# Exploiting Program Guides for Contextualisation

Lotte Belice Baltussen
Netherlands Institute for Sound and Vision
Hilversum, The Netherlands
Email: lbbaltussen@beeldengeluid.nl

Themistoklis Karavellas
Netherlands Institute for Sound and Vision
Hilversum, The Netherlands
Email: tkaravellas@beeldengeluid.nl

Roeland J.F. Ordelman
University of Twente,
Enschede, The Netherlands
Email: roeland.ordelman@utwente.nl

*Abstract*—Archives of cultural heritage organisations typically consist of collections in various formats (e.g. photos, video, texts) that are inherently related. Often, such disconnected collections represent value in itself but effectuating links between 'core' and 'context' collection items in various levels of granularity could result in a 'one-plus-one-makes-three' scenario both from a contextualisation perspective (public presentations, research) and access perspective. A key issue is the identification of contextual objects that can be associated with objects in the core collections, or the other way around. Traditionally, such associations have been created manually. For most organizations however, this approach does not scale. In this paper, we describe a case in which a semi-automatic approach was employed to create contextual links between television broadcast schedules in program guides (context collection) and the programs in the archive (core collection) of a large audiovisual heritage organisation.

*Index Terms*—*Contextualisation, Cultural Heritage, Data Analysis, Linked Collections*

## I. INTRODUCTION

The collections of heritage organisations such as libraries, museums and archives typically consist of collections in various formats that are inherently related. For instance, a National Library possesses books as a physical object and digital representation, but possibly also private collections of the authors of these books such as letters, notes or even objects. An art museum may have paintings, but also the original pamphlets used for previous exhibitions featuring these paintings [1]. An audiovisual archive might have a collection of broadcast video, and also television program guides or photos taken during the production of a program on the set. Although each collection represent value in itself, making explicit links between these 'core' and 'context' collection items in various levels of granularity increases their value both from a contextualization perspective (public presentations, research) and access perspective. Access scenarios are for example the use of information from contextual sources to enrich content descriptions of core data to improve 'searchability', or providing new entry points to either core or context collections via the establishment of explicit (hyper)links between items from the collections. For instance, the program descriptions in program guides can be used as *metadata* for audiovisual broadcast content, or end-users could be pointed towards audiovisual sources while browsing the collection of program guides.

A key issue, either from a contextualization perspective or an access perspective, is the identification of contextual media objects that can be associated with media objects in the core collections, or the other way around. Traditionally, such associations have been created manually by media archivists or collection specialists. For most organizations however, this approach does not scale. For large-scale contextual interlinking a (semi-) automatic approach is required: algorithms detect relations between items and suggest links that are validated manually by humans (e.g. archivists).

We consider a use case of an audiovisual archive with a large collection of broadcast video data, and a digitized collection of program guides representing four decades of television history. The guides contain TV and radio broadcast schedules, articles, pictures, advertisements and other elements that provide insight in the cultural-historical context of the broadcast videos. The use case aims at the automatic identification of program guide material that can be linked to the broadcasts in the archive, and by doing so, to contextualise the core collection with context data. Practically, the program guide material will be used as metadata to improve 'searchability' of archive video. Moreover, via an online program guide browser, direct access from the guides to related archival content will be provided.

In this paper, we describe how our work relates to other developments in the Digital Humanities domain, our approach for extracting and matching data from OCR files to the Netherlands Institute for Sound and Vision's (NISV) av-archive and our lessons-learned.

## II. RELATED WORK

In the LODLAM community[1] (Linked Open Data in Libraries, Archives and Museums) [2], [3] many activities are undertaken to connect collections from heritage institutions to both internal and external information. Linked Open Data technologies are used to support the process of creating connections. For example, the *Free Your Metadata*[2] initiative provides easy-to-use tools for making metadata part of the Linked Data cloud and organizes hands-on workshops.

