

AAA4LLL – Acquisition, Annotation, Augmentation for Lively Language Learning

Bartholomäus Wloka  

University of Vienna, Centre for Translation Studies, Vienna, Austria

Werner Winiwarter  

University of Vienna, CSLEARN – Educational Technologies, Vienna, Austria

Abstract

In this paper we describe a method for enhancing the process of studying Japanese by a user-centered approach. This approach includes three parts: an innovative way of acquiring learning material from topic seeds, multifaceted sentence analysis to present sentence annotations, and the browser-integrated augmentation of perusing Wikipedia pages of special interest for the learner. This may result in new topic seeds to yield additional learning content, thus repeating the cycle.

2012 ACM Subject Classification Information systems → Browsers; Computing methodologies → Lexical semantics; Applied computing → E-learning

Keywords and phrases Web-based language learning, augmented browsing, natural language annotation, corpus alignment, Japanese computing, semantic representation

Digital Object Identifier 10.4230/OASlcs.LDK.2021.29

1 Introduction

As most of us know, learning a foreign language, unless it is done in early childhood, is a challenging task, demanding motivation, patience, and last but not least a good teacher or learning method. The endeavour becomes even more challenging when the language differs greatly from the ones we already know. We address this issue in this paper by proposing a self-directed, contextualized learning method for English speakers to learn Japanese.

Apart from the stark difference of the writing system, Japanese has a fairly unique style of grammar, heavily dependent on postpositions, the tendency to omit personal pronouns, and several registers of politeness, which are often expressed by entirely different verb forms. The Japanese writing style is not only different, but also uses a combination of the syllabary *kana*, which comes in two forms, *hiragana* and *katakana*, and a large collection of logographic characters, which are called *kanji*. Each of these pictograms has several possible readings and meanings and in most cases a complicated decomposition into smaller building blocks. Apart from that, their pronunciation, meaning, and grammatical function can be heavily modified by the embedding context. They can also be combined to form new compound expressions and terminology. Japanese children are taught kanji throughout their high school education with 80 to 200 kanji per school year [13]. This slow and gradual process of acquiring this writing system allows for a strong foundation and complex compound expressions and terms can be easily learned step by step.

Clearly, adult learners of Japanese do not have this luxury and the complex characters, grammar, and spoken Japanese have to be learned at the same time. To put things in perspective, there are far more characters in everyday use and kanji are learned throughout the entire adult life in Japan. Standard dictionaries include more than 10,000 characters and over 50,000 words built from these characters. This means that learning methods need to streamline the process as much as possible just to assist the learner in catching up with this life long learning process of a native Japanese speaker.



© Bartholomäus Wloka and Werner Winiwarter;
licensed under Creative Commons License CC-BY 4.0

3rd Conference on Language, Data and Knowledge (LDK 2021).

Editors: Dagmar Gromann, Gilles Sérasset, Thierry Declerck, John P. McCrae, Jorge Gracia, Julia Bosque-Gil, Fernando Bobillo, and Barbara Heinisch; Article No. 29; pp. 29:1–29:15

OpenAccess Series in Informatics



Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

It is clear, that a learning environment needs to be interesting and engaging in order to increase the students' motivation. We find that this is best achieved by granting the learners the most possible freedom in selecting the material, while leading them towards understanding by offering the best possible decomposition of difficult concepts and their explanation. After presenting a possible translation to the learners, we deconstruct the Japanese sentence and enhance it with lexical, syntactic, conceptual, and relational annotation. The additional information is presented in a visually appealing way using colors and images to improve the overview of the structure. The learners can, of course, adjust the level of detail to their current language skill. We call these two parts *Acquisition* and *Augmentation of Learning Contexts*.

We first discuss the relevant related work in Sect. 2. We then describe the technical details of the implemented proof-of-concept framework in Sect. 3. In Sect. 4 we explain the automatic preparation of the learning material starting from seed topics chosen by the learner. In Sect. 5 we provide a detailed intuitive example of the augmented display presented to the learner. We summarize in Sect. 6 with a discussion, and our plans on how to extend and evaluate our approach in an in-class setting.

2 Related Work

The idea we build on in this paper is that motivation, the ability to choose the content, a keen interest in the subject matter, and a multimodal and multilayered environment are key to success when learning a new language. This is discussed extensively in [15]. A different approach to the same idea is shown in the incidental learning technique in [17], where the students are presented with information while browsing on-the-fly and in an unintrusive way. While browsing content of interest in the language they learn, like articles about their hobbies or daily news, the students are supported with facts about the text; a very elegant method which, however, clearly requires a relatively high level of skill in the foreign language, hence is reserved for advanced learners.

We extend on this concept in this paper on several levels, including the opportunity to use it even earlier in the language learning process. We do this by letting the learner decide the context, while still maintaining a feasible didactic structure and sensible levels of difficulty of the content. Research such as [14] shows the effectiveness of this approach. The interviews carried out in [6] further support these findings from the subjective point of view of the learners.

The component that makes this possible is the pre-selection of context, for which we use a bilingual alignment method of Wikipedia content, based on a metric obtained by matching of lexical units. This method and a discussion about how well Wikipedia is suited as a source of parallel content for the Japanese-English language pair can be found in [18].

Research in alignment and harvesting of bilingual material has been very active in the last decades, especially due to the importance of the application of such data in machine translation and other language technologies requiring training data. We would like to mention the largest of these approaches including the Japanese-English language pair, *JParaCrawl* [10], resulting in a collection of 10 million Japanese-English sentence pairs. The *WikiMatrix* project produced a multilingual collection with a large amount of parallel language data in 85 languages. This was done by using LASER sentence embeddings [1]. With such a volume of data, the quality varies greatly between the language pairs and is comparatively low for Japanese-English. These two large scale approaches are representative of the current paradigm of multilingual data collection, namely the brute-force black box approach. The

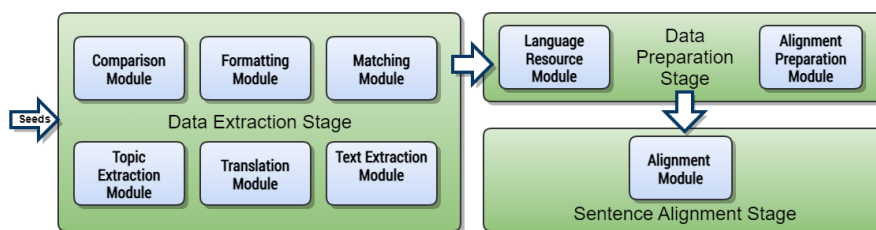
quality and huge computational requirements problems – an issue raised well in [3] – aside, the lack of transparency or a quality score makes it difficult to trace a path to the result and to judge the fitness of aligned data for language learning.

