



UNIVERSITY OF AMSTERDAM

## UvA-DARE (Digital Academic Repository)

### Media Suite

*Unlocking Archives for Mixed Media Scholarly Research*

Noordegraaf, J.J.

[Link to publication](#)

*Citation for published version (APA):*

Noordegraaf, J. J. (2018). Media Suite: Unlocking Archives for Mixed Media Scholarly Research. 21-25.

#### General rights

It is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), other than for strictly personal, individual use, unless the work is under an open content license (like Creative Commons).

#### Disclaimer/Complaints regulations

If you believe that digital publication of certain material infringes any of your rights or (privacy) interests, please let the Library know, stating your reasons. In case of a legitimate complaint, the Library will make the material inaccessible and/or remove it from the website. Please Ask the Library: <https://uba.uva.nl/en/contact>, or a letter to: Library of the University of Amsterdam, Secretariat, Singel 425, 1012 WP Amsterdam, The Netherlands. You will be contacted as soon as possible.

UvA-DARE is a service provided by the library of the University of Amsterdam (<http://dare.uva.nl>)

# **CLARIN Annual Conference 2018**

## **PROCEEDINGS**

Edited by

Inguna Skadiņa, Maria Eskevich

8-10 October 2018

Pisa, Italy

# Programme Committee

## Chair:

- Inguna Skadiņa, Institute of Mathematics and Computer Science, University of Latvia & Tilde (LV)

## Members:

- Lars Borin, Språkbanken, University of Gothenburg (SE)
- António Branco, Universidade de Lisboa (PT)
- Koenraad De Smedt, University of Bergen (NO)
- Griet Depoorter, Institute for the Dutch Language (NL/Vlanders)
- Jens Edlund, KTH Royal Institute of Technology (SE)
- Tomaž Erjavec, Dept. of Knowledge Technologies, Jožef Stefan Institute (SI)
- Francesca Frontini, University of Montpellier (FR)
- Eva Hajičová, Charles University (CZ)
- Erhard Hinrichs, University of Tübingen (DE)
- Nicolas Larrousse, Huma-Num (FR)
- Krister Lindén, University of Helsinki (FI)
- Bente Maegaard, University of Copenhagen (DK)
- Karlheinz Mörth, Institute for Corpus Linguistics and Text Technology, Austrian Academy of Sciences (AT)
- Monica Monachini, Institute of Computational Linguistics "A. Zampolli" (IT)
- Costanza Navarretta, University of Copenhagen (DK)
- Jan Odijk, Utrecht University (NL)
- Maciej Piasecki, Wrocław University of Science and Technology (PL)
- Stelios Piperidis, Institute for Language and Speech Processing (ILSP), Athena Research Center (EL)
- Kiril Simov, IICT, Bulgarian Academy of Sciences (BG)
- Marko Tadić, University of Zagreb (HR)
- Jurgita Vaičėnienė, Vytautas Magnus University (LT)
- Tamás Váradi, Research Institute for Linguistics, Hungarian Academy of Sciences (HU)
- Kadri Vider, University of Tartu (EE)
- Martin Wynne, University of Oxford (UK)

## Reviewers:

- Ilze Auziņa, LV
- Bob Boelhouwer, NL
- Daan Broeder, NL
- Silvia Calamai, IT
- Roberts Dargis, LV
- Daniël de Kok, DE
- Riccardo Del Gratta, IT
- Christoph Draxler, DE
- Dimitrios Galanis, GR
- Maria Gavrilidou, GR
- Luís Gomes, PT
- Normunds Grūzītis, LV
- Jan Hajič, CZ
- Marie Hinrichs, DE
- Pavel Ircing, CZ
- Mateja Jemec Tomazin, SI
- Neeme Kahusk, EE
- Fahad Khan, IT
- Alexander König, IT
- Jakub Mlynar, CZ
- Jiří Mírovský, CZ
- Marcin Oleksy, PL
- Petya Osenova, BG
- Haris Papageorgiou, GR
- Hannes Pirker, AT
- Marcin Pol, PL
- Valeria Quochi, IT
- João Rodrigues, PT
- Ewa Rudnicka, PL
- Irene Russo, IT
- João Silva, PT
- Egon W. Stemle, IT
- Pavel Stranak, CZ
- Thorsten Trippel, DE
- Vincent Vandeghinste, BE
- Jernej Vičič, SI
- Jan Wieczorek, PL
- Tanja Wissik, AT
- Daniel Zeman, CZ
- Claus Zinn, DE
- Jerneja Žganec Gros, SI

