

1 International Journal of Cooperative Information Systems
2 Vol. 21, No. 4 (2012) 1–21
3 © World Scientific Publishing Company
4 DOI: 10.1142/S0218843013500019



5 **INFORMATION OVERLAP IN MULTILINGUAL**
6 **WIKIPEDIA AND SUMMARIZATION**

7 ELENA FILATOVA
8 *Computer and Information Sciences Department*
9 *Fordham University, 441 East Fordham Road*
10 *Bronx NY 10458, USA*
11 *filatova@cis.fordham.edu*

12 Received 12 November 2011
13 Accepted 2 October 2012
14 Published

15 Wikipedia is used as a training corpus for many information selection tasks: summa-
16 rization, question-answering, etc. The information presented in Wikipedia articles as
17 well as the order in which this information is presented, is treated as the gold stan-
18 dard and is used for improving the quality of information selection systems. However,
19 the Wikipedia articles corresponding to the same entry (person, location, event, etc.)
20 written in different languages have substantial differences regarding what information is
21 included in these articles. In this paper we analyze the regularities of information overlap
22 among the articles about the same Wikipedia entry written in different languages: some
23 information facts are covered in the Wikipedia articles in many languages, while others
24 are covered only in a few languages. We introduce a hypothesis that the structure of
25 this information overlap is similar to the information overlap structure (pyramid model)
26 used in summarization evaluation, as well as the information overlap/repetition struc-
27 ture used to identify important information for multidocument summarization. We prove
28 the correctness of our hypothesis by building a summarization system according to the
29 presented information overlap hypothesis. This system summarizes English Wikipedia
30 articles given the articles about the same Wikipedia entries written in other languages.
31 To evaluate the quality of the created summaries, we use Amazon Mechanical Turk as
32 the source of human subjects who can reliably judge the quality of the created text. We
33 also compare the summaries generated according to the information overlap hypothesis
34 against the lead line baseline which is considered to be the most reliable way to generate
35 summaries of Wikipedia articles. The summarization experiment proves the correctness
36 of the introduced multilingual Wikipedia information overlap hypothesis.

37 *Keywords:* Summarization; crowdsourcing; information overlap; multilingual wikipedia.

38 **1. Introduction**

39 Many information selection applications utilize multiple descriptions of the same
40 object and often rely on information repetition in these multiple descriptions. Sum-
41 marization applications use multiple descriptions of the input entry (person, event,
etc.) to extract the most important information about the input entry. Answer

2 *E. Filatova*

1 validation applications use the information extracted from different documents to
2 identify which of the suggested answers is the correct one, etc.

3 In addition to processing multiple descriptions about the same entry, most infor-
4 mation selection systems require the gold standard against which their output can
5 be compared. Typically, the gold standard is created by human annotators: correct
6 answers to questions, model summaries, etc. Creating the gold standard is often
7 expensive and time-consuming as it requires a lot of hours of work conducted by
8 trained human annotators.

9 Wikipedia^a provides articles about people, events, locations, etc. in many lan-
10 guages. It is constantly monitored by Wikipedia editors to ensure that Wikipedia
11 articles comply with the Wikipedia standards and is up-to-date. Many informa-
12 tion selection systems utilize these Wikipedia features assuming that the informa-
13 tion presented in Wikipedia is correct. Some systems use Wikipedia content as
14 the source of reliable and up-to-date information, others utilize the structure of
15 Wikipedia articles (see Sec. 2).

16 In this work, we raise an issue that certain Wikipedia peculiarities should be
17 taken into consideration by the researchers that use Wikipedia as the source of
18 the gold standard for information selection. Namely, most Wikipedia entries have
19 articles in several languages. These articles are not translations of a Wikipedia arti-
20 cle from one language into other languages. Rather, Wikipedia articles in different
21 languages are independently created by different users. Thus, these documents are
22 not *identical* but rather can be treated as *different descriptions* of the same entry.
23 At the same time, despite all the differences, the Wikipedia articles about the same
24 entry written in different languages obviously have a certain degree of information
25 overlap. Some information facts are repeated in the articles in many languages while
26 others are covered in the articles in only a handful of languages.

27 Given the fact that Wikipedia articles about the same entry written in differ-
28 ent languages contain different sets of facts, the question arises whether or not
29 Wikipedia can be used as the gold standard for information selection tasks. The
30 issue of the trustworthiness of the Wikipedia information is currently studied for
31 monolingual Wikipedias. Obviously, this issue becomes even more complex when
32 multilingual Wikipedia is analyzed.

33 We compare descriptions of Wikipedia entries written in different languages and
34 investigate their information overlap pattern. We introduce a hypothesis regarding
35 the pattern of the information overlap in multilingual Wikipedia. According to our
36 hypothesis: information overlap in Wikipedia articles about the same entry written
37 in different languages corresponds well to the pyramid summarization evaluation
38 model.^{1,2}

39 To show the correctness of this hypothesis we create a summarization system
40 that follows our information overlap hypothesis. As most extractive summariza-
41 tion systems work with sentences as the summary generation units, in our work

^a<http://en.wikipedia.org/wiki/Wikipedia>

1 we use a sentence as a proxy for an information fact: the more articles in different
2 languages contain a certain sentence (information fact) — the higher the prior-
3 ity for this sentence (information fact) to be included into the output summary.
4 The summarization system that follows our multilingual Wikipedia information
5 overlap hypothesis is close in nature to the multidocument summarization systems
6 that rely on information repetition to rank information facts according to their
7 importance.^{3,4}

8 The summarization system created for our hypothesis testing combines both
9 single- and multidocument summarization aspects: it summarizes a single docu-
10 ment (English Wikipedia article about a certain entry) given multiple reference
11 documents (Wikipedia articles about the same entry written in different languages).
12 The evaluation experiment shows that the summarization system built according
13 to our multilingual Wikipedia information overlap hypothesis has a high level of
14 user satisfaction. We believe this result does not only prove the correctness of our
15 hypothesis, but also helps the understanding of the combined value of multilingual
16 Wikipedia entry descriptions. On the one hand, in order to get a complete pic-
17 ture about a Wikipedia entry, the articles in all languages should be combined.^b
18 On the other hand, the structure of the multilingual Wikipedia information over-
19 lap can be used for a variety of NLP tasks, including summarization, information
20 trustworthiness estimation, etc.

21 The rest of the paper is structured as follows. In Sec. 2, we describe related
22 work. In Sec. 3, we describe in detail our hypothesis about the nature of the infor-
23 mation overlap in multilingual Wikipedia. In Sec. 4, we describe our corpus, the
24 summarization-based experiments that we use to prove our hypothesis regarding
25 the structure of the information overlap in multilingual Wikipedia, and discuss the
26 results of this experiment. In Sec. 5, we outline the avenues for future research.