For the lexical and syntactic annotation of the sentences we use the de facto standard *universal dependencies* [11], whereas we rely on the *abstract meaning representation* (AMR) [2] for the semantic representation. Recently, there has been a renewed interest on different approaches for meaning representation (for a recent shared task see [12]). We have chosen AMR over other competing approaches because we judged it most suitable for language learning purposes.

3 System Architecture

We have implemented a language learning environment in which the users can study Japanese Wikipedia pages based on topics that really interest them. The learners select one or several topics as seeds. The following architecture turns them into a topic-specific collection of parallel sentences.

The architecture is divided into three stages, and each stage is further divided into modules. The modular approach allows for flexibility and ease of maintenance. Figure 1 gives an overview of this architecture.



■ **Figure 1** Stages with their corresponding modules.

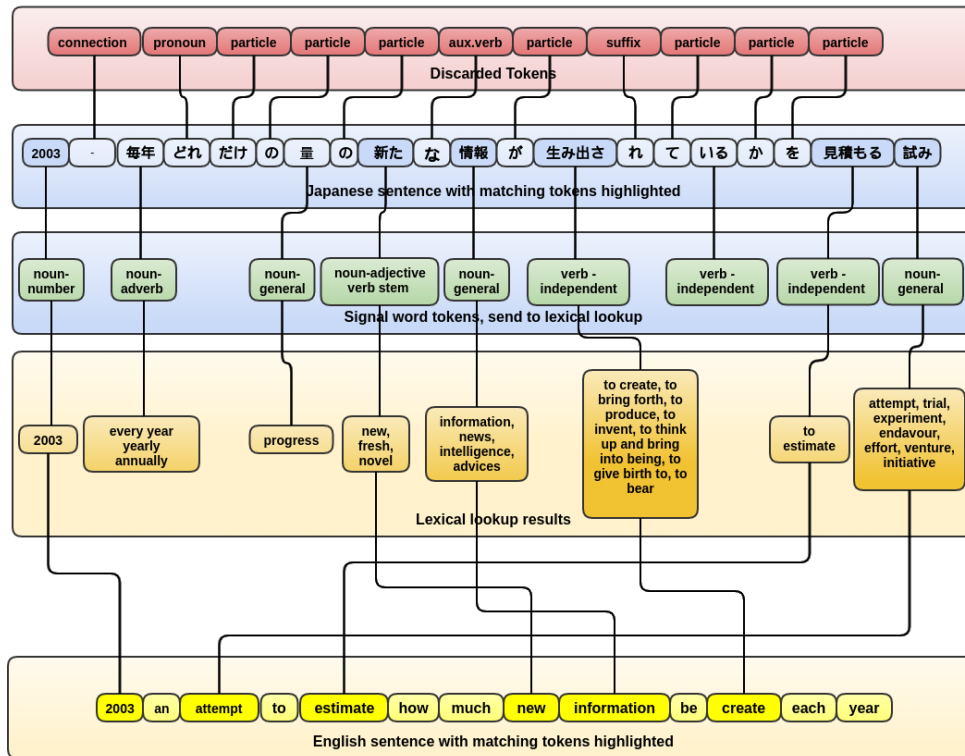
The flow of the data through this pipeline of stages and modules starts at the top with the *Data Extraction Stage*. This stage processes the seed input and extracts text data from Wikipedia accordingly. This is done by processing each English Wikipedia page of a seed topic and finding all links to further pages within this article. This process is repeated for the according Japanese pages. This is taken as the first measure of similarity between the contents. A discussion and preliminary results for the question of how much of the content between Japanese and English articles is comparable can be found in [18]. We use a threshold value to adjust the degree of similarity vs. the volume of candidate data.

Since we traverse the pages recursively and define each link as the starting point for the next iteration we obtain text that is in some way related to the initial seed. Naturally, the semantic distance increases with the number of links from the initial topic, which can also be adjusted as needed. Once we have collected topics, which we deem as good candidates, we extract the text from the Wikipedia pages.

In the *Data Preparation Stage* we prepare the text for further processing, i.e., the alignment. We segment the English text on a sentence and word level, lemmatize it, and determine the part-of-speech. We use dictionary resources from *EDRDG*¹. We segment and annotate the Japanese text with *MeCab* [8], which we later also use for our *lexical annotation*.

¹ <http://www.edrdg.org/>

Once the data is prepared, we align the bilingual input in the *Sentence Alignment Stage*. With the assumption that at least a good portion of the data has translation equivalents – depending on the above-mentioned threshold value – we traverse the entire pre-processed data set for matches. Selected lexical units in the Japanese sentence are examined for potential alignment indicators. An example of this process is shown in Fig. 2.



■ **Figure 2** Alignment example.

On top of the schematic depiction of this alignment process we see the tokenized Japanese sentence. Above the sentence is a row with discarded PoS tokens, below are the tokens taken into account for alignment. The reason we discard some of the tokens is that they contribute less to the alignment process. A detailed discussion can be found in [19].

The remaining tokens are being looked up in the lexical resources, i.e. bilingual dictionaries. It is important to mention that we retrieve each possible lexical equivalent of the words, as shown in the line with the lexical lookup results in Fig. 2.

Translations of sentences often vary in style and several differently sounding sentences might convey the same content. We consider these variations with our alignment method that takes into account all possible synonyms. This helps us to identify several – often stylistically varying – candidates, which is particularly interesting in a learning context.

In the process of matching the individual parts we compute an alignment score based on the number of matches normalized by the sentence length. We use this score, which indicates the alignment quality, hence the lexical similarity between the Japanese and the English sentence, to sort by alignment confidence. The output, i.e. the set of best alignment candidates, is the input to the annotated presentation.

For the purpose of presenting the annotated view to the learners we have designed a Web-based client-server architecture using *augmented browsing* technology to enhance the Web documents with event handlers at the client to retrieve annotations from the learning server and display the information in a comprehensible way. This is realized with the *WebExtensions API cross-browser technology*² and the *jQuery*³ and *jQueryUI*⁴ libraries.