## **CLARIN 2018 submissions, review process and acceptance**

- Call for abstracts: 17 January 2018, 28 February 2018
- Submission deadline: 30 April 2018
- 77 submissions in total were received and reviewed (three reviews per submission)
- Face-to-face PC meeting in Wroclaw: 21-22 June 2018
- Notifications to authors: 2 July 2018
- 44 accepted submissions: 21 oral presentations, 23 posters/demos

More details can be found at <https://www.clarin.eu/event/2018/clarin-annual-conference-2018-pisa-italy>.

# Table of Contents

## **Thematic Session: Multimedia, Multimodality, Speech**

### *EXMARaLDA meets WebAnno*

Steffen Remus, Hanna Hedeland, Anne Ferger, Kristin Bührig and Chris Biemann ..... 1

### *Human-human, human-machine communication: on the HuComTech multimodal corpus*

Laszlo Hunyadi, Tamás Váradi, István Szekrényes, György Kovács, Hermina Kiss and Karolina Takács  
6

### *Oral History and Linguistic Analysis. A Study in Digital and Contemporary European History*

Florentina Armaselu, Elena Danescu and François Klein ..... 11

### *The Acorformed Coprus: Investigating Multimodality in Human-Human and Human-Virtual Patient Interactions*

Magalie Ochs, Philippe Blache, Grégoire Montcheuil, Jean-Marie Pergandi, Roxane Bertrand, Jorane Saubesty, Daniel Francon and Daniel Mestre ..... 16

### *Media Suite: Unlocking Archives for Mixed Media Scholarly Research*

Roeland Ordelman, Liliana Melgar, Carlos Martinez-Ortiz, Julia Noordegraaf and Jaap Blom .. 21

## **Parallel Session 1: CLARIN in Relation to Other Infrastructures and Projects**

### *Using Linked Data Techniques for Creating an IsiXhosa Lexical Resource - a Collaborative Approach*

Thomas Eckart, Bettina Klimek, Sonja Bosch and Dirk Goldhahn ..... 26

### *A Platform for Language Teaching and Research (PLT&R)*

Maria Stambolieva, Valentina Ivanova and Mariyana Raykova ..... 30

### *Curating and Analyzing Oral History Collections*

Cord Pagenstecher ..... 34

## **Parallel Session 2: CLARIN Knowledge Infrastructure, Legal Issues and Dissemination**

### *New exceptions for Text and Data Mining and their possible impact on the CLARIN infrastructure*

Pawel Kamocki, Erik Ketzan, Julia Wildgans and Andreas Witt ..... 39

### *Processing personal data without the consent of the data subject for the development and use of language resources*

Aleksei Kelli, Krister Lindén, Kadri Vider, Pawel Kamocki, Ramūnas Birštonas, Silvia Calamai, Chiara Kolletzek, Penny Labropoulou and Maria Gavrilidou ..... 43

### *Toward a CLARIN Data Protection Code of Conduct*

Pawel Kamocki, Erik Ketzan, Julia Wildgans and Andreas Witt ..... 49

### **Parallel Session 3: Use of the CLARIN infrastructure**

<i>From Language Learning Platform to Infrastructure for Research on Language Learning</i> David Alfter, Lars Borin, Ildikó Pilán, Therese Lindström Tiedemann and Elena Volodina . . . . .	53
<i>Bulgarian Language Technology for Digital Humanities: a focus on the Culture of Giving for Education</i> Kiril Simov and Petya Osenova . . . . .	57
<i>Multilayer Corpus and Toolchain for Full-Stack NLU in Latvian</i> Normunds Grūzītis and Artūrs Znotiņš . . . . .	61
<i>(Re-)Constructing "public debates" with CLARIAH MediaSuite tools in print and audiovisual media</i> Berrie van der Molen, Jasmijn van Gorp and Toine Pieters . . . . .	66
<i>Improving Access to Time-Based Media through Crowdsourcing and CL Tools: WGBH Educational Foundation and the American Archive of Public Broadcasting</i> Karen Cariani and Casey Davis-Kaufman . . . . .	66

