

Location, Creator and Membership-related Knowledge Acquisition from Wikipedia-based Information-rich Taxonomy for ConceptNet Expansion

Marek Krawczyk, Rafal Rzepka, Kenji Araki

Graduate School of Information Science and Technology, Hokkaido University, Sapporo, Japan
{marek,rzepka,araki}@ist.hokudai.ac.jp

Abstract

In this paper we present a method for extracting IsA assertions (hyponymy relations), AtLocation assertions (informing of the location of an object or place), LocatedNear assertions (informing of neighboring locations), CreatedBy assertions (informing of the creator of an object) and MemberOf assertions (informing of group membership) automatically from Japanese Wikipedia XML dump files. These assertions would be suitable for introduction to the Japanese part of the ConceptNet common sense knowledge ontology. We use the Hyponymy extraction tool v1.0, which analyzes definition, category and hierarchy structures of Wikipedia articles to extract IsA assertions and produce an information-rich taxonomy. From this taxonomy we extract additional information, in this case AtLocation, LocatedNear, CreatedBy and MemberOf types of assertions, using our original method. The presented experiments prove that we achieved our research goal on a large scale: both methods produce satisfactory results, and we were able to acquire 5,866,680 IsA assertions with 96.0% reliability, 131,760 AtLocation assertion pairs with 93.5% reliability, 6,217 LocatedNear assertion pairs with 98.5% reliability, 270,230 CreatedBy assertion pairs with 78.5% reliability and 21,053 MemberOf assertions with 87.0% reliability. Our method surpassed the baseline system in terms of both precision and the number of acquired assertions.

1 Introduction

In the information society of today it becomes increasingly difficult to have a creative and enjoyable conversation without a broad background knowledge. People like to discuss concepts and ideas, and this requires a mixture of commonsense and general knowledge spanning through a broad spectrum of topics. If we want machines to understand such conversations and interact with people using natural language, we need to equip them with large-scale general knowledge bases that will allow them to do so. A few examples of such bases include Cyc [Lenat, 1995], YAGO

[Suchanek *et al.*, 2007] and ConceptNet [Liu and Singh, 2004b]. In this paper we will focus on ConceptNet, a knowledge representation project that provides a large semantic graph describing general human knowledge. We have chosen ConceptNet as it captures a wide range of common sense concepts and relations, and its simple semantic network structure makes it easy to use and manipulate [Liu and Singh, 2004a]. ConceptNet was designed to contain knowledge collected by the Open Mind Common Sense project's website [Singh *et al.*, 2002], as well as knowledge from similar websites and online word games which automatically collect general knowledge in several languages. The current goal of ConceptNet is to expand the knowledge base with data mined from Wiktionary¹ and Wikipedia². This open-source knowledge base is used for many applications such as topic-gisting [Speer *et al.*, 2010], affect-sensing [Cambria *et al.*, 2010b], dialog systems [Korner and Brumm, 2009], daily activities recognition [Ullberg *et al.*, 2010], social media analysis [Cambria *et al.*, 2010a] and handwriting recognition [Wang *et al.*, 2013]. It is also applied to open-domain sentiment analysis as an integral element of a common and common sense knowledge core, which is then transformed into more compact multidimensional vector space [Cambria *et al.*, 2014]. Manual expansion of the knowledge base would be a long and labor-intensive process. For example, nadya.jp³, an online project that aims to gather knowledge by using a game with a purpose [Nakahara and Yamada, 2011], since its launch in 2010 has been able to introduce a little over 43,500 entries to ConceptNet. It is therefore evident that we need to develop automatic methods to gather new data.