The language learning server is implemented in *SWI-Prolog*⁵. It is not only an excellent choice for natural language processing tasks but also offers a scalable Web server solution as well as libraries for the efficient handling of huge XML and RDF files. The annotation data is transferred via XMLHttpRequests in JSON and assembled at the client in popup divs.

The only external software that we use is the Japanese dependency parser *CaboCha* [7]. We take its output as starting point for analyzing and annotating a sentence in four steps: the lexical, syntactic, conceptual, and relational level. We use several lexicosemantic resources: the dictionary files from EDRDG, *WordNet* [9], and datasets from *DBpedia*⁶.

For *lexical annotation*, we take the output of CaboCha, which includes the above-mentioned part-of-speech and morphological analyzer MeCab. The latter uses a fine-grained hierarchical tagset with up to four levels and additional conjugation types and forms. We map this tagset to *universal POS tags*⁷. Since MeCab follows a rather extreme segmentation strategy, which we found quite unsuitable for educational purposes, we use our lexical resources to merge adjacent tokens to achieve a more compact and comprehensible presentation. For the *kanji cards* in our display, we include *ideographic description sequence* data⁸ from the *CHISE* project⁹. We also show images as visual clues for kanji, which were hand-collected from Wikipedia pages.

CaboCha transforms a sentence into a sequence of *segments*, which are linked through *dependency patterns*. As CaboCha does not output any syntactic relation names, we had to add the appropriate *universal syntactic relation* names for displaying the *syntactic annotation*. We also had to arrange the relations vertically into several rows to offer an appealing visual representation.

The *conceptual annotation* is based on XML frame files from the *OntoNotes* project available from *LDC*¹⁰, and *AMR resource lists*¹¹. We extend the AMR approach by mapping words also to Wikipedia pages through DBpedia disambiguation links and to WordNet synsets, whenever we cannot find a suitable frame. We also display short abstracts and thumbnails for Wikipedia pages, again retrieved from DBpedia. For disambiguation, we rely on contextual and distributional data from the current sentence and its English equivalent, the Wikipedia page, and the collected topic-specific corpus.

As a final step, we add *roles* to the display, using again data from the OntoNotes frame files to offer a *relational annotation* to round off our augmented view of the linguistic and semantic properties of a Japanese sentence. In Sect. 5, we go through a detailed example, which illustrates the individual annotation levels as perceived from the perspective of the language student.

² <https://developer.mozilla.org/en-US/docs/Mozilla/Add-ons/WebExtensions>

³ <https://jquery.com>

⁴ <https://jqueryui.com>

⁵ <https://www.swi-prolog.org/>

⁶ <https://wiki.dbpedia.org/>

⁷ <https://universaldependencies.org/guidelines.html>

⁸ <https://github.com/cjkvi/cjkvi-ids>

⁹ <https://www.chise.org/>

¹⁰ <https://www ldc.upenn.edu/>

¹¹ <https://amr.isi.edu/download.html>

The use of SWI-Prolog for implementation has the big advantage that all the data is stored and accessed in a declarative way as Prolog *fact files*, which can be easily customized and reconsulted dynamically. We also developed several visual interfaces and editors for expert users in previous research efforts (see [16, 17]).

4 Acquisition of Learning Contexts

The initial step in our learning environment is the selection and acquisition of the learning material. Since we firmly believe that the best way to learn is to examine interesting content while having the difficulty of the material custom shaped to the language level, rather than reading boring texts that clearly only aim at conveying vocabulary or a grammatical concept, we let the user choose their context freely. The initial step is to find engaging content on English Wikipedia. A topic or several related topics of interest then become the seed(s) towards a collection of example sentences to start the learning journey, while discovering the desired information about a certain topic on Wikipedia.

The sentences are selected according to the description in Sect. 3. Thanks to the efficiency of the alignment algorithm, the learners can extract new information about any topic within minutes or a few hours, depending on the size of the dataset. Naturally, the exact time of the sentence alignment depends on the desired size of the example sentence corpus. Table 1 shows an example of collecting 805 example sentences for one seed topic.

■ **Table 1** Runtime example for a small dataset.

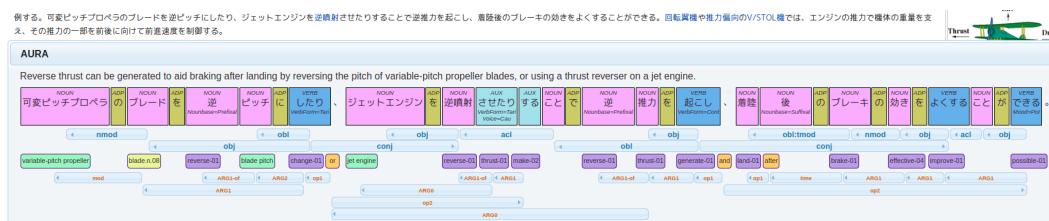
Module	Time	Output
Data Extraction Stage		
Topic Extraction	30m36s	en: 801087 articles, ja: 56736 articles
Text Extraction	18.1s	en: 2037 lines, ja: 2072 lines en: 85,804 tokens, ja: 75,393 tokens
Data Preparation Stage		
Alignment Preparation	59m45s	en: 3510 sentences, ja: 805 sentences
Sentence Alignment Stage		
Alignment	24m50s	805 aligned sentences

It is important to mention that the runtime efficiency increases with the number of alignments run on the learner's computer, since frequently occurring topic equivalents are stored locally and do not have to be looked up repeatedly in the Wikipedia database. This and other runtime tweaks are described in detail in [19]. After the output is generated, the learners can select which and how many of the best scoring sentences they want to examine. With the sentences they have chosen, the users then continue their learning experience with the Augmentation part of our language learning solution.

5 Augmentation of Learning Contexts

In this section we present the augmented information displayed to the user while working with Japanese Wikipedia pages as study material. We discuss each annotation level in a cleanly separated subsection. Throughout this section, we use one sentence as running example, taken from the Japanese Wikipedia page on *thrust*¹². Figure 3 shows the complete pop-up div, which appears when the student clicks on the sentence in the Web page.

