

Analyzing learner language: towards a flexible natural language processing architecture for intelligent language tutors

Luiz Amaral^a, Detmar Meurers^{b*} and Ramon Ziai^b

^a*Department of Languages, Literatures and Cultures, University of Massachusetts, Amherst, MA, USA;* ^b*Seminar für Sprachwissenschaft, Universität Tübingen, Tübingen, Germany*

Intelligent language tutoring systems (ILTS) typically analyze learner input to diagnose learner language properties and provide individualized feedback. Despite a long history of ILTS research, such systems are virtually absent from real-life foreign language teaching (FLT). Taking a step toward more closely linking ILTS research to real-life FLT, in this article we investigate the connection between FLT activity design and the system architecture of an ILT system. We argue that a demand-driven, annotation-based natural language processing (NLP) architecture is well-suited to handle the demands posed by the heterogeneous learner input which results when supporting a wider range of FLT activity types. We illustrate how the unstructured information management architecture (UIMA) can be used in an ILTS, thereby connecting the specific needs of activities in foreign language teaching to the current research and development of NLP architectures in general. Making the conceptual issues concrete, we discuss the design and realization of a UIMA-based reimplementation of the NLP in the TAGARELA system, an intelligent web-based tutoring system supporting the teaching and learning of Portuguese.

Keywords: intelligent language tutoring systems (ILTS); intelligent computer-assisted language learning (ICALL); natural language processing (NLP); unstructured information management architecture (UIMA); demand-driven annotation-based architecture; individualized feedback

1. Introduction

In the context of computer-assisted language learning, intelligent language tutoring systems (ILTS) provide individualized feedback to learners working on activities. ILTS may also individually adjust the sequencing of instruction. Typically the focus of the analysis is on form errors made by the learner, even though in principle feedback can also target aspects of meaning or highlight correctly used forms.

While for some restricted exercises it is possible to anticipate all potential learner input and intended system responses, for most types of language learning activities such a direct mapping between potential learner input and feedback is not feasible (cf. Nagata, 2009). Instead, it is necessary to abstract from the specific string to more general classes by automatically analyzing the learner input using algorithms and

*Corresponding author. Email: detmar.meurers@uni-tuebingen.de

resources from Natural Language Processing (NLP). Generation of feedback can then be based on the information obtained through such NLP analysis.

Common to most current ILTS is that the NLP modules are integrated into a pipeline architecture (cf., e.g., Delmonte, 2003; Heift, 2003; Levin & Evans, 1995; Nagata, 2002; Rypa & Feuerman, 1995). The system calls the NLP modules in a pre-defined order, transforming one data structure into another and terminating when specific conditions are met, e.g., when the learner response matches a pre-stored target response, or when spell checking fails.

Such a pipeline architecture works well as long as the system deals with learner input from activity types that are uniform with respect to the required NLP processing. For example, in the E-Tutor (Heift, 2003) and Robo-Sensei (Nagata, 2002) the learner answers consist of single sentences and the lexical material to be used by the learner is constrained explicitly by listing the stems or implicitly by eliciting the student answers through translation (cf., Amaral & Meurers, 2011, section. 3.1). The NLP diagnosis and feedback generated apparently is the same for all activities.

A uniform pipeline architecture becomes problematic, however, when trying to integrate a wider range of activity types resulting in learner input of a heterogeneous nature, which potentially should also be evaluated using various criteria. Based on our experience from creating a range of activities for TAGARELA (<http://purl.org/ical/tagarela>), a web-based ILTS for the instruction of Portuguese as a foreign language, in this article we thus argue for a more flexible, demand-driven architecture for this type of system.

In such an architecture, the use and sequencing of the different NLP modules is triggered by demands for particular information based on the activity models for the different activity types. Each NLP module enriches the input with annotations until all information required to evaluate the learner's performance to provide feedback on a particular activity is present. In other words, the fixed algorithmic pipeline is replaced by whatever processing sequence is needed to obtain the particular information that is required by the activity model to provide feedback for a given exercise.

Relating this point to the broader NLP context, an ILTS can be viewed as an instance of an application required to deal with heterogeneous input and different information needs based on that input. Our approach thus is reminiscent of current approaches to information extraction, where, e.g., IBM's OmniFind makes use of the Unstructured Information Management Architecture (UIMA, Ferrucci & Lally, 2004) to obtain a range of annotations depending on the specific information needs.

In this article, we show how such an architecture can be realized for an ILTS. We discuss our reimplementations of the NLP in the TAGARELA system based on the UIMA architecture and showcase the benefits of this demand-driven, annotation-based architecture.