There are already working projects, such as NELL [Carlson *et al.*, 2010] or KNEXT [Schubert, 2002], that aim to extract semantic assertions from unstructured text data found on the Internet. Alternatively, we could use information present in the existing semi-structured sources and transfer it into a knowledge base. As a considerable amount of human validation has already been involved in the process of creating such sources, the reliability of information gathered

¹A multilingual, web-based free content dictionary
<https://www.wiktionary.org/>

²A free-access, free content Internet encyclopedia
<https://www.wikipedia.org/>

³<http://nadya.jp/>

in this way would be considerably higher. Wikipedia is probably the best example of an open-source, large-scale information pool covering an extremely broad range of human knowledge. DBpedia project, apart from the previously-mentioned YAGO, also aims to transfer knowledge gathered in Wikipedia into a more formalized, digitally processable form [Mendes *et al.*, 2011]. English part of DBpedia has already been merged to ConceptNet. However, the Japanese part has not been transferred yet, leaving this part of the knowledge base at the size of roughly 1/10th of the English language domain. The problem with using the DBpedia repository is that the information gathering algorithms used to prepare the knowledge base were designed for multilingual input processing and therefore introduce a considerable amount of noise. It is therefore vital to widen the scope of the Japanese part independently, as the knowledge gathered in ConceptNet is in large part language-specific.

This paper elaborates on the efforts of [Krawczyk *et al.*, 2015]. We extended the scope of acquired assertions and explored the possibilities of deriving common sense knowledge from instance-related information triplets.

2 Hyponymy relation as IsA relation

To extract hyponymy relation pairs from Wikipedia's XML dump files we utilize the Hyponymy extraction tool v1.0⁴. The tool has been developed specifically to process Japanese language entries. It consists of four modules, three of which deal with extraction of hyponymy pairs from different parts of Wikipedia content: definition, category and hierarchy structures [Sumida and Torisawa, 2008]. The program uses the Pecco library⁵ (SVM-like machine learning tool) to assess the plausibility level of the extracted hyponymy relation pairs and boost the precision and recall of the system [Sumida *et al.*, 2008]. The extracted hyponymy pairs may be transferred to ConceptNet as two concepts related to each other by IsA relationship (Table 1 lists examples of the extracted pairs). According to [Yamada *et al.*, 2010] these pairs are not informative enough to be useful for NLP tasks such as Question Answering; however they do fall into the scope of ConceptNet, a domain representing common sense and general knowledge. They are simple enough not to interfere with the ConceptNet's usage flexibility, yet informative enough to introduce new and valuable input to the knowledge base.

3 Extracting other relations

The last, fourth module of the Hyponymy extraction tool v1.0 generates intermediate concepts of hyponymy relations using the output of the first three modules [Yamada *et al.*, 2010]. The tool executes the following procedure: first it acquires basic hyponymy relations from Wikipedia using the method proposed by [Sumida *et al.*, 2008]. Next, it augments each acquired hypernym with the title of the Wikipedia article

⁴<https://alaginrc.nict.go.jp/hyponymy/>

⁵<http://www.tkl.iis.u-tokyo.ac.jp/ynaga/pecco/>

⁶All Japanese language phrases are transliterated and written in italics.

Table 1: Examples of extracted 'IsA' relationship pairs.

Hypernym	Hyponym
<i>kouen</i> ⁶ (park)	<i>Motomiya-kouen</i> (Motomiya Park)
<i>kougu</i> (tool)	<i>baisu</i> (vice)
<i>Werudaa Bureemen-no senshu</i> (Werder Bremen player)	Klaus Allofs
<i>Nihon-no SF shouseitsu</i> (Japanese SF novel)	<i>Maikai Suikoden</i> (Hell's Water Margin)

from which the basic hyponymy relation was extracted and consolidates the basic hypernym with the newly generated augmented hypernym (so-called 'T-INTER'). Finally, it generates an additional intermediate concept ('G-INTER') by generalizing the enriched hypernym. As a result, it acquires four-level, information-rich hyponymy relations.

Examples of augmented hyponymy relations include: *tojo-jinbutsu* (character) – *SF eiga no tojo-jinbutsu* (character of SF movie) – *WALL-E no tojo-jinbutsu* (character of WALL-E) – M.O; *seihin* (product) – *kigyō no seihin* (product of a company) – *Silicon Graphics no seihin* (product of Silicon Graphics, Inc.) – IRIS Crimson; *sakuhin* (work) – *America no shosestu-ka no sakuhin* (work of American novelist) – *J.D. Salinger no sakuhin* (work of J.D. Salinger) – *A boy in France*; *machi* (town) – *England no shu no machi* (town in a county in England) – *East Sussex no machi* (town in East Sussex) – Uckfield. As we can see from the examples, the generated augmented hypernyms are too specific to be incorporated into ConceptNet directly. However some additional information about their corresponding hyponyms may be extracted from them, such as information concerning location, neighboring locations, creator, membership and so on. Knowledge about location, creator and membership may be directly transferred into ConceptNet through already built-in *AtLocation*, *LocatedNear*, *CreatedBy* and *MemberOf* relations. It should be noted that according to the ConceptNet documentation⁷ the *CreatedBy* relation relates to processes, however inspection of the existing *CreatedBy* assertions show that they include creations and their authors as well. The remaining part of the acquired information related to the hyponyms may be represented by a more general *RelatedTo* relation.