### **Parallel Session 4: Design and construction of the CLARIN infrastructure**

<i>Discovering software resources in CLARIN</i> Jan Odijk . . . . .	72
<i>Towards a protocol for the curation and dissemination of vulnerable people archives</i> Silvia Calamai, Chiara Kolletzek and Aleksei Kelli . . . . .	77
<i>Versioning with Persistent Identifiers</i> Martin Matthiesen and Ute Dieckmann . . . . .	82
<i>Interoperability of Second Language Resources and Tools</i> Elena Volodina, Maarten Janssen, Therese Lindström Tiedemann, Nives Mikelic Preradovic, Silje Karin Ragnhildstveit, Kari Tenfjord and Koenraad de Smedt . . . . .	86
<i>Tweak Your CMDI Forms to the Max</i> Rob Zeeman and Menzo Windhouwer . . . . .	91

### **Poster session**

<i>CLARIN Data Management Activities in the PARTHENOS Context</i> Marnix van Berchum and Thorsten Trippel . . . . .	95
<i>Integrating language resources in two OCR engines to improve processing of historical Swedish text</i> Dana Dannélls and Leif-Jöran Olsson . . . . .	100
<i>Looking for hidden speech archives in Italian institutions</i> Vincenzo Galatà and Silvia Calamai . . . . .	104
<i>Setting up the PORTULAN / CLARIN centre</i> Luís Gomes, Frederico Apolónia, Ruben Branco, João Silva and António Branco . . . . .	108
<i>LaMachine: A meta-distribution for NLP software</i> Maarten van Gompel and Iris Hendrickx . . . . .	112
<i>XML-TEI-URS: using a TEI format for annotated linguistic resources</i> Loïc Grobol, Frédéric Landragin and Serge Heiden . . . . .	116
<i>Visible Vowels: a Tool for the Visualization of Vowel Variation</i> Wilbert Heeringa and Hans Van de Velde . . . . .	120
<i>ELEXIS - European lexicographic infrastructure</i> Milos Jakubicek, Iztok Kosem, Simon Krek, Sussi Olsen and Bolette Sandford Pedersen . . . . .	124
<i>Sustaining the Southern Dutch Dialects: the Dictionary of the Southern Dutch Dialects (DSDD) as a case study for CLARIN and DARIAH</i>	

Van Keymeulen Jacques, Sally Chambers, Veronique De Tier, Jesse de Does, Katrien Depuydt, Tanneke Schoonheim, Roxane Vandenberghe and Lien Hellebaut	128
<i>SweCLARIN – Infrastructure for Processing Transcribed Speech</i>	
Dimitrios Kokkinakis, Kristina Lundholm Fors and Charalambos Themistokleous	133
<i>TalkBankDB: A Comprehensive Data Analysis Interface to TalkBank</i>	
John Kowalski and Brian MacWhinney	137
<i>L2 learner corpus survey – Towards improved verifiability, reproducibility and inspiration in learner corpus research</i>	
Therese Lindström Tiedemann, Jakob Lenardič and Darja Fišer	142
<i>DGT-UD: a Parallel 23-language Parsebank</i>	
Nikola Ljubešić and Tomaž Erjavec	147
<i>DI-ÖSS - Building a digital infrastructure in South Tyrol</i>	
Verena Lyding, Alexander König, Elisa Gorgaini and Lionel Nicolas	151
<i>Linked Open Data and the Enrichment of Digital Editions: the Contribution of CLARIN to the Digital Classics</i>	
Monica Monachini, Francesca Frontini, Anika Nicolosi and Fahad Khan	155
<i>How to use DameSRL: A framework for deep multilingual semantic role labeling.</i>	
Quynh Ngoc Thi Do, Artuur Leeuwenberg, Geert Heyman and Marie-Francine Moens	159
<i>Speech Recognition and Scholarly Research: Usability and Sustainability</i>	
Roeland Ordelman and Arjan van Hessen	163
<i>Towards TICCLAT, the next level in Text-Induced Corpus Correction</i>	
Martin Reynaert, Maarten van Gompel, Ko van der Sloot and Antal van den Bosch	169
<i>SenSALDO: a Swedish Sentiment Lexicon for the SWE-CLARIN Toolbox</i>	
Jacobo Rouces, Lars Borin, Nina Tahmasebi and Stian Rødven Eide	173
<i>Error Coding of Second-Language Learner Texts Based on Mostly Automatic Alignment of Parallel Corpora</i>	
Dan Rosén, Mats Wirén and Elena Volodina	177
<i>Using Apache Spark on Hadoop Clusters as Backend for WebLicht Processing Pipelines</i>	
Soheila Sahami, Thomas Eckart and Gerhard Heyer	181
<i>UWebASR – Web-based ASR engine for Czech and Slovak</i>	
Jan Švec, Martin Bulín, Aleš Pražák and Pavel Ircing	186
<i>Pictograph Translation Technologies for People with Limited Literacy</i>	
Vincent Vandeghinste, Leen Sevens and Ineke Schuurman	190