2. TAGARELA and the analysis it needs to perform

2.1. The NLP modules and what they are used for

The TAGARELA system makes use of a number of NLP modules. The form analysis includes a tokenizer which takes into account specifics of Portuguese such as cliticization, contractions, and abbreviations. Full-form lexical lookup returns all analyses based on the CURUPIRA lexicon (Martins, Nunes, & Hasegawa, 2003), which generally provides multiple analyses for each token. Finite state disambiguation rules are used to narrow down this lexical information based on where in a

sentence the token appears. This is similar in spirit to Constraint Grammar (Bick, 2000, 2004; Karlsson, Voutilainen, Heikkilä, & Anttila, 1995), where local disambiguation rules can be specified without the need to fully describe the language to be analyzed. A bottom-up chart parser is used to check agreement, case and some global well-formedness conditions. To analyze the meaning of the learner input, shallow semantic matching strategies between the student's input and target answers are used, in line with the Content Assessment Module proposed by Bailey and Meurers (2008).

2.2. Analyzing Portuguese learner data

To identify the types of errors the system has to handle, we collected a corpus of approximately 10,000 words from written assignments of students in an introductory Portuguese course at the college level, and we created a taxonomy of expected errors. Among the most common classes of errors in our corpus were *spelling* (24%), *agreement* (16%), *missing word* (12%), *extra word* (7.5%), and *word choice* (3.2%). Because beginners deal with very restricted types of constructions, it is not surprising that the most common error classes are *spelling* and *agreement*. The difficulties English speaking students encounter with subject–verb and nominal agreement in Romance languages are well known by instructors and researchers (cf., e.g., Koike & Klee, 2003; Montrul, 2004). Our corpus confirmed this and the insights gained by the analysis of the student data provided important empirical guidance for the creation of the NLP tools.

Besides confirming the most common errors in beginner performance, our corpus helped us define the error taxonomy to be used with our student population. Error taxonomies are an essential component of ILTS development. The classification of errors is used by the system to select the NLP tools needed to identify particular error types in student input, and it is important for prioritizing and formulating the appropriate feedback message. The error taxonomies used in ILTS development are mostly driven by the activities offered by the system and the specific properties of the target language.¹ For the more general task of creating error taxonomies for the annotation of learner corpora, a range of annotation schemes have been proposed (Díaz-Negrillo & Fernández-Domínguez, 2006), yet no consensus has been reached on which annotation scheme includes the necessary distinctions and supports identification with sufficiently high agreement (Meurers, 2009).

For the TAGARELA development, we thus developed our error taxonomy around the types of errors which arise for the beginning learners of Portuguese, as exemplified in the corpus of learner answers we collected.

The most common *agreement errors* we observed are between determiners and nouns, such as the gender agreement error in (1), followed subject–verb agreement errors in person or number as illustrated by (2). TAGARELA's parser is used to identify such errors.

- (1) Eu vou **na** cinema.
I go to the *fem*cinema_{*masc*}
- (2) Eu **trabalha** no jornal.
I_{*.sg*} work_{*3.sg*} at the newspaper

The system uses the shallow semantic matching modules to deal with errors classified as *missing words*, *extra words*, and *word choice*. Missing words range from typical

function words, such as the prepositions in (3), to lexical heads, such as the missing verb in (4); in the examples, the missing word is shown in bold in square brackets.

- (3) Nós começamos [**a**] falar com eles.
we started to speak to them
- (4) Eu [**tenho**] muito trabalho. Tchau! Obrigada!
I have much work goodbye thank you

The most common cases of *extra words* were extra articles, such as in (5). We also found cases of extra prepositions, complementizers, and pronouns, such as the clitic ‘se’ in (6); in the examples, the extra word is shown in bold.

- (5) Vocês **os** dois sempre querem a sobremesa.
you the two always want the dessert
Both of you always want dessert.
- (6) Eu me chamo-**se** John.
I myself call-oneself John
My name is John.

Word choice are errors that have their origin in false cognates or in false translations chosen from bilingual dictionaries. In example (7), the student translated ‘have a drink’ literally into ‘*ter* uma bebida’, even though in Portuguese the correct expression is ‘*tomar* uma bebida’ (to take a drink).

- (7) Eu pretendo ir ao clube e **ter** uma bebida.
I intend to go to the club and have a drink

The resulting error taxonomy used by TAGARELA consists of seven general groups of errors: *non-words*, *orthography*, *agreement*, *missing concept*, *extra concept*, *word order*, and *word choice*. Some of the error groups are subdivided further; for example, *agreement* is divided into *subject–verb* for person and number, and *adjective–noun*, *determiner–noun* and *subject–predicative* for number and gender. As shown in Table 1, the error types are grouped into meaning-related and form-related errors, depending on whether the error impacts the meaning or the well-formedness of the learner utterance.

2.3. Providing feedback

Once errors are diagnosed in the student input, a module called Feedback Manager decides on the error message to report, generates the message, and displays it to the

Table 1. Error taxonomy.

Meaning-related	Form-related
Word choice	Word order
Missing concept	Agreement
Extra concept	Non-words
Missing and extra concept	Punctuation
	Capitalization

Figure 1. Example feedback provided by TAGARELA.

student. An example of a feedback message can be seen in Figure 1, where the student made a word choice error, among others. The feedback message contrasts the infinitive of the word used by the learner with the infinitival form of the correct word choice.