The procedure of acquiring additional information is presented in Figure 1 and exemplified in Figure 2. First (Step 1), we scan the G-INTER using our handcrafted primary rule base in search of tags referring to locations, creators or members, for example {city}⁸, {district}, {cartoonist}, {writer}, {member} and so on. In the case of acquiring *LocatedNear* pairs, we confirm that the basic hypernym contains a marker

⁷<https://github.com/commonsense/conceptnet5/wiki/Relations>

⁸Curly brackets were used to mark the tags' representations.

indicating physical proximity (such as the Chinese character meaning 'neighboring'). Next (Step 2), we filter the basic hyponym through a secondary rule base to exclude items that would introduce noise. For example, we can extract information about the birthplaces of famous people; however this does not mean that we can build an *ATLocation* kind of relationship between the person and his or her birthplace. If so, hyponyms indicating people are excluded from the analysis of location. When analysing *LocatedNear* pairs we filter out ambiguous items. If the basic hyponym is positively assessed by the secondary rule base, then (Step 3) we assume that the phrase acquired by deleting the basic hyponym from the *G-INTER* is a valid location, creator or member tag. Using the example from Figure 2, we check that 'adjacent municipality' is a valid tag to describe a nearby location. In the next stage (Step 4) we compare the validated location, creator or member tag with the content of the *T-INTER*. This way, using the previous example, we can extract the knowledge that the municipality we refer to is *Tomi-shi*. Finally (Step 5), we join the newly acquired information to the base hyponym with a proper relationship tag to extract a new relation, for example *Komoro-shi-LocatedNear-Tomi-shi*.

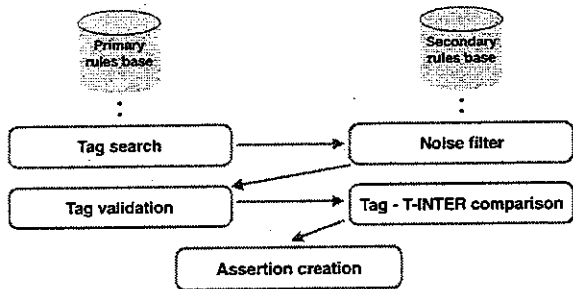


Figure 1: Flowchart of our proposed method.

The method mainly depends on the number and nature of introduced rules to both the primary and secondary rule bases to extract the required information effectively. At this stage we used 58 primary rules and 16 secondary rules, which allowed us to extract assertions concerning location, neighboring locations, creators and members. The manually crafted rules have been created using heuristics after analysis of the input data. We chose this kind of approach because the information units contain Chinese or Japanese characters indicating a type of location, a city, province, school, creator or a member. We use the rules to detect these characters, and this way we are able to obtain the named entities referring to locations, creators and members. Due to the qualities of the Japanese language's writing system these rules are often very simple, containing a single character, but are still effective for detecting the language units we want to extract. For example, the secondary rules used for detecting people include the suffix '~sha', which describes different professions. For English such a shortcut would be harder to apply, and therefore person detection would require a much larger rule base covering a long list of names of professions and appropriate suffixes (like '~er', '~or' or '~ist').