¹²<https://ja.wikipedia.org/wiki/%E6%8E%A8%E5%8A%9B>



■ **Figure 3** Example of augmentation.

Unfortunately, the sentence is rather long and complicated, which is fairly typical for sentences found on Japanese Wikipedia pages. Luckily, our acquisition step (see Sect. 4) equips us with a reasonably good translation from the corresponding English Wikipedia page to significantly facilitate the challenging task for the language learner to make any sense of this text. To offer sufficient resolution for readability, we divide the presentation of the sentence into three parts in the following subsections. Since Japanese grammar is exclusively head-final and strongly left-branching with abundant use of postpositional particles, we consequently go through the sentence from right to left.

5.1 Lexical Annotation

We annotate the sentence using *universal POS tags* for the individual tokens, which are also visually emphasized through different colors. In our example sentence, these are: *blue* for verbs, *light blue* for auxiliaries, *pink* for nouns, and *olive* for adpositions.

Whenever necessary, we add *universal features*. In some cases, we also decided to use *language-specific features* and *values*. All these choices can be easily adjusted by expert users to accommodate their personal preferences.

The user can click on each token to open a popup div with further information. This includes the English glosses retrieved from our lexical resources and *kanji cards* for the kanji that are part of the word. The kanji cards include the *radical* number, *on’yomi* readings (in uppercase), *kun’yomi* readings, and English glosses. Radicals are 214 special kanji that are used as components of other kanji. One of them is always singled out as the radical of a kanji to look up the character in a kanji dictionary. On’yomi readings descend from approximations of original Chinese pronunciations whereas kun’yomi readings are based on pronunciations of native words approximating the meaning of the character when it was introduced. Finally, we display an image to offer a visual clue for memorizing the kanji. The *ideographic description sequence* defines the spatial structure of the kanji based on simpler components, the radical of the kanji is highlighted in *red* in this sequence, other radical components in *orange*.

In Fig. 4, we show the information for the noun 着陸 (*chakuriku*). As can be seen, the correct contextual pronunciations of the kanji are indicated in pink in the lists of possible readings.

The right part of the sentence contains the following lexical tokens for content words:

- the verb できる (*dekiru*), which is actually the *potential* form of the verb する (*suru*) “to do”, therefore, meaning “to be able to do”;
- the noun こと (*koto*) “thing”, which just nominalizes the preceding clause;
- the verb よくする (*yokusuru*) “to improve”;
- the noun 効き (*kiki*) “effectiveness”, derived from the *continuative* form of the verb *kiku*;

landing; alighting; touch down

着 123 𠄎 𠄎 𠄎 𠄎 𠄎
CHAKU; JAKU
ki.ru; -gi; ki.seru; -ki.se;
tsu.ku; tsu.keru
don; arrive; wear; counter for
suits of clothing

陸 170 𠄎 𠄎
RIKU; ROKU
oka
land; six

■ **Figure 4** Example of lexical analysis – right part.

- the loanword ブレーキ (*burēki*) “brake”;
- the suffix 後 (*go*) “after”, which is used like a *temporal* adposition; and
- the noun 着陸 (*chakuriku*) shown in the popup div.

Since adpositions and auxiliaries mainly serve as syntactic function words, we discuss them later in Sect. 5.2. Figure 5 shows the middle part of the sentence with the following lexical units:

- the continuative form 起こし (*okashi*) of the verb 起こす (*okasu*), which here means “to generate”, the continative form is used like the conjunction “and” to loosely connect the two clauses;
- the noun 推力 (*suiryoku*) “thrust”;
- the prefix 逆 (*gyaku*) “reverse”, shown in the popup div;
- the noun 逆噴射 (*gyakufunsha*) “reverse thrust”; and
- the loanword ジェットエンジン (*jettoenjin*) “jet engine”.

Finally, the left part of the sentence, as shown in Fig. 6, consists of the following lexical units:

- the *tari*-form したり (*shitari*) of the verb *suru*, which is used like the conjunction “or” to connect several exemplars with the pattern *-tari ...-tari suru*;
- the loanword ピッチ (*pitchi*), meaning here “blade pitch” or “angle”;
- the loanword ブレード (*burēdo*) “blade”; and
- the compound noun (part loanword) 可変ピッチプロペラ (*kahenpitchipuropera*) “variable pitch propeller”, for which the details are shown in the popup div.

5.2 Syntactic Annotation

For the syntactic annotation, we add *universal syntactic relations* to the display. However, to improve the comprehensiveness of the visual representation, we omit obvious relations between adjacent tokens. As mentioned in Sect. 3, we use the output of the Japanese dependency parser *CaboCha* for this purpose, but enhance it with syntactic relation names. Figure 7 shows the dependencies for the right part of our example sentence:

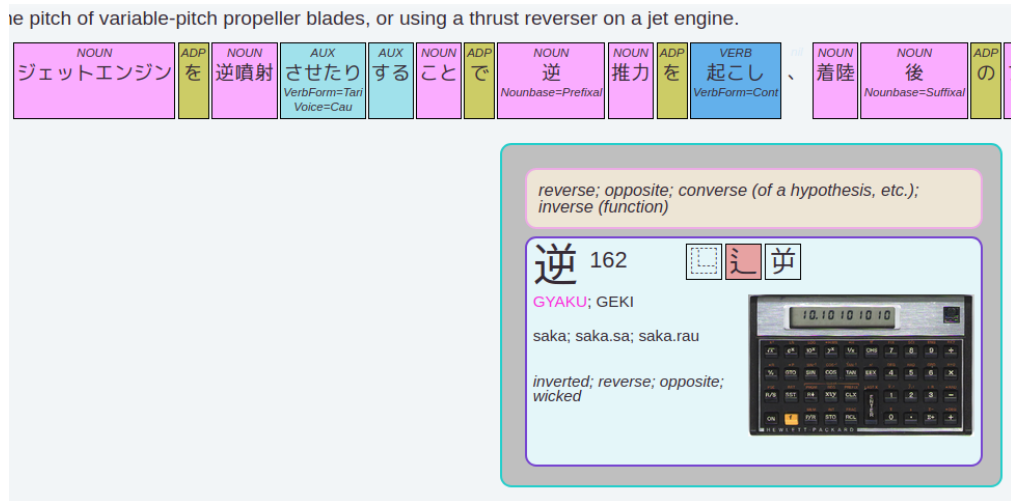


Figure 5 Example of lexical analysis – middle part.