## Media Suite: Unlocking Archives for Mixed Media Scholarly Research

**Roeland Ordelman**

Netherlands Institute for Sound and Vision  
University of Twente  
The Netherlands  
rordelman@beeldengeluid.nl

**Liliana Melgar**

Department of Media Studies  
University of Amsterdam  
The Netherlands  
melgar@uva.nl

**Carlos Martinez-Ortiz**

Netherlands eScience Center  
Amsterdam  
The Netherlands  
c.martinez@esciencecenter.nl

**Julia Noordegraaf**

Department of Media Studies  
University of Amsterdam  
The Netherlands  
J.J.noordegraaf@uva.nl

### Abstract

This paper discusses the rationale behind the development of a research environment –the Media Suite– in a sustainable, dynamic, multi-institutional infrastructure that supports mixed media scholarly research with large multimedia data collections, serving media scholars and digital humanists in general.

### 1 Introduction

In some domains of scholarly research, the focus is on the creation of new data collections. In astronomy for instance, new collections of astronomical observations are made publicly available on a regular basis. In other domains such as Media Studies research focuses on data collections maintained at cultural heritage institutions, archives, libraries, and knowledge institutions. However, especially when audiovisual media are concerned, access to, and use of these collections is often restricted due to intellectual property rights (IPR) or privacy issues (e.g., with respect to recorded interviews). Moreover, individual institutions often do not have the technical infrastructure in place to serve basic scholarly needs with respect to search, exploration and inspection of individual items (play-out, viewing). Therefore, scholars either fall back on collections that are openly available or spend considerable amounts of time in *onsite* visits to archives for consulting data collections. Data collections at these institutes can be regarded as “locked”, or at least hard to use for scholarly research.

To unlock these “institutional” collections and let scholars take advantage of the sheer quantity and richness of these data sets, we are developing an infrastructure for *online* scholarly exploration of collections that are distributed across various “institutional” content owners. Specifically, we focus on *audiovisual* data collections and related *mixed-media* sources, such as radio and television broadcasts, film, oral history interviews but also (news)paper archives, film posters and eyewitness reports. The *Media Suite* serves as the online portal to the infrastructure where first of all, content and metadata can be explored, browsed, compared, and stored in personal collections. In addition, the Media Suite provides a workspace for working with mixed media collections, providing tools for manual and automatic annotation, visualization, analysis, and sharing.