27 2. Related Work

28 Wikipedia has become a popular corpus for training and testing information selec-
29 tion systems.

30 Ahn *et al.*⁷ use Wikipedia “both as a source of answers to factoid questions and
31 as an importance model”. Buscaldi *et al.*⁸ use Wikipedia for answer verification;
32 they also use “Wikipedia *categories* in order to determine a set of patterns that
33 should fit with the expected answer”. Ko *et al.*⁹ use Wikipedia for measuring answer
34 relevance and as a source of synonyms. In Refs. 7 and 8 monolingual (English and
35 Spanish, respectively) Wikipedias are used. Ko *et al.*⁹ use English, Chinese and
36 Japanese Wikipedias. However, these three corpora are used independently and
37 thus, one can say that Ref. 9 uses three versions of monolingual Wikipedia. In
38 these papers, Wikipedia articles are treated as the gold standard: the systems do
39 not question the correctness of the information presented in Wikipedia.

^bThe preliminary experiments and findings for this paper have been published in Refs. 5 and 6.

4 *E. Filatova*

1 Baidysy *et al.*¹⁰ use the structure of Wikipedia articles describing people to learn
2 the order in which facts should be presented in a biography. Their results show
3 that the summarization system that generates biographies following the structure
4 of Wikipedia articles describing people “significantly outperforms all systems that
5 participated in DUC 2004, according to the ROUGE-L metric, and is preferred by
6 human subjects”.

7 Many successful multidocument summarization systems heavily rely on infor-
8 mation repetition to produce good quality summaries.^{19–21} These systems use the
9 idea that an important information fact that should be included into the final sum-
10 mary is likely to be repeated in many input documents. The idea of the parallel
11 between information importance and its use in the summaries is the foundation of
12 the pyramid-based summarization evaluation system.² The pyramid-based summa-
13 rization evaluation approach relies on the model summaries created by humans: the
14 more humans included a certain information fact or summary content unit (SCU)
15 into the summary — the more important for a summarization system to include
16 this information fact into its output.

17 Nelken *et al.*¹¹ use article revision history for collection and correction of lexical
18 errors; for training sentence compression algorithms; and for discerning lexicalized
19 models.

20 The multilingual aspect of Wikipedia is used for a variety of text processing
21 tasks. The potential of Wikipedia as the source of parallel corpora is investigated
22 in Ref. 12. The authors construct an English–Dutch parallel corpus and describe
23 two ways of looking for similar sentences in Wikipedia pages (using matching trans-
24 lations and hyperlinks). Multilingual Wikipedia is used to annotate a large corpus
25 of text with Named Entity tags in Ref. 13; to facilitate cross-language IR in Ref. 14,
26 and to perform cross-lingual QA in Ref. 15.

27 The described applications do not question the trustworthiness of the informa-
28 tion presented in Wikipedia. In a separate line of research, approaches are developed
29 to rate the trustworthiness of Wikipedia information. These approaches, however,
30 deal with monolingual Wikipedias.

31 One approach to compute the trustworthiness of Wikipedia information is to
32 use the information from its edit history. Wikipedia content trustworthiness can
33 be estimated by using a combination of the amount of the content revised and the
34 author’s reputation performing this revision Ref. 16. Another way to use edit history
35 to estimate information trustworthiness is to treat Wikipedia article editing as a
36 dynamic process and to use a dynamic Bayesian network trust model that utilizes
37 rich revision information in Wikipedia.¹⁷ Wikipedia trustworthiness can also be
38 estimated by using an additional tab (*Trust tab*).¹⁸

39 The research closest to ours is presented in Ref. 22, where the main goal is to
40 use self-supervised learning to align and/or create new Wikipedia infoboxes across
41 four languages (English, Spanish, French, German). Wikipedia infoboxes contain
42 a small number of facts about Wikipedia entries in a semi-structured format. In
43 our work, however, we deal with plain text and disregard any structured data such

1 a infoboxes, tables, etc. Although, we and Adar *et al.* analyze different types of
2 Wikipedia information, we arrive to similar conclusions: the most trusted informa-
3 tion is repeated in the Wikipedia articles in different languages. At the same time,
4 no single article can be considered as the complete source of information about a
5 Wikipedia entry.

6 **3. Information Coverage in Multilingual Wikipedia**

7 Wikipedia is a prominent example of collaborative knowledge generation and shar-
8 ing. People contributing to Wikipedia articles' writing and verification are driven
9 by a variety of reasons, including career, social, ideology, etc.²³ All Wikipedia arti-
10 cles are monitored by the Wikipedia community, which ensures the control of the
11 quality and content of the articles.

12 Currently, Wikipedia has articles in more than 200 languages. The language that
13 contains the largest number of articles is English,^{12,24} but the size of non-English
14 Wikipedia is growing fast and represents a rich corpus.^c

15 Most current information selection applications that use Wikipedia as their
16 corpus analyze Wikipedia articles written only in one language and assume that
17 a Wikipedia article in any language is a reliable source of information. However,
18 according to our observations, Wikipedia articles about the same entry (person,
19 location, event, etc.) written in different languages frequently cover different sets of
20 facts. Studying these differences can boost the development of various NLP appli-
21 cations (i.e. summarization, QA, new information detection, machine translation,
22 etc.). According to our Wikipedia analysis, there are two major sources of differ-
23 ences in the articles about the same Wikipedia entry written in different languages:

- 24 • the amount of the information covered by the Wikipedia articles (the length of
25 the Wikipedia articles);
- 26 • the choice of the information covered by the Wikipedia articles.

27 **3.1. Wikipedia article length difference**

28 The lengths of the Wikipedia articles about the same entry written in different
29 languages vary. Perhaps, the amount of information in Wikipedia articles depends
30 on how important a particular Wikipedia entry is for the community of people
31 speaking in a particular language. In this work, the length of a Wikipedia article is
32 measured in sentences used in the text description of a Wikipedia entry.

33 For example, baseball is popular in the USA, Latin America, and Japan but it
34 is not widely spread in, for example, Europe or Africa. Thus, the articles about a
35 legendary baseball player *Babe Ruth* exist in 26 languages.^d The longest and the

^chttp://meta.wikimedia.org/wiki/List_of_Wikipedias

^dWikipedia is changing constantly. The Wikipedia examples and data analyzed in this paper were collected on February 10, 2009, between 14:00 and 21:00 PST.

6 *E. Filatova*

1 most detailed articles are the English and Japanese ones. The article about *Babe Ruth*
2 *Ruth* in Croatian has one sentence, in Lithuanian — two sentences, in Finnish —
3 six sentences. These short articles only list a few general biographical facts such as:
4 date of birth, death; the fact that *Babe Ruth* was a legendary baseball player.

5 It is likely that the facts from the Lithuanian, Croatian, and Finnish articles
6 about *Babe Ruth* (e.g. general biography facts) will be listed in a summary of the
7 English language Wikipedia article about him. Moreover, all the sentences from
8 the Lithuanian, Croatian, and Finnish articles about *Babe Ruth* have corresponding
9 sentence in the English article about him and neither Lithuanian, nor Croatian, nor
10 Finnish articles have information that is not covered in the English article.