As part of the CATCH Programme, together with the Intelligent Systems Lab Amsterdam at the University of Amsterdam, the Centre for Television in Transition at Utrecht University, we developed the tool CoMeRDa (Contextualizing Media Research Data[3]). It is an aggregated system that allows searching full-text across diverse data collections including television catalogue metadata, newspaper articles, photo collections, program guides and a specialized wiki. The interface combines information from the archives' core collection (catalogue info) with its context collections. In CoMeRDa, it is possible to use a record as a search query across the

---

[1]http://lodlam.net/
[2]http://freeyourmetadata.org/
[3]http://vps46235.public.cloudvps.com/bridge/tools/the-comerda-tool/

collections, but no explicit Linked Open Data or other linking technologies were used to make relations between records and collections explicit. [4]

Within the scope of the BBC Genome project [4] 350,000 pages of Radio Times program guides from 1922-2009 were digitised. Optical Character Recognition (OCR) software was used to automatically convert the characters on the digitised pages into machine-encoded text. This OCR process not only makes it possible to provide full-text search, but also to segment the information on the pages into distinct chunks. By analysing the mark-up and layout of the program schedules, the information about programs aired on television and radio were automatically turned into a massive database, listing a total of almost 4.5 million BBC program records. These records contain information on the program title, description, date, channel, cast, airing time, related to the specific page of the guide the information came from.

The project described in this paper took the approaches described above one step further by linking the extracted data of the context collection of the guides to the core av-collection of the archive, and explores how to optimise the quality of the extracted output and detected links.

## III. APPROACH

### A. Data collection

Our corpus consists of 330 guides of six public broadcasters[5], spanning four decades: 1951-1986. We didn't scan all guides from this period, but used intervals of five years, and for each year scanned two months worth of issues. For example, for 1951 we scanned the guides of all six broadcasters from January and September, for 1956 we did the same for the February and October, and so on. In total, 25,000 scanned pages were part of our corpus.

For each page, the OCR software Abbyy FineReader 10 was used to convert the image scans into a computer readable format (XML). We used a parsing script[6] developed specifically for the project to segment the XMLs of the broadcast schedules of TV programs of public broadcasters, in order to create a database of broadcasts that spans four decades of Dutch television history. For each record in the database (representing one program) we tried to capture the following structured information: date, start time of the broadcast, end time of the broadcast, channel, broadcaster name, title and description. In a next step, we mapped the extracted titles with those in the AV catalogue of NISV, in order to make the match between specific programs in the archive and the program guide information about them. This results in linking the previously analogue, paper contextual collections and the core av-archive.

### B. Parsing the OCR

There is more information in the guides than TV schedules, such as advertisements, articles and photos. Although these additional data provide great contextual information as well, we focused only on extracting descriptions regarding Dutch public TV programs. This process starts with identifying the correct information to parse by analysing the OCR output captured in a raw OCR file. In this raw file, different parameters are used to described each detected character: for example the coordinate of the character on the page, confidence score, font type, font size, the position of the character within a word, and if it is bold or italic. Our data extraction algorithm groups separate characters into words. A group of words form a line and a group of lines form a block. The block of text is analysed to determine whether it contains TV program schedule information.

A block of text is checked for elements that typically describe TV programs, such as temporal elements representing start and end times, a date, and broadcaster and TV channel names from a predefined list. If temporal elements and a title and description are found, we assume that the block contains information about one single program. The first temporal element in this block is considered to be the start time. The rest of the text is parsed to detect the other elements: end time, broadcaster name, and the title and program description. The channel name and date are usually found outside of the text block, at the top of a page and are extracted separately. (see Fig. 1). Following this approach, each text block describing an individual TV program populates an XML record. Each attribute gets its own place within this record, including the page of the guide on which the text block can be found. After all pages in a guide are parsed and all detected TV program text blocks are examined, the extracted records are saved in an XML file representing the entire guide.