3. Demands on the NLP architecture

3.1. Handling a range of activity types

We mentioned in section 1 that one of the motivations for having a demand-driven architecture in an ICALL system is to adjust the processing of the input and the feedback messages to different types of activities. In this section, we describe the different types of activities supported by the TAGARELA system, and present some of their specifications. In section 4.2, we illustrate how TAGARELA uses that information to process the student input and provide feedback.

TAGARELA includes six types of activities for beginning learners of Portuguese at the university level: reading, listening, description, rephrasing, vocabulary, and fill-in-the-blanks. The activity types represent different tasks that have to be performed by the learner. As shown in Table 2, activity specifications directly affect: (a) the nature of the student input; (b) the type of NLP processing that is necessary to handle such input; and (c) the nature of the feedback message that should be generated.

TAGARELA was designed to be an electronic workbook that could incorporate activity types commonly found in current textbooks. Each activity type was designed to trigger specific types of answers and make students practice certain skills and language patterns. For example, while *fill-in-the-blank* activities are meant to practice specific verbal and nominal morphology, *listening* and *reading* activities were designed to make students practice the respective skills while they answer questions

Table 2. Activity types and their properties.

Activity	Input	NLP	Feedback on
FIB	Words	Spelling, lexicon	Word form, missing word
Vocabulary	Phrase	<i>same as FIB</i>	Missing + extra word, word form
Rephrasing	Sentence	<i>same as FIB</i> + parsing	Form-related errors
Description	Sentence	<i>same as rephrasing</i>	Missing word, agreement
Reading	Sentence	<i>same as rephrasing</i> + content	Meaning-related errors
Listening	Sentence	<i>same as reading</i>	Meaning-related errors

that are more content-oriented. It is important to notice that this is a design choice specific to TAGARELA, and not a necessary demand or restriction imposed by activity types. In other words, in general nothing prevents the use of listening activities to practice verbal morphology or the development of dictation skills, it is just that in our system, listening activities are not used for those purposes.

In section 4, we will see that by specifying the pedagogical goals and the properties of input for each activity type, we are able to better filter the error analysis and select more appropriate feedback messages that can focus on the pedagogical requirements of any given activity.

3.1.1. *Learner input properties*

Reading, listening, description, and rephrasing require the learner to produce a full sentence. The target answer for vocabulary activities is usually a noun phrase, while the fill-in-the-blank exercises are typically answered with one word per blank.

3.1.2. *Processing requirements*

Because of the different expected input types, the NLP modules required by each activity can vary. Fill-in-the-blanks only require the spell-checker and a simple matching mechanism, while for all other types of activities the input processing usually starts with tokenization and lexical look-up, and may end up requiring a full syntactic analysis of the input sentence. Even within the same type of activity, the properties that need to be identified by the NLP modules can differ: some reading exercises target forms explicitly given in the text, others require more semantic analysis or inferences.

To give a concrete example, reading comprehension questions in TAGARELA all take full sentences as answers, yet are heterogeneous in terms of processing and feedback. In the most simple reading comprehension task based on *text identity*, the student is required to identify an answer explicitly given in the text. For such a task, content analysis can often be successfully performed using simple string matching. In the second type, corresponding to a more general *information extraction* task, the student is required to extract information which is given in the text but not as a contiguous sequence. While full string matching with simple edit distance measures will not be sufficient for such tasks, shallow content analysis using token matching can still be successful. Finally, for reading comprehension tasks requiring the student to draw inferences based on a text, content analysis requires deeper analysis to compare concepts and relations in the learner answer with those in the target answer.

Besides the need for specific NLP analysis for different activity types, there is also the issue of the reliability of such analysis. For example, the output of some matching techniques can be more reliably used to generate feedback messages with activities such as rephrasing and description than with activities like reading and listening. This happens because the elicitation techniques used by activities like rephrasing and description more effectively constrain the possible variation in student responses than the *wh*-questions used in reading and listening comprehension, which often allow for multiple possible answers.

One can ask at this point why one does not always perform all NLP analyses for every activity. Essentially the answer can be boiled down to ‘don’t guess what you know’. The more we know about the linguistic properties, the types of variation, and the potential errors that the NLP needs to detect, the more specific information we can diagnose with higher reliability.

Naturally, the more activity types and NLP resources are present in a system, the more beneficial the kind of architecture presented here becomes – an issue we return to in section 4.2.

3.1.3. *Feedback messages*

Information about specific activity types also impacts the feedback messages that can be displayed to the student. Feedback messages for reading, listening, and description activities should prioritize meaning over form. For these types of activities, meaning-related errors are displayed first whenever multiple errors are diagnosed in TAGARELA. Feedback messages for rephrasing activities, on the other hand, can focus on syntactic errors at the sentence level, while vocabulary and fill-in-the-blanks activities tend to require feedback messages that target specific misuses of lexical items or morphemes.

3.2. *Combining information from different NLP modules*

A flexible, annotation-based ICALL architecture is also motivated by the need to support interleaving of contributions from different modules. For instance, tokenization can already resolve some part-of-speech ambiguities. Take, for example, the Portuguese token *a*, which can be a preposition (*to*), a pronoun (*her*, clitic direct object), or an article (*the*, feminine singular). But when the token *a* arises in the analysis of a contraction, such as in (8), the correct part-of-speech can unambiguously be determined.