11 **3.2. Choice of information**

12 At the same time, there exist Wikipedia entries whose shorter articles contain infor-
13 mation facts that are not covered by longer articles. For example, the Wikipedia
14 entry about *Isadora Duncan* has articles in 47 languages. The lengths of these arti-
15 cles vary greatly: from almost 150 sentences in the English language article to four
16 sentences in Danish. *Isadora Duncan* was an American-born dancer who was very
17 popular in Europe and was married to a Russian poet, *Sergey Esenin*. Certain facts
18 (i.e. major biography dates) about *Isadora Duncan* are repeated in the articles in
19 every language. However, shorter articles are not necessarily summaries of longer
20 articles. For example, the article in Russian is almost four times shorter than the
21 article in English, however, it contains information that is not covered in the article
22 written in English. The same can be noted about the articles in French and Spanish.

23 **4. Testing the Hypothesis about the Multilingual Wikipedia**

24 **Information Overlap Structure: Summarization Experiment**

25 Despite the differences in length and information coverage described in Sec. 3, it is
26 obvious that Wikipedia articles about the same entry written in different languages
27 have a substantial amount of information that is common for articles in all (or in
28 most) languages.

29 In our work, we suggest the hypothesis that the information overlap in multilin-
30 gual Wikipedia corresponds well to the pyramid summarization evaluation model as
31 well as the multidocument summarization information importance model. To prove
32 this hypothesis we build a summarization system based on the presented multi-
33 lingual Wikipedia information overlap hypothesis. We use this system to generate
34 summaries of the English Wikipedia articles given the articles about the same entry
35 in other languages.

36 **4.1. Summarization and wikipedia: Information overlap and**

37 **information position**

38 Currently, information repetition (redundancy) is successfully used for both sum-
mary generation and summarization evaluation. The state-of-the-art (automatic

1 and manual) summarization evaluation approaches compare the automatically gen-
2 erated summaries against several model summaries. Such models are typically
3 created manually. The more model summaries contain a certain information fact —
4 the greater value it gets in the automatically generated summary.^{1,2,25}

5 Many current successful multidocument summarization systems rely on the
6 observation that the information facts that are good candidates to be included
7 into the summary are likely to be repeated in many input documents.^{3,4} At the
8 same time, in the field of single document summarization there exists a consensus
9 that taking the first several sentences is typically the best way to summarize a
10 news article. This single news article summarization approach is often used as an
11 evaluation baseline, called the *lead line* baseline, and is very difficult to beat.²⁶ The
12 strength of this baseline is due to the rules that are typically applied by journalists
13 to structure their news articles.

14 In contrast to news summarization, the lead section of Wikipedia articles can
15 be used not as a baseline but rather as a human generated model summary.
16 Wikipedia articles have a well-defined structure. Even though authors can write
17 Wikipedia articles disregarding the established conventions, the crowdsourcing
18 nature of Wikipedia leads toward the enforcement of the Wikipedia Manual of
19 Style^e for all Wikipedia articles. This phenomenon is successfully used by NLP sys-
20 tems to learn the typical structure of Wikipedia articles in general²⁷ and biograph-
21 ical Wikipedia articles in particular.¹⁰ Specifically, the lead section^f of Wikipedia
22 articles “serves both as an introduction to the article and as a summary of its most
23 important aspects. The lead should be able to stand alone as a concise overview
24 of the article”. In our experiment, we use this peculiarity of the Wikipedia article
25 structure as the source of human generated model summaries.

26 In the summarization experiment that we set up to prove the correctness of our
27 multilingual Wikipedia information overlap structure hypothesis we use the infor-
28 mation repetition degree in the articles about the same Wikipedia entry written in
29 different languages to identify those sentences from the English articles about this
30 entry that are important. In addition, we rely on the typical structure of Wikipedia
31 articles and output the identified important sentences in the order they appear in
32 the English article. We use a sentence as a proxy for an information fact because
33 most existing extractive multidocument summarizers use sentences as main gener-
34 ation units.

35 4.2. Data set

36 For our experiment, we use the list of people created for the Task 5 of DUC 2004:
37 biography generation task (48 people).^g First, we downloaded from Wikipedia all
38 the articles in all the languages corresponding to each person from the DUC 2004

^ehttp://en.wikipedia.org/wiki/Wikipedia:Manual_of_Style

^f[http://en.wikipedia.org/wiki/Wikipedia:Manual_of_Style_\(lead_section\)](http://en.wikipedia.org/wiki/Wikipedia:Manual_of_Style_(lead_section))

^g<http://duc.nist.gov/duc2004/tasks.html/>

1 evaluation. We used Wikitext, the text that is used by Wikipedia authors and edi-
2 tors. Wikitext complies with the wiki markup language and can be processed by
3 the Wikimedia content manager system into HTML which can then be viewed in a
4 browser. This is the text that can be obtained through Wikipedia dumps.^h For our
5 experiment, we removed from the wikitext all the markup tags and tabular infor-
6 mation (e.g. infoboxes and tables) and kept only plain text. There is no commonly
7 accepted standard wikitext language, thus our final text had a certain amount of
8 noise.

9 For each Wikipedia entry (i.e. DUC 2004 person) we downloaded the corre-
10 sponding articles in all the languages, including Esperanto, Latin, etc. We used
11 the name of a person to find the article about this person in English and then we
12 followed the links from the left side panel of the Wikipedia page template to get
13 the articles in other languages. To facilitate the comparison of the articles written
14 in different languages we used the Google machine translation toolⁱ to translate the
15 downloaded articles into English. The number of languages covered currently by
16 the Google translation system (41) is less than the number of languages used in
17 Wikipedia (265). However, the language distribution in the collected corpus corre-
18 sponds well to the language distribution in Wikipedia and the collected Wikipedia
19 subset can be considered a representative sample.⁵

20 Five people from the DUC 2004 set had only English Wikipedia articles: *Paul*
21 *Coverdell*, *Susan McDougal*, *Henry Lyons*, *Jerri Nielsen*, *Willie Brown*. Thus, they
22 were excluded from the analysis. The person whose Wikipedia entry had articles
23 in the most languages (86) was *Kofi Annan*. On average, a Wikipedia entry for a
24 DUC 2004 person had articles in 25.35 languages. The article in English was not
25 always the longest article: in 17 cases the longest article of a Wikipedia entry for a
26 DUC 2004 person was in a language other than English.

27 The Google machine translation system does not have publicly available eval-
28 uation numbers for all the language pairs. However, none of the freely available
29 research systems cover as many languages as the Google machine translation system
30 does. Given our information overlap hypothesis and the design of our summariza-
31 tion experiment we opt to use the Google machine translation system. It is more
32 important for us to analyze Wikipedia articles in as many languages as possible,
33 even though their translations into English might be poor, rather than have only a
34 few high quality translations into English.