As our experiments revealed, extracting creator information is more complex and creates some challenges. While extracting location and member-related information, the introduced rules may be simple and straightforward. In the case of creators, the rules not only have to cover the qualities of the writing system, but also take into consideration the importance of particular roles while creating a given piece of work. For example our annotators indicated that a number of professionals taking part in the creation of films may not be considered as the creators of these films. Actors, actresses and voice actors, even if they make a great contribution to the work, should not be labeled as its creators. Further experiments showed that similarly animators, animation directors, sound directors, and storyboard creators, according to the annotators, do not qualify to be included in the common sense *CreatedBy* assertions. The question whether all these roles should be indeed excluded from the creator category is open to discussion. If we changed our perspective and considered that not only one person or role is to be credited as the creator of a given piece of work, then we could assess some of these roles as correct in the *CreatedBy* assertions. The problem of different opinions on this matter would however remain. As the algorithm bases on keywords, it is unable to distinguish, for example, between director and sound director. Such distinction would be possible if we employed an additional, concept-based knowledge base.

In future we would like to investigate the possibility of combining heuristics with automated rule discovery methods in order to achieve higher precision and recall. The number and reliability level of the data acquired with our method is presented in the Evaluation section.

4 Evaluation

To verify the reliability level declared by Sumida [Sumida *et al.*, 2008] and evaluate our proposed method for obtaining additional relations we used the 2014-11-04 version of the Japanese Wikipedia dump data. We ran the definition, category and hierarchy modules of the Hyponymy extraction tool v1.0 at 93% precision rate using the biggest available training set, and obtained 6,014,194 hyponym-hyponym pairs. The number of unique hyponymy pairs was 5,866,680, which indicates that 147,514 pairs have been extracted by more than one module. The 93% reliability level declared by the authors of the method has been verified by three human annotators, whose task was to evaluate a sample of the data and decide whether the extracted pairs a) represent a correct hyponymy relation, b) represent related concepts, but not in a hyponymy relation, or c) represent unrelated concepts. The annotators assigned 1, 0.5 and 0 points respectively to 300 randomly selected assertions. We decided to assign 0.5 points to related concepts as they may be used to create correct assertions (see Future Work section). If two or more annotators assessed an item as belonging to one category, their decision was regarded as the evaluation output. In cases where their decisions varied (which happened 10 times), the first author decided the score. The procedure follows

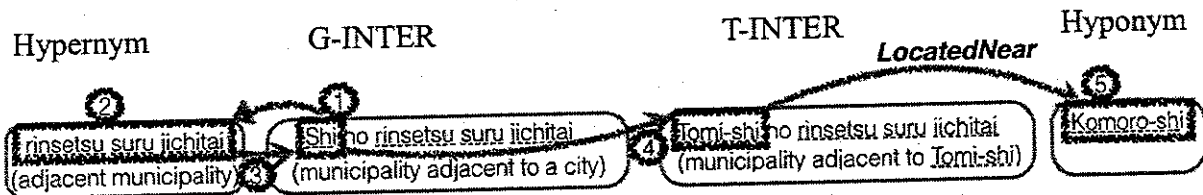


Figure 2: Procedure of our proposed method exemplified on the extracted relation.

a modified Sumida *et al.* [Sumida *et al.*, 2008] evaluation method.

Table 2 presents the evaluation results. 283 pairs were assessed as representing a correct hyponymy relation, 10 pairs as related concepts, but not in a hyponymy relation and 7 as unrelated concepts. This results in 96.0% precision value of the tested sample, which surpasses the 93% declared by Sumida *et al.* The level of overall agreement between annotators was 86.9%, and the Kappa value⁹ was 0.80, which indicates that the annotation judgement was in substantial agreement [Randolph, 2005].

Table 2: Evaluation results for IsA relations.

Correct hyponymy	Related concepts	Unrelated concepts	Precision	Total number of pairs
0.943 (283/300)	0.033 (10/300)	0.023 (7/300)	0.960	5,866,680

We obtained 2,738,211 basic hypernym-G-INTER-T-INTER-basic hyponym sets by running the fourth 'extended' module of the Hyponymy extraction tool v1.0 on the same Wikipedia dump data. By applying our method for extracting additional information, we were able to produce 131,760 pairs representing AtLocation relation, 6,217 pairs representing LocatedNear relation, 270,230 pairs representing CreatedBy relation and 21,053 pairs representing MemberOf relation. For comparison, nadya.jp, the baseline system, has provided only 8,706 AtLocation relations and no LocatedNear, CreatedBy or MemberOf relations in four years of its operation. In the case of AtLocation pairs, we evaluated 100 pairs¹⁰ randomly selected from our method's output and 100 pairs randomly selected from nadya.jp's AtLocation assertions [Nakahara and Yamada, 2011]. While evaluating LocatedNear, CreatedBy and MemberOf relations, a comparison with the baseline was not possible, as ConceptNet 5.3 does not yet contain any LocatedNear, CreatedBy or MemberOf pairs in its Japanese language section. These assertions were

⁹To measure the agreement level between judges, we used Randolph's free marginal multirater kappa instead of Fleiss' fixed-marginal multirater kappa, due to high agreement low kappa paradox.