- (8) a. $da = de + a_{article}$
 b. $vê-la = ver + a_{cliticpronoun}$
 c. $â = a_{preposition} + a_{article}$

The tokenizer thus should already be able to assign the part of speech in such cases. To support this, in an annotation-based architecture the same data structure is accessed by all NLP modules and can be enriched monotonically. This means that throughout processing, information is generally added to the input, not replaced or removed.

Another example for several NLP modules contributing to the same representation is part-of-speech disambiguation. For the lexical ambiguity which arises at the

point of lexical lookup, disambiguation can be based on two distinct sources of information. On the one hand, the Constraint Grammar-like disambiguation rules introduced in section 2.1 attempt to use information about the local context to reduce or eliminate the ambiguity. On the other hand, the frequency of a given tag as specified in the CURUPIRA lexicon is used for disambiguation of any remaining cases, and one can readily imagine the integration of more complex statistical ambiguity resolution modules.

3.3. Combining information from different sources: learner input, activity model, learner model

In addition to the information obtained from the learner input, information about the learner and the activity performed can also play an important role and thus needs to be integrated into the ILTS architecture.

3.3.1. Integrating information from the activity model

The issue essentially brings us back to the very beginning of the article, where we stated that for some constrained exercise types it is possible to anticipate and hard-code for each potential input the corresponding feedback, whereas for others it is necessary to abstract and generalize using NLP techniques. Important for us here is that even for exercise types leading to a wider range of well-formed and ill-formed input, it is possible to hand-specify certain cases and the feedback to be provided for them. This insight can be useful to avoid costly NLP steps (either to avoid them completely, or to only run them once and then cache the result), or to manually provide information which cannot be reliably identified or identified at all using available NLP techniques. It is thus an important requirement for an ILTS processing architecture to support the flexible integration of the NLP analysis of the learner input with hand-encoded information for specific learner inputs provided as part of the activity models.

Interestingly, one can also see this flexible integration of static information from the activity model with the dynamic identification of information obtained by analyzing the learner input as presenting us with a continuum of systems stretching all the way from traditional CALL systems, where all learner input and the feedback to be provided for it needs to be anticipated and hard-coded, to an ICALL system without explicit activity information (such as a system providing feedback on free form essays), where all feedback is provided based exclusively on information derived by processing the learner input using general NLP resources.

3.3.2. Integrating information from the learner model

Information about the learner, their typical language use and errors, and the strategies they use to perform a particular activity can be crucial for disambiguating the learner input (cf. Amaral & Meurers, 2008). It is thus necessary to obtain an architecture which flexibly integrates information from the learner model and activity model with the NLP analysis of the learner input, both in terms of the processing itself and in terms of combining the output of the different sources of information into a single annotated representation of the learner input.

4. Realization of the architecture in UIMA

Having motivated the use of a demand-driven, annotation-based NLP framework for ILTS, in this section we connect the conceptual issues to a concrete NLP architecture. We introduce the Unstructured Information Management architecture (UIMA) and describe how we used its features to realize the envisaged ILTS architecture.

4.1. What UIMA offers

UIMA (Ferrucci & Lally, 2004) is a general framework for the automatic annotation of unstructured data, such as audio or text. In practice, UIMA is almost exclusively used for NLP applications. Central to UIMA's design is the idea of storing analysis results in a shared repository, the Common Analysis System (CAS, cf. Götz & Suhre, 2004). The CAS stores the text to be annotated separately from the annotations that pertain to specific parts of the text, similar to stand-off XML annotation in current linguistic corpus annotation schemes (cf., e.g., Ide et al., 2000). Annotations can be seen as enrichment layers that add information. This contrasts with traditional NLP architectures transforming the input from one representation to another, depending on what each module in such a pipeline expects. In UIMA, the so-called *Annotators* provide information by adding new annotations to the CAS. Subsequent modules can then retrieve previously added annotations and add new ones as needed.

In connection with the central data repository, UIMA introduces the idea of a global *type system*, where the types of annotation to be stored are defined. Annotations are described in terms of typed feature structures, where feature values can be of any type. For example, a type *Token* might have a feature *tag* with value of type *String*, which stores the token's part-of-speech tag. An advantage of type systems is that they enable meaning-related access; type definitions are explicitly established for the whole application, so every Annotator knows which information can be found where and no data structure checking is necessary.

4.2. Using UIMA to implement an ILTS architecture

4.2.1. Type system

Our type system is organized around three main units to which all other information is attached. *Tokens* are the smallest unit, to which the result of spell checking and lexical lookup for part-of-speech and morphological information is annotated.

Phrases represent partial or complete parse trees, built of tokens and other phrases. For this purpose, they have a list of daughter nodes that represent their subtrees. For convenience, one can also store a reference to the parent node of a phrase. Figure 2 exemplifies the structure for the *Phrase* 'menina bonita' ('beautiful girl'). Phrases can be annotated with the result of daughter agreement, such as subject-verb agreement.

