Report from Dagstuhl Seminar 14302

# Digital Palaeography: New Machines and Old Texts

**Edited by**

# Tal Hassner[1], Robert Sablatnig[2], Dominique Stutzmann[3], and Ségolène Tarte[4]

1   **Open University of Israel – Raanana, IL**, `hassner@openu.ac.il`
2   **TU Wien, AT**, `sab@caa.tuwien.ac.at`
3   **Institut de Recherche et d'Histoire des Textes (CNRS) – Paris, FR**, `dominique.stutzmann@irht.cnrs.fr`
4   **University of Oxford, GB**, `segolene.tarte@classics.ox.ac.uk`

──── **Abstract** ────

This report documents the program and the outcomes of Dagstuhl Seminar 14302 "Digital Palaeography: New Machines and Old Texts", which focused on the interaction of Palaeography and computerized tools developed in Computer Vision for the analysis of digital images. This seminar intertwined research reports from the most advanced teams in the field and interdisciplinary discussions on the potentials and limitations of future research and the establishment of a community of practice in Digital Palaeography. It resulted in new research directions in the Computer Sciences and new research strategies in Palaeography and in a better understanding of how to conduct interdisciplinary research across all the fields of expertise involved in Digital Palaeography.

## 1   Executive Summary

*Dominique Stutzmann*
*Ségolène Tarte*

Digital Palaeography emerged as a research community in the late 2000s. Following a successful Dagstuhl Perspectives Workshop on Computation and Palaeography (12382)[1], this seminar focused on the interaction of Palaeography and computerized tools developed in Computer Vision for the analysis of digital images. Given the present techniques developed to enhance damaged documents, optical text recognition or computer-assisted transcription, identification and categorisation of scripts and scribes, the current technical challenge is

---

[1]  http://dx.doi.org/10.4230/DagMan.2.1.14

to develop "new machines", i. e. efficient solutions for palaeographic tasks, and to provide scholars with quantitative evidence towards palaeographical arguments, even beyond the reading of "old texts" (ancient, medieval and early modern documents), which is of interest to the industry, to the wider public, and to the broad community of genealogists.

The identified core issue was to create the conditions of a fluid and seamless communication between Humanities and Computer Sciences scholars in order to advance research in Palaeography, Manuscript Studies and History, on the one hand, and in Computer Vision, Semantic Technologies, Image Processing, and Human Computer Interaction (HCI) systems on the other hand. Indeed, researchers must articulate their respective systems of proof, in order to produce efficient systems that present palaeographical data quickly and easily, and in a way that scholars can understand, evaluate, and trust. To establish fruitful collaborations, it is thus essential to address the "black box" issue, to make a better use of the outreach potential offered by computerized technologies to enrich palaeographical knowledge, and to facilitate the sharing of both the CS and palaeographical methodologies.

This seminar was able to shed light onto two major evolutions between 2012 and 2014; these notable shifts are to do with interdisciplinary communication and with access to "black box" expertise. On the one hand, the notion of "communication" or "bridging the gap" (as expressed by seminar 14301, which took place in conjunction with our own seminar) has become more specific in that issues and problems are now better identified, understood, and expressed. While the two-fold expression "digital palaeography" might lead one to believe that the communication involves only two sorts of actors, it has been expressed in ways clearer than ever that Digital Palaeography as a field is much more complex than a simplistic adjunction of Computer Sciences and Palaeography; indeed CS research, engineering and software development, support and service, linguistics, palaeography, art history, and cultural heritage institutions (Galleries, Libraries, Archives, and Museums – GLAM) all form part of the Digital Palaeography research arena. Good communication requires correct identification of the roles and competence of each actor, and a well-balanced project has to associate/include/foresee the participation of the other actors. It is for example important to clarify that palaeographers are not responsible for copyright or image quality provided by GLAM institutions, in the same way as CS researcher are not responsible for designing interfaces. Within each community, a better understanding of methods and interests of the actors of the other communities is needed to find the right partners (e. g.: keyword spotting is not alignment; writer identification is not script classification). On the other hand, the "black box" issue seems to have been addressed by most teams through the introduction or increase of interactivity of the software tools they presented; interactivity was used not only as a means to produce clear and convincing results, but also to overcome the shortcomings of strictly automatic approaches. In this sense, the reintroduction of "the human into the loop" (or "the use of the users") is part of a process allowing a better understanding on both sides. The "human in the loop" can and should be integrated at all stages, and, even if this need is not always perceived, it is crucial that substantial efforts be dedicated to making implicit assumptions or knowledge explicit. Special attention should be given to avoid the development of tools relying on tautological approaches where tools or datasets incorporate expectations as an underlying (and often implicit) model. In this regard, one cannot overestimate that an unclear result is as important for historians as a clear-cut clustering. In the middle, the "human" gives feedback on preliminary results, enables the enhancement and improvement of the model, as well as creates ground-truth. The display of intermediary results and the integration of user feedback within the process are a welcome solution offered by the latest developments. Likewise, palaeographers have developed new

strategies, in their ways of formulating tool requirements or expressing requirements for which they can evaluate the results themselves, regardless of the software being an opaque black-box (P. Stokes, D. Stutzmann, M. Lawo with B. Gottfried).

Overall, this seminar seems to have operated a paradigm shift from black-box issues to trust issues, in the sense that when we first identified black-box issues, we focussed on "computational black boxes", when "human black boxes" are in fact just as problematic. Instead of focussing on computational black-boxes as an issue, we were able to formulate that the important endeavour is that of establishing trust in the respective methodological approaches to the research questions of the research domains. This trust in methodologies is usually mediated by human interactions ("humans in the loop" again!), and the ways in which scholars are able to share an intuitive understanding of their respective expertises with non-experts.

It hence follows that a new (technical) challenge arises, consisting in the creation and implementation of an integrated software tool, web service suite, or environment that would allow users to access and work with extant datasets and tools. The impetus to take up this challenge resides as much in the Humanities as it does in the Computer Sciences. By aggregating the multiple, isolated, specific tools developed by CS researchers through a common access point, digital humanists would support the development of better evaluation metrics and promote a wider use of CS technologies among more traditional Humanities scholars, who could thus become more aware of the existing tools, more autonomous (i. e. less dependant on CS researchers) and thereby empowered. As a reciprocal positive effect, CS researchers could more easily validate their results and gain access to a wider range of annotated datasets. This challenge is also naturally related to trending key concepts such as "interoperability" and "open access". It furthermore engages with the question of the nature of success metrics in the Humanities, where a successful tool is not only the one giving the best results, it is also one enjoying wide acceptance and a large number of users. Improving ergonomics is mandatory, to put the user in the middle and to accumulate a consistent critical mass of annotations (both as feedback and ground-truth).