Fig. 1. Example of a TV broadcast schedule of four progammes, containing the start times, titles and descriptions. The names of the cast also mentioned. In this case, the date and channel name are found elsewhere on the page

In total, 990,000 distinct broadcasts were identified within the corpus of 330 guides containing 25,000 pages. It needs to be noted that the bulk of these broadcasts overlap, as guides –originally published by different broadcasters– represent the same years and weeks. Also, this set contains programs that were broadcast more than once (reruns). A sample set was manually checked for accuracy of the extracted elements, which showed that 91.6 percent of the title and description elements are (almost) complete and correct and that in almost all cases the start time of a program was extracted accurately. The elements channel, broadcaster and date were however only parsed correctly in around 50-55 percent of the cases.

---

[4] http://genome.ch.bbc.co.uk/

[5] In the Netherlands, the public broadcast system consists of many different broadcasters, each targeting a societal group (e.g. Protestants, Catholics). We scanned guides from the following broadcasters: AVRO,KRO,NCRV,TROS,VARA,VPRO

[6] https://github.com/beeldengeluid/tvguide-segmenter

## C. Main parsing issues

We encountered a number issues which made it difficult to parse information correctly. The main issues are:

- Incorrect OCR output: OCR-output typically has flaws. Due to ink bleeding, rips and tears in the paper, ink fading or poor print quality some characters in the OCR output were incorrectly detected.

- Text flow order: the guides present TV schedules in different ways. Some guides present them in neatly defined vertical columns per broadcasters, whereas others use a horizontal presentation. This makes it hard to identify to which channel a text block belongs and to identify when the text in a text block ends on the bottom of a column and continues at the top of the next column (see left page in Fig. 2).

- Layout identification: each guide has its own graphic design and way to arrange elements on a page. Each guide has its own way of presenting information: some print the start and end times of programmes in Italic, some in bold, others always print program titles in bold and font size 16, in some guides the date is presented at the top left of a page and so on. Usually, this styling of elements differs per broadcaster, per year.

- Page selection for segmentation: we only wanted to parse TV program schedules from Dutch public broadcasters. To this end, we had to create specific lists of TV broadcaster and channel names in order to identify the correct pages to parse.

- Date parsing: in general, the day of the program schedule is mentioned at the top of the page. However, many formats are used to describe the date e.g. "Friday 13 May" or "12 Jan". This makes it very hard to convert the date –if detected by the parsing script at all– to a unified format.

The most important action undertaken in the project to improve parsing accuracy was creating so-called configuration files. These files convey information about the text flow order and layout of information of a certain guide, e.g. the location of the date, the font size of titles and which elements were presented in bold. Configuration files were manually created for each guide per year and provided to the parsing script. Although it took a couple of days to create these configuration files, the accuracy output of the script improved dramatically, in some cases resulting in the detection of three times as many individual program blocks on a single page.

## D. Matching parsed results with television programs

After the TV program schedules were parsed, links were identified between the extracted program guide records and the corresponding programs in NISV's catalogue. This was done by creating a list of program titles in the catalogue corresponding to the date range of each parsed program guide. So, if a guide represents information from 8 January 1961 to 14 January 1961, only programs in the catalogue from this date range are considered in the matching process. When a title was detected in the description of a parsed record, this was recorded as a match. This process is however not 100 percent accurate:

- The date extracted from the guide was not always correct or not detected at all.

- The matching script also allowed partial matching. This means that a title like "Jeugdjournaal" (Children's News) also matches the title "Journaal" (News).

- Reruns of programs appear in the program schedule of a guide, but in the NISV archive generally only the first, original broadcasts are stored. This can result in a false positive match when a program is rerun within the date range of a single guide.

