

Automatic Annotation Suggestions for Audiovisual Archives: Evaluation Aspects

LUIT GAZENDAM AND CHRISTIAN WARTENA

Novay, Enschede, The Netherlands

VÉRONIQUE MALAISÉ AND GUUS SCHREIBER

Department of Computer Science, Vrije Universiteit Amsterdam, The Netherlands

ANNEMIEKE DE JONG

Netherlands Institute for Sound and Vision, Hilversum, The Netherlands

HENNIE BRUGMAN

Max Planck Institute for Psycholinguistics, Nijmegen, The Netherlands

In the context of large and ever growing archives, generating annotation suggestions automatically from textual resources related to the documents to be archived is an interesting option in theory. It could save a lot of work in the time consuming and expensive task of manual annotation and it could help cataloguers attain a higher inter-annotator agreement. However, some questions arise in practice: what is the quality of the automatically produced annotations? How do they compare with manual annotations and with the requirements for annotation that were defined in the archive? If different from the manual annotations, are the automatic annotations wrong? In the CHOICE project, partially hosted at the Netherlands Institute for Sound and Vision, the Dutch public archive for audiovisual broadcasts, we automatically generate annotation suggestions for cataloguers. In this paper, we define three types of evaluation of these annotation suggestions: (1) a classic and strict evaluation measure expressing the overlap between automatically generated keywords and the manual annotations, (2) a loosened evaluation measure for which semantically very similar annotations are also considered as relevant matches, and (3) an in-use evaluation of the usefulness of manual versus automatic annotations in the context of serendipitous

browsing. During serendipitous browsing, the annotations (manual or automatic) are used to retrieve and visualize semantically related documents.

KEYWORDS Extraction, Evaluation, Automatic annotation, Audiovisual archives, Semantic evaluation, Semantic browsing

Context

The Netherlands Institute for Sound and Vision is in charge of archiving publicly broadcasted TV and radio programmes in the Netherlands. Two years ago, the audiovisual production and archiving environment changed from analogue to digital data. This effectively quadrupled the inflow of archival material and, as such, the amount of work for cataloguers.¹

Sound and Vision faces the challenge to create a durable continuous access to the daily increasing collections with the same number of cataloguers (40 people). The manual annotation is the bottleneck in the archiving process: it may take a cataloguer up to three times the length of a TV-programme to annotate it manually, depending on the genre (news item, game-show, documentary). During annotation, cataloguers often consult and use available contextual information such as TV-guide synopses, official TV-programmes web site texts, and subtitles.

The annotation process follows strict guidelines. All catalogue descriptions conform to a metadata scheme called iMMiX. The iMMiX metadata model is an adaptation for 'audiovisual' catalogue data of the FRBR data model² which was developed in 1998 by the International Federation of Library Associations (IFLA).

The iMMiX metadata model captures four important aspects of a broadcast:

1. information content (who, what, when, where, why and how, includes keywords, organizations, locations)
2. audiovisual content (What can be seen or heard? Includes descriptions like *close-up*)
3. formal data (e.g. intellectual property rights)
4. document management data (e.g. document ID).

Choices for some of the iMMiX fields (subject, location, persons, etc.) are restricted to a controlled vocabulary named GTAA. GTAA is a Dutch acronym for 'Common Thesaurus [for] Audiovisual Archives' and contains about 160,000 terms, organized in six facets: **Locations, People, Names (of organizations, events, etc.), Makers, Genres and Subjects**. This latest facet contains 3800 keywords and 21,000 relations between the keywords belonging to the ISO-2788 defined relationships of Broader Term, Narrower Term, Related Term and Use/Use for. It also contains linguistic information such as preferred textual representations of keywords and non-preferred representations. Each keyword on average has 1 broader, 1 narrower and 3.5 related terms. Cataloguers are instructed to select keywords that describe the programme as a whole, are specific and allow good retrieval.

Automatic annotation suggestions

The annotation of audiovisual material is one of the core tasks of the archive. At the moment this is a manual process. The archive considers the inclusion of automatic techniques in this process promising, but risky. Therefore a human-in-the-loop is required during innovation: automatically generated keywords need to be suggested to cataloguers who will perform an *a-posteriori* validation.

The CHOICE project investigates how to automatically suggest GTAA keywords to cataloguers during their annotation task. We assume that by applying Natural Language Processing and Semantic Web techniques to the contextual documents (e.g. TV guide texts describing the broadcast), reasonable annotation suggestions can be generated. These context documents differ from the audiovisual content of the programme, which is problematic in theory: how can an automatic process analysing context documents derive good keywords when cataloguers performing the same task also inspect the original audiovisual material? In practice, however, this different nature of context documents can also be an advantage: often they summarize the content of the programme which makes it easier to find summarizing keywords. After a brief overview of automatic annotation tools and platforms proposed in the literature (*See Related work section*), we introduce our processing pipeline.

The suggestions are intended to increase a cataloguer's working speed and consistency. Typical measures of inter-cataloguer consistency range from 13 to 77 per cent (with an average of 44 per cent) when a controlled vocabulary is used (Leininger 2000). The topology of disagreement shows that a portion of these differences is small semantic differences. This disagreement can be problematic when manual annotations serve as a gold standard for the evaluation of our automatic annotation suggestions. Nevertheless, the manual annotations are our best baseline for evaluation.