## **2**  **Table of Contents**

## 3    Overview of Talks

### 3.1    Interdisciplinary Approach to the Study of Tibetan Manuscripts and Xylographs: The State of the Art and Future Prospects

*Orna Almogi (Universität Hamburg, DE)*

From the point of view of a student of the history of ideas who is primarily interested in the intellectual culture, intellectual history, philosophy, and religion of any given civilization, past and present, it is assumed that there is no effective way to gain a nuanced and well-founded knowledge of them without a profound knowledge of the pertinent languages and without extensively exploring the diverse indigenous textual sources. Despite significant progress that has been made during the past decades in this regard in the field of Classical Tibetan Studies, a relatively new discipline, scholars have barely managed to scratch the surface of the enormously vast, diverse, and rich textual material that has come down to us in the form of manuscripts and xylographs produced by the Tibetan civilization over the centuries. Recent decades have witnessed a significant increase in the accessibility of old Tibetan (mainly Buddhist) texts produced and transmitted from the seventh century until the present. These new discoveries of old primary textual material have no doubt significant implications in the field, posing new challenges and at the same time offering fascinating new opportunities for Tibetologists. However, this tremendous increase in the accessibility of hitherto inaccessible and unexplored textual material some of it fragmentary and often no longer in its original place of deposit but scattered over various libraries around the world heighten the desire to refine existing research tools and seek new ones that are more efficient and more powerful for investigating this material and the ideas transmitted therein. In my presentation I presented the state of affairs in the field of Tibetan textual studies, briefly discussing the major difficulties Tibetologists face in dealing with the large and diverse textual material, and finally described three computerized tools aiming at facilitating Tibetan textual studies that are currently in development.

### 3.2    Encoding Scribe Variability

*Vincent Christlein (Universität Erlangen-Nürnberg, DE)*

Like faces or speech, handwritten text can serve as a biometric identifier. This talk gives an overview of recent methods in scribe identification and verification. Scribe identification methods can be divided into two categories: allograph based methods and textual based ones. Although textual based methods are easier to interpret, the best results so far were achieved by allograph based approaches. One such approach is based on GMM supervectors. This method is compared against other allograph based methods on contemporary datasets such as the ICDAR 2013 competition set and the CVL dataset showing TOP-1 accuracy

of more than 97%. Finally, the method has been applied on a set of datum lines of high medieval papal charters. Background artifacts reduce the accuracy of the classification, thus a word based approach built on GMM supervectors, which reduces the error by a large margin, was developed. This also reveals the limit of current datasets which consist of too few scribes and are too clean in contrast to historical documents. However, in general writer identification / verification methods perform very well, especially when they are applied on contemporary documents, and can thus reduce the effort of large-scale identification / verification drastically.

## 3.3 Algorithmic Paleography

*Nachum Dershowitz (Tel Aviv University, IL)*

Modern algorithms can help in many tasks of interest to scholars of the humanities and, in particular, in the analysis of old manuscripts and texts. We describe ongoing research in the application of methods developed in the fields of computer vision, bioinformatics, and machine learning to endeavors such as the paleographic analysis of manuscripts, finding documents in the same hand, searching within images, and tracing fibers in papyri. Our examples include the Dead Sea Scrolls, the Cairo Genizah, and the Tibetan Buddhist corpus.

## 3.4 Appearance Modeling for Handwriting Recognition

*Gernot Fink (TU Dortmund, DE)*

In this presentation I give an overview of appearance modeling techniques for offline handwriting recognition, i.e., the recognition of handwriting from document images.

I first present the traditional techniques that have been proposed for the recognition of isolated characters and follow a classical pattern recognition pipeline (cf. e.g. [1, Chap. 10–11]).

Then I focus on the recognition of cursive script, where segmentation-based approaches fail due to the very nature of cursive writing and the high variability of the data. Therefore, so-called segmentation-free methods have been proposed, the most well-known being based on hidden Markov models (HMMs) (cf. [2, 5]). I present the general architecture of an HMM-based handwriting recognition system and introduce the sliding-window approach that is essential for converting images of handwritten script into sequences of feature vectors that can be modeled by HMMs. Afterwards, I describe how structured recognition models can be built based on elementary modeling units. For these mostly the characters of the respective script are used but there also exist approaches where context-dependent characters or sub-character units are applied.

In addition to modeling approaches for handwriting recognition I briefly present how script appearance is represented in today's handwriting retrieval systems that are based

on query-by-example word spotting techniques (cf. [3]). In this field image descriptors based on gradient statistics are used for building holistic models of individual query words following the Bag-of-Features (BoF) principle (cf. [4]). In order to improve the performance beyond basic BoF-based word-spotting systems, the BoF principle can be combined with the sequential statistical modeling provided by HMMs. These BoF HMMs today deliver excellent handwriting retrieval performance [7, 6].

From these considerations it can be concluded that impressive results can be achieved for problems with large, annotated training data sets. Language constraints can be described well statistically, but the training of such models for non-contemporary data remains an open problem. A further challenge is that special attention to character appearance is almost exclusively achieved via preprocessing and feature extraction and there exist no principled approaches for sharing of structural cues between character models. It is especially unclear how to transfer such "appearance knowledge" to different writing styles, from printed to handwritten material, or to an entirely new type of script.

Therefore, from a Pattern Recognition viewpoint it appears to be especially interesting to automatically extract script-specific information from example data, to exploit semi-supervised learning strategies, i. e., to learn appearance models from a few labeled and a huge number of unlabeled samples, and to systematically transfer or adapt appearance models to new tasks. With respect to applications in paleographic research it will be important to involve paleographic experts as humans-in-the-loop such that automatic pattern recognition methods rather provide assistance than try to compute necessarily imperfect finalized solutions.

### References

**1**    David Doermann and Karl Tombre, editors. *Handbook of Document Image Processing and Recognition.* Springer, London, 2014.

**2**    Gernot A. Fink. *Markov Models for Pattern Recognition, From Theory to Applications.* Advances in Computer Vision and Pattern Recognition. Springer, London, 2 edition, 2014.

**3**    Josep Lladós, Marçal Rusiñol, Alicia Fornés, David Fernández, and Anjan Dutta. On the influence of word representations for handwritten word spotting in historical documents. *Int. J. Pattern Recognition and Artificial Intelligence*, 26(5), 2012.