The third type of annotation, *AnalysisResults*, does not refer to a particular portion of the input string but rather to the input as a whole. It is a repository of information which can be directly used by the diagnosis part of the system, such as the result of global agreement checking and the content assessment module introduced in section 2.1.

A notable aspect of the *Token* type is that it can associate an underlying form. For example, as we saw in section 3.2, certain Portuguese words such as contractions are syntactically complex, i.e., they can be thought of as consisting of two underlying tokens. The contraction ‘da’ we saw in (8a) is analyzed as consisting of the preposition ‘de’ (of) and the determiner ‘a’ (the). During the analysis, modules such as the parser need to refer to ‘de’ and ‘a’ separately if phrase structure rules are to apply properly. However, as discussed in Amaral and Meurers (2009), feedback to the learner needs to be given in terms of the surface representation, in this case ‘da’. To address the need to encode both perspectives, a *Token* can store a list of underlying *Tokens* under DEEPFORM as illustrated in Figure 3.

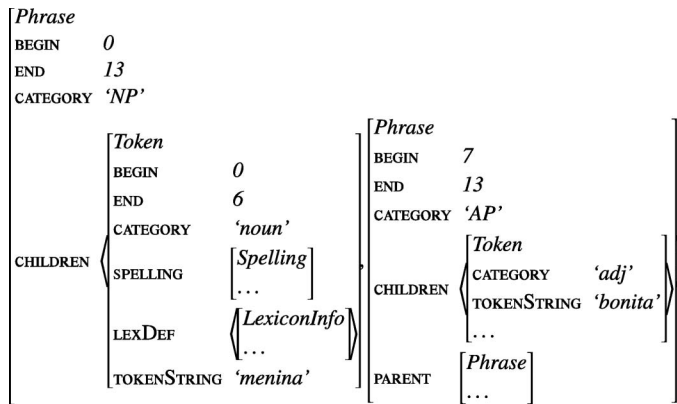


Figure 2. Typed feature structure representation of the phrase *menina bonita*.

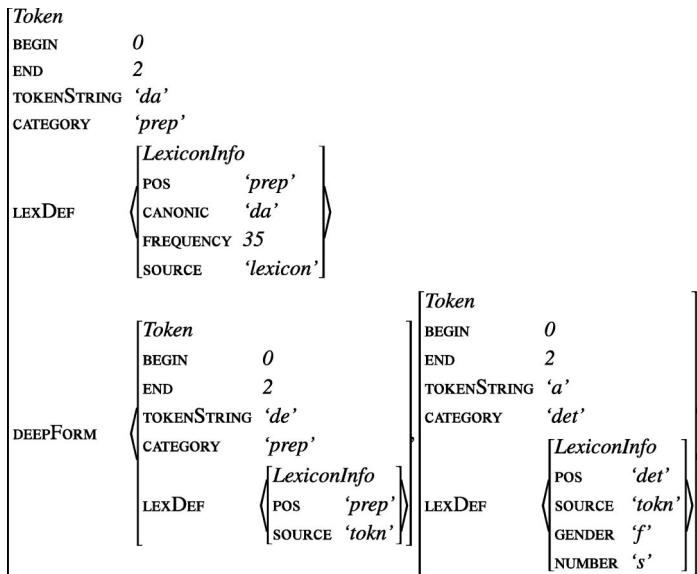


Figure 3. Typed feature structure representation of the contraction *da* (‘of the’).

4.2.2. *Multiple views for learner and target answers*

Similar to other ILTS, TAGARELA uses target answers as a reference for analysis of learner answers. Both target and learner answers need to be processed (tokenized, part-of-speech tagged, etc.) independently, but certain analysis steps need to make reference to both of the annotated representations. For example, the content assessment modules need to map tokens in the learner answer to the ones in the target answer. Consequently, we need a data structure to encode such mappings. To avoid doubly encoded information, it should also provide access to the annotated analysis results for both learner and target inputs. UIMA provides an adequate solution called *multiple views*. In contrast to a regular CAS, a multi-view CAS can hold more than one text. This allows us to analyze learner and target answers independently and to create mappings using the representations of either part.

Consider, for example, the ill-formed learner input ‘Ele se chamo Carlos’ (‘His name is Carlos’), where the learner made a wordform error in the verb ‘chama’, using first-person ‘chamo’ instead. This can be encoded through a mapping from the learner token ‘chamo’ to the target token ‘chama’. The mapping also makes the token annotations, such as the morphological information accessible, allowing the system to pinpoint the exact nature of the mismatch when giving feedback to the learner.

4.2.3. *Input and output specifications for analysis modules*

We argued in section 3.1 that different activities require different NLP analysis. So how can the different activity demands flexibly be translated into processing strategies? Our solution to this problem makes use of the explicitness enforced by UIMA type systems. Given that types have to be declared before any analysis is done, they can be used as a means of specifying analysis requirements. For example, a fill-in-the-blank activity can require *Tokens* with a specified `LEXDEF` feature, which means that lexical lookup must be done in addition to tokenization.