The ultimate goal is to (i) enable distant reading (Schulz, 2011), that is, identifying patterns or new research questions in all aggregated collections, (ii) facilitate close reading: the detailed examination of individual items (e.g., videos) in a collection or parts of these items (e.g., video segments) during search and scholarly interpretation, and (iii) make sure that the “scholarly primitives” (Unsworth, 2000;

---

This work is licenced under a Creative Commons Attribution 4.0 International Licence. Licence details: <http://creativecommons.org/licenses/by/4.0/>

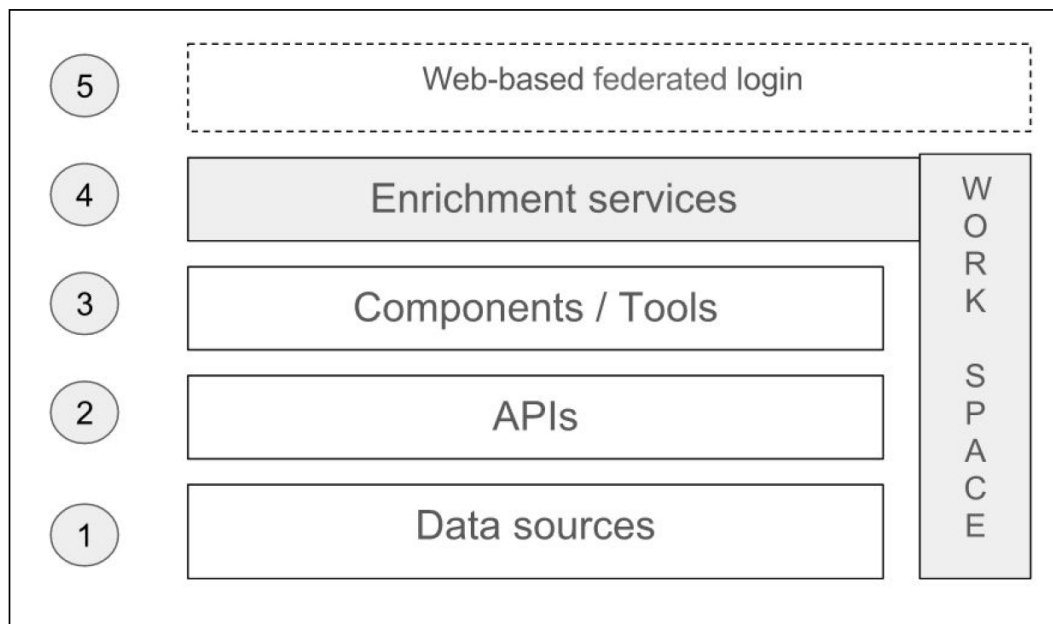


Figure 1: The building blocks of the CLARIAH Media Suite

Blanke and Hedges, 2013), basic activities common to research across humanities disciplines, are well supported.

### 1.1 Challenges

Questions however are: How to facilitate “close reading” when the media objects cannot be accessed because of copyright issues? How to enable “distant reading” when metadata is sparse, or diverse, and incomplete? How to cater to the needs of scholars with specific research questions and methods in the context of an infrastructure that has to be generic enough to be feasible? How to enable scholars to work with collections from different institutes using the same tools, when these collections are “locked”? How to enable scholars that are computer literate to work directly with the data or to deploy private content analysis tools such as computer vision or sentiment analysis?

The approach of the CLARIAH Media Suite to tackle these challenges is to provide mechanisms that enable researchers to work with tools and aggregated data *within* the closed environment of the infrastructure sealed with a federated authentication mechanism (SURFConext<sup>1</sup>) that currently only serves scholars with a university account in the Netherlands, but that soon will be expanded to the CLARIN federation. Also, the so called ‘homeless users’ that do not have an account with an academic institution, will eventually have the opportunity to request for a login. We refer to this approach as to “bringing the tools to the data”, as opposed to “bringing the data to the tools”.

Figure 1 shows the main elements that constitute the Media Suite research environment. Below we discuss shortly each of these elements.

## 2 Data Sources – Data Governance

Institutional collection maintainers have internal data governance processes to ensure that data assets are formally managed. One important aspect covered by governance processes is licensing: who has

<sup>1</sup><https://www.surf.nl/en/services-and-products/surfconext/index.html>

permission to access the data. However, data governance with respect to external processes –loosely defined as being part of an ‘infrastructure’– is typically not accounted for. This means that key data governance areas such as availability (e.g., metadata can be harvested), usability (e.g., source data can be viewed), integrity (e.g., protocols are in place to handle duplication and enrichment), and security (e.g., provenance information is maintained), need to be (re)organized or (re)considered, formalized and supported by the Media Suite and the emerging infrastructure in which it is embedded.

