

Workshop on Biographical Linked Data and Prosopography

22 Jan, 2016, University of Oxford

Goals of the Workshop

The general goal of the workshop is to sort out and discuss Digital Humanities research problems underlying the Reassembling the Republic of Letters EU COST action (RRL), related to [Work Group 2: People and Networks](#). The focus in this meeting is especially on modeling biographical data and using it for prosopographical research. Concrete research problems and use cases are solicited for discussion within the RRL community.

Important topics for discussion in the workshop include e.g. the following:

- What prosopographical data models to use in RRL; are the foundations secure?
- How we are going to fill that data model with data? When to do it manually, semi-automatically, and automatically?
- What content will be used for populating the data model?
- What are the biographical and prosopographical reference sources most relevant for researching the republic of letters?
- Which of these reference sources are available for open access on the Web? Can all of these be provided to the users of the cataloging systems?
- How to select and prioritize the most valuable sources?
- To which of these resources might our infrastructure wish to add new entries., e.g. new people records for VIAF or Getty ULAN, or new names forms for CERL?
- From which of these resources might we wish to extract data automatically and how?
- How would we go about negotiating collaborations of this kind with external dataset publishers, regarding using and enriching the data?

We solicit **proposals for pilot applications** where prosopographical research problems could be addressed in novel ways using the linked data approach, semantic computing, knowledge discovery, and visualizations.

To create a basis for such research and discussions, the programme presents different approaches to modeling historical people, networks, and events using ontologies. Also modeling and representing biographical data is in focus. These presentations are given by invited scholars that have been involved in related works.

The talks also present initial results in designing and building a collaborative Linked Data infrastructure for RRL. The idea is that this infrastructure could be piloted within the project by interested participants, and it could make a basis for an anticipated Horizon 2020 project proposal, discussed in more detail in the final session.

The participants are also invited to send a proposal for a short talk to the Programme Committee. The proposal should explain clearly, why this presentation and topic for discussions is valuable to the goals of the workshop. Indicate the desired length of the presentation in your submission.

Presentation proposals should be sent in PDF format to the workshop chair Eero Hyvönen eero.hyvonen@aalto.fi by email.

Program

Time	Presentation
09:00-09:10	Opening of the Workshop Eero Hyvönen: Goals and Contents of the Workshop PDF
09:10-10:40	Aims and Research Questions (chair Eero Hyvönen) Howard Hotson: Conceptual fundamentals of the EMLO prosopographical data model PDF Robin Buning: Prosopographical Research Questions PDF Thomas Wallnig: Historiographical Research Questions and Preparations for the Warsaw conference Discussion
10:40-11:10	Coffee Break
11:10-12:30	The EMLO Data Model: Is it fit for use? (chair Thomas Wallnig) Jetze Touber and Mikkel Jensen: Using the EMLO prosopographical data model, means for data input PDF Martin Hadley: Visualizing Prosopographical Data based on Correspondences PDF Discussion, next steps towards releasing the data model
12:30-13:30	Lunch
13:30-15:00	Towards Linked Data (chair Arno Bosse) Eetu Mäkelä: EMLO as a Linked Data Service PDF Jouni Tuominen: Dynamic Ontology Service for Historical Persons and Places Based on Crowdsourcing PDF Antske Fokkens: Lessons Learned from BiographyNet PDF Discussion, defining an infrastructure for prosopographical research
15:00-15:30	Coffee Break
15:30-16:00	Future Prospects, Topics, and Partners (chair Howard Hotson) Sandra Toffolo: Prosopography in Network Presentations of WG2 members and their collaborators: Thea Lindquist: Early Stuart Diplomatic Service: Prosopography and Networks PDF Discussion
16:00-17:00	Discussion: Future Steps for WG2 and Horizon 2020 (chair Howard Hotson) Next steps regarding research questions, data model, tools, and partners. Assignment of tasks and deadlines for WG2 and Horizon 2020.
	Dinner

Participation

The workshop is open to all members of the RRL EU Cost project and their collaborators.

Venue

University of Oxford, [St. Anne's College](#):

St Anne's College
Woodstock Road
Oxford, OX2 6HS, UK
More information about [getting there](#).