**4**    Stephen O'Hara and Bruce A. Draper. Introduction to the bag of features paradigm for image classification and retrieval. *Computing Research Repository*, arXiv:1101.3354v1, 2011.

**5**    Thomas Plötz and Gernot A. Fink. *Markov Models for Handwriting Recognition.* SpringerBriefs in Computer Science. Springer, 2011.

**6**    Leonard Rothacker, Marcal Rusinol, and Gernot A. Fink. Bag-of-features HMMs for segmentation-free word spotting in handwritten documents. In *Proc. Int. Conf. on Document Analysis and Recognition*, Washington DC, USA, 2013.

**7**    Leonard Rothacker, Szilard Vajda, and Gernot A. Fink. Bag-of-features representations for offline handwriting recognition applied to Arabic script. In *Proc. Int. Conf. on Frontiers in Handwriting Recognition*, Bari, Italy, 2012.

## 3.5    Separating glyphs of handwritings with Diptychon

*Björn Gottfried (Universität Bremen, DE)*

My presentation is about a transdisciplinary project in the context of digital palaeography in which methods are developed in order to support palaeographers in comparing handwritings. It is supported by the German Research Foundation, DFG, under grant number GO 2023/4-1 (LA 3066), LA 3007/1-1.

As one important objective the separation of handwritings into their constituent glyphs is discussed and motivated as follows:

- Separated glyphs allow the search for strings in the original document, showing the context of specific glyphs,
- facilitate the character-wise comparison of handwritings, and
- enable the characterisation of the specificities of single glyph images.

Though being generally very difficult and sometimes even impossible, the extraction of single glyphs is challenging but not impossible. An interactive human-machine methodology enables the extraction of single glyphs by combining both the precision and efficiency of the computer as well as the expertise and flexibility of the user. An example of an automatic method is provided in [1].

The methodology has been applied to different handwritings between the 9th and 18th centuries and depends on the specific characteristics of each handwriting. The interaction effort to correct imperfect suggestions provided by the computer lies in the average around 2 seconds per glyph and ranges between 0.6 and 1.4 operations per glyph.

### References
**1**    Jan-Hendrik Worch, Mathias Lawo, Björn Gottfried. *Glyph spotting for mediaeval handwritings by template matching.* ACM, Springfield, Paris, France, September 4–7, 2012.

## 3.6    Deciphering and Mapping the Socio-Cultural Landscape of 12th Century Jerusalem: Texts, Artifacts and Digital Tools

*Anna Gutgarts-Weinberger (The Hebrew University of Jerusalem, IL) and Iris Shagrir (The Open University of Israel – Raanana, IL)*

Research on the urban layout of medieval Jerusalem has traditionally been based on the integration of written descriptions and archaeological investigation. Especially privileged in this context were the monumental buildings whose architecture was both described in detail and has been in many cases, still visible on the ground. In the Crusader period in Jerusalem realities changed. First, there was an upsurge in documentation regarding the life in the city. This was produced by various institutions and agents operating in the newly established Christian capital in the Levant. Secondly, we have information not only about the big monuments but also on private buildings, urban zoning, the commercial areas and religious endowments. From this wealth of documents we aim to produce a detailed database which will allow for qualitative and quantitative analysis of the urban configuration and

topographical layout, in a manner that has never been performed before. We aim to use this database to study the social, cultural and perhaps economic development and hopefully clarify further the distribution of various sectors of the population within the city.

**Method outline:**   The project presented here aims at reconstructing and analyzing chronologically and spatially the development of medieval Jerusalem, 11th- 13th centuries, based on an analysis of the entire corpus of legal, historical, descriptive and religious documents pertaining to sites and events in Jerusalem of the crusader period. The plan is to juxtapose textual and archaeological data, derived from excavations conducted over recent decades. To date, no integrative study of medieval Jerusalem exists. The combination of documentary and archaeological data is expected to enable a comprehensive spatial-temporal reconstruction and analysis of the plan, topography, property-ownership and urban development of Jerusalem over this period. The study aims at assimilating up-to-date insights from the Digital Humanities, in order to create an integrative record of spatially positioned archaeological and topographical data captured and represented on a Geographical Information System (GIS), with carefully categorized text-based historical analysis. The project promises to yield results that will greatly augment our understanding of the history of the Holy City, and generate new questions and further research. Considering the different nature and number of the available sources, the main challenge in the construction of the database lies in the conversion and standardization of historical and archaeological sources into data that can be collated and analyzed from a chronological as well as spatial perspective. This can be demonstrated on the documents pertaining to the city during the period in question. These documents record transactions involving exchanges of properties in and around the city of Jerusalem, conducted among various agents. In order to isolate and trace multiple strands of information, the documents were collected and organized according to their chronological order and the geographic information they hold. They were then broken down into multiple subcategories according to several main thematic clusters, among which are agency, institutional association, property details and connections to other documents. This deconstruction of the documents into their primary elements is designed to accommodate for multifaceted cross-sectioning of the data, allowing an examination and analysis of correlations between multiple clusters of information, thus incorporating both chronological and spatial evolution. This type of analysis yields a detailed and dynamic representation of the underlying mechanisms responsible for the changes that occurred in the cityscape throughout the 12th century. It also reflects the balance and relationship between socio-economic functions and the urban setting they inhabited, helping deciphering and better understanding Frankish Jerusalem's urban fabric.

**Sample issues/challenges for DH:**   Developing software tools that support the process of interpretation and digital tools to complement the human expertise in actions such as:

- Cross-referencing narrative and archeological data.
- Representation of static vs. dynamic data.
- Representation of discrete objects vs. abstractions.
- Codifying and calibrating non-specific property descriptions.
- Automatic identification of different name variants.
- Isolation, classification and analysis of transactions, and statistical significance.

## 3.7    Positioning computational tools

*Tal Hassner (The Open University of Israel – Raanana, IL)*

The conclusions of the Schloss Dagstuhl – Leibniz Center for Informatics, Perspective Workshop on "Computation and Palaeography: Potentials and Limits", 2012, expressed in its subsequent manifesto [1], listed a number of crucial points of concern regarding the collaboration between computer scientists and palaeographers. In my talk, I focus on two of these, namely data availability and its significance to the development and training of computerized systems; and the so-called "black-box" issue, relating to the need of palaeography scholars to have more understanding and interaction with their computerized tools. Taking as an example the specific task of transcript alignment, I attempt to draw a taxonomy of available computerized tools, based on the data required to train them versus the amount of interaction they require of the scholar. The key question raised is where in this taxonomy would an ideal computerized palaeographic tool be positioned, in order for it to be both realistic in its prerequisite data and effective in its capabilities?