35 It must be noted that there exist summarization systems that specifically deal
36 with the multilingual summarization task. Evans *et al.*²⁸ generate English sum-
37 maries from English, Russian, and Japanese sources. Methods for identifying sim-
38 ilarities and differences across English and Arabic news articles are studied in
39 Ref. 29. Litvak *et al.*³⁰ work with two monolingual as well as with a bilingual
40 English–Hebrew corpora. There exist many other systems that produce high quality

^h<http://download.wikimedia.org/>

ⁱ<http://translate.google.com/>

1 summaries given input documents written in a specific set of languages. However,
2 the main goal of our paper is to test the hypothesis that information overlap in the
3 multilingual Wikipedia has a pyramid-like structure and thus, we need to use many
4 more languages that are currently handled by the most successful multilingual mul-
5 tidocument summarizers. Also, the output summary of our system is a summary
6 of a single document (Wikipedia article written in English) based on the informa-
7 tion overlap among the articles about this entry currently existing in multilingual
8 Wikipedia.

9 **4.3. Data processing tools**

10 After the Wikipedia articles about the DUC 2004 set were collected and translated,
11 we divided these articles into sentences using the LingPipe sentence chunker.³¹
12 For each DUC 2004 person we compared an article about this person in English
13 against the articles about this person in the other languages that were handled by
14 the Google translation system. We counted articles in how many languages had
15 sentences corresponding to the sentences in the articles in English. To identify
16 matching sentences we used the LingPipe string matching tool based on TF/IDF
17 distance which “is based on vector similarity (using the cosine measure of angular
18 similarity) over dampened and discriminatively weighted term frequencies. [...] two
19 strings are more similar if they contain many of the same tokens with the same
20 relative number of occurrences of each. Tokens are weighted more heavily if they
21 occur in few documents”.³¹ Before the similarity measure is applied to sentences, we
22 eliminate all the stop words from the sentences. For our experiment, each sentence
23 was treated as a separate document. The IDF value was computed based on the two
24 articles under consideration (one — in English, the other one — translation into
25 English). We used three similarity thresholds: 0.5, 0.35, 0.2.

26 **4.4. What was measured**

27 To evaluate how much information is repeated in the articles about the same per-
28 son in different languages, we measure the similarity of the person’s description in
29 English and in other languages. Since each sentence is treated as a separate docu-
30 ment, the number of tokens (words) for comparison is rather small. Thus, for the
31 0.5 similarity threshold, the sentences marked as similar are almost identical. The
32 0.35 and 0.2 thresholds allow us to search for nonidentical sentences that still have
33 a substantial word overlap.

34 Our hypothesis is that those information facts (sentences) that are mentioned in
35 the articles about a person in different languages fit well the pyramid summarization
36 model. For example, if we are to summarize an article about a person from the
37 English Wikipedia: first, we should add to the summary those sentences that have
38 their counterparts in the most number of articles about this person in the languages
39 other than English. Sentences added on this step correspond to the top level of the
40 pyramid and represent the most important part of the English article about the

Table 1. Algorithm outline.

| Algorithm | |
|-----------|--|
| 1. | Submit the person's name to Wikipedia |
| 2. | Get Wikipedia articles about this person in all languages |
| 3. | Remove nonplain text information from the articles |
| 4. | For all the languages handled by the Google MT, translate the articles from these languages into English |
| 5. | Break English articles (the original English article and the translations into English from other languages) into sentences |
| 6. | Identify what sentences from the original English-language article have similar sentences in the articles about the same person written in other languages |
| 7. | Rank all the sentences from the original English-language article according to the number of languages that have similar sentences |
| 8. | If several sentences are placed on the same level (have the same rank), list these sentences in the order they appear in the original English-language Wikipedia article |
| 9. | Use the top three levels from the above ranking |

1 person under analysis. If the length of the summary is not exhausted, then on the
 2 next step, we add to the summary the sentences that appear in the next most
 3 number of languages, and so on. Thus, we can place sentences on different levels of
 4 the pyramid, with the top level being populated by the sentences that appear in
 5 the most languages and the bottom level having sentences that appear in the least
 6 number of languages. For our experiments, we used sentences from the top three
 7 levels of this pyramid. All of the sentences added to the summary should appear in
 8 at least two languages other than English. Table 1 has a schematic outline of the
 9 described algorithm.

10 **4.5. Output example and experiment discussion**

11 Table 2 presents three-level summaries for the English Wikipedia article about *Gene*
 12 *Autry*. Wikipedia has articles about *Gene Autry* in 11 languages: in English and in
 13 ten other languages, each of which can be translated into English by the Google
 14 Translation system.

15 When using a similarity of 0.5 we get one sentence from the English article that
 16 has counterparts in at least two other languages (here, in three other languages).
 17 This sentence's ID is 0: it is the first, introductory sentence of the English Wikipedia
 18 article about *Gene Autry*.

19 When using a similarity 0.35 we get a summary consisting of six sentences.
 20 The sentence on the top level is the same sentence that was listed in the previous
 21 summary. However, having a more permitting similarity threshold, this sentence
 22 was mapped to similar sentences in seven languages, rather than in three. The
 23 next level consists of the sentences from the English article that were mapped to
 24 sentences in three other languages. Sentences on the third level were mapped to
 25 sentences from two other languages. It is interesting to notice that the sentences
 26 included in the summary are coming from different parts of the document that has
 27 88 sentences.

Table 2. Three-level summaries for Gene Autry (#: the summary level; *Lang.*: number of languages that contain a similar sentence; *ID*: the position of the sentence in the English article; *Text*: the sentence itself).

| # | Lang. | ID | Text |
|-----------------|-------|----|--|
| Similarity 0.5 | | | |
| 1 | 3 | 0 | Orvon Gene Autry (September 29, 1907 — October 2, 1998) was an American performer, who gained fame as the Singing Cowboy on the radio, in movies and on television |
| Similarity 0.35 | | | |
| 1 | 7 | 0 | Orvon Gene Autry (September 29, 1907 — October 2, 1998) was an American performer, who gained fame as “The Singing Cowboy” on the radio, in movies and on television |
| 2 | 3 | 1 | Autry, the grandson of a Methodist preacher, was born near Tioga, Texas. |
| | | 13 | His first hit was in 1932 with “That Silver-Haired Daddy of Mine”, a duet with fellow railroad man, Jimmy Long |
| 3 | 2 | 3 | After leaving high school in 1925, Autry worked as a telegrapher for the St. Louis-San Francisco Railway |
| | | 14 | Autry also sang the classic Ray Whitley hit “Back in the Saddle Again”, as well as many Christmas songs including “Santa Claus Is Coming to Town”, his own composition “Here Comes Santa Claus”, “Frosty the Snowman”, and arguably his biggest hit “Rudolph the Red-Nosed Reindeer” |
| | | 72 | Gene Autry died of lymphoma at age 91 at his home in Studio City, California and is interred in the Forest Lawn, Hollywood Hills Cemetery in Los Angeles, California |
| Similarity 0.2 | | | |
| 1 | 7 | 0 | Orvon Gene Autry (September 29, 1907 — October 2, 1998) was an American performer, who gained fame as “The Singing Cowboy” on the radio, in movies and on television |
| 2 | 6 | 1 | Autry, the grandson of a Methodist preacher, was born near Tioga, Texas |
| | | 73 | His death on October 2, 1998 came nearly three months after the death of another celebrated cowboy of the silver screen, radio, and TV, Roy Rogers |
| 3 | 5 | 21 | From 1940 to 1956, Autry had a huge hit with a weekly radio show on CBS, “Gene Autry’s Melody Ranch”. His horse, Champion, also had a radio-TV series “The Adventures of Champion” |
| | | 72 | Gene Autry died of lymphoma at age 91 at his home in Studio City, California and is interred in the Forest Lawn, Hollywood Hills Cemetery in Los Angeles, California |

1 The summary created using the 0.2 threshold contains the introductory sentence
2 as well as sentences not included in the summaries for the 0.5 and 0.35 similarity
3 threshold.