### **3 APIs – Sustainable development**

A digital infrastructure should use existing protocols, conventions, and standards. Besides obtaining data by harvesting using the OAI-PMH protocol, or using application programming interfaces (APIs), the functionalities have been organized in a modular approach, which includes (Martinez-Ortiz et al., 2017):

- Components that use API’s to perform specific tasks.
- Tools that incorporate a number of components in a tool.

### **4 Components/Tools – User-friendly interaction design**

Developing new tools “from scratch” for every research question would be a very inefficient (and costly!) endeavour. The digital infrastructure should provide tools that are suitable both for common scholarly tasks and for specific tasks required by each discipline. However, the digital humanities community incorporates a wide diversity of scholars with different research questions, methods, and levels of expertise in working with information processing techniques and technologies. We address this challenge by (i) focusing on the similarities in research methods from different disciplines (de Jong et al., 2011; Melgar Estrada and Koolen, 2018), (ii) analyzing tools that support qualitative methods (Melgar et al., 2017), and (iii) working with scholars as co-developers in the process. The resulting functionalities are built in a modular (lego) approach that supports both flexible software development of components and user-friendly interaction with assembled tools.

### **5 Work Space – Working with audio-visual content and private data**

In addition to IPR and privacy restrictions, access to the audiovisual content in the Media Suite is also limited due to its nature; consisting of pixels (video) and samples (audio) and hopefully some manually generated metadata or subtitles (text). Typically, scholars want to search audiovisual data using (key)words that may be ‘hidden’ (encoded) in the pixels or the samples. This is called the semantic gap (Smeulders et al., 2000) that needs to be “bridged” by decoding the information in the pixels and the samples to semantic representations, e.g., a verbatim transcription of the speech or labels of visual concepts in the video (a car, a face, the Eiffel Tower), that can be matched with the keywords from the scholars. These semantic representations can be generated manually or, especially when data collections are large, automatically using automatic speech recognition (ASR) or computer vision technology. The generation of semantic representations is addressed in different ways. On the one hand, tools such as ASR are regarded as ‘must have’ components in an infrastructure focusing on fine-grained access. We are implementing an automatic speech recognition service that resides within the CLARIAH infrastructure that can handle requests from the infrastructure itself (e.g., bulk processing of collections, possibly activated by a scholar with an interest in a specific data set), but also requests from individual scholars that want to process their private collections. On the other hand, supporting manual annotation is key for interpretation in scholarly contexts. The Media Suite aims to support the generation of both ways of semantic representations in complementary ways via information workflows centred around a “Work Space” (see Figure 2) that has the following functionalities:

- Storing individual items from different “institutional” collections resulting in a private, virtual, multimedia, research collection.
- Storing private session data such as queries and filtering options.