To reduce the shortcomings of an evaluation based on a strict string-based comparison: *classical evaluation*, we propose a second type of evaluation: *semantic evaluation*. In a third evaluation, we then investigate the potential value of automatically generated keywords. These can bring *new types* of search or archival behaviour, that cannot be evaluated against current practices. We designed a new kind of archive access named *serendipitous browsing* to test this potential value. During serendipitous browsing, overlapping annotations are used to link documents. These links can be used to browse through the archive. Each link departing from a document will highlight other aspects of the original document and each document reached will have new aspects which were not part of the original. Serendipitous browsing will thus contextualize documents resulting in new interpretations and allow for the discovery of new and unexpected documents related to the original.

Related work

The tools and architectures that have been implemented for generating semantic annotations based on ontologies or other concept-based representations of a controlled vocabulary can be roughly categorized into:

- tools for manual annotation: an interface providing help for a human to insert semantic annotations in a text
- tools for semi-automatic annotation: a system providing help and automatic suggestions for the human annotation
- tools for automatic annotation: a system providing annotation suggestions, possibly to be validated or modified *a posteriori*.

Tools like Annotea (Kahan and Koivunen 2001) and SHOE (Heflin and Hendler 2000) provide environments for *manually* assigning annotations to documents; we aim at automatically suggesting them in our project, to ease some of the annotation burden.

The second category of tools proposes annotation suggestions after a learning process. They are represented by tools such as Amilcare (Ciravegna and Wilks 2003) and T-Rex (Iria 2005), that learn rules at annotation time in order to provide the annotator with suggestions. They are both based on the GATE platform (Cunningham *et al.* 2002), a generic Natural Language Processing platform that implements simple Named Entity recognition modules and a rule language to define specific patterns to expand on simple string recognition. These interactive annotation tools are designed to work on the same texts the cataloguers are annotating, but in our situation, cataloguers annotate audiovisual programmes and not the textual context documents themselves. Therefore, tools from the third category were considered the most relevant.

We opted for the semantic annotation performed by tools that generate them without human interaction. A typical example of this third type of tools is the KIM platform (Kiryakov *et al.* 2005); the MnM tool (Vargas-Vera *et al.* 2002) is mixed, providing both semi-automatic and automatic annotations. Although they can be adapted to different domains or use cases, the adaptation requires a lot of work, and in the case of KIM, the upper level of the ontology cannot be changed. The MnM system integrates an ontology editor with an information extraction pipeline, and this is also the approach that we decided to follow in our project, but we used GATE for this purpose, because of its openness and adaptability.

Annotation and ranking pipeline in the CHOICE-project

Our approach to suggesting keywords automatically to annotate TV-programmes is based on information extraction techniques, applied to textual resources describing the TV-programme's content, such as TV-guide texts or web-site texts. Our system transforms these texts into a suggestion list of thesaurus keywords. The system comprises three parts:

1. *A text annotator.* The text annotator tags occurrences of thesaurus keywords in the texts. GATE (Cunningham *et al.* 2002) and its plug-in Apolda (Wartena *et al.* 2007) implement this process
2. *TF.IDF computation.* TF.IDF is a statistical measure which expresses how distinctive a keyword is for a document compared to the rest of the collection

3. A *cluster-and-rank process* which uses the thesaurus relations to improve upon the TF.IDF ranked list.

The TF.IDF of a keyword depends on the frequency of a keyword in the document, divided by the number of documents in the collection in which the keyword appears (e.g. if the keyword *economy* appears a lot in one document, this is relevant, unless all documents in the collection have *economy* in it). We use TF.IDF as baseline which we try to beat with ranking algorithms. In the next subsection, we elaborate upon the cluster-and-rank process.

Cluster-and-rank process

The keywords tagged in the context documents of a TV-programme are sometimes related to each other by thesaurus relationships. Together, the keywords and the relations form a graph. In Figure 1, one can see the graph for a text containing the keywords *Ministers*, *Government*, *Civil servants*, *Soldiers* and *Armed forces*. A direct connection exists between *Ministers* and *Government*. To increase the connectedness of our graph, we also included indirect relations (in which an intermediate keyword connects two found keywords). An indirect connection exists between *Ministers* and *Soldiers* with *Professions* as an intermediate term. *Professions* did not appear in the original text.

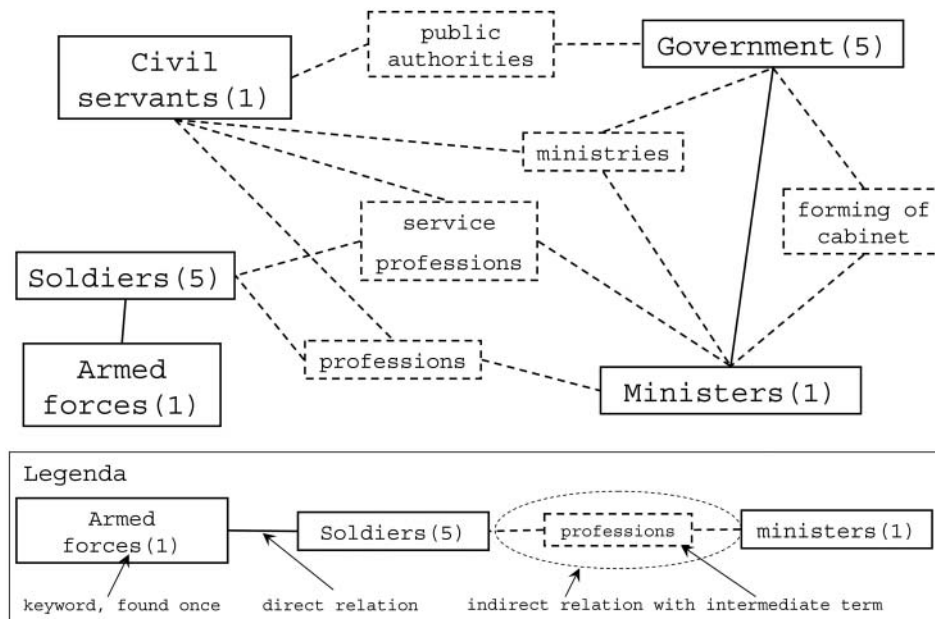


FIGURE 1 Relations found between a set of keywords.