4 Despite the fact that for our experiment we chose the set of people used for the
5 DUC 2004 biography generation task, we could not use the DUC 2004 model sum-
6 maries for our evaluation. These models were created using the DUC 2004 corpus,
7 while in our experiments we used a subset of multilingual Wikipedia. Moreover,
8 Wikipedia entry articles about the DUC 2004 people had dramatic updates since

12 *E. Filatova*

1 2004. For example, *Jörg Haider* died of injuries from a car crash on October 11,
2 2008 and this information is included into our three-level summaries.

3 In the experiment described in this paper, we analyze only the sentences from the
4 English text that appear in at least two other languages. Thus, using the DUC 2004
5 set we created:

- 6 • for the similarity score of 0.2
 - 7 — **one**-level summaries for two people^j;
 - 8 — **two**-level summaries for two people^k;
 - 9 — **three**-level summaries for 38 people.^l
- 10 • for the similarity score of 0.35
 - 11 — **one**-level summaries for three people;
 - 12 — **two**-level summaries for four people;
 - 13 — **three**-level summaries for 35 people.
- 14 • for the similarity score of 0.5
 - 15 — **one**-level summaries for five people;
 - 16 — **two**-level summaries for eight people;
 - 17 — **three**-level summaries for 25 people.

18 From the above numbers it is clear that not all the people under analysis have
19 articles in two languages other than English with sentences that are similar to the
20 sentences in the English article. Given that on the sentence level the similarity of 0.5
21 corresponds to almost identical sentences, we fail to generate summaries for several
22 people.

23 The above results should be read as follows: given the similarity of 0.2, we can
24 create summaries for 42 people: 38 summaries have sentences on all three levels,
25 two summaries have sentences only on the first two levels, and the remaining two
26 summaries have sentences only on the first level. The length of the created sum-
27 maries is measured in sentences. Table 3 presents information about the average
28 number of sentences on each level.

Table 3. Average number of sentences for each level.

| Similarity Measure | Level One | Level Two | Level Three |
|--------------------|-----------|-----------|-------------|
| 0.2 | 1 | 2.2 | 2.42 |
| 0.35 | 1 | 1.6 | 2 |
| 0.5 | 1 | 7 | 6 |

^j**One**- level summary: only one level (Level 1) is filled in.

^k**Two**- level summary: two levels (Level 1 and 2) are filled in.

^l**Three**-level summary: all three levels (Level 1, 2 and 3) are filled in.

4.6. Model summaries

To prove the correctness of our hypothesis, in addition to the summaries generated following the algorithm outlined in Table 1, we create model lead line summaries for the English Wikipedia articles about the DUC 2004 people. We use three sets of model summaries created for each of the similarity measures (0.2, 0.35, 0.5). For a given person, similarity measure, and summary level, we extract from the English Wikipedia article about the person under analysis as many sentences as are added to the respective summary level by our summarization system based on the multilingual Wikipedia information overlap structure hypothesis.

As described in Sec. 4.2, to analyze the information overlap in Multilingual Wikipedia we use pre-processed Wikitext. To generate model summaries we do not use this pre-processed text, rather we manually extract the first sentences from the Wikipedia articles about the DUC 2004 people written in English. As described in Sec. 4.1, the lead section of a Wikipedia article can be used as the summary of this article. However, according to the above procedure, the model summaries can contain only a part of the lead section.

Table 4 presents the model (lead section) three-level summaries for the English Wikipedia article about *Gene Autry*. The number of sentences on each level is the same as the number of sentences on the respective levels of summaries presented in Table 2.

4.7. Evaluation procedure

To prove the validity of our hypothesis, we evaluate the quality of the summaries produced by our system that exploits information overall in multilingual Wikipedia and compare it to the quality of the model summaries. We used Amazon Mechanical Turk^m as the source of human subjects who can reliably evaluate certain NLP tasks.³² All in all we wanted to evaluate 258 outputs: six summaries for each of the 43 people from the DUC 2004 set, three automatically generated summaries (one — for each similarity measure), and three corresponding model summaries. For each of the 258 outputs we recruited five human annotators. The annotators were provided with: the name of the person; link to the English-language Wikipedia article about this person; three-level summary of this Wikipedia article. The annotators did not know whether they got an automatically generated or a model summary. We asked our human annotators to answer the following questions:

- Do you agree that the sentences listed on Level 1 are a good summary of the Wikipedia article about *Person* (assume, the number of sentences in the summary cannot exceed the number of sentences listed on Level 1)?
- Assume that the summary of the Wikipedia article about *Person* can have as many sentences as listed on Level 1 and Level 2 combined. Do you agree that the sentences listed on Level 1 and Level 2 are a good summary?

^m<http://www.mturk.com>

14 *E. Filatova*

Table 4. Three-level summaries for Gene Autry consisting of the lead sentences. The number of sentences on each level is the same as the number of sentences on the respective levels of summaries presented in Table 2.

| # | ID | Text |
|-----------------|----|--|
| Similarity 0.5 | | |
| 1 | 0 | Orvon Gene Autry (September 29, 1907 — October 2, 1998) was an American performer, who gained fame as “The Singing Cowboy” on the radio, in movies and on television |
| Similarity 0.35 | | |
| 1 | 0 | Orvon Gene Autry (September 29, 1907 — October 2, 1998) was an American performer, who gained fame as “The Singing Cowboy” on the radio, in movies and on television |
| 2 | 1 | Autry, the grandson of a Methodist preacher, was born near Tioga, Texas |
| | 2 | His parents, Delbert Autry and Elnora Ozment, moved to Oklahoma in the 1920s |
| 3 | 3 | After leaving high school in 1925, Autry worked as a telegrapher for the St. Louis-San Francisco Railway |
| | 4 | Talent with the guitar and his voice led to performing at local dances |
| | 5 | After an encouraging chance encounter with Will Rogers, he began performing on local radio in 1928 as “Oklahoma’s Yodeling Cowboy” |
| Similarity 0.2 | | |
| 1 | 0 | Orvon Gene Autry (September 29, 1907 — October 2, 1998) was an American performer, who gained fame as “The Singing Cowboy” on the radio, in movies and on television |
| 2 | 1 | Autry, the grandson of a Methodist preacher, was born near Tioga, Texas |
| | 2 | His parents, Delbert Autry and Elnora Ozment, moved to Oklahoma in the 1920s |
| 3 | 3 | After leaving high school in 1925, Autry worked as a telegrapher for the St. Louis-San Francisco Railway |
| | 4 | Talent with the guitar and his voice led to performing at local dances |

- 1 • Assume that the summary of the Wikipedia article about *Person* can have as
 2 many sentences as listed on Level 1, Level 2, and Level 3 combined. Do you agree
 3 that the sentences listed on Level 1, Level 2, and Level 3 are a good summary?