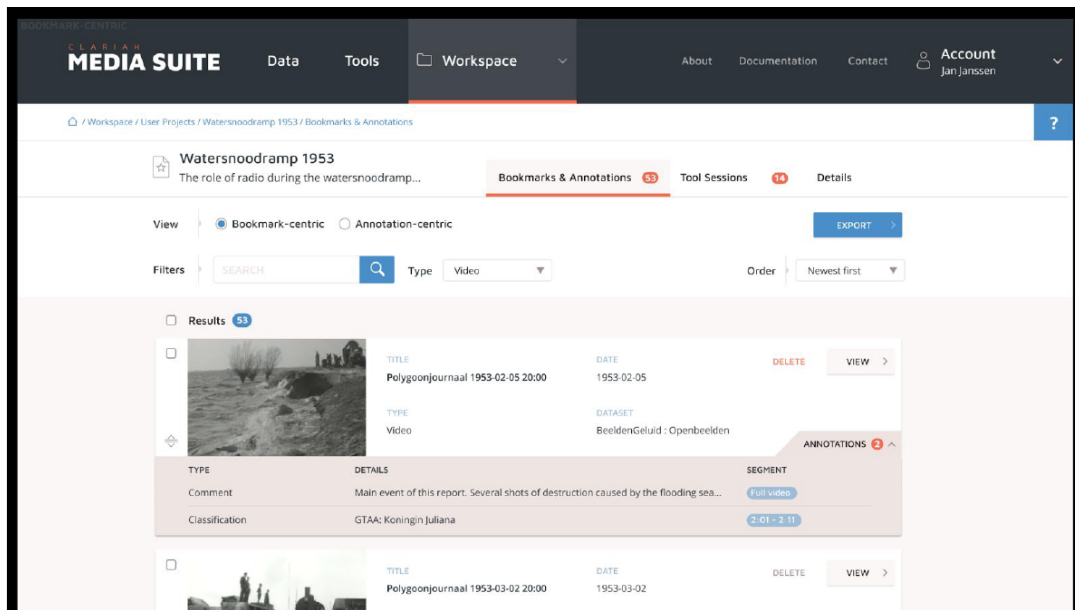


Figure 2: The CLARIAH Media Suite's Workspace

- Uploading private data and perform enrichment services to these data (e.g., speech recognition)
- Running private code on data collections in the infrastructure for creating data visualization (e.g., Jupyter Notebooks).

## 6 Conclusion and future work

We described the challenges found in building an infrastructure that satisfies the needs of humanities scholars working with audio-visual media and contextual collections. We choose the approach of building a research environment that adheres to infrastructural requirements while at the same time being flexible and user-friendly. In order to develop this environment in a sustainable way, that can be used and developed further after the project's lifetime, we need to carefully align the requirements of scholars with the context of the ecosystem the Media Suite needs to live in: an ICT infrastructure hosted and maintained by multiple institutions that in turn, adheres to a diverse set of institutional requirements with respect to, for instance, data access permissions and software development and maintenance. In order to have this infrastructure it is required that it is generic enough to cater for the general needs of every group that we have identified, while at the same time it incorporates flexible functionality capable of addressing very specialistic research questions. The Media Suite is currently functional and used by scholars doing actual research projects and will be developed further, e.g., by incorporating additional data sources (e.g., social media data), increasing metadata granularity (e.g., adding computer vision or emotion recognition), adding advanced annotation tools, and supporting missing data visualization (data critique) for heterogeneous datasets.

## References

Tobias Blanke and Mark Hedges. 2013. Scholarly primitives: Building institutional infrastructure for humanities e-science. *Future Generation Computer Systems*, 29(2):654–661.

- Franciska M.G. de Jong, Roeland J.F. Ordelman, and Stef Scagliola, 2011. *Audio-visual Collections and the User Needs of Scholars in the Humanities: a Case for Co-Development*, pages –. Centre for Language Technology, Copenhagen, 11. eemcs-eprint-20868.
- Carlos Martinez-Ortiz, Roeland Ordelman, Marijn Koolen, Julia Noordegraaf, Liliana Melgar, Lora Aroyo, Jaap Blom, Victor de Boer, Willem Melder, Jasmijn van Gorp, Eva Baaren, Kaspar Beelen, Norah Karrouche, Oana Inel, Rosita Kiewik, Themis Karavellas, and Thomas Poell. 2017. From tools to “recipes”: Building a media suite within the dutch digital humanities infrastructure clariah. DHBenelux.
- Liliana Melgar, Marijn Koolen, Hugo Huurdeman, and Jaap Blom. 2017. A process model of scholarly media annotation. In *Proceedings of the 2017 Conference on Conference Human Information Interaction and Retrieval, CHIIR '17*, pages 305–308, New York, NY, USA. ACM.
- Liliana Melgar Estrada and Marijn Koolen. 2018. Audiovisual media annotation using qualitative data analysis software: A comparative analysis. *The Qualitative Report*, 23(13):40–60.
- Kathryn Schulz. 2011. What is distant reading. *The New York Times*, 24.
- Arnold WM Smeulders, Marcel Worring, Simone Santini, Amarnath Gupta, and Ramesh Jain. 2000. Content-based image retrieval at the end of the early years. *IEEE Transactions on pattern analysis and machine intelligence*, 22(12):1349–1380.
- John Unsworth. 2000. Scholarly primitives: What methods do humanities researchers have in common, and how might our tools reflect this. In *Symposium on Humanities Computing: Formal Methods, Experimental Practice. King's College, London*, volume 13, pages 5–00.