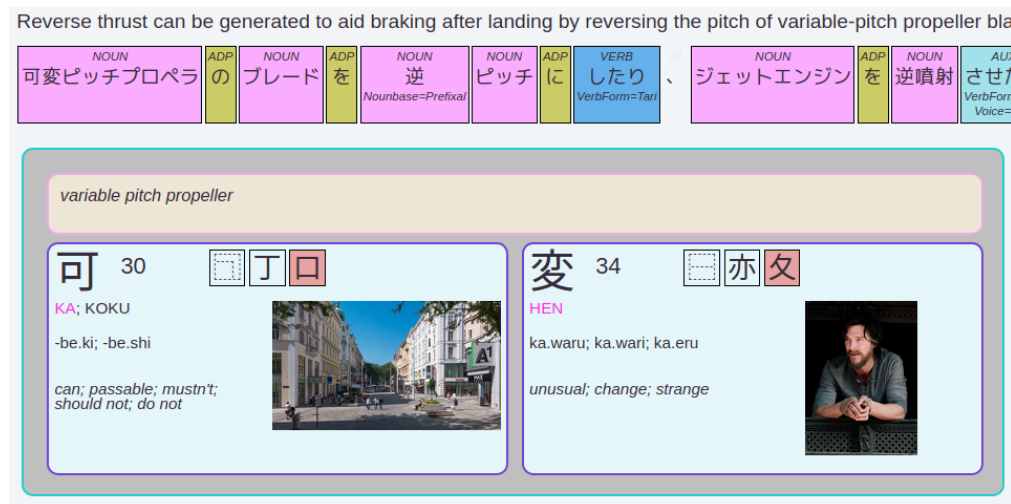


Figure 6 Example of lexical analysis – left part.

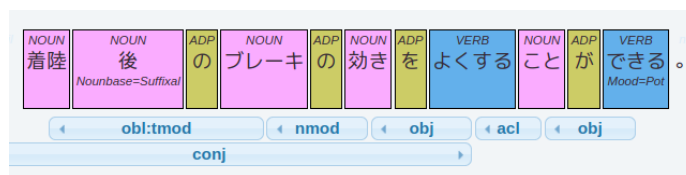
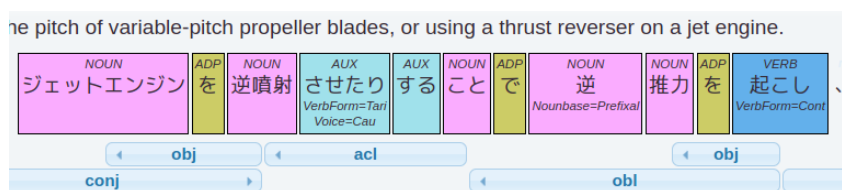


Figure 7 Example of syntactic annotation – right part.

- the *object* **obj** relation between *dekiru* and *koto*, indicating the “thing” that “can be done”, the postposition *が* (*ga*) usually marks the subject, however, in this context the direct object of the verb;
- the **acl** relation between *koto* and *yokusuru*: as mentioned before, this *adnominal clause* relation nominalizes the preceding clause, resulting in “improvement”, however, in combination with *koto ga dekiru* this effect is somehow canceled out as it just means “it is possible to improve”;

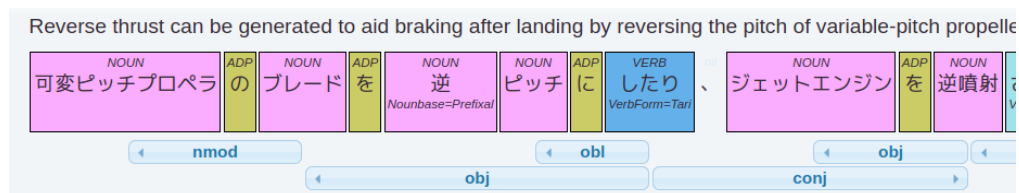
- the **obj** relation between *yokusuru* and *kiki* tells us that the “effectiveness can be improved”, here we can see the usual direct object marker を (o);
- the *nominal modifier* **nmod** relation between *kiki* and *burēki* with the corresponding adposition の (no) means that the former is an attribute or genitive complement, i.e. “the effectiveness of the brakes” or “the effectiveness of braking”;
- because of the special meaning “after” of the suffix *go*, we have a *temporal modifier* **obl:tmod** relation between *burēki* and *chakuriku*: “braking after landing”;
- finally, there is a *conjunct* **conj** relation between this clause and the preceding clause shown in Fig. 8, due to the continuative form *okashi*, as mentioned before.



■ **Figure 8** Example of syntactic annotation – middle part.

The middle part of the sentence, as shown in Fig. 8, contains the following relations:

- *suiryoku* is the direct object of *okashi*, i.e. we “generate thrust”;
- there is an *oblique nominal* **obl** relation between *okashi* and *koto*, the extremely polysemous adposition で (*de*) indicates here the means by which we generate the thrust;
- the two auxiliaries *させたり* (*sasetari*) and *する* (*suru*) combine with the preceding noun *gyakufunsha* and verbalize it so we end up with something like “reverse thrusting”, again this is undone by the **acl** relation with the noun *koto*,
- the **obj** relation to “jet engine” shows that we “reverse thrust the jet engine”, actually, here the *causative* form of the verb is used, so it literally means “make reverse thrust the jet engine”;
- ultimately, we have again a **conj** relation to the left part of the sentence, thanks to the already mentioned *-tari ...-tari suru* construction.



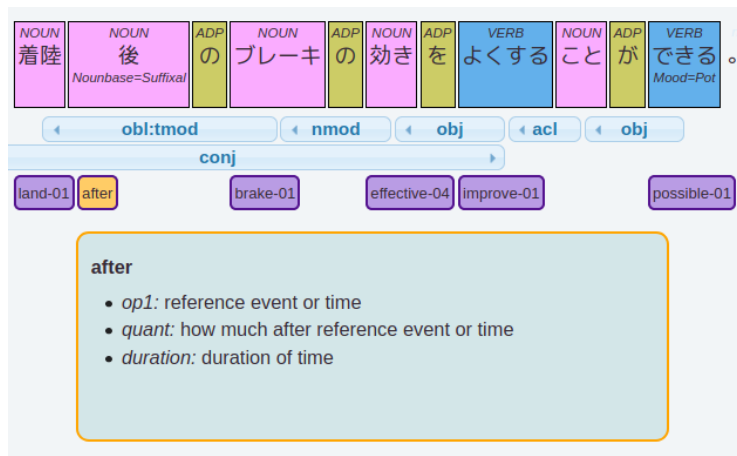
■ **Figure 9** Example of syntactic annotation – left part.