4 If the summary did not have Level 2 and/or Level 3 sentences, the annotator
 5 was asked to skip answering the corresponding questions. Table 5 shows the per-
 6 centage of summaries that are considered *good* on each level by the majority of the
 7 annotators (at least three out of five) for both model summaries and summaries
 8 that are created relying on the hypothesis about the pyramid-structure of the infor-
 9 mation overlap in multilingual Wikipedia. It is possible to see that the scores of the
 10 system generated summaries based on our hypothesis are lower than the scores for
 11 the human generated models. However, we would like to emphasize that:

- 12 • Lead section for Wikipedia articles are human generated summaries that are
 13 constantly checked and approved by Wikipedia editors;
 14 • Lead section model summaries used in our experiment are created by manually
 15 extracting the first sentences from Wikipedia articles, while system-generated

Table 5. The percentage of the DUC 2004 people whose summaries are considered good by the majority of MTurk annotators.

| Similarity Measure | System or Model | Three-Level Summary | Two-Level Summary | One-Level Summary |
|--------------------|-----------------|---------------------|-------------------|-------------------|
| 0.2 | model | 97.62% | 100% | 97.37% |
| | system | 76.19% | 87.5% | 84.21% |
| 0.35 | model | 92.90% | 87.18% | 80.00% |
| | system | 81.00% | 71.80% | 71.43% |
| 0.5 | model | 94.74% | 90.91% | 96.00% |
| | system | 60.53% | 57.58% | 80.00% |

1 summaries consist of the sentences that are extracted from the automatically
2 processed text from the *Edit* Wikipedia tab.

3 We believe that despite the fact that the summaries produced by our summa-
4 rization system get lower scores than the model summaries, the numbers presented
5 in Table 5 can be used as a proof of the validity of our multilingual Wikipedia infor-
6 mation overlap structure hypothesis. In all the cases the summaries built according
7 to this hypothesis are considered reasonably good by the majority of the human
8 annotators. In Sec. 4.8, we present the error analysis that shows that many of
9 the automatically-generated summaries' weak points can be avoided by using more
10 sophisticated text processing tools.

11 4.8. Error analysis

12 As described in Sec. 4.2, we use Wikitext to analyze the information overlap in
13 Multilingual Wikipedia. Wikitext is the text that is used by Wikipedia authors and
14 editors. There is no commonly accepted standard Wikitext language, thus, our final
15 text had a certain amount of noise. According to our analysis, substantial number
16 of summaries that were judged as being *bad* contained sentences that were affected
17 by the pre-processing of the Wikitext data.

18 After analyzing the sentences included into the system-generated summaries,
19 we detected four types of problems that could cause the negative rating of the pro-
20 duced summaries: broken sentences; weak content; poor context; redundant infor-
21 mation. Table 6 summarizes our error analysis and contains the frequency of each
22 problem type.

23 4.8.1. Error type 1: Broken sentences

24 According to our analysis, several summaries contain truncated sentences due to
25 incorrect sentence chunking. Let us analyze the following sentence from the English-
26 language Wikipedia article about *Paul Wellstone*:

27 [Paul David Wellstone (July 21, 1944 - October 25, 2002) was a two-
28 term U.S.] [Senator from the U.S. state of Minnesota and member of

Table 6. Error analysis.

| Level # (# of sentences) | Broken Sentence | Weak Content | Poor Context | Redundant Sentence |
|-----------------------------|--------------------|-----------------|-----------------|-----------------------|
| Similarity 0.5 | | | | |
| Level 1 (38) | 3 | 6 | 2 | 0 |
| Level 2 (56) | 17 | 14 | 3 | 0 |
| Level 3 (50) | 7 | 21 | 8 | 1 |
| Similarity 0.35 | | | | |
| Level 1 (40) | 3 | 4 | 0 | 0 |
| Level 2 (57) | 7 | 13 | 1 | 1 |
| Level 3 (55) | 5 | 15 | 4 | 5 |
| Similarity 0.2 | | | | |
| Level 1 (42) | 4 | 3 | 1 | 0 |
| Level 2 (87) | 17 | 9 | 10 | 9 |
| Level 3 (92) | 25 | 13 | 5 | 1 |

1 the Democratic-Farmer-Labor Party, which is affiliated with the national
2 Democratic Party.]

3 This sentence was broken into two sentences (identified above by the square brack-
4 ets), and only the first portion of the sentence was added to the output summary.
5 Despite the fact that this portion contains important biographical information, it
6 cannot be used as a stand-alone sentence. This portion of the sentence was used as
7 the only sentences on Level 1 and thus, obviously, the summary consisting only of
8 one level was judged as *bad* by the majority of the annotators.

9 We call such sentences *broken* as they do not flow comprehensibly because a
10 period used for an abbreviation is treated as a full stop. This results in sentences
11 ending at a person's middle initial, title, or any acronym involving periods. Con-
12 sequently, a piece of text that we erroneously treat as a sentence is only a part of
13 a sentence: typically, either the beginning or the end of a sentence. The pieces of
14 text either finish abruptly or begin awkwardly. Our corpus has a few cases where
15 such an awkward continuity of a single sentence goes on to take up to three or four
16 text pieces that are added to the output as separate sentences, thus increasing the
17 number of broken sentences in the output.

18 4.8.2. Error type 2: Weak content

19 Determining the quality of a sentence's content is subjective, since different users
20 could be searching for different details.^{1,2} When a number of sentences included
21 into the summary is dramatically smaller than the article itself, this issue becomes
22 even more profound. For example, an article can contain a wide variety of details
23 on many topics, however, each sentence contains information only on one of the
24 topics. Thus, if the number of sentences is considerably smaller than the number
25 of topics covered in the article, it is obvious that many topics will not be covered
26 in the summary.

1 For example, one of the MTurk annotators who rated the summary about *Gene*
2 *Autry* (Table 2: similarity measure of 0.35) notes that

3 “The summary includes too many song references and fails to mention
4 Autry’s hall of fame induction”.