The meeting will take place in the Seminar Room 9; you can find it using St Anne's [site map](#).

Travel Expenses

RRL EU Cost project will support, using the EU COST guidelines, the travel expenses of presenters and WG2 members. Since there is only limited travel budget available, the programme committee will decide on funding allocation, if there are more participants than can be fully funded.

Accommodation

The participants are expected to find and book accommodation by themselves. In addition to hotels, some colleges offer accommodation. This [website](#) provides a central point for finding rooms.

More Information

For more information, please contact workshop chair Eero Hyvönen, eero.hyvonen@aalto.fi, tel. +358 50 384 1618.

Programme Committee

- Eero Hyvönen (chair)
- Ruth Ahnert
- Sebastian Ahnert
- Arno Bosse
- Howard Hotson
- Tanya Gray Jones
- Eetu Mäkelä
- Thomas Wallnig

Participants

Dominic Oldman (British Museum)
Miranda Lewis (Cultures of Knowledge, Oxford)
Eero Hyvönen (Aalto, Semantic Computing Group)
Howard Hotson (Cultures of Knowledge, Oxford)
Eetu Mäkelä (Aalto, Semantic Computing Group)
Thomas Wallnig (Vienna)
Sandra Toffolo (Tours)
Robin Buning (Cultures of Knowledge, Oxford)
Alexander Butterworth (East Anglia)
Philip Carter (ODNB, Oxford)
Jo Payne (ODNB & ADBN, Oxford)
Mikkel Munthe Jensen (European Institute, Florence)

Jetze Touber (Utrecht)
Tanya Gray Jones (Bodleian Library, Oxford)
Christoph Kudella (SUB, Göttingen)
Martin Hadley (IT Services, Oxford)
Antske Fokkens (Biography.Net, Amsterdam)
Arno Bosse (Cultures of Knowledge, Oxford)
Jouni Tuominen (Aalto, Semantic Computing Group)

Eero Hyvönen, Opening Remarks

Eero opened the meeting

- introduction of the participants
- goals of the workshop
- research questions and topics
- programme

Topic I: Aims and Research Questions

Howard Hotson

- context and background of the RRL project, “letters only serve their function by being scattered”
- needs for technological solutions for prosopographical research
- disambiguation, understanding networks & content
- balancing between simplicity and adaptability
- Basic epistolary data model comparatively simple; prosopography is harder
- need to balance friendliness with detail, durability with flexibility
- some of the data model's problems (of data linking) can be overcome by moving from the spreadsheet form of data ingestion into webform
- H2020 can only fund part of what we want to do

DO: Are we interested in existing data sets?

HH: Yes, ODNB, CERL obvious. Working w/others who are working w/structured biographical data. Their way may be better. And for merging, linking and pooling.

DO: King's College. John Bradley, soon to retire. Have assembled prosopographies. Should contact.

EH: Can create web pages for collecting repositories of biographical information; participants can send own references to be added onto the page

ML: A lot of EMLO's contributors have chronicles, lists.. event streams tucked away, and could offer them.

DO: When you start mixing in data from other sources; how to get sustainability from a mixed model

HH: We're hoping that we can base this on a set of 'core' fields, like the six we've selected for EMLO. Core metadata fields w/different (broader) rights from fields adding on top of this. One of our objectives for this meeting.. is this a feasible model?

CK: for many correspondence corpora, esp. from the 16th and 17th century, only the location of dispatch of the individual letters is known. prosopographical data on the correspondents could help to semi-automatically establish where letters have been sent to, which in turn is important for geospatial visualizations ('scope' of a network etc); more in general: library and information science has decades of experience in modeling biographical data and generating unique identifiers (cf 'authority records') - currently there is a trend toward an internationalization of these records, see for example VIAF; it would be important for the action as a whole to tap into the expertise of library and information science in this regards