The cluster-and-rank component uses the graph to create a (re)ranked list as output. We implemented three algorithms that build ranked lists from this graph: the well known algorithm named PageRank (Brin and Page 1998) (uses only graph information), our own method called CARROT (also uses TF.IDF information), and our second method which is called Mixed (also uses TF.IDF information and the whole graph of the thesaurus as additional information).

CARROT

CARROT (Malaisé *et al.* 2007) stands for Cluster And Rank Related Ontology concepts or Thesaurus terms. It combines the local connectedness of a keyword and the TF.IDF score. The only graph property CARROT uses is the local connectedness of a keyword. It creates four groups each having the same local connectedness (group 1: both direct and indirect connections (*Soldiers, Government, Ministers*), group 2: only direct connections (*Armed forces*), group 3: only indirect connections (*Civil servants*) and group 4: no connections). Each group is sorted on the TF.IDF values.

PageRank

PageRank (Brin and Page 1998) is used to determine the centrality of items in a network. One way to understand the working of PageRank is by imagining activation spreading through a network. The initial (e.g. TF.IDF) activation spreads itself equally via each available relation to other nodes in the network. It then spreads again via the relations of the network, some back to the original starting nodes and some further. In the end, on each node in the network, a dynamic equilibrium will be reached (each moment the same activation that leaves the node is also fed onto the node from other nodes; dynamic equilibrium). This equilibrium is no longer dependent on the starting activation, only on the network structure. The activation on each node corresponds to the PageRank score and expresses its importance.

In research similar to our own by Wang *et al.* (2007), PageRank was used to determine the most central WordNet keywords in scientific articles. They compared PageRank with TF.IDF and showed that PageRank suggested much better keywords.

PageRank is performed upon the same cluster as CARROT, but the PageRank algorithm also assigns PageRank scores to the intermediate terms so these are included in the suggestion list (also include the dashed terms of Figure 1).

Mixed algorithm using general keyword importance

For the Mixed algorithm, we wanted to retain some of the relevancy information conveyed by TF.IDF while performing the spreading of activation. We start with the TF.IDF activation and only spread it around with the official PageRank formula during three iterations. At that moment, some influence of the original TF.IDF is still present and at the same time, some activation accumulates at the central nodes in the network. This PageRank at $t=3$ is multiplied with the general importance of the keywords. The idea behind the weighting with keyword importance is that we want to favour

keywords which are considered more important in general. The way we determine the general importance of keywords is by PageRanking the GTAA as a whole. We assume that the modelling of the GTAA reflects the importance of the keywords: topics which are considered important according to the GTAA makers from Sound and Vision are modelled with many keywords and many relations. The five keywords with the highest GTAA PageRank are *businesses, buildings, people, sports, animals*. They are on average connected to 60 other keywords and are central nodes in the thesaurus. Their importance is reflected in the catalogue: on average, these central keywords are used once per 42 documents. The five keywords with the lowest GTAA PageRank are *lynchings, audiotapes, holography, autumn, spring* (each having one relation). The makers of the GTAA do not consider these important enough for the GTAA as a whole to model more relations. This lower importance is reflected in their usage: on average, they appear once in every 9900 documents.

Experimental set-up

Upon a test corpus, we perform three evaluations in two experiments. In our first experiment, we generate keyword suggestions from context documents with the four different settings of our pipeline and we evaluate these against manually assigned keywords. We evaluate these resulting lists of suggestions in two different ways: classically and semantically.

Our first evaluation is a classic Precision/Recall evaluation, inherited from the information extraction world. Given the reality of inter-annotator disagreement however, we questioned beforehand whether this classic evaluation methodology was appropriate for the task of suggesting keywords in the archival domain.

The second evaluation introduces a measure of semantic overlap between the Automatic Annotations and the target against which we evaluate them: the manual annotations of the TV-programmes. This setting is still biased towards current annotation practices and does not show another dimension: what can Automatic Annotations bring in the context of possible *new applications*?

In order to evaluate the possibilities in terms of new practices in archives, we tuned a second experiment, which underlines the possible value of Automatic Annotations and Manual Annotations in the context of a particular search through an archive: serendipitous browsing. With it, we test the value of the manual annotations and the CARROT keyword annotation suggestions for retrieving semantically related documents. By doing so, we feed an idea from the Semantic Web (inherited from Semantic Browsing (Faaborg and Lagoze 2003; Hildebrand 2008)) back into the archival world to bring new solutions to their core task: find relevant information/documents in large archives. Although the value of this idea needs to be tested, it reminds Sound and Vision's customer service of the loose search performed by users by flipping through a physical card-tray. The arrangement of physical cards in trays on one topic made it possible to browse for strong, semi or loosely

related documents. This option was lost when the archives' access with card trays was replaced by computers.