The final left part of the sentence in Fig. 9 only contains three additional relations:

- there is an **obl** relation between *shitari* and “pitch”, together with the again very polysemous adposition に (*ni*) this indicates here that we change something to a new state, i.e. “reversed pitch”;
- what we change is expressed by the **obj** relation, namely the “blade”; and
- to clarify matters through an **nmod** relation, it is the “blade” of a “variable pitch propeller”.

5.3 Conceptual Annotation

At the third level of annotation we map content words to concepts within the semantic representation framework AMR. In AMR we can distinguish between dedicated *AMR frames* like the one shown in Fig. 10 (*orange*) and *OntoNotes frames* as displayed in Fig. 11 (*purple*). As usual, the popup divs can be inspected by just clicking on the concept names. The frames provide a definition and several *roles* (see Sect. 5.4). For the right part of the sentence, we can see the following mappings:



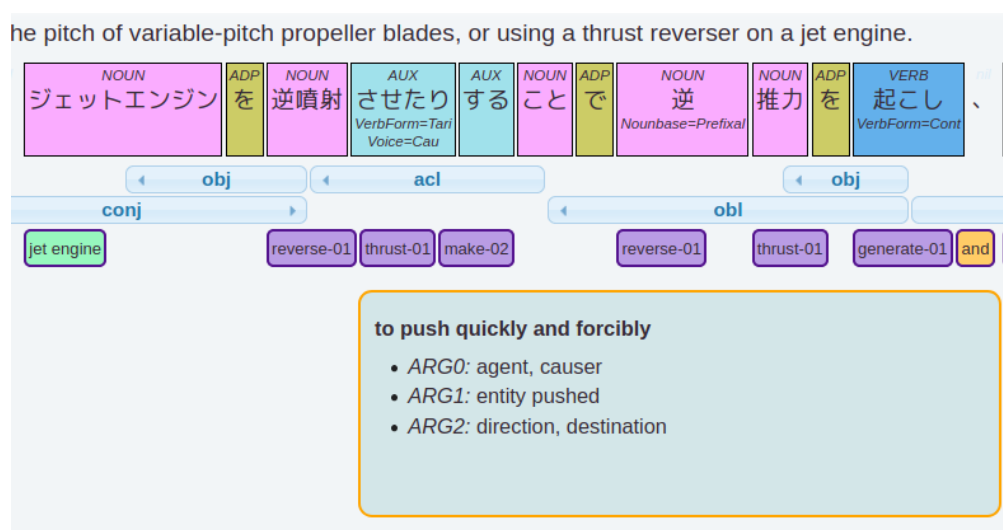
■ **Figure 10** Example of conceptual annotation – right part.

- `dekiru` ⇒ `possible-01`: “likely or able to be/occur”,
- `yokusuru` ⇒ `improve-01`: “make better”,
- `kiki` ⇒ `effective-04`: “cause an effect, successful in creating a desired effect”,
- `brēki` ⇒ `brake-01`: “slow a car via brakes”,
- `go` ⇒ `after`,
- `chakuriku` ⇒ `land-01`: “bring to land, from water or air”.

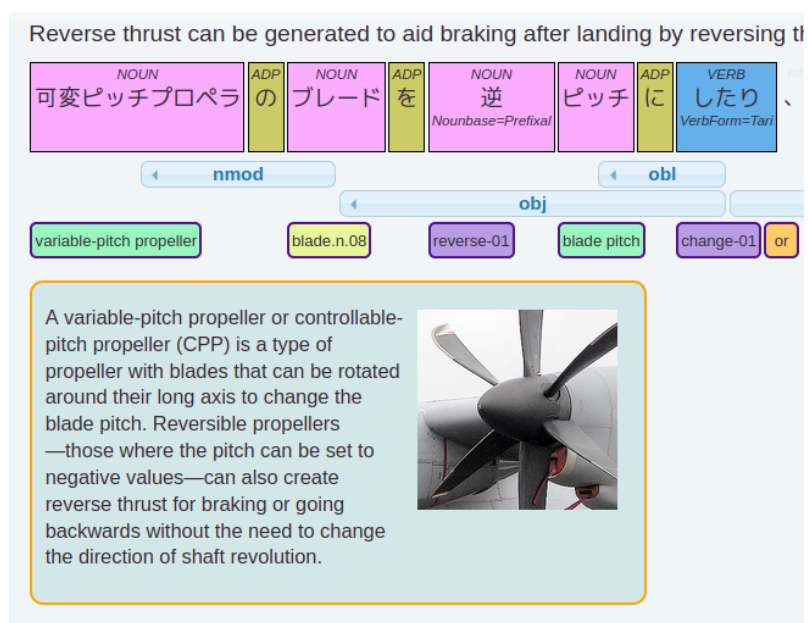
We use the glosses from the lexical annotation to retrieve possible frames, and contextual and distributional data to disambiguate among likely candidates (see Sect. 3). As can be seen in Fig. 11, conjunctions are mapped in AMR to special **and** and **or** frames, in addition, there are the following mappings to OntoNotes frames for the middle part:

- `okoshi` ⇒ `generate-01`: “create”,
- `suiryoku` ⇒ `thrust-01`: “to push quickly and forcibly”,
- `gyaku` ⇒ `reverse-01`: “turn around, change direction”,
- `gyakufunsha sasetari suru` ⇒ `reverse-01`, `thrust-01`, and `make-02`: “cause (to be)”.

The last entry is an example of assigning several concepts to one position in the sentence: the first two concepts correspond to the noun `gyakufunsha`, the last one is represented by the auxiliaries `sasetari suru`, however, since auxiliaries are not annotated as dependencies by CaboCha, we also map the third concept to the noun. In Fig. 11 there is also one mapping of a word to the corresponding Wikipedia page: `jettoenjin` ⇒ “jet engine” (*green*). Whenever we cannot map a word to a frame, we try to find a Wikipedia page representing the concept. If a user clicks on such a concept name, we display the short abstract and thumbnail retrieved via DBpedia (see Fig. 12). We use mainly the *DBpedia disambiguation links* data to identify



■ **Figure 11** Example of conceptual annotation – middle part.