Binding activities to certain annotation types has the drawback that it only takes into account the NLP analysis – but the ultimate goal of an ILTS is to give useful feedback. We therefore added explicit error types to the system in order to fill in the missing link, as discussed at the beginning of section 2.2. The strategy can be described as top-down: pedagogical intervention opportunities are expressed as targeted error types, i.e., the particular errors a specific activity focuses on. These error types are mapped to required annotation types that are then annotated by the relevant analysis modules. The system thus knows what errors to focus on in processing and what to prioritize in the feedback for a given activity model. We elaborate on this process with examples in section 4.3.

In addition to formulating the overall analysis requirements, annotation types are also used to express dependencies between individual analysis modules. The lexicon lookup, for example, needs to work on tokenized input data. Hence, the input specification for the lexicon module is required to contain at least the type *Token*. The output specification then additionally contains the feature `LEXDEF`, whose value is a set of lexical entries.

4.3. *Handling a range of activity types*

We can now turn to describing in more detail, how the variation in TAGARELA's activities outlined in section 3.1 can be handled by our architecture. As mentioned in the previous section, our architecture dynamically adapts processing and feedback to the activity using the error types specified in the activity model.

The strategy employed by the first version of TAGARELA (Amaral, 2007) was to always prefer feedback on meaning over feedback on form. In our UIMA-based reimplementations of TAGARELA, the strategies are specified as part of each activity model to be able to adapt the strategy to the activity. We motivate the choice of different, activity-dependent strategies by describing two activity types, reading comprehension and rephrasing, with sample learner input and generated feedback.

4.3.1. *Reading comprehension*

For a reading comprehension activity, the original strategy of preferring feedback on meaning over feedback on form makes good sense because the questions aim to test the learner's understanding of a given text. They are not designed to test certain grammatical constructions or morphological characteristics. Hence, the activity model for reading comprehension activities explicitly states that feedback should concentrate on any kind of meaning-related error that can be detected. Concretely, such errors can be inappropriate lexical choices or missing concepts, among others. Let us consider an example task from TAGARELA where the text describes a woman named Patricia. One of the questions is 'Quantos anos ela tem?' ('How old is she?'). A possible learner answer is given in example (9):

- (9) Ela é quinze **ano**.
she is fifteen year

There are two errors in this learner sentence. First, age here is expressed using the verb 'ser' ('to be') instead of the correct verb 'ter' ('to have') used for this purpose in Portuguese, so the learner made a lexical choice error. Second, 'ano' ('year') should be in its plural form 'anos' ('years') in order to agree with 'quinze' ('fifteen') in number. Both errors are detected by the system. However, since meaning-related feedback is preferred according to the activity model, TAGARELA selects the wrong lexical choice as the more important error and responds with the error message we already saw in Figure 1, saying 'I am not expecting the verb "ser" for this answer. Try using "ter" instead.' Once the learner has changed the verb and re-submitted his answer, the system reports the remaining agreement error.

4.3.2. *Rephrasing*

TAGARELA's original 'meaning-over-form' strategy works fine for content-oriented activities such as answering reading comprehension questions, but for activities such as rephrasing the content is explicitly given. For those activities, the focus is on circumscribing a given sentence in a different way, using particular lexical material. Consequently, the activity model for rephrasing activities in the reimplemented TAGARELA system states that form errors should be the main target of feedback.

Figure 4. Form-focused feedback on a Rephrasing activity.

One such activity requires the learner to rephrase the sentence ‘Eu sou americano’ (‘I am American’) using the expression ‘Estados Unidos’ (‘United States’). The intended correct target is ‘Eu sou dos Estados Unidos’ (‘I am from the United States’). Consider the erroneous answer provided by a learner in (10):

- (10) Eu sou **das** Estados.
 I am from the *fem* States *masc*

There are two different errors. First, the learner forgot to include part of the proper name, the modifier ‘Unidos’ (‘United’), in the rephrased sentence, which according to TAGARELA is a meaning-related error. Second, ‘das’ (‘from the’) has the wrong gender, yielding a form-related error. The particular activity model instructs TAGARELA to prioritize the form-related error, resulting in the message shown in Figure 4. Once this form error is fixed, TAGARELA provides feedback on the lexical content error.

In summary, these two cases are meant to illustrate that due to the demand-driven architecture employed in the reimplemented TAGARELA, the activity designer is given some control over the behavior of the system in case of multiple learner errors. By encoding targeted error types into the activity model, the system can be told what feedback to prioritize. And as outlined in section 4.2.3, these error types are also used to guide the system’s NLP.

5. Conclusion

In this article, we discussed what we believe is needed to develop ILT systems capable of integrating a wider range of FLT activity types, resulting in learner input of a heterogeneous nature and evaluated on different criteria. We argued that in place of the fixed, uniform NLP pipelines employed by current ILTS, it is beneficial to view the NLP modules as enriching the input with annotations, giving the activity model control over which annotations can or must be provided in a demand-driven architecture. We showed how such an architecture can be realized in the UIMA architecture developed for general NLP analysis of unstructured information.