Source material

Our corpus consists of 258 broadcasted TV-documentaries. In total, 80 per cent of these broadcasts belonged to three series of TV-programmes: *Andere Tijden*, which is a series of Dutch historical documentaries, *Beeldenstorm*, which is a series of art documentaries presented by Henk van Os, the former director of the Rijksmuseum, and *Dokwerk*, which is a series of historical political documentaries. Each broadcast is associated with one or more texts from the broadcasters' web site (we name these *context documents*) and one manual catalogue description made by Sound and Vision. The 258 TV-broadcasts are associated with 362 context documents. The length of the context documents varied between 25 words and 7000 words with an average of 1000 words.

Catalogue descriptions

Each TV-broadcast in our corpus has a catalogue description. These catalogue descriptions contain keywords which were assigned manually by cataloguers from Sound and Vision. The catalogue descriptions on average contain 5.7 keywords with a standard deviation of 3.2 keywords. The minimum number of terms is 1, the maximum is 15. These keywords are the ground truth against which we evaluate the TF.IDF baseline and the three ranking algorithms in the next two experiments.

Classical evaluation

We want to measure the quality of the automatically derived keywords. For this purpose, we compare the automatic annotations with the existing manual annotations. The standard way of evaluating our systems output against manual annotation is with the information retrieval measures of *precision* and *recall* (Salton and McGill 1983). **Precision** is defined as the number of relevant keywords suggested by our system for one TV-programme divided by the total number of keywords that are given by our system for that programme, and **recall** is defined as the number of relevant keywords suggested by our system for one TV-programme divided by the total number of existing relevant keywords for that TV-programme (which should have been suggested for that TV-programme). Often precision and recall are inversely related, so it is possible to increase one at the cost of reducing the other. For this reason, they are often combined into a single measure, such as the balanced F-measure, which is the weighted harmonic mean of precision and recall.

Given the fact that our system produces ranked lists, we can look at average precision and recall for different top parts of our list: *precision@5* and *precision@10* express, respectively, the precision of the first 5 and the first 10 suggestions. For the suggestion of keywords to cataloguers, only these top terms are important: a cataloguer will only read a limited number of

suggestions. The cataloguer will stop when the suggestions are good (he has read enough good suggestions so he is satisfied (Simon 1957)) and stop when the suggestions are poor (he is not expecting reasonable suggestions anymore).

Classical evaluation of the results

Table 1 shows the classic evaluation for our four ranking algorithms.

The first observation we make is that only the PageRank setting is considerably worse than the others. This is probably attributable to the fact that PageRank lacks the ability to incorporate any relevancy information from the TF.IDF scores. The performance of PageRank in the experiment of Wang *et al.* (2007) makes this result unexpected.

A second observation is that the Mixed model starts out as a very poor, but that it catches up with the better settings such as the TF.IDF baseline and CARROT. The TF.IDF seems best, but this difference is not statistically significant (at $p < 0.05$).

The final observation is the big jump in F-score between @1 and @3 for all methods. This is interesting as it tells us that one suggestion just cannot contain that much information and that lists with 3 or 5 suggestions are better.

Discussion

Medelyan and Witten (2006) conducted an experiment similar to ours. They automatically derived keywords from the Agrovoc thesaurus (containing 16,600 preferred terms) for FAO documents (Food and Agriculture Organization of the United Nations). Their results show similar low numbers of around 0.20 for precision, recall and F-score. Their best method KEA++ reached the best F-score@5 of 0.187 with a precision@5 of 0.205 and a recall@5 of 0.197. Given that their documents are on average 17 times longer than ours (which helps for retrieving good keywords) but that their number of possible keywords is five times as large too (which makes it harder to pick the right

TABLE 1
CLASSICAL EVALUATION OF OUR RESULTS

Precision	@1	@3	@5	@10
Baseline: TF.IDF	0.38	0.30	0.23	0.16
CARROT	0.39	0.28	0.22	0.15
PageRank	0.19	0.17	0.14	0.11
Mixed	0.23	0.21	0.19	0.15
Recall				
Baseline: TF.IDF	0.08	0.18	0.23	0.31
CARROT	0.08	0.15	0.21	0.27
PageRank	0.04	0.09	0.13	0.20
Mixed	0.05	0.12	0.18	0.28
F-score				
Baseline: TF.IDF	0.13	0.22	0.23	0.21
CARROT	0.13	0.20	0.21	0.20
PageRank	0.07	0.12	0.14	0.14
Mixed	0.08	0.16	0.19	0.20

keyword), we can only state that our best methods produce reasonable results.

Inspection of individual suggestion lists reveals a mismatch between our sense of quality of the suggestions and the classic evaluation: many good suggestions do not contribute at all to the precision and recall numbers. To give an example, the first six CARROT suggestions for TV-programme *Andere Tijden 04-09-2000* are *Jews, camps, deportations, interrogations, trains* and *boys*. The topic of this TV-programme was the Dutch deportation camp of Westerbork from which Jews were deported to concentration camps in the Second World War. The manually assigned keywords were *deportations, persecution of Jews, history* and *concentration camps*. According to the classic evaluation, however, only the suggestion of *deportations* is correct. However, most of the other keywords do convey valuable information. When we look at the relations of these suggested keywords in the GTAA, we see that *camps* is the broader term for *concentration camps* and that *Jews* is related to *persecution of Jews*. These thesaurus relations are used during semantic evaluation.