As a potential answer, I provide the recently developed OCR-Free transcript alignment system [2]. This system directly matches the pixels in an image of a historical text with those of a synthetic image created from the transcript for the purpose. This, rather than attempting to recognize individual letters in the manuscript image using optical character recognition (OCR). It therefore does not require manual labeling or pre-segmentation of letters nor massive training data required to learn particular alphabets and characteristics of scribal hands. I visualize the output of this system and discuss the ways in which it may be manipulated by the scholar in order to quickly and effectively correct for alignment errors. I conclude with suggesting future work, discussing how such corrections can potentially be used to learn, on the fly, the particular characteristics of the manuscript at hand, and improve alignment from one line of text to the next.

### References
1   Hassner, T., Rehbein, M., Stokes, P.A., Wolf, L.: Computation and Palaeography: Potentials and Limits (Dagstuhl Perspectives Workshop 12382). Dagstuhl Manifestos **2** (2013)
2   Hassner, T., Wolf, L., Dershowitz, N.: OCR-free transcript alignment. In: Document Analysis and Recognition (ICDAR), 2013 12th International Conference on, IEEE (2013), pp. 1310–1314

### 3.8 DIVADIA & HisDoc 2.0 Approaches at the University of Fribourg to Digital Paleography

*Marcus Liwicki (DFKI – Kaiserslautern, DE)*

In this article we present DIVADIA, a toolkit for labeling medieval documents and the HisDoc and HisDoc 2.0 projects on Document Image Analysis (DIA) funded by the Swiss National Science Foundation (SNSF). At the University of Fribourg, we conceptualize a workspace comprising methods for input and presentation of humanists' research on historical documents. The underlying architecture of the workspace consists of three modules concerned with Item Description, Content Representation, and Research Data. Each of the modules provides computational methods for semi-automatic processing of document images, transcriptions, annotations, and research data. DIVADIA is ongoing research at DIVA research group at the University of Fribourg. In its current state it provides Document Image Analysis (DIA) methods for layout analysis, script analysis, and text recognition of historical documents. The methods build on the concept of incremental learning and provide users with semi-automatic labeling of document parts, such as text, images, and initials. The future goal is to provide means for labelling, annotating, searching, browsing, viewing, and comparing documents as well as presenting research data in adequate visualizations. In the HisDoc projects we perform research on textual heritage preservation. HisDoc aimed at layout and textual content analysis of historical documents, i. e., focusing on philological studies. HisDoc 2.0 will take the approach a step further: it will be dedicated to paleographical studies and incorporate semantic domain knowledge automatically extracted from existing document databases into DIA methods in order to facilitate large-scale processing. As such we will investigate the yet missing ingredients for automatic large-scale analysis of historical documents, and how to make the results useful for historians. While concentrating on medieval manuscripts, we intend to develop methods easily adaptable to other kinds of documents and scripts. During the discussion we presented the current stage of the DIVA-HisDB which will contain larger amounts of annotated historical images with difficult layouts. Every year during the ICDAR and ICFHR conferences we will publish new data along with a benchmark competition. An interesting discussion point raised during the seminar was the presentation of document processing results; as developer of document enhancement methods we should make it clear that the output of the enhancement method (e. g., binarization) is a processed image and not a direct photograph of the original document. The main reason for that is that each data processing step introduces derivations from the original image and might also introduce errors. In the worst case a paleographer investigating only a processed image without being aware of the processing steps might draw conclusions which would not have been drawn when investigating the original physical document.

## 3.9 Word spotting in historical manuscripts. The "Five Centuries of Marriages" project

*Josep Llados (Autonomus University of Barcelona, ES)*

Search centered at people is very important in historical research, including historical demography, people trajectories reconstruction and genealogical research. Queries about a person and his/her connections to other people allow to get a picture of a historical context: a person's life, an event, a location at some period of time. For this purpose, scholars use documents like birth, marriage, or census records.

From a technical point of view, word spotting plays a central role in searching among historical people records. Word spotting is the process of retrieving all instances of a queried keyword from a digital library of document images. We have proposed different word spotting approaches for historical manuscript retrieval. In particular, we have evaluated the performance within the EU-ERC project Five Centuries of Marriages (5CofM), which consists in the analysis of marriage license records from the Barcelona Cathedral.

We have made some contributions in context-aware word spotting. Usually word spotting is built based solely on the statistics of local terms. The use of correlative semantic labels between codewords adds more discriminability in the process. Three levels of context can be defined in a word spotting scenario. First, the joint occurrence of words in a given image segment. Second, the geometric context involving a language model regarding to the relative 1D or 2D position of objects. Third, the semantic context defined by the topic of the document. A number of document collections convey an underlying structure.

We take advantage of the structure to boost the search of words, with a joint search of the query word and its context.

## 3.10 Modern Technologies for Manuscript Research

*Robert Sablatnig (TU Wien, AT)*

Manuscript analysis and reconstruction has long been solely the domain of philologists who had to cope with complex tasks without the aid of specialized tools. Technical scientists were only engaged in recording and conservation of valuable objects. In recent years, however, interdisciplinary work has constantly gained importance, concentrating not on a few special tasks only, like the development of OCR software, but comprising an increasing amount of relevant interdisciplinary fields like material analysis and document reconstruction. It may be expected that in the long run the decipherment, study and edition of such sources will predominantly be done based on digital images. This relieves the originals, makes their investigation independent of the place of preservation and permits a lossless storage of the

contents. Additionally, a more precise and less time-consuming investigation of manuscripts through automatic image analysis is made possible. Especially for information invisible to the human eye, spectral imaging methods are applied in order to visualize lost content. Digital cameras sensitive to an extended spectral band are used to produce multi-spectral images which (in combination with digital image processing) allow enhancing the readability of "hidden" texts and an automated investigation of structure and content of the manuscripts.

In order to acquire manuscripts in libraries, a system is needed that is easily portable, robust, and permits quick handling and fast imaging. Thus, we combined a Nikon D2Xs RGB camera to obtain conventional color images and a Hamamatsu C9300-124 high resolution camera with a spectral response from Ultra-Violet (UV) to Near-Infra-Red (NIR, 330 to 1000 nm) and a resolution of 4000x2672 pixels. The lighting system consists of two LED panels with 13 narrow spectral bands. Additionally, four white light LED panels are used for the RGB photographs, since LED lighting does not impose additional heat radiation on the manuscript.