¹⁰We adjusted the number of evaluated pairs to balance the proportion between the total number of pairs and the test sample.

therefore evaluated independently. The evaluation procedure follows the previously applied one: 1 point being applied to correct AtLocation, LocatedNear, CreatedBy or MemberOf assertions, 0.5 point to related concepts, but not in the evaluated relation, and 0 points to unrelated concepts. In 15 cases the annotators' evaluation was inconsistent, and therefore the first author decided the score.

Table 3 shows the evaluation results of our AtLocation pairs generation method in comparison with the baseline system. 88 pairs generated by our method were evaluated as representing a correct AtLocation relation, 11 pairs as related concepts, but not in an AtLocation relation, and 1 as unrelated concepts. This results in a 93.5% precision value. In the case of the baseline system, 64 pairs were evaluated as correct AtLocation assertions, 20 as related concepts, but not in an AtLocation relation, and 16 as unrelated concepts. The precision value for the baseline system is 74.0%. The level of overall agreement between annotators was 73.6% and the Kappa value was 0.60, which indicates that the annotation judgment was in moderate agreement. Examples of the extracted AtLocation assertions are presented in Table 4.

Table 3: Evaluation results for AtLocation relations in comparison with the nadya.jp baseline.

	Correct AtLocation	Related concepts	Unrelated concepts	Precision	Total number of pairs
Proposed	0.880 (88/100)	0.110 (11/100)	0.010 (1/100)	0.935	131,760
Baseline	0.640 (64/100)	0.200 (20/100)	0.160 (16/100)	0.740	8,706

$p < 0.001$, t -score = 4.6291

Table 5 contains the evaluation result of the generated LocatedNear relations. 97 pairs were evaluated as correct LocatedNear pairs, 3 as related concepts and none as unrelated concepts, which results in 98.5% precision. The level of overall agreement between annotators was 86.6% and the Kappa value was 0.80, which indicates that the annotation judgment was in substantial agreement. Examples of the extracted LocatedNear assertions are presented in Table 6.

Table 7 contains the evaluation result of the generated CreatedBy relations. 60 pairs were evaluated as correct CreatedBy pairs, 37 as related concepts and 3 as unrelated con-

Table 4: Examples of generated AtLocation assertions.

<i>Tomato Ginkou</i> (Tomato Bank)	AtLocation	<i>Okayama-shi</i> (Okayama city)
<i>Otao hoikuen</i> (Outao nursery)	AtLocation	<i>Sakai-shi</i> (Sakai city)
<i>Sandifukku</i> (Sandy Hook)	AtLocation	<i>Eriotto-gun</i> (Elliott County)
<i>Hoteru Kadoya</i> (Kadoya Hotel)	AtLocation	<i>Tochigi-shi</i> (Tochigi city)

Table 5: Evaluation results for LocatedNear relations

Correct Located-Near	Related concepts	Unrelated concepts	Precision	Total number of pairs
0.970 (97/100)	0.030 (3/100)	0.000 (0/100)	0.985	6,217

Table 6: Examples of generated LocatedNear assertions.

<i>Ougoe-machi</i> (Ogoe city)	LocatedNear	<i>Ono-machi</i> (Ono city)
<i>Iseri-gawa</i> (Iseri river)	LocatedNear	<i>Konoha-gawa</i> Konoha river
Daiting	LocatedNear	Monheim
<i>Kumotori-yama</i> (Mount Kumotori)	LocatedNear	<i>Karamatsuo-yama</i> (Mount Karamatsuo)

cepts, which results in 78.5% precision. The level of overall agreement between annotators was 71.6% and the Kappa value was 0.57, which indicates that the annotation judgment was in moderate agreement. Examples of the extracted CreatedBy assertions are presented in Table 8.