Semantic evaluation

The classic type of evaluation takes place on the basis of exact match or *terminological consistency* (Iivonen 1995). We argue that this exact type of evaluation does not measure the quality of our suggestions well. We want keywords which present a semantic similarity with the manually assigned keywords to be counted as correct too. This is good enough for the task of suggesting keywords and it tackles part of the problem of the inter-annotator disagreement. This semantic match is known as *conceptual consistency* (Iivonen 1995).

Medelyan and Witten (2006) describe a practical implementation of evaluation against conceptual consistency instead of terminological consistency. They use the relations in a thesaurus as a measure for conceptual consistency. The conceptually consistent terms are all terms which are within a certain number of thesaurus relationships from the target term. In their experiment, Medelyan and Witten consider all terms reachable in two relations to be conceptually consistent (given their task and thesaurus). We chose to consider all terms within one thesaurus relationships to be conceptually consistent. This choice for one relationship is not purely motivated by the structure of our thesaurus, as it also would allow two steps of distance, but we face the risk of interaction between semantically based ranking methods (which use thesaurus relations) and the semantic evaluation methodology (which also uses thesaurus relations).

Results

We semantically evaluated the four settings against the manually assigned keywords and the results are presented in Table 2.

In this table, we see two things. First, we observe from the F-scores that the Mixed setting is the best setting, but only @5 and @10. Its better F-score is

TABLE 2
SEMANTIC EVALUATION OF OUR RESULTS

Precision	@1	@3	@5	@10
Baseline: TF.IDF	0.50	0.43	0.37	0.30
CARROT	0.53	0.45	0.40	0.32
PageRank	0.47	0.40	0.36	0.30
Mixed	0.52	0.46	0.42	0.36
Recall				
Baseline: TF.IDF	0.16	0.32	0.40	0.54
CARROT	0.17	0.28	0.36	0.48
PageRank	0.14	0.30	0.38	0.51
Mixed	0.16	0.31	0.40	0.53
F-score				
Baseline: TF.IDF	0.24	0.37	0.39	0.38
CARROT	0.25	0.35	0.38	0.39
PageRank	0.22	0.34	0.37	0.38
Mixed	0.24	0.37	0.41	0.43

only statistically significant @10. The PageRank setting is again the worst setting, however, it is only significantly poorer than Mixed @5 and @10. The second observation is the difference in behaviour with respect to precision and recall of the different methods. The Mixed model is good in precision, but normal in recall. CARROT is poor in recall and slightly better in precision.

When we compare Tables 1 and 2, we see a big improvement in performance. This not unexpected as the semantic evaluation effectively lowers the number of possible classes. We also see that the Mixed and the PageRank setting improved much more than the other methods. Now we will look at the results qualitatively.

Qualitative analysis

A qualitative analysis of the lists generated by the four different settings can give us some more insight into the value of the four ranking algorithms and into a possible interaction between semantic ranking methods and the semantic evaluation: does a setting score well during the semantic evaluation because it is just a good setting, or because the evaluation prefers semantically connected keywords and the semantic settings (PageRank, CARROT and Mixed) happen to suggest these. The TV-documentary *Andere Tijden: Mining accident at Marcinelle* is chosen for illustration.

Sound and Visions' catalogue describes this programme as follows: *Episode of the weekly programme Andere Tijden. In this episode a mining accident in the fifties of last century in Belgium is addressed. In this mining accident, many Italian foreign workers died during a fire.* The first 12 ranks generated by our four settings are displayed in Table 3. The cataloguer attached the keywords *history, disasters, coalmines, miners and foreign employees* to this programme. The catalogue keywords are not ranked (all are equally correct).

The keywords in Small Caps are exact matches with the catalogue keywords. The keywords in **bold** are conceptually consistent and the keywords in *italics* are wrong.

TABLE 3
SUGGESTED TERMS FOR ANDERE TIJDEN 2003-11-11: MINING DISASTER AT MARCINELLE

Rank	TF. IDF	CARROT	PageRank	Mixed	Catalogue
1	<i>mines</i>	MINERS	<i>mines</i>	mining	history
2	MINERS	DISASTERS	mining	MINERS	disasters
3	DISASTERS	<i>fire</i>	COALMINES	COALMINES	coalmines
4	<i>fire</i>	FORGN EMPL.	<i>publications</i>	DISASTERS	miners
5	<i>cables</i>	<i>fathers</i>	<i>human body</i>	accidents	forgn empl.
6	FORGN EMPL.	<i>corpses</i>	<i>buildings</i>	blue-collar workers	
7	<i>fathers</i>	coal	<i>art</i>	coal	
8	<i>corpses</i>	<i>mothers</i>	MINERS	<i>mines</i>	
9	coal	<i>firemen</i>	accidents	fires	
10	<i>safety</i>	fires	<i>families</i>	<i>families</i>	
11	<i>governments</i>	immigrants	mining accidents	lignite	
12	<i>mothers</i>	immigration	DISASTERS	goldiggers	

From the table, we make four observations. First, we see that each list contains exactly three correct suggestions. In the TF.IDF and CARROT settings, the keywords *miners*, *disasters* and *foreign employees* are in the list. The PageRank and the Mixed settings have *miners* and *disasters* too, but they have *coalmines* as a third. Both the TF.IDF and the CARROT settings have many wrong suggestions in the list. The suggestion *mines* which is at the top of the TF.IDF list, is wrong as it means an *under water bomb* in the GTAA. CARROT did not have this suggestion in the first group so it correctly is lower on the list. It also had *cables*, *safety* and *government* in a lower group.