5 Another example of what can be considered as weak content can be found
6 in the summary about *Hugo Chávez* constructed using the similarity measure
7 of 0.35:

8 • Level 1

9 (1) Early life (1954–1992) Chávez was born on July 28, 1954 in the town of
10 Sabaneta, Barinas.

11 • Level 2

12 (1) After a two-year imprisonment, Chávez was pardoned by President Rafael
13 Caldera in 1994.

14 • Level 3

15 (1) Chávez went on to win the 1998 presidential election on December 6, 1998
16 with 56% of the votes.

17 The annotators’ comments submitted for this summary are listed below. They
18 show that though the summary contains useful information, this information is
19 not enough for a summary to be considered good. Perhaps, in future work this
20 obstacle can be overcome by allowing more sentences to be included in the output
21 summary.

- 22 • “There is no context between the three sentences. It does not say what he is
23 famous or notable for. He was imprisoned. . . for what? He won the presidency. . .
24 of which country?”
25 • “Very thin summary”.

26 Here is how the MTurk annotators rated the above summary about *Hugo*
27 *Chávez*:

- 28 • Level 1: Good - 1; Bad - 4;
29 • Level 2: Good - 0; Bad - 5;
30 • Level 3: Good - 1; Bad - 4.

31 4.8.3. *Error type 3: Poor context*

32 The problem of poor context is related to the problem of weak content. In sev-
33 eral cases the problem of poor context can be eliminated by simple re-ordering of
34 sentences.

35 Below is the summary output produced by our system with the similarity of
36 0.2 for the articles about *Thabo Mbeki*. Only three sentences are included into the
summary (one for each level).

18 *E. Filatova*

1 • Level 1

2 (1) In November 2008, “The New York Times” reported that due to Thabo
3 Mbeki’s rejection of scientific consensus on AIDS and his embrace of AIDS
4 denialism, an estimated 365,000 people perished in South Africa.

5 • Level 2

6 (1) His father was a stalwart of the African National Congress (ANC) and the
7 South African Communist Party.

8 • Level 3

9 (1) Thabo Mbeki was the executive face of government in South Africa from 1994.

10 Three of the five annotators left comments concerning the quality of the output for
11 the article about *Thabo Mbeki*:

12 • “The summaries don’t work well in their starts. Start with presidential, then
13 notable aids epidemic then father figure”.

14 • “The sentence listed in Level 3 is the only one which provides detail on who Thabo
15 Mbeki is. The Level 1 and Level 2 summary describes the subject’s actions, belief
16 and relations”.

17 • “I think the information would better be presented in reverse, with the sentence
18 from Level 3 at the start, followed by the sentence from Level 2, and concluding
19 with the sentence from Level 1”.

20 Here is how MTurk annotators rated the above summary about *Thabo Mbeki*:

21 • Level 1: Good - 1; Bad - 4;

22 • Level 2: Good - 2; Bad - 3;

23 • Level 3: Good - 3; Bad - 2.

24 Obviously, the sentences included into the summary contain information worth
25 summarizing; however, their placement relative to the rest of the output sentences
26 rendered it incomprehensible. This usually occurs when a sentence is chosen from
27 the center of an article’s content, rather than being the first sentence of an article’s
28 new section.

29 According to the annotators’ comments, a simple re-ordering of the levels could
30 have improved the rating given to the summary. We believe that in our future
31 experiments this problem could be fixed by using the technique developed for iden-
32 tifying the most likely order of information in a Wikipedia article that is described
33 in Ref. 27.

34 4.8.4. *Error type 4: Redundant information*

35 In our system we do not have any control for the redundancy of the information.
36 According to Table 6, information overlap in the sentences within the same article
37 does not happen often. In our future experiments we plan on using of the algorithms
38 developed specifically to avoid redundancy in summarization systems.^{4,33}

4.8.5. Error analysis: Conclusion

Overall, the quality of the summaries generated according to our information overlap hypothesis is lower than the quality of the model summaries. However, many of the automatically generated summaries' errors can be eliminated by using more sophisticated text processing tools. Thus, we believe that the quality of the automatically generated summaries is sufficient to proof the correctness of our hypothesis.

5. Conclusion and Future Work

In this paper we study the structure of information overlap in multilingual Wikipedia. We introduce and prove the hypothesis that the structure of this information overlap is similar to the information overlap structure (pyramid model) used in summarization evaluation, as well as the information overlap/repetition structure used to identify important information for multidocument summarization.

We believe that the understanding of the information overlap structure in multilingual Wikipedia can be used in a variety of ways. For example, it can be used for automatic generation of short, place-holder entry descriptions in new languages. Second, the pyramid structure of information overlap can be used for judging the trustworthiness of information facts mentioned in Wikipedia articles: the more important (and the more trusted) the information fact for a particular Wikipedia entry is, the higher the chances that this information is covered (repeated) in the articles in many languages.

While in this paper we focus on studying the structure of information overlap, in our future work, we are also interested in investigating information asymmetries in multilingual Wikipedia. Highlighting differences in Wikipedia articles about the same entry written in different languages can be used for suggesting automatic updates of Wikipedia articles, information tracking, information trustworthiness evaluation, etc. Information asymmetries is an interesting topic to study given that the choice of information can reveal the authors attitude toward the Wikipedia entry, especially in controversial topics.

References

1. S. Teufel and H. V. Halteren, Evaluating information content by factoid analysis: Human annotation and stability, in *Proc. 42th Annual Meeting of the Association for Computational Linguistics (ACL 2004)*, Barcelona, Spain, 2004, pp. 419–426.
2. A. Nenkova, R. Passonneau and K. McKeown, The pyramid method: Incorporating human content selection variation in summarization evaluation, *ACM Trans. Speech Lang. Process.* **4**(2) (2007) 4.
3. E. Filatova and V. Hatzivassiloglou, A formal model for information selection in multi-sentence text extraction, in *Proc. 20th Int. Conf. Computational Linguistics (COLING 2004)*, Geneva, Switzerland, 2004.
4. D. Radev, H. Jing, M. Styś and D. Tam, Centroid-based summarization of multiple documents, *Inf. Process. Manage.* **40**(6) (2004) 919–938.