The UIMA-based reimplementation of TAGARELA (Ziai, 2009) offers the full functionality of the original system described in Amaral (2007). Additionally, it flexibly adjusts feedback both to the activity and to the learner, based on properties of the respective models. Using the popular UIMA framework, exchangeability of individual NLP components is made easier, which also makes porting the general architecture to other languages simpler. A system for Spanish is already being planned. The source code of the UIMA components will also be made available under an open-source license to encourage collaborative development with other researchers. This aspect of using UIMA as a standard architecture is in line with the recent argument of Wood (2008) for such standardization to facilitate compatibility and reusability of resources in ICALL development.

Making use of the increased modularity, we also plan to extend the system with the full student model proposed by Amaral and Meurers (2008). In the current version, task strategies such as scanning a text are not associated with activities yet and thus cannot be included in the student model. The new modular architecture should enable us to implement these ideas with little technical overhead.

Finally, let us mention a striking parallel between the issue of bridging between the complex FLT needs and the ILTS architecture discussed in this article and the complex business demands handled in current service oriented architectures (SOA). More specifically, the adaptivity and dynamic configuration of the NLP processing sequence informed by the activity model and driven by the feedback needs that we have argued for in this article seems to bear an interesting resemblance to the Case Handling approach which has been proposed in the context of business process management, for which Weske, van der Aalst, and Verbeek (2004, pp. 3f) argue: 'While straight-through processing strives for more automation, case handling addresses the problem that many processes are much too variable or too complex to capture in a process diagram (van der Aalst & Berens, 2001). One way to do this is to make workflows data-driven rather than process-driven and allow for authorizations to skip or undo activities'.

Acknowledgments

We would like to thank the organizers and the audience of the Symposium on NLP in CALL at EUROCALL 2007, where the ideas worked out in this article were first presented. We are also grateful to the two anonymous journal reviewers for their useful feedback and to Mat Schulze for providing insightful comments on the article. Finally, we discovered the relation to case handling in Service Oriented Architectures thanks to a presentation by Christoph F. Strnadl and the helpful pointers he provided.

Note

1. See, e.g., the error taxonomy for the BRIDGE system in Weinberg, Garman, Martin, and Merlo (1995), or the one for the ALICE-chan system in Levin and Evans (1995).

Notes on contributors

Luiz Amaral (<http://people.umass.edu/amaral/>) is an assistant professor of Hispanic Linguistics at the University of Massachusetts Amherst. His research focuses on second language acquisition, intelligent computer-assisted language learning (ICALL), and syntax. He designed the ICALL system TAGARELA (<http://purl.org/icall/tagarela>) together with Detmar Meurers. Research issues there involve activity design, student models, and natural language processing of learner language for intelligent language tutoring systems.

Detmar Meurers (<http://purl.org/dm>) is a professor of computational linguistics at the University of Tübingen and an adjunct associate professor at the Ohio State University. In 2003, he started the OSU ICALL research group (<http://purl.org/net/icall>) to explore the interface of natural language processing, real life foreign language teaching, and second language acquisition research. His research in this domain investigates intelligent web-based workbooks (<http://purl.org/icall/tagarela>), supporting a range of meaning-based activities through automatic content assessment (<http://purl.org/icall/comic>), learner modeling, automatic visual input enhancement (<http://purl.org/net/werti>), task-based learner corpora (<http://purl.org/icall/welcome>), and automatic analysis and annotation of learner corpora. In support of interdisciplinary collaboration as a crucial component of sustainable ICALL research, he co-organized several workshops on the 'Interfaces of ICALL' (<http://purl.org/net/iicall>) and the 'Automatic Analysis of Learner Language' (<http://purl.org/calico/aall08.html> and <http://purl.org/calico/aall09.html>).

Ramon Ziai (<http://www.sfs.uni-tuebingen.de/~rzi>) is a researcher and PhD candidate at the Collaborative Research Center 833 at the University of Tübingen. His main research interest and background is in computational linguistics, with his current project focusing on shallow semantic analysis and the processing of potentially ill-formed input. The work presented in this article builds on his MA thesis. He has also published work on dependency parsing of learner language, the SACODEYL corpus search tool, the WELCOME platform for distributed collection of structured learner corpus data, integrating parallel analysis modules to evaluate the meaning of answers to reading comprehension questions, and enhancing authentic web pages for language learners.