The PageRank starts with three reasonable suggestions, but then from rank 4 until 7 gives very general suggestions. It favours suggestions that are very connected (and thus very general). The semantics of these suggestions is too general (not specific enough), which is often the case with the PageRank suggestions. The following keywords appear among the top 10 in many of PageRank's suggestion lists: *publications*, *buildings*, *businesses*, *transportation*, *human body* and *professions*. If we were to judge keywords within two relations as correct as Medelyan and Witten (2006) did, we would sometimes evaluate these general terms as correct.

The Mixed setting has a nice trade-off between general and specific suggestions. It has some of the general suggestions like *mining* and *blue collar workers* which were introduced by PageRank, but it also has suggestions specific enough to match the level of the usual manual annotations. Furthermore, it has many more of the conceptually consistent suggestions in its list, not directly in the beginning, but further down the list. It does not generate more direct hits (Table 1), but more semantic matches as Table 2 shows. Mixed gives more closely related suggestions.

Serendipitous browsing

After inspection of several lists of automatically derived keywords suggestions, we discovered that they contained four types. To illustrate the four types, we again use the TV-programme *Andere Tijden 04-09-2000* about the Dutch concentration camp Westerbork. The suggestion lists contain:

1. main topic descriptors, e.g. *Jews, camps*
2. keywords related to the main topic, e.g. *interrogations*
3. subtopic descriptors, e.g. *trains*
4. wrong suggestions, e.g. *boys*.

The value of the first and non-value of the fourth type are clear. This second and third type would not be chosen by cataloguers to index a programme, but they do convey interesting aspects of the programme. Our lists of annotation suggestions contain exact suggestions, semantically related suggestions, subtopics and wrong suggestions. Lists belonging to two different broadcasts can contain the same keyword suggestions. This overlap can be used to link the broadcasts. Overlapping lists of annotation suggestions, although imprecise, might be a good measure of relatedness between two broadcasts. In the same manner, overlapping manual annotations can relate two documents.

The value for users of these relations between documents can be great: to be able to browse through the archives, discover unsuspected relationships, thus creating new interpretations. It can create an accidental discovery or a moment of serendipity.

Experimental set-up for serendipitous browsing

We tested the value of the manual annotations and automatic annotations for serendipitous browsing with an experiment. During this experiment, we created a table for our corpus, both for the manual annotations and for the automatic annotations in which we store the overlap between documents. From both tables, we selected the 10 pairs with the biggest overlap. So we are cherry picking, but we did this for a reason. Our corpus contains only 258 programmes, which represents only a small fraction of the entire catalogue of over one million documents. For the entire catalogue, we would get much better results. The best matches in our corpus give a better idea of what the method would mean for the entire catalogue.

For the automatic annotations, the pairs had between 13 and 5 overlapping keywords. For the manual annotations, these pairs had between 9 and 4 overlapping keywords. For each document in the top pairs, we selected its four closest neighbours. This means that for each document A, we have the five documents X1–X5 which have the highest number of overlapping keywords with document A. The first pair, A–X1 is one of the 10 best pairs of either the manual annotations or the automatic annotations. The pair X1–A appears a second time as the first pair in the list of the five best pairs for document X1. The overlapping keywords for each pair represent the semantics of the link between the two documents.

In our list of results, we identify three types of pairs:

1. Document X1 has a semantic overlap with document A
2. X1 and A are two context documents of the same TV-programme
3. Documents X1 and A are parts of one TV-programme which was broadcasted in multiple episodes.

When pairs had a semantic overlap, we judged the similarity between the two documents on a five-point Likert scale (Likert 1932): Strongly disagree, Disagree, Neither agree nor disagree, Agree, Strongly Agree.

Results

Of all links in our top 10 list with manual annotations, five linked broadcasts with a semantic overlap, four linked two parts of a TV-programme which was broadcasted in multiple episodes and one link was an error in the database.

Of all links in our top 10 list with automatic annotations, four linked broadcasts with a semantic overlap, four linked two parts of a TV-programme which was broadcasted in multiple episodes, one linked two context documents associated to the same broadcast and one link was an error in the database.

For each document in the top 10 list, we also inspected the next four with the most overlapping keywords. Of these 100 links, most were semantic links (83 for the automatic annotations and 86 for the manual annotations).

In Figure 2, we display all Likert scale judgements of the semantic links. It seems that the average quality of the semantic links is not very high: the average tends slightly more to *neutral* than to *poor* for both sets. Given the small size of our corpus, this is not very unexpected. It contains too few documents to generate many very good links. Still, both the automatic annotations and the manual annotations have 21 judgements in the *good* or *very good* group. So with both annotations, we could find some quite interesting links between documents for this small corpus. They do generate very different results however. Only eight of the pairs appear in both sets (8 out of 100), i.e. eight pairs were linked both via the manual annotations and the automatic annotations. Six of these are parts of one TV-programme which was broadcasted in multiple episodes. Both their catalogue descriptions and their context documents were very similar.

