evaluation of error mark-up in a learner corpus of Czech]*. Doctoral dissertation, Charles University, Faculty of Arts, Prague.
- Štindlová, B., Rosen, A., Hana, J. and Škodová, S. (2012). CzeSL – An error tagged corpus of Czech as a second language. In P. Pežik (ed.) *PALC 2011 – Practical Applications in Language and Computers (Łódź 13–15 April 2011), Łódź Studies in Language*. Peter Lang.
- Tenfjord, K., Hagen, J. E., and Johansen, H. (2009). Norsk andrespråkskorpus (ASK) – design og metodiske forutsetninger. *NOA norsk som andrespråk*, 25(1): 52–81.
- Tono, Y. (2003). Learner corpora: design, development and applications. In *Proceedings of the Corpus Linguistics 2003 Conference*. Lancaster, United Kingdom, 800–809.