■ **Figure 12** Example of conceptual annotation – left part.

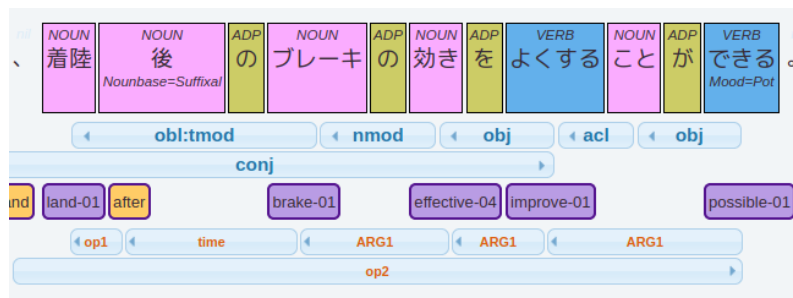
ambiguous words and select the Wikipedia page representing the correct word sense. Finally, if there is no existing Wikipedia page, we use *WordNet* as backup, again retrieving the correct *synset* (yellow) using word sense disambiguation based on contextual data and relational information derived from WordNet. If we click on a WordNet concept, the synset definition is displayed. This results in the following mappings for the left part of the sentence:

- shitari ⇒ change-01: “transform”,
- pitchi ⇒ blade pitch,
- burēdo ⇒ blade.n.08: “flat surface that rotates and pushes through air or water”,
- kahenpitchipuropera ⇒ variable-pitch propeller.

As can be seen, we can successfully narrow down the senses of the polysemous word “pitch” by following the disambiguation link to the correct Wikipedia page on `blade pitch`. In the case of “blade”, this is not possible, because there is only a Wikipedia page for the sense “sharp cutting part, for instance of a weapon or tool”.

5.4 Relational Annotation

With the relational annotation level, we complete the picture by adding semantic roles to the display to offer a semantic representation of the meaning of the sentence within the AMR framework. As can be seen in Fig. 13, we use a visual representation similar to that of universal dependencies in the syntactic annotation.



■ **Figure 13** Example of relational annotation – right part.

Whenever possible, we use *core roles*, defined in the OntoNotes frames (`ARG0`, `ARG1`, ...). In addition, AMR offers an inventory of *non-core roles*, e.g. `time` in Fig. 13 indicates the time when the braking occurs. The roles `op1`, `op2`, ... are special roles only used in AMR frames. Therefore, we have the following roles for the right part of the sentence:

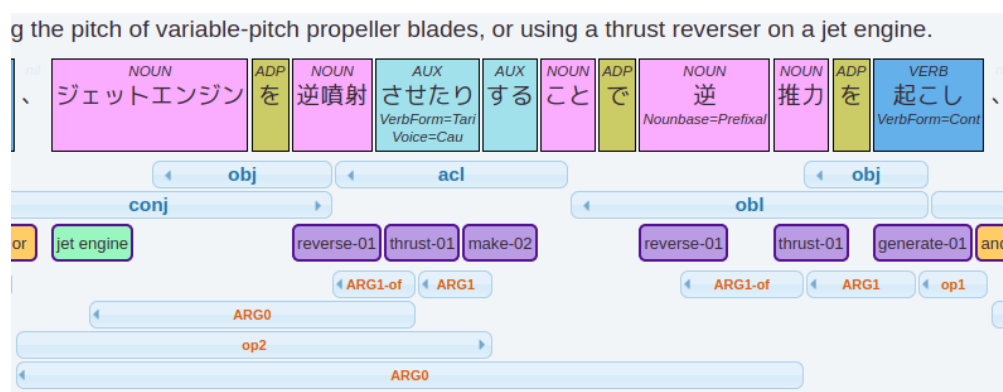
- `possible-01` $\xrightarrow{\text{ARG1}}$ `improve-01`: `improve-01` is the “thing that is possible”,
- `improve-01` $\xrightarrow{\text{ARG1}}$ `effective-04`: `effective-04` is the “thing improving”,
- `effective-04` $\xrightarrow{\text{ARG1}}$ `brake-01`: `brake-01` is the “domain in which arg0 (cause) is effective; outcome effected”.

We only indicate roles for which there is explicit evidence in the Japanese sentence. Since Japanese omits many details that are usually expressed in other languages at least through anaphora (a phenomenon also known as *zero anaphora* [5]), it is not often necessary to use the variable mechanism of AMR to refer to antecedents. The middle part of the sentence (see Fig. 14) contains the following roles, it also shows the use of *inverse roles* to *re-focus* the AMR representation:

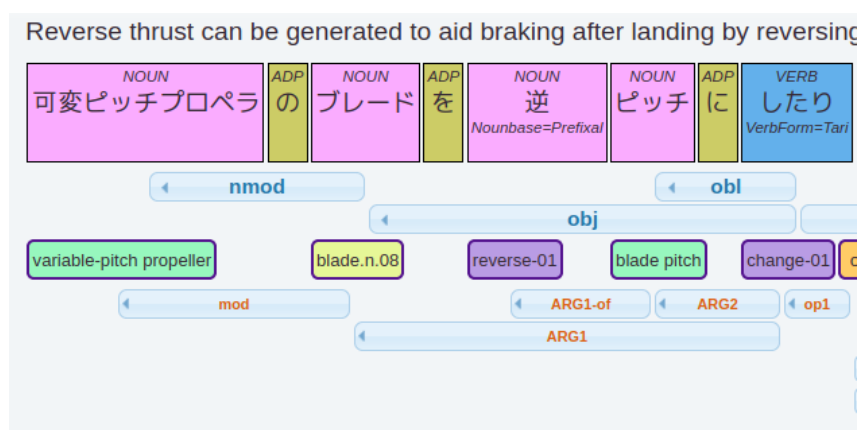
- `generate-01` $\xrightarrow{\text{ARG1}}$ `thrust-01`: `thrust-01` is the “thing created”,
- `thrust-01` $\xrightarrow{\text{ARG0}}$ `or`: the whole `-tari ...-tari suru` construct is the “agent, causer”,
- `thrust-01` $\xrightarrow{\text{ARG1-of}}$ `thrust-01`: this is an inverse role indicating that `thrust-01` is the “thing turning around”,
- `make-02` $\xrightarrow{\text{ARG1}}$ `thrust-01`: `thrust-01` is the “impelled action/ predication”,
- `thrust-01` $\xrightarrow{\text{ARG0}}$ `jet engine`: in this case, `jet engine` is the “agent, causer”.