Table 7: Evaluation results for CreatedBy relations.

Correct CreatedBy	Related concepts	Unrelated concepts	Precision	Total number of pairs
0.600 (60/100)	0.370 (37/100)	0.030 (3/100)	0.785	270,230

The analysis of the relatively low precision score of the assessed CreatedBy assertions revealed the following: in 24 cases it was the annotators' opinion that actors, voice actors, animators, storyboard creators or sound directors cannot be considered as creators of works they contribute to. Although it would be valid to include such persons in the RelatedTo kind of relationship with the work they helped to create, defining them as creators would go against common sense. This is a valid observation and it will be taken into consider-

Table 8: Examples of generated CreatedBy assertions.

Dark Horse	CreatedBy	George Harrison
<i>Kaze</i> (Wind)	CreatedBy	<i>Kubota Koutarou</i>
<i>Manuke-na Oukami</i> (Sheep Wrecked)	CreatedBy	Michael Lah
The Point of View	CreatedBy	Alan Crosland

ation when re-designing and expanding the rule base for the next version of the algorithm. There were also cases of assertions assessed as invalid due to errors passed from the output of the Hyponymy extraction tool to the proposed method. Table 9 contains examples of assertions that were assessed as erroneous by the annotators.

Table 9: Examples of erroneous CreatedBy assertions.

Road 88	CreatedBy	<i>Tomita Yasuko</i> (actress)
<i>Kaiketsu Zorori</i> (Incredible Zorori)	CreatedBy	<i>Yamada Etsuji</i> (sound director)
<i>Kishin Douji Zenki</i> (Zenki)	CreatedBy	<i>Hayashi Akemi</i> (animator)
Human (incomplete name error)	CreatedBy	Nicholson Baker

Table 10 contains the evaluation result of the generated MemberOf relations. 76 pairs were evaluated as correct MemberOf pairs, 22 as related concepts and 2 as unrelated concepts, which results in 87.0% precision. The level of overall agreement between annotators was 80.6% and the Kappa value was 0.71, which indicates that the annotation judgment was in substantial agreement. Examples of the extracted MemberOf assertions are presented in Table 11.

Table 10: Evaluation results for MemberOf relations.

Correct MemberOf	Related concepts	Unrelated concepts	Precision	Total number of pairs
0.760 (76/100)	0.220 (22/100)	0.020 (2/100)	0.870	21,053

In the 13 cases the annotators decided that the generated MemberOf assertion refer to the former member of relative group, and therefore assigned it as the related concepts. The question whether these pairs should be considered as representing concepts in MemberOf relation is currently under discussion. If we would consider that the status of a member, once granted, is not temporary, then the precision

Table 11: Examples of generated MemberOf assertions.

Henning Schmitz	MemberOf	<i>Kurafutowaaku</i> (Kraftwerk)
Dir.F	MemberOf	<i>Suiyoubi no Kanpanera</i> (Wednesday Canpanella)
Oono Satoshi	MemberOf	<i>Arashi</i>
Nishimura Akihiro	MemberOf	<i>Nikkan Giin Renmei</i> (Japan-Korea Parliamentarians' Union)

rate of the tested sample would be higher, reaching 93.5%.

The results show that IsA relation pairs generated by the definition, category and hierarchy of the Hyponymy extraction tool v1.0, as well as AtLocation, LocatedNear and MemberOf relation pairs extracted by our proposed method may be incorporated into ConceptNet. Considering the number of the newly acquired assertions as well as reliability of the data in comparison with the resources already present in the knowledge base, such operation would be beneficial for ConceptNet. CreatedBy relation pairs could also be added after the revision of introduced rules and a substantial increase of the precision rate.