A multi-spectral representation of the page (one object in multiple spectral ranges) acquired in this manner is the basis for our subsequent analyses like image enhancement, since this data representation holds a great potential for increasing the readability of historic texts, especially if the manuscripts are (partially) damaged and consequently hard to read. The readability enhancement is based on a combination of spatial and spectral information of the multivariate image data, a so called Multivariate Spatial Correlation (MSC). The benefit of this method is the possibility to specifically consider individual text regions in document images. Additionally, Independent Component Analysis (ICA), Principal Component Analysis (PCA), and Fisher Linear Discriminate Analysis (LDA) have been successfully applied in order to reduce the dimension of the multispectral scan and for the separation and enhancement of diverse writings. Since LDA is a supervised dimension reduction tool, it is necessary to label a subset of multispectral data. For this purpose, a semi-automated label generation step was developed, which is based on an automated detection of text lines. Thus, the approach is not only based on spectral information – like PCA and ICA – but also on spatial information. A qualitative analysis shows, that the LDA based dimension reduction gains better performance, compared to unsupervised techniques.

Another interesting aspect when working with manuscripts is the automatic identification of authors based on their scribes. We investigated scribe identification on the example of historical Slavonic manuscripts. The quality of these documents is partially degraded by faded-out ink or varying background. The writer identification method used is based on textual features, which are described with Scale Invariant Feature Transform (SIFT) features. A visual vocabulary is used for the description of handwriting characteristics, whereby the features are clustered using a Gaussian Mixture Model and employing the Fisher kernel. The writer identification approach is originally designed for grayscale images of modern handwritings. But contrary to modern documents, the historical manuscripts are partially corrupted by background clutter and water stains. As a result, SIFT features are also found on the background. Since the method shows also good results on binarized images of modern handwritings, the approach was additionally applied on binarized images of the ancient writings. Experiments show that this preprocessing step leads to a significant performance increase: The identification rate on binarized images is 98.9%, compared to an identification rate of 87.6% gained on grayscale images.

### References

**1**    Fabian Hollaus and Melanie Gau and Robert Sablatnig„ *Enhancement of Multispectral Images of Degraded Documents by Employing Spatial Information.* Proc. of 12th International

Conference on Document Analysis and Recognition (ICDAR 2013). 2013, pp. 145–149

**2**    Fabian Hollaus and Melanie Gau and Robert Sablatnig, *Acquisition and Enhancement of Multispectral Images of Ancient Manuscripts*, Proc. of 11th Culture and Computer Science Conference, 2013, ed. Sieck, J., Franken-Wendelstorf, R.

## 3.11  tranScriptorium

*Joan Andreu Sanchez Peiro (Polytechnic University of Valencia, ES)*

TranScriptorium (http://www.transcriptorium.eu) [1] aims to develop innovative, efficient and cost-effective solutions for the indexing, search and full transcription of historical handwritten document images, using modern, holistic Handwritten Text Recognition (HTR) technology.

tranScriptorium will turn HTR technology into a mature technology by addressing the following objectives:

1. Enhancing HTR technology for efficient transcription.
   Departing from state-of-the-art HTR approaches, tranScriptorium will capitalize on interactive-predictive techniques for effective and user-friendly computer-assisted transcription.

2. Bringing the HTR technology to users.
   Expected users of the HTR technology belong mainly to two groups: a) individual researchers with experience in handwritten documents transcription interested in transcribing specific documents. b) volunteers which collaborate in large transcription projects.

3. Integrating the HTR results in public web portals.
   The HTR technology will become a support in the digitization of the handwritten materials. The outcomes of the tranScriptorium tools will be attached to the published handwritten document images. This includes not only full, correct transcriptions, but also partially correct transcription and other kinds of automatically produced metadata, useful for indexing and searching.

### References
**1**    J. A. Sanchez and G. Mühlberger and B. Gatos and P. Schofield and K. Depuydt, R. M. Davis and E. Vidal and J. de Does. *tranScriptorium: a European Project on Handwritten Text Recognition.* ACM Symp. on Document Engineering DOCENG, 2013, pp. 227–228.

## 3.12  Text Classification and Medieval Literary Genres

*Wendy Scase (University of Birmingham, GB)*

This presentation reported on investigation of problems in text classification that are experienced in the creation and querying of large corpora of texts and images. Humanists know that genre information is relevant to the palaeographical analysis of documents. For

example, the genre of a text can influence the scribe's choice of script. A legal document may be written in a cursive script, reflecting the need for speed and economy in the production of the document, whereas a bible will often be written in a formal script that requires the scribe to create letters from many small, careful strokes. This choice may reflect the aspirations of the patron (to save his soul, or display his wealth), the need for the scribe to conceal his identity (where the manuscript may be considered heretical) and so on. So classification by genre is relevant to the interpretation of material in corpora, especially where the corpora are produced from many different parent resources (e.g. archives of documents, collections of literary texts, chronicles). Classification of genres from the medieval point of view is however still little understood. A further problem occurs when resources from different modern genres are federated in a resource (e.g. dictionaries, catalogues, full-text transcriptions). The user needs to know the genre of the text retrieved to interpret it accurately. Manuscripts Online (www.manuscriptsonline.org) is an experiment with federating resources relating to medieval; British texts was used to illustrate these problems and some partial solutions. More work needs to be done. The final part of the presentation reported on work towards the expansion and further enhancement of a corpus reported on at Dagstuhl Perspectives Workshop 12382. The Vernon manuscript scribe's text and image corpus (Bodleian Library, MS Eng. Poet.a.1) has been increased with the digitisation of the Simeon manuscript (British Library, Addit. MS 22283), also partly copied by the Vernon scribe. Many research questions could be explored if the images could be provided with aligned transcription. The presentation proposed that the existing files of the Vernon manuscript project could be harnessed to create a training set that would permit semi-automated labelling of the images of the Simeon manuscript.

## 3.13   Describing Handwriting – Again

*Peter A. Stokes (King's College – London, GB)*

When considering the identification of characters and scripts, two important aspects that were identified in the 2012 Dagstuhl Perspectives Workshop on computing and palaeography are ontologies and mid-level features [1]. This paper focussed on those two aspects, partly in deliberate contrast to the highly computational approach that most studies in the field have taken to date.