The total *match potential* was $6,462$. In other words, for the dates of the parsed guides, there are $6,462$ individual television programs in NISV's archive. Of these programmes, $2,611$ were matched. Because we scanned six different guides from the same weeks, in many cases more than one match was found per program, resulting in a total of $6,224$ matches (see Table I). This amount may seem low, considering that $990,000$ records were found. However, there are large gaps in the NISV archive for the period of the parsed guides (330 guides from 1951-1986). The reason is that it only became possible to archive all broadcasts in the 2000s, when fully digital workflows were implemented. Before, tapes were sent back and forth between broadcasters and the archive, and a selection policy was employed as it was too costly to archive everything. Moreover, sometimes other programs besides Dutch public television programs were parsed by our script, such as radio broadcasts and programs from foreign channels.

| Matches per program | Occurrences | Total matches |
|---|---|---|
| 1 | 1,147 | 1,147 |
| 2 | 432 | 864 |
| 3 | 288 | 864 |
| 4 | 422 | 1,688 |
| 5 | 275 | 1,375 |
| 6 | 47 | 282 |
| **TOTAL** | **2,611** | **6,220** |

TABLE I.    Number of matches per program, ranging from a match from 1 guide to 6 guides

A representative sample of the matched records was checked manually for accuracy. We found that $89.2$ percent of the matches are accurate, and that for $6.6$ percent of the inaccurate matches the right program was matched, but not the right episode. The latter was due to the date being incorrect or not present in the parsed output. This means that in $4.2$ percent of the cases the match was completely inaccurate. However, the match to the guide itself is always correct, since the broadcast date of the program always falls within the date range of the guide.

## E. Extracting thesaurus terms

For each parsed title and description element within an identified text block, Named Entity Recognition was used to extract terms from the archive's thesaurus[7], which contains terms relevant to the NISV archive. For $63.8$ percent of the $6,224$ total matches, one or more terms were extracted. 73

---

[7] http://openskos.beeldengeluid.nl/api

Fig. 2. Program guide text flow order and layout examples

percent of these terms are person names, and 17 percent are geographical names. A sample set was analysed and showed that 30 percent of these extracted terms match the terms that were manually added by archivists to the records in the NISV catalogue. This means that the remaining 70 percent adds new information to these records, mainly concerning the names of cast and crew.

## IV. FUTURE WORK

The project was finalised by importing the matched records into the NISV catalogue. However, many catalogue systems, including the one of NISV are not well-suited yet for storing these links. For instance, it is not possible to click on through to a program guide from the record of matched TV program. Catalogue systems would require significant changes in order to truly facilitate storing contextual links and for machines to access them through APIs. Also, applications need to be developed that allow end users to fully explore these links.

Furthermore, the parsing script, configuration files and matching scripts can be improved to increase accuracy of the output in various stages of the workflow. This is not just needed to improve results for the set we have used for this pilot contextualisation project, but especially when scaling up. Besides the 330 guides used for the pilot, the archive has scanned 3,300 more guides from six Dutch public broadcasters from 1951-2014, covering 250,000 pages containing the complete overview of what was broadcast on Dutch television and radio.

Finally, we would like to explore the contextualisation possibilities further, by creating links between the program guides and the archive's large photo collections and (linked) data initiatives like Wikidata[8] to come full circle and connect not only our internal collections, but to link to and from relevant external sources as well.

---

8http://www.wikidata.org/

## REFERENCES

[1] G. Auffret and B. Bachimont, "Audiovisual cultural heritage: from tv and radio archiving to hypermedia publishing," in *Research and Advanced Technology for Digital Libraries*. Springer, 1999, pp. 58–75.

[2] D. W. Oard, A. S. Levi, R. L. Punzalan, and R. Warren, "Bridging communities of practice: Emerging technologies for content-centered linking," in *Museums and the Web*, 2014.

[3] C. A. Harper, "Linked open data in libraries, archives and museums (lodlam): current trends, tools & techniques, and the role of vendors," 2013.

[4] M. Bron, J. van Gorp, F. Nack, M. de Rijke, A. Vishneuski, and S. de Leeuw, "A subjunctive exploratory search interface to support media studies researchers," in *Proceedings of the 35th International ACM SIGIR Conference on Research and Development in Information Retrieval*, ser. SIGIR '12. New York, NY, USA: ACM, 2012, pp. 425–434. [Online]. Available: http://doi.acm.org/10.1145/2348283.2348342