5 Generalizing over assertions

Wikipedia contains a lot of information related to instances of certain concepts, such as Salvador Dali as an instance of an artist. Filling up ConceptNet with instances is a valid task, as it is very hard to establish the boundaries of common sense knowledge – facts that are obvious to one group of people overlap to a large proportion with the knowledge of another group, but there is always a discrepancy. This issue raises a question: would it be possible to come to more general conclusions on the basis of the numerous instances? In order to solve this problem we created and performed an initial test of the following method: we took each of the additional information lists (representing LocatedAt, LocatedNear, CreatedBy and MemberOf relations) and analyzed each assertion one by one. For both concepts in the assertion we found their hypernyms in the generated IsA relations list. Next, we generated assertions representing all possible combinations between concept A's hypernyms and concept B's hypernyms. We repeated the process for all assertions in the additional information list and calculated the generated hypernym assertions' occurrence frequency. As predicted, the assertions with the highest occurrence frequency represent general, common sense observations. This is true for AtLocation, CreatedBy and MemberOf lists, but it is not the case when processing the LocatedNear list, because of the relatively low number of LocatedNear assertions. It became apparent that the higher number of initial assertions increases the probability of generating meaningful general assertions. See Table 12 for the examples of generated general assertions. The procedure requires further development in terms of the method for frequency calculations and automatic filtering of non-general assertions.

Table 12: Examples of generated general assertions.

<i>toshi oyobi machi</i> (city and town)	AtLocation	<i>gun</i> (province)
<i>shougakkou</i> (elementary school)	AtLocation	<i>machi</i> (city)
<i>sakuhin</i> (work)	CreatedBy	<i>zonmei jinbutsu</i> (living person)
<i>shutsuen sakuhin</i> (performance art)	CreatedBy	<i>bunkajin</i> (cultural figure)
<i>zonmei jinbutsu</i> (living person)	MemberOf	<i>Nihon no kashu gruupu</i> (Japanese singer group)
<i>owarai geinin</i> (comedian)	MemberOf	<i>manzaishi</i> (comic duo)

6 Conclusion

In this paper we presented a method for automatic acquisition of common sense knowledge triplets from the Japanese Wikipedia. It allowed us to mine IsA, AtLocation, LocatedNear, CreatedBy and MemberOf assertions with precision estimated at the levels of 96.0%, 93.5%, 98.5%, 78.5% and 87.0% respectively. We also demonstrated the possibility of formulating common sense assertions on the basis of generated instances data. As the Japanese part of the current ConceptNet 5.3 consists of 1,071,046 assertions, a contribution of 6,295,940 new assertions would be significant. It would mean an almost sixfold increase and could potentially make ConceptNet applicable to many Japanese language analysis problems. Moreover, as Wikipedia is a constantly expanding source, we could acquire more assertions simply by applying our method to the updated Wikipedia XML dump files.

7 Future work

In order to extend the functionality of our proposed method, we intend to update the primary and secondary rules, which would allow the system to increase its precision and the scope of extracted information. We would also like to explore the possibility of using a machine learning algorithm for automatic rule generation combined with the already present heuristics. Such a combination could potentially be more effective in increasing precision and recall, as well as finding new rules to extract even more relations.

We also plan to create an interface for the evaluation of the method's output by Japanese native speakers, which would allow us to utilize the pairs representing related concepts.

References

- [Cambria *et al.*, 2010a] Erik Cambria, Amir Hussain, Catherine Havasi, and Chris Eckl. Sentic computing: exploitation of common sense for the development of emotion-sensitive systems. In *Development of Multi-*