Qualitative inspection

When we look at examples of semantic overlap, we see very interesting results. We see for example that *Andere Tijden 2004-01-06* and *Andere Tijden*

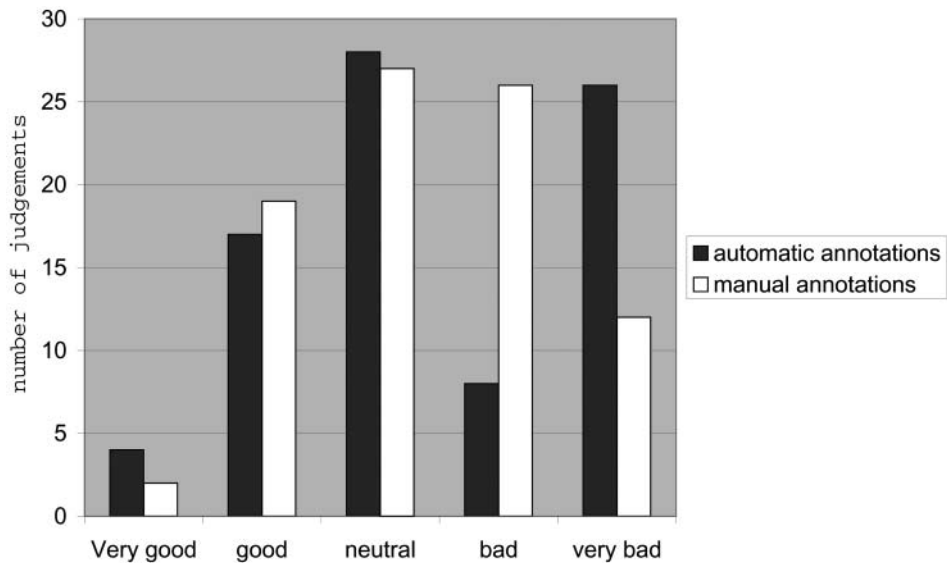


FIGURE 2 Judgements of the semantic links on a five-point scale.

2004-12-07 get paired by the automatic annotations. The second programme incorporated much of the content of the first programme. According to the catalogue description, the topic of the first programme is: “*the first Bilderberg-Conference which was held in 1954 under the presidency of Prince Bernhard*”. The topic of the second programme is: “*the role Prince Bernard played in the international circuit of politicians, soldiers and businessmen, especially his presidency of the international Bilderberg-meeting and his friendship with journalist Martin van Amerongen*”. This second programme was broadcasted just after the death of Prince Bernard and incorporated much of the first programme’s material. The catalogue description does not mention the relation between the programmes and the catalogue descriptions do not show a large overlap in terms of manual keywords. We managed to relate these documents because the original makers adapted a context document of the first programme and associated it with the second programme. The automatic annotations derived from the original and the adapted context document show a large overlap. The manual annotations have only one overlapping keyword. The first programme was indexed with the keywords *history, post-war rebuilding, secrecy, foreign policy, anti-Americanism, anti-communism*. The second programme was indexed with the keywords *history, conferences, politicians, entrepreneurs*. This difference is not only the result of the difference in the programme. It serves as an example of inter-annotator differences within the archives of Sound and Vision.

Sometimes one document is semantically very similar for multiple different programmes (e.g. *Andere Tijden 2004-11-23 Rushdie affaire* has five overlapping manual keywords with *Andere Tijden 2003-09-30 Khomeiny* and four with *Andere Tijden 2005-02-01 The arrival of the mosque*). Some broadcasts address topics relevant for many others. So there is a tendency to cluster around quintessential documents. For a collection, these characteristic documents may be very interesting starting points for visualization and navigation.

Discussion

Serendipitous browsing was created as a new way to evaluate the perceived value of the automatic annotations. We were not able to capture this value in the evaluation against manual annotations, neither in the exact evaluation nor in the semantic evaluation. However, the information specialists from Sound and Vision appreciated the new use of automatic techniques in a practical archive setting. In particular, the automatic linking of documents, whether it is done on the basis of manual annotations or automatic annotations, appears valuable and reminds of usages of the archive with the former physical card system. This linking of documents cannot be performed by hand (i.e., by human cataloguers) and lies outside the scope of the current archiving. An interesting result is the similar value for semantic browsing of automatic annotations compared to manual annotations: both sets of annotations generated the same amount of ‘good’ and ‘very good’ relations and, on average, both relations were judged with the same score. This suggests that, although the automatic annotations are not as precise as the manual annotations, for semantic browsing purposes, they have the same value.

Discussion and perspectives

We set out to evaluate in three ways the value of automatic annotation suggestions for the audiovisual archive of Sound and Vision. The classic precision/recall evaluation showed that the baseline formed by TF.IDF ranking is the best ranking method. For the task of keyword suggestion within an archive however, this evaluation is too strict. The loosened semantic precision/recall measure showed that, instead of the TF.IDF ranking, the Mixed model performed best. As the Mixed model starts out poorest followed by the TF.IDF ranking, this result was only significant for the group of 10 first suggestions. Manual inspection showed that the Mixed model tended to suggest more general terms. The third evaluation of manual and automatic annotations was in the serendipitous browsing experiment. This showed that the manual annotations and the automatic annotations have the same value for finding interesting related documents. With this experiment, we only use the CARROT suggestions, so we are not able to differentiate ranking methods.

When we combine these three evaluation results and add to this the limited inter-annotator agreement, it becomes hard to see how manual annotations can serve as a *gold* standard. It is, however, the only material which we have. The question is how to evaluate against this resource and how to interpret the relevance of the outcome. As a first step, it is good to apply semantic evaluation. A second step which we are working on is a user evaluation of our keyword suggestions by cataloguers from Sound and Vision. This user study is meant to produce a human validation of the interest of the keyword suggestions for annotation and to obtain a deeper understanding of evaluation of our automatic keyword suggestion system.