References

- Amaral, L. (2007). *Designing intelligent language tutoring systems. Integrating natural language processing technology into foreign language teaching* (PhD thesis). The Ohio State University.
- Amaral, L., & Meurers, D. (2008). From recording linguistic competence to supporting inferences about language acquisition in context: Extending the conceptualization of student models for intelligent computer-assisted language learning. *Computer Assisted Language Learning*, 21(4), 323–338.
- Amaral, L., & Meurers, D. (2009). Little things with big effects: On the identification and interpretation of tokens for error diagnosis in ICALL. *CALICO Journal*, 27(1), 580–591.
- Amaral, L., & Meurers, D. (2011). On using intelligent computer-assisted language learning in real-life foreign language teaching and learning. *ReCALL*, 27(1), 4–24.
- Bailey, S., & Meurers, D. (2008). Diagnosing meaning errors in short answers to reading comprehension questions. *Proceedings of the 3rd Workshop on Innovative Use of NLP for Building Educational Applications, ACL*. Retrieved from <http://aclweb.org/anthology-new/W08-0913>.
- Bick, E. (2000). *The parsing system "Palavras": Automatic grammatical analysis of Portuguese in a constraint grammar framework*. Aarhus: Aarhus University Press.
- Bick, E. (2004). PaNoLa: Integrating constraint grammar and CALL. In H. Holmboe (Ed.), *Nordic language technology (Yearbook 2003)* (pp. 183–190). Copenhagen: Museum Tusulanum.
- Delmonte, R. (2003). Linguistic knowledge and reasoning for error diagnosis and feedback generation. *CALICO Journal*, 20(3), 513–532.
- Díaz-Negrillo, A., & Fernández-Domínguez, J. (2006). Error tagging systems for learner corpora. *Revista Española de Lingüística Aplicada (RESLA)*, 19, 83–102.
- Ferrucci, D., & Lally, A. (2004). UIMA: An architectural approach to unstructured information processing in the corporate research environment. *Natural Language Engineering*, 10(3–4), 327–348.
- Götz, T., & Suhre, O. (2004). Design and implementation of the UIMA common analysis system. *IBM Systems Journal*, 43(3), 476–489.
- Heift, T. (2003). Multiple learner errors and meaningful feedback: A challenge for ICALL systems. *CALICO Journal*, 20(3), 533–548.
- Heift, T., & Schulze, M. (2007). *Errors and intelligence in computer-assisted language learning: Parsers and pedagogues*. New York: Routledge.

- Ide, N., Bonhomme, P., & Romary, L. (2000). XCES: An XML-based encoding standard for linguistic corpora. *Proceedings of the 2nd International Conference on Language Resources and Evaluation* (pp. 825–830). Retrieved from <http://www.cs.vassar.edu/~ide/papers/xces-lrec00.pdf>.
- Karlssoon, F., Voutilainen, A., Heikkilä, J., & Anttila, A., (Eds.). (1995). *Constraint grammar: A language-independent system for parsing unrestricted text*. Berlin and New York: Mouton de Gruyter.
- Koike, D., & Klee, C. (2003). *Lingüística aplicada: Adquisición del español como segunda lengua*. New York, NY: John Wiley and Sons, Inc.
- Levin, L.S., & Evans, D.A. (1995). ALICE-chan: A case study in ICALL theory and practice. In V. Holland, J. Kaplan, & M. Sams (Eds.), *Intelligent language tutors. Theory shaping technology* (pp. 77–98). New Jersey: Lawrence Erlbaum Associates, Inc.
- Martins, R., Nunes, G., & Hasegawa, R. (2003). Curupira: A functional parser for Brazilian Portuguese. *Computational Processing of the Portuguese Language, 6th International Workshop, PROPOR. LNCS, 2721, 195*. Berlin, Heidelberg: Springer.
- Meurers, W.D. (2009). On the automatic analysis of learner language: Introduction to the special issue. *CALICO Journal, 26*(3), 469–473.
- Montrul, S. (2004). *The acquisition of Spanish*. Amsterdam, Philadelphia: John Benjamins Publishing Company.
- Nagata, N. (2002). BANZAI: An application of natural language processing to web based language learning. *CALICO Journal, 19*(3), 583–599.
- Nagata, N. (2009). Robo-Sensei's NLP-based error detection and feedback generation. *CALICO Journal, 26*(3), 562–579.
- Rypa, M., & Feuerman, K. (1995). CALLE: An exploratory environment for foreign language learning. In V. Holland, J. Kaplan, & M. Sams (Eds.), *Intelligent language tutors. Theory shaping technology* (pp. 55–76). New Jersey: Lawrence Erlbaum Associates, Inc.
- van der Aalst, W.M.P., & Berens, P.J.S. (2001). Beyond workflow management: Product-driven case handling. In S. Ellis, T. Rodden, & I. Zigurs (Eds.), *International ACM SIGGROUP Conference on Supporting Group Work (GROUP 2001)* (pp. 42–51). New York: ACM Press.
- Weinberg, A., Garman, J., Martin, J., & Merlo, P. (1995). A principle-based parser for foreign language tutoring in German and Arabic. In V. Holland, J. Kaplan, & M. Sams (Eds.), *Intelligent language tutors. Theory shaping technology* (pp. 23–44). New Jersey: Lawrence Erlbaum Associates, Inc.
- Weske, M., van der Aalst, W.M.P., & Verbeek, H.M.W. (2004). Advances in business process management. *Data & Knowledge Engineering, 50*, 1–8.
- Wood, P. (2008). Developing ICALL tools using GATE. *Computer Assisted Language Learning, 21*(4), 383–392.
- Ziai, R. (2009). *A flexible annotation-based architecture for intelligent language tutoring systems* (Master's thesis). Universität Tübingen. www.sfs.uni-tuebingen.de/~rziai/papers/Ziai-09.pdf.