To this end, problems not only of terminology but also of conceptual ambiguity and imprecision in palaeography were introduced. The ontology developed for the DigiPal project was briefly presented as a response to this, including the way that it has been used in practice for describing writing in the Latin and Hebrew alphabets as well as for decoration [2, 3]; initial work has also been done to use it for Greek and Latin inscriptions and cursive Latin script. The ontology was presented not as an ideal solution but rather as a pragmatic one that has proven useful in a variety of circumstances, and as a starting-point to a very difficult problem with many challenges that still remain.

The second part of the talk considered possible mid-level features, presenting a selection of potential characteristics of handwriting that are relevant to palaeographers and that seem to this author to be relatively easily amenable to computational analysis but which seem not

to have been considered in practice. These included 'stabbing' strokes (perhaps indicating a scribe accustomed to writing on wax), 'equilibrium' (the regularity or otherwise of strokes, perhaps a sign of fluency, experience, forgery or imitation), and the effective visualisation of these particularly in the context of other factors such as the codicological structure of the book. As an aside, DigiPal's RESTful API was also introduced as a potential source of annotated images for the training of computer vision systems.

None of these methods or approaches is necessarily appropriate for writer identification, but they suggest other directions in which computer vision might be taken and which perhaps are more pertinent to research in medieval manuscripts than some of the work done to date.

**References**
**1**    T. Hassner, M. Rehbein, P. A. Stokes, L. Wolf (eds). Computation and Palaeography: Potentials and Limits. *Dagstuhl Manifestos* 2(1):14–35, 2013. DOI: 10.4230/DagMan.2.1.14
**2**    *DigiPal: Digital Resource and Database of Palaeography, Manuscript Studies and Diplomatic.* London: 2011–14. http://www.digipal.eu/
**3**    P. A. Stokes, S. Brookes, G. Noël, G. Buomprisco, D. Matos and M. Watson. The DigiPal Framework for Script and Image. *Digital Humanities 2014 Book of Abstracts* (Lausanne, 2014), pp. 541–3. http://dharchive.org/paper/DH2014/Poster-193.xml

## 3.14    Bridging the gap between Digital Palaeography and Computational Humanities

*Dominique Stutzmann (Institut de Recherche et d'Histoire des Textes (CNRS) – Paris, FR)*

As part of the common introduction to seminars 14301 (Computational Humanities – bridging the gap between Computer Science and Digital Humanities) and 14302 (Digital Palaeography New Machines and Old Texts), the first paper presented the specific field of the Digital Humanities devoted to the history of scripts, aka "digital palaeography" and why it is of interest even for textual scholars. Texts are transmitted through signs; signs are transmitted through shapes; the shapes for each sign evolve and are perceived for their meaning and in their historical context. Moreover, scripts convey a particular meaning for themselves, as do the *litterae elongatae* and diplomatic script in a diploma of Charlemagne, referring to imperial *litterae caelestes* and supporting the claim of a new Empire, while the same emperor, on the other hand, could support the Caroline script, named after him. Issues for the palaeographer encompass the history of script, cultural history, writer identification, dating and assigning a place of origin for any written sample. As demonstrated by the examples from the transmission of Cicero's works, textual scholarship need to envision the materiality of the transmitted text (not least for classical texts for which there are only medieval witnesses) and digital palaeography addresses the notions of text through image, layout and shape, through their materiality, their history, origin and provenance of the witnesses, through their cultural significance. Digital Palaeography means: how to use computers to help the humanities identifying the relevant historical phenomena, to identify interscript, interscribal, intra-script and intra-scribal variations as well as cultural and textual relevant features. Some bridges with Computational Humanities are obvious: Keyword Spotting and retrieval is similar to indexing techniques; Handwritten Text Recognition is

linked to scholarly editing textual transmission, ideas and their reception. In the issues raised by 14301 are mentioned the kind of results and the transfer to other fields (methodology and applicability), the difficulties in cross-disciplinary collaboration, the human resources and communication, the variability and quality of data, the evaluation and ground-truth. Demonstration and proof in the Humanities and Computer Science, or the measure of success supposes a unique ground-truth, which does not always exist, while the result of a calculation generally represents only an additional clue in the complex reality. All these issues, as well as the crucial notion of reciprocal uncertainties have been addressed in the perspective workshop 12832. Indeed, the four core issues identified issues in 2012 were "Communication and roles in the interdisciplinary interplay, the notions of black box and meaning of calculation", the evaluation of and need for "quality and quantity" in the data from the humanities, and the new audiences (with correlations in interoperability, rights managements and engaging with other communities). These issues are now to be addressed by Digital Palaeographers on a technical and epistemological level, but are also common to all fields in the Digital Humanities and should appeal for a more intense dialogue.

## 3.15 Digital Palaeography. Text-Image Alignment and Script/Scribal Variability (ANR ORIFLAMMS / Cap Digital)

*Dominique Stutzmann (Institut de Recherche et d'Histoire des Textes (CNRS) – Paris, FR)*

**License** ![cc] Creative Commons BY 3.0 Unported license
          © Dominique Stutzmann
**Joint work of** Stutzmann, Dominique; Lavrentiev, Alexei; Kermorvant, Christopher; Bluche, Théodore; Leydier, Yann; Ceccherini, Irene; Eglin, Véronique; Vincent, Nicole; Debiais, Vincent; Treffort, Cécile; Ingrand-Varenne, Estelle; Smith, Marc
**URL** http://oriflamms.hypotheses.org/
**URL** http://www.agence-nationale-recherche.fr/projet-anr/?tx_lwmsuivibilan_pi2[CODE]=ANR-12-CORP-0010