As future work, we plan to experiment with the suggestion of keywords based on automatic speech transcripts from the broadcasts and compare the results with the output generated from the context documents presented in this paper.

The interdisciplinary circle in this paper has come to a close: the practical archive setting forced us to change the classical way of evaluation and adopt novel ways of evaluation of our keyword suggestion system. However, the changed view on the evaluation came back to the archive in the form of serendipitous browsing, which is perceived as a very interesting and probably valuable option for the daily archive. Even more interesting are our changed views: the problematic nature of evaluation changes the way we perceive information extraction and the archive, and gives a radical new view on the future of archiving: this foresees that it will encompass 80 per cent automatic annotation and 20 per cent manual annotation. Furthermore, the thinking on automatic annotation will generate new ideas for interacting with the archive.

Our research follows a storyline often seen in the humanities, but uncommon for the sciences: instead of finding an improved solution to a known problem, as is common in the sciences, we obtained an almost Socratic understanding of evaluation: we now know that we have a very limited understanding of evaluation and are only starting to grasp the vastness of its problematic nature: we found problems and wonderment, as is common in the humanities.

Notes

- ¹ For information on users of Sound and Vision's television broadcast archive, see Hollink *et al.* (2009) in this issue.
- ² Functional Requirements for Bibliographical Records, <http://archive.ifa.org/VII/s13/frbr/frbr.pdf> (last accessed 06/03/09).

Bibliography

- Brin, Sergey, and Lawrence Page. 1998. The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems* 30(1–7): 107–17.
- Ciravegna, Fabio, and Yorick Wilks. 2003. Designing adaptive information extraction for the semantic web in amilcare. In *Annotations for the semantic web*, ed. S. Handschuh and S. Staab, Volume 1, 112–127. Amsterdam: IOS Press.
- Cunningham, H., D. Maynard, K. Bontcheva, and V. Tablan. 2002. GATE: A framework and graphical development environment for robust NLP tools and applications. *Proceedings of the 40th Anniversary Meeting of the ACL*.
- Faaborg, Alexander, and Carl Lagoze. 2003. Semantic browsing. *ECDL*, 70–81.
- Heflin, J., and J. Hendler. 2000. Searching the web with shoe. *Proceedings of the AAAI-2000 Workshop on AI for Web Search*. Montenegro: Budva.
- Hildebrand, M. 2008. Interactive exploration of heterogeneous cultural heritage collections. In *The Semantic Web — ISWC 2008*, Volume 5318 of *Lecture Notes in Computer Science*, 483–98.
- Hollink, Laura, Bouke Huurnink, Michiel van Liempt, Johan Oomen, Annemieke de Jong, Maarten de Rijke, Guus Schreiber, and Arnold Smeulders. 2009. A multidisciplinary approach to unlocking television broadcast archives. *Interdisciplinary Science Reviews*, this issue.
- Iivonen, M. 1995. Consistency in the selection of search concepts and search terms. *Information Processing and Management* 31(2): 173–190 (March–April).
- Iria, Jos. 2005. T-Rex: A flexible relation extraction framework. *Proceedings of the Eighth Annual CLUIK Research Colloquium*. Manchester.
- Kahan, J., and M.-R. Koivunen. 2001. Annotea: an open RDF infrastructure for shared web annotations. *World Wide Web* 6:23–32.
- Kiryakov A., B. Popov, I. Terziev, D. Manov, and D. Ognyanoff. 2005. Semantic annotation, indexing, and retrieval. *Journal of Web Semantics* 2(1): 49–79.
- Leininger, K. 2000. Inter-indexer consistency in PsycINFO. *Journal of Librarianship and Information Science* 32(1): 4–8.
- Likert, Rensis. 1932. A technique for the measurement of attitudes. *Archives of Psychology* (140): 155.
- Malaisé, Véronique, Luit Gazendam, and Hennie Brugman. 2007. Disambiguating automatic semantic annotation based on a thesaurus structure. *14e conférence sur le Traitement Automatique des Langues Naturelles (TALN)*.
- Medelyan, Olena, and Ian H. Witten. 2006. Thesaurus-based index term extraction for agricultural documents.
- Salton, G., and M.J. McGill. 1983. *Introduction to modern information retrieval*. McGraw-Hill.
- Simon, Herbert. 1957. *Models of man*. New York: Wiley.
- Vargas-Vera M., Enrico Motta, John Domingue, Mattia Lanzoni, Arthur Stutt and Fabio Ciravegna. 2002. MnM: Ontology driven semi-automatic and automatic support for semantic markup. In *Proceedings of the Thirteenth International Conference on Knowledge Engineering and Management (EKAW-2002)*. Spain: Siguenza.
- Wang, Jinghua, Jianyi Liu, and Cong Wang. 2007. Keyword extraction based on PageRank. *Advances in Knowledge Discovery and Data Mining* 4426/2007: 857–864.
- Wartena, Christian, Rogier Brussee, Luit Gazendam, and Wolf Huijsen. 2007. Apolda: A practical tool for semantic annotation. *The Fourth International Workshop on Text-based Information Retrieval (TIR 2007)*. Germany: Regensburg.

Notes on Contributors

Correspondence to: Luit Gazendam, Novay, Postbus 589, 7500 AN Enschede, The Netherlands.

Email: luit.gazendam@novay.nl

Copyright of *Interdisciplinary Science Reviews* is the property of Maney Publishing and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.