20 E. Filatova

- 1 5. E. Filatova, Directions for exploiting asymmetries in multilingual Wikipedia, in *Proc. NAACL Workshop on Cross Lingual Information Access*, Boulder, CO, USA, 2009.
- 2
- 3 6. E. Filatova, Multilingual Wikipedia, summarization, and information trustworthiness,
- 4 in *Proc. NAACL Workshop on Information Access in a Multilingual World*, Boston,
- 5 MA, USA, 2009.
- 6 7. D. Ahn, V. Jijkoun, G. Mishne, K. Müller, M. de Rijke and S. Schlobach, Using
- 7 Wikipedia at the TREC QA track, in *Proc. Text REtrieval Conf. (TREC 2004)*,
- 8 2004.
- 9 8. D. Buscaldi and P. Rosso, Mining knowledge from Wikipedia for the question answer-
- 10 ing task, in *Proc. Fifth Int. Conf. Language Resources and Evaluation (LREC 2006)*,
- 11 Genoa, Italy, 2006.
- 12 9. J. Ko, T. Mitamura and E. Nyberg, Language-independent probabilistic answer rank-
- 13 ing for multilingual question answering, in *Proc. 45th Annual Meeting of the Associ-*
- 14 *ation for Computational Linguistics (ACL 2007)*, Prague, Czech Republic, 2007.
- 15 10. F. Baidys, J. Hirschberg and E. Filatova, An unsupervised approach to biography
- 16 production using Wikipedia, in *Proc. 46th Annual Meeting of the Association for*
- 17 *Computational Linguistics (ACL 2008)*, Columbus, OH, USA, 2008.
- 18 11. R. Nelken and E. Yamangil, Mining Wikipedia's article revision history for training
- 19 computational linguistics algorithms, in *Proc. AAAI Workshop on Wikipedia and*
- 20 *Artificial Intelligence: An Evolving Synergy*, Chicago, IL, USA, 2008.
- 21 12. S. F. Adafre and M. de Rijke, Finding similar sentences across multiple languages
- 22 in Wikipedia, in *Proc. Conf. European Chapter of the Association for Computational*
- 23 *Linguistics, Workshop on New Text — Wikis and Blogs and other Dynamic Text*
- 24 *Sources*, Trento, Italy, 2006.
- 25 13. A. Richman and P. Schone, Mining Wiki resources for multilingual named entity
- 26 recognition, in *Proc. 46th Annual Meeting of the Association for Computational Lin-*
- 27 *guistics (ACL 2008)*, Columbus, OH, USA, 2008.
- 28 14. P. Schönhofen, A. Benczúr, I. Bíró and K. Csalogány, Performing cross-language
- 29 retrieval with Wikipedia, in *Proc. Working Notes for the CLEF 2007 Workshop*,
- 30 Budapest, Hungary, 2007.
- 31 15. S. Ferrández, A. Toral, Ó. Ferrández, A. Ferrández and R. Munoz, Applying
- 32 Wikipedia's multilingual knowledge to cross-lingual Question Answering, *Lecture*
- 33 *Notes in Computer Science (LNCS): Natural Language Processing and Information*
- 34 *Systems*, Vol. 4592, (2007), pp. 352–363.
- 35 16. B. T. Adler, K. Chatterjee, L. de Alfaro, M. Faella, I. Pye and V. Raman, Assign-
- 36 ing trust to Wikipedia content, in *Proc. Int. Symp. Wikis (WikiSym 2008)*, Porto,
- 37 Portugal, 2008.
- 38 17. H. Zeng, M. Alhossaini, L. Ding, R. Fikes and D. McGuinness, Computing trust
- 39 from revision history, in *Proc. Int. Conf. Privacy, Security and Trust (PST 2006)*,
- 40 Markham, Canada, 2006.
- 41 18. D. McGuinness, H. Zeng, P. P. da Silva, L. Ding, D. Narayanan and M. Bhaowal,
- 42 Investigations into trust for collaborative information repositories: A Wikipedia case
- 43 study, in *Proc. Workshop on the Models of Trust for the Web (MTW 2006)*, Edin-
- 44 burgh, UK, 2006.
- 45 19. K. McKeown, J. Klavans, V. Hatzivassiloglou, R. Barzilay and E. Eskin, Towards
- 46 multidocument summarization by reformulation: Progress and prospects, in *Proc.*
- 47 *16th National Conf. American Association for Artificial Intelligence (AAAI 1999)*,
- 48 Orlando, Florida, 1999, pp. 453–460.
- 49 20. R. Barzilay, N. Elhadad and K. McKeown, Inferring strategies for sentence ordering
- 50 in multidocument news summarization, *J. Artif. Intell. Res.* **17** (2002) 35–55.

- 1 21. S. Harabagiu, F. Lăcătușu and S. Maiorano, Multi-document summaries based on
2 semantic redundancy, in *Proc. 14th Florida AI Conf. (FLAIRS 2003)*, St. Augustine,
3 Florida, 2003.
- 4 22. E. Adar, M. Skinner and D. Weld, Information arbitrage in multi-lingual Wikipedia,
5 in *Proc. Second ACM Int. Conf. Web Search and Data Mining*, Barcelona, Spain,
6 2009.
- 7 23. O. Nov, What motivates Wikipedians? *Commun. ACM* **50**(11) (2007) 60–64.
- 8 24. L. Denoyer and P. Gallinari, The Wikipedia XML Corpus, SIGIR Forum.
- 9 25. C.-Y. Lin and E. Hovy, Automatic evaluation of summaries using N-gram co-
10 occurrence statistics, in *Proc. Language Technology Conf. (HLT-NAACL 2003)*,
11 Edmonton, Canada, 2003.
- 12 26. P. Over and W. Liggett, Introduction to DUC-2002: An intrinsic evaluation of generic
13 news text summarization systems, in *Proc. Workshop on Automatic Summarization*
14 (*DUC 2002*), Philadelphia, PA, USA, 2002.
- 15 27. C. Sauper and R. Barzilay, Automatically generating Wikipedia articles: A structure-
16 aware approach, in *Proc. 47th Annual Meeting of the Association for Computational*
17 *Linguistics (ACL 2009)*, Singapore, 2009.
- 18 28. D. Evans, J. Klavans and K. McKeown, Columbia Newsblaster: Multilingual news
19 summarization on the web, in *Demonstration Papers at HLT-NAACL 2004, HLT-*
20 *NAACL-Demonstrations '04, Association for Computational Linguistics*, Strouds-
21 burg, PA, USA, 2004, pp. 1–4.
- 22 29. D. Evans and K. McKeown, Identifying similarities and differences across English and
23 Arabic news, in *Int. Conf. Intelligence Analysis*, McLean, VA, USA, 2005.
- 24 30. M. Litvak, M. Last and M. Friedman, A new approach to improving multilingual sum-
25 marization using a genetic algorithm, in *Proc. 48th Annual Meeting of the Association*
26 *for Computational Linguistics, Association for Computational Linguistics*, Uppsala,
27 Sweden, 2010, pp. 927–936.
- 28 31. Alias-i, Lingpipe 3.7.0. (accessed January 19, 2009), <http://alias-i.com/lingpipe>
29 (2009).
- 30 32. V. Sheng, F. Provost and P. Ipeirotis, Get another label? Improving data quality
31 and data mining using multiple, noisy labelers, in *Proc. Fourteenth ACM SIGKDD*
32 *Int. Conf. Knowledge Discovery and Data Mining (KDD 2008)*, Las Vegas, NV, USA,
33 2008.
- 34 33. J. Carbonell and J. Goldstein, The use of MMR, diversity-based reranking for reorder-
35 ing documents and producing summaries, in *Proc. 21st Annual Int. ACM SIGIR*
36 *Conf. Research and Development in Information Retrieval (SIGIR 1998)*, New York,
37 NY, USA, 1998, pp. 335–336.