Medieval scripts are a challenge to historical analysis, as for describing and representing the graphical evidence, analyzing and clustering letter forms and their features through Computer Vision and analyzing historical phenomena. The ANR funded research project ORIFLAMMS (Ontology Research, Image Feature, Letterform Analysis on Multilingual Medieval Scripts, 2013-2016) gathers seven partners from the Humanities and Computer Science (IRHT = Institut de Recherche et d'Histoire des Textes, CNRS; CESCM = Centre d'Études Supérieures de Civilisation Médiévale; École Nationale des Chartes; ICAR = Interactions Corpus Apprentissages Représentations, École Normale Supérieure de Lyon, for the Humanities; A2iA; LIRIS = Laboratoire d'InfoRmatique en Image et Systèmes d'information, INSA Lyon; LIPADE = Laboratoire d'Informatique de Paris Descartes, for Computer Science). It aims at studying the coherence and variability of graphical systems, according to their language, level of formality, support, genre, date and place, as well as creating an ontology of medieval signs, through the alignment of text and images, by extracting letterforms, abbreviations and signs, then perform pattern similarity analysis, and enhance the results with computational linguistics and paleographical analysis. In order to achieve representative results, several core corpora have been identified (charters, books, books of charters such as cartularies and registers, inscriptions). The research is based on XML-TEI compliant editions and compels to deepening our understanding of scribal systems and forms [1, 2]. As part of this research, a software has been developed in order to easily visualize and validate the text-image alignment. The latter is produced by two different systems developed in this project: the first one without prior knowledge [3], the second one

with GMM and DNN with very good results. By now, two large data sets have been aligned: *Queste du Graal* including 130 pages, 10700 lines, more than 115'000 words and 400'300 characters; *Fontenay* including 104 pages, 1341 lines, more than 22'200 words and 99'900 characters. This is a major first step. With the following corpuses, this research contributes to both Humanities (letterform identification, historical semiotics) and Computer Science (Handwriting recognition), with the core idea of not reinventing the wheel, but using former research, computer and human brain at their maximal capacities.

**References**
**1**     D. Stutzmann. Paléographie statistique pour décrire, identifier, dater ... Normaliser pour coopèrer et aller plus loin?  In *Kodikologie und Paläographie im digitalen Zeitalter 2 – Codicology and Palaeography in the Digital Age 2*, Norderstedt, 2010, pp. 247–277
**2**     D. Stutzmann. Ontologie des formes et encodage des textes manuscrits médiévaux. Le projet ORIFLAMMS, In *Document numérique*, 16/3 (2013):81–95. DOI: 10.3166/DN.16.3.69-79.
**3**     Y. Leydier, V. Eglin, S. Bres, D. Stutzmann. Learning-free text-image alignment for medieval manuscripts. In *Proc. Int. Conf. on Frontiers in Handwriting Recognition*, Crete, Greece, 2014.

## 3.16 Digital Images of Ancient Textual Artefacts: Connecting Computational Processing and Cognitive Processes

*Ségolène Tarte (University of Oxford, GB)*

Drawing on examples from palaeographical scholarship rooted in Classics and in Assyriology, this talk will give an overview of how it might be possible to connect computational processing and cognitive processes. As a preamble, considering the type of material that palaeographers (be they Classicists, Mediaevalists, or Assyriologists) work from, I will argue that an image of an ancient textual artefact is a digital avatar of the textual artefact. In digital palaeography, these images are an absolute prerequisite, but it is crucial to be aware that as avatars they are already part of the interpretative workflow that transforms the data (the textual artefact) into knowledge and meaning. Digital avatars are interpretative; they express a certain form of presence of the textual artefact, they are contingent on the act of digitization and they have an expected performative value [1]. All those implicit aspects that participate in the act of knowledge creation coexist with the intuitive strategies that scholars develop to carry out their task. I will present three such strategies identified through ethnographic studies of Classicists and Assyriologists at work [2]. Establishing a correspondence between these ethnographic observations and cognitive processes (as identified in the cognitive sciences literature), I will show examples of how these cognitive processes influenced and supported the choice of computational processing made by the scholars. Namely: embodied cognition and an awareness of the materiality of a papyrus suggested modelling it as a roll to justify the repositioning of a fragment; kinaesthetic facilitation was supported through digital tracing of the text of another artefact, thereby supporting the establishment of the connection

between the text as a shape and the text as a meaning; depth perception through monocular parallax motion was supported for yet another artefact by the digitization process, allow to interactively relight the artefact. These examples are vivid illustrations of the fact that understanding scholars' cognitive involvement have the exciting potential to facilitate the seamless integration of the use of computational tools within the research workflow whilst at the same time supporting embodied sense-making practices.

### References

**1** Tarte, Ségolène M. The Digital Existence of Words and Pictures: The Case of the Artemidorus Papyrus In *Historia 3:61*, pp. 325–336 (+bibliog. pp. 357-61; fig. pp 363-5), 2012.

**2** Tarte, Ségolène, Interpreting Textual Artefacts: Cognitive Insights into Expert Practices In: *Proc. of the Digital Humanities Congress 2012*, Ed. Clare Mills, Michael Pidd, and Esther Ward, Sheffield: HRI Online Publications, Studies in the Digital Humanities, 2014.

## 3.17 Text classification

*Nicole Vincent (Paris Descartes University, FR)*

Classification, and text classification, has to be done with respect to some objectives. These objectives are varying according to the field of interest possibly being medical, security or palaeography. Some questions are rising, such as: Do you have some ground truth available defining the classes and their number? One point is the definition of features. But how to choose them? Choose many to have a large amount of information. Not too many because of dimensionality problem and because the aim is to decrease complexity. What about feature selection? What about learning? What may be the criteria to choose features: have they to be understandable? Should they be local or global, addressing details? Should they be invariant towards different factors? What about the process? Defined by the expert, blind based on computer science theory, based on pixels or features or primitives, involving an interaction with the user? 4 examples of text classification are presented. They have been developed in the GRAPHEM project funded by French National Research Agency:

- One involving a decomposition of writing that models the way the drawing is done
- One based on the statistical analysis of the writing contour
- One trying to be the automated version of an expert palaeographer
- One base on the statistical analysis of some low level patterns