Finally, Fig. 15 displays the semantic roles for the remaining right part of the sentence:

- `change-01` $\xrightarrow{\text{ARG1}}$ `blade.n.08`: `blade.n.08` is the “thing changing”,



■ Figure 14 Example of relational annotation – middle part.



■ Figure 15 Example of relational annotation – left part.

- $\text{change-01} \xrightarrow{\text{ARG2}} \text{blade pitch}$: blade pitch is the “end state”,
- $\text{blade.n.08} \xrightarrow{\text{mod}} \text{variable-pitch propeller}$: this non-core role tells us that the latter is a *modifier* of blade.n.08.

6 Conclusion

We have presented a Lively Language Learning solution that enables the student to explore customized material in a dynamic way through Acquisition, Annotation, and Augmentation (AAA4LLL). We have described how the users can choose their learning context by selecting seed topics, finding appropriate translations of interesting sentences from Wikipedia with the help of a transparent and traceable alignment technique, and inspecting these sentences by studying the individual parts, enriched with lexical, syntactic, conceptual, and relational annotation. The learning process can then be repeated by selecting new or additional topics.

As future work we will evaluate our learning solution in a classroom setting, for which we will involve graduate level language students. We will assess both system performance and learning outcomes by additionally employing novel evaluation approaches, such as learner centered development as described in [4]. Finally, we are planning to release our learning environment as open software together with instructive demos and an extensive documentation of the annotation formats.

References

- 1 Mikel Artetxe and Holger Schwenk. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. *Transactions ACL*, 7:597–610, 2019. doi:10.1162/tac1_a_00288/43523.
- 2 Laura Banarescu et al. Abstract meaning representation for sembanking. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 178–186. ACL, 2013. URL: <https://www.aclweb.org/anthology/W13-2322/>.
- 3 Marta Bañón et al. ParaCrawl: Web-scale acquisition of parallel corpora. In *Proceedings of the 58th Annual Meeting of the ACL*, pages 4555–4567. ACL, 2020. doi:10.18653/v1/2020.acl-main.417.
- 4 Hendrik Heuer and Daniel Buschek. Methods for the design and evaluation of HCI+NLP systems. *arXiv*, 2102.13461 [cs.CL], 2021. arXiv:2102.13461.
- 5 Ryu Iida, Kentaro Torisawa, Jong-Hoon Oh, Canasai Kruengkrai, and Julien Kloetzer. Intra-sentential subject zero anaphora resolution using multi-column convolutional neural network. In *Proceedings EMNLP 2016*, pages 1244–1254. ACL, 2016. doi:10.18653/v1/d16-1132.
- 6 Maki Kubota. Post study abroad investigation of kanji knowledge in Japanese as a second language learners. *System*, 69:143–152, 2017. doi:10.1016/j.system.2017.07.006.
- 7 Taku Kudo and Yuji Matsumoto. Fast methods for kernel-based text analysis. In *Proceedings of the 41st Annual Meeting of the ACL*, pages 24–31. ACL, 2003. doi:10.3115/1075096.1075100.
- 8 Taku Kudo, Kaoru Yamamoto, and Yuji Matsumoto. Applying conditional random fields to Japanese morphological analysis. In *Proceedings EMNLP 2004, ACL 2004*, pages 230–237. ACL, 2004. URL: <https://www.aclweb.org/anthology/W04-3230/>.
- 9 George A. Miller. WordNet: A lexical database for English. *Commun. ACM*, 38(11):39–41, 1995. doi:10.1145/219717.219748.
- 10 Makoto Morishita, Jun Suzuki, and Masaaki Nagata. JParaCrawl: A large scale web-based English-Japanese parallel corpus. In *Proceedings of LREC 2020*, pages 3603–3609. ELRA, 2020.
- 11 Joakim Nivre et al. Universal Dependencies v2: An evergrowing multilingual treebank collection. In *Proceedings of LREC 2020*, pages 4034–4043. European Language Resources Association, 2020. URL: <https://www.aclweb.org/anthology/2020.lrec-1.497/>.
- 12 Stephan Oepen et al. MRP 2020: The second shared task on cross-framework and cross-lingual meaning representation parsing. In *Proceedings of the CoNLL 2020 Shared Task: Cross-Framework Meaning Representation Parsing*, pages 1–22. ACL, 2020. doi:10.18653/v1/2020.conll-shared.1.
- 13 Simon Paxton. Kanji matters in a multilingual Japan. *The Journal of Rikkyo University Language Center*, 42:29–41, 2019.
- 14 Harald Wahl and Werner Winiwarter. A technological overview of an intelligent integrated computer-assisted language learning (iiCALL) environment. In *Proceedings of World Conference on Educational Multimedia, Hypermedia and Telecommunications*. AACE, 2011.
- 15 Werner Winiwarter. Mastering Japanese through augmented browsing. In *Proceedings of iiWAS 2013, iiWAS '13*, pages 179–188. ACM, 2013.
- 16 Werner Winiwarter. JAMRED: a Japanese Abstract Meaning Representation Editor. In *Proceedings of iiWAS 2015*, pages 11:1–11:5. ACM, 2015. doi:10.1145/2837185.2837246.
- 17 Werner Winiwarter. JILL: Japanese Incidental Language Learning. In *Proceedings of iiWAS 2015*, pages 9:1–9:9. ACM, 2015. doi:10.1145/2837185.2837191.
- 18 Bartholomäus Wloka. Identifying bilingual topics in Wikipedia for efficient parallel corpus extraction and building domain-specific glossaries for the Japanese-English language pair. In *Proceedings of LREC 2018*. ELRA, 2018.
- 19 Bartholomäus Wloka. *Automated Creation of Domain-Specific Bilingual Corpora for Machine Translation, focusing on Dissimilar Language Pairs*. PhD thesis, University of Vienna, 2020.