- modal Interfaces: Active Listening and Synchrony*, pages 148–156. Springer, 2010.
- [Cambria et al., 2010b] Erik Cambria, Amir Hussain, Catherine Havasi, and Chris Eckl. SenticSpace: visualizing opinions and sentiments in a multi-dimensional vector space. In *Knowledge-Based and Intelligent Information and Engineering Systems*, pages 385–393. Springer, 2010.
- [Cambria et al., 2014] Erik Cambria, Yangqiu Song, Haixun Wang, and Newton Howard. Semantic multidimensional scaling for open-domain sentiment analysis. *Intelligent Systems, IEEE*, 29(2):44–51, 2014.
- [Carlson et al., 2010] Andrew Carlson, Justin Betteridge, Bryan Kisiel, Burr Settles, Estevam R Hruschka Jr, and Tom M Mitchell. Toward an architecture for never-ending language learning. In *AAAI*, volume 5, page 3, 2010.
- [Korner and Brumm, 2009] Sven J Korner and Torben Brumm. Resi—a natural language specification improver. In *Semantic Computing, 2009. ICSC'09. IEEE International Conference on*, pages 1–8. IEEE, 2009.
- [Krawczyk et al., 2015] Marek Krawczyk, Rafal Rzepka, and Kenji Araki. Extracting conceptnet knowledge triplets from japanese wikipedia. In *Proceedings of the 21st Annual Meeting of The Association for Natural Language Processing*, pages 1052–1055, 2015.
- [Lenat, 1995] Douglas B Lenat. Cyc: A large-scale investment in knowledge infrastructure. *Communications of the ACM*, 38(11):33–38, 1995.
- [Liu and Singh, 2004a] Hugo Liu and Push Singh. Commonsense reasoning in and over natural language. In *Knowledge-based intelligent information and engineering systems*, pages 293–306. Springer, 2004.
- [Liu and Singh, 2004b] Hugo Liu and Push Singh. Conceptnet—a practical commonsense reasoning tool-kit. *BT technology journal*, 22(4):211–226, 2004.
- [Mendes et al., 2011] Pablo N Mendes, Max Jakob, Andrés García-Silva, and Christian Bizer. Dbpedia spotlight: shedding light on the web of documents. In *Proceedings of the 7th International Conference on Semantic Systems*, pages 1–8. ACM, 2011.
- [Nakahara and Yamada, 2011] Kazuhiro Nakahara and Shigeo Yamada. Development and evaluation of a web-based game for common-sense knowledge acquisition in japan. In *Unisys Technology Review no. 107*, pages 295–305. 2011.
- [Randolph, 2005] Justus J Randolph. Free-marginal multi-rater kappa (multirater k [free]): An alternative to fleiss' fixed-marginal multirater kappa. *Online Submission*, 2005.
- [Schubert, 2002] Lenhart Schubert. Can we derive general world knowledge from texts? In *Proceedings of the second international conference on Human Language Technology Research*, pages 94–97. Morgan Kaufmann Publishers Inc., 2002.
- [Singh et al., 2002] Push Singh, Thomas Lin, Erik T Mueller, Grace Lim, Travell Perkins, and Wan Li Zhu. Open mind common sense: Knowledge acquisition from the general public. In *On the Move to Meaningful Internet Systems 2002: CoopIS, DOA, and ODBASE*, pages 1223–1237. Springer, 2002.
- [Speer et al., 2010] Robert H Speer, Catherine Havasi, K Nichole Treadway, and Henry Lieberman. Finding your way in a multi-dimensional semantic space with luminoso. In *Proceedings of the 15th international conference on Intelligent user interfaces*, pages 385–388. ACM, 2010.
- [Suchanek et al., 2007] Fabian M Suchanek, Gjergji Kasneci, and Gerhard Weikum. Yago: a core of semantic knowledge. In *Proceedings of the 16th international conference on World Wide Web*, pages 697–706. ACM, 2007.
- [Sumida and Torisawa, 2008] Asuka Sumida and Kentaro Torisawa. Hacking wikipedia for hyponymy relation acquisition. In *IJCNLP*, volume 8, pages 883–888. Citeseer, 2008.
- [Sumida et al., 2008] Asuka Sumida, Naoki Yoshinaga, and Kentaro Torisawa. Boosting precision and recall of hyponymy relation acquisition from hierarchical layouts in wikipedia. In *LREC*, 2008.
- [Ullberg et al., 2010] Jonas Ullberg, Silvia Coradeschi, and Federico Pecora. On-line adl recognition with prior knowledge. In *Proceedings of the 2010 conference on STAIRS 2010: Proceedings of the Fifth Starting AI Researchers' Symposium*, pages 354–366. IOS press, 2010.
- [Wang et al., 2013] Qiu-Feng Wang, Erik Cambria, Cheng-Lin Liu, and Amir Hussain. Common sense knowledge for handwritten chinese text recognition. *Cognitive Computation*, 5(2):234–242, 2013.
- [Yamada et al., 2010] Ichiro Yamada, Chikara Hashimoto, Jong-Hoon Oh, Kentaro Torisawa, Kow Kuroda, Stijn De Saeger, Masaaki Tsuchida, and Junichi Kazama. Generating information-rich taxonomy from wikipedia. In *Universal Communication Symposium (IUCS), 2010 4th International*, pages 97–104. IEEE, 2010.

25