## 3.18   Diplomatics and Digital Palaeography

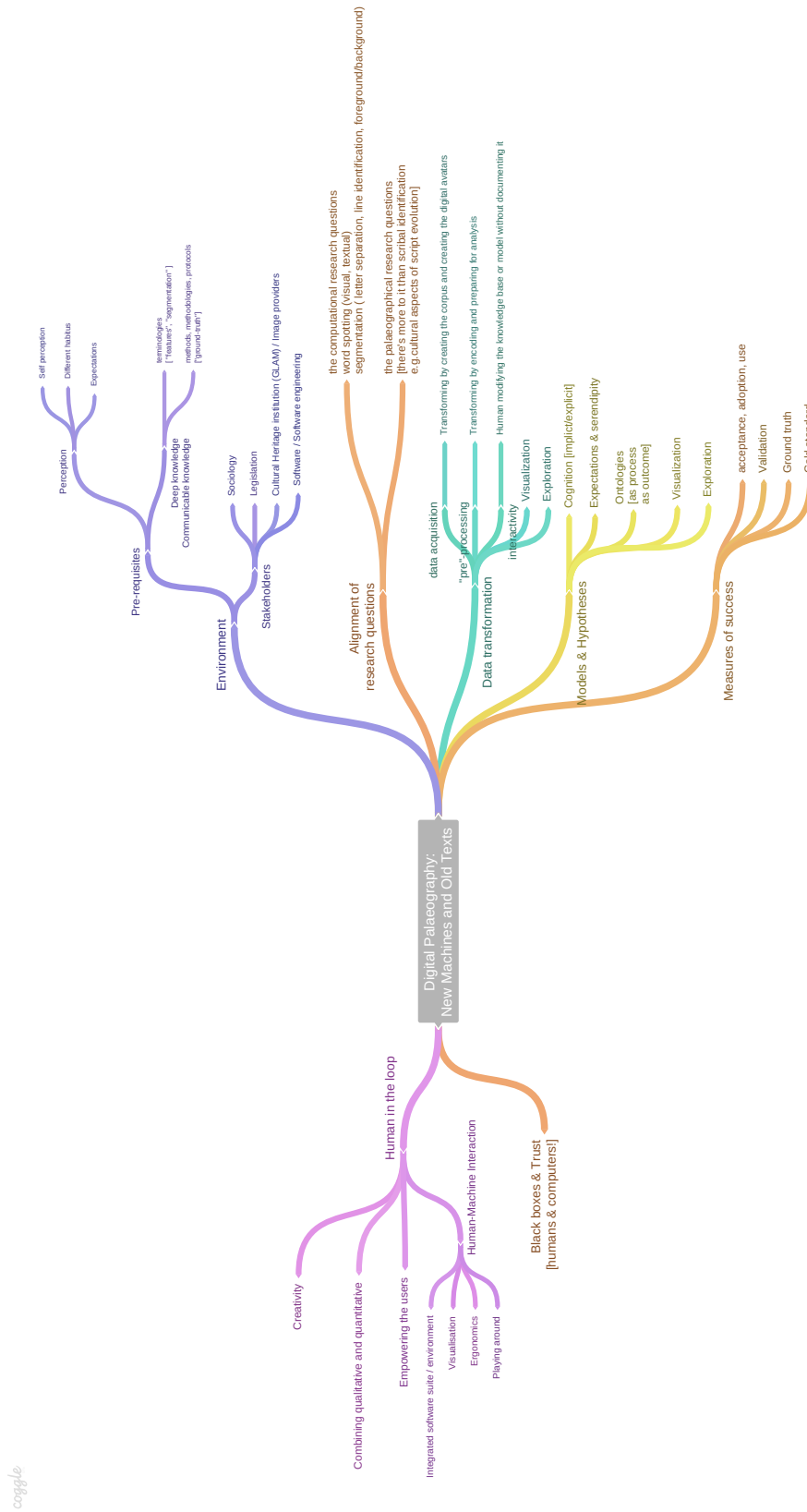*Georg Vogeler (Karl-Franzens-Universität Graz, AT)*

Manuscripts documenting a single legal act usually authenticated by special means are an important source for the history of the middle ages and the early modern period. They are subject of the research field of "diplomatics", which includes skills in philology, sphragistics, chronology and certainly palaeography. The paper gave an overview on the issues of image based digital methods in diplomatics and their applicability to digital palaeography, and addressed the following questions: what *diplomatics* can contribute to *Digital Palaeography*? Should/Can we build an integrated "Virtual Research Environment" for digital palaeography?

- Digital Palaeography applied to diplomatic sources confronts new challenges in comparison with literary manuscripts, since charters are short, very numerous, formulaic. However, there is usually substantial context information and metadata (date, place).
- Diplomatic writings document the history of Latin script in a specific manner (multiple hands / multiple scripts on one document; "documentary writing style", functional scripts, chancery scripts vs. notarial hands, stylistic influences between book and diplomatic scripts)
- Large digital charter collections and diplomatic databases like http://monasterium. net/ offer new possibilities for research in the field of digital palaeography (discovering imitations, forgeries, copies; identifying "writing landscapes")
- The recently started project "Illuminated Charters" (http://illuminierte-urkunden.uni-graz. at) demonstrates how legal instruments may be considered by their value for art history. It allowed to discuss the basic functionalities for an VRE to be used in the project and the role of controlled vocabularies/formal ontologies in these contexts.

The paper demonstrated that legal documents ("charters", "legal instruments") are a rich source for experiments with digital methods on historical sources as they convey a large data set with relatively precise historical metadata (date, place, partially even writer) and suggested to work on the definition of interfaces and standards to reuse software tools in a web based palaeographic tool chain, also in order to build "trust" under a cognitive aspect (how does the tool shape the perception of the task?).

## 4    A Graphical Representation of the Discussed Subjects

The mind map in Fig. 1 presents on overview of the subjects that where broached during the seminar. Each item and sub-item represents an area in which substantial efforts might be concentrated in the future to further research in computational palaeography.

**Figure 1** Overview of the themes and issues discussed during the seminar.

## Participants

Orna Almogi
Universität Hamburg, DE

Vincent Christlein
Univ. Erlangen-Nürnberg, DE

Nachum Dershowitz
Tel Aviv University, IL

Véronique Eglin
INRIA / INSA – Lyon, FR

Jihad El-Sana
Ben Gurion University –
Beer Sheva, IL

Gernot Fink
TU Dortmund, DE

Björn Gottfried
Universität Bremen, DE

Anna Gutgarts-Weinberger
The Hebrew University of
Jerusalem, IL

Tal Hassner
The Open University of Israel –
Raanana, IL

Rolf Ingold
University of Fribourg, CH

Noga Levy
Tel Aviv University, IL

Marcus Liwicki
DFKI – Kaiserslautern, DE

Josep Lladós
Autonomus University of
Barcelona, ES

Frederike Neuber
Karl-Franzens-Univ. Graz, AT

Jean-Marc Ogier
University of La Rochelle, FR

Robert Sablatnig
TU Wien, AT

Joan Andreu Sanchez Peiro
Polytechnic University of
Valencia, ES

Wendy Scase
University of Birmingham, GB

Iris Shagrir
The Open University of Israel –
Raanana, IL

Peter A. Stokes
King's College – London, GB

Dominique Stutzmann
Institut de Recherche et
d'Histoire des Textes (CNRS) –
Paris, FR

Ségolène Tarte
University of Oxford, GB

Nicole Vincent
Paris Descartes University, FR

Georg Vogeler
Karl-Franzens-Univ. Graz, AT