Robin Buning

- collective biography
- prosopography can be applied on a variety of historical and sociological research
- concerned with identifying patterns of activity and connectivity in large quantities of data
- not the same as social network analysis
- method to identify shared characteristics and interests, and connections within a group that are hidden or unclear
- method to analyze the structure and changing composition of a group, and the shifting roles of individuals or subgroups
- the information that can be captured must be person-centred and concern external features of people: gender, age, family relationships, ethnic background, religious denomination, social class, personal relationships, education, occupations and political faction
- the captured data can act as a biographical database that can be searched, filtered and visualized
- research case Samuel Hartlib
- data problem: >4250 letters by c.325 correspondents during his life in England, but no letter data from his earlier life
- research questions related to:
 - geographical origin
 - family relationships
 - contact history
 - location history and institutions attended
 - educational and occupational history
- classification of research questions:
 - ethnic, religious and socio-economic background of members of a social group with respect to their characteristics and distinctions
 - origin, composition and structure of a network, and on a larger scale how networks intersect
 - contextualization of specific individuals and subgroups within a larger group

CK: We have to keep in mind that social network analysis was developed in sociology and has later been used by historians (1970s, early in 2000s), historians have always talked about networks, but rarely used social network analysis methods

CK: attribute data: not a link between two nodes + relational information, you can visualize these (person related to place), prosopographical information having these 2 distinct layers, mixed model

CK: term "event streams" has to be defined, c has used term dynamic data; if i query my prosopographical data to see how the composition of the network has evolved, e.g., from religious denomination perspective, e.g., in 16th century, you need a data model that has a temporal dimension included

EH: what is the ultimate goal of prosopography? meta-level, deeper questions

HH: we need tools for various research questions; historians try to understand past in its messy complexity (vs. sociological needs); data model with input tools + analysis and visualization tools, that variety of disciplines can use; understanding process of historical change (vs. in sociology, social behaviour of groups of people)

Thomas Wallnig

- creation of biographical dictionaries = canon formation within the RRL
- when discussing prosopography we should also consider historiography in a broader sense (philosophy, history of science, history of institutions, etc.)
- anthony grafton, correct to remind us of 4 line 'biographies', with a person's virtues, originally from St. Hieronymus (De viris illustribus), with roots back to the old testament

- ODNB entries still carry a faint echo of this, but at the same time represent a national paradigm of the 19th / 20th centuries: “great men of a nation” - actually, when working with those data we learn more about the 19th century than about the people themselves
- we can draw from past canonical forms, but must also draw on present forms
- Leibniz was not born a mathematician, but what events could help mark the transition?
- Ambiguities and anachronistic vocabulary . e.g. “class”, is a 19th century category; “Bohemia”, is it part of “Germany” (yes: Holy Roman Empire; no: not German)? etc.
- 20th century Ethnicity *can* be projected back, but we have to be very self-aware, and bear in mind that there are contemporary identities and categories of description
- Pro Domo database, future project w/Dagmar Mroziak, comparing Jesuits and Benedictines; e.g. in the 19th century, as the state provides better employment perspectives, what effect does this have on the makeup of the religious orders?
- Point/recommendation: use *also* historical terminology to describe past lives, statuses, disciplines, institutions/places

DO: interesting contrast between the presentations, some with very definite ideas about what you are recording, vs. surprise by what one encounters. crucial element in deciding how to model, or not to model at all.

EH: What kind of data models are used in museums?

DO: most museums have biographical databases linked to their objects; how to model or not model

DO: Paper by [John Bradley, KCL], discussing getting flexibility into very fixed models, and failing. His progression over 10 years in this area is interesting to follow. ‘Factoids’ — how do you atomize propositions you can’t predict and model in advance?

DO: one of my recent topics has been how to atomize propositions in datasets, to be searchable; fixed SQL systems are not good for these flexible needs

Topic II: Is the EMLO data model fit for use?

Jetze Touber

- STSTM on experimenting how the EMLO prosopographical data model works, with focus on visualizing networks of correspondence
- observations:
 - What kind of tool is it.. in order to manage expectations, what are its aims?
 - different user needs for capturing data
 - basic metadata birth, death
 - social exchanges
 - Cupers’ profile consisted largely of social events
 - building biographies (crowdsourcing)
 - persons and events
 - prosopography is about persons, but the model is event-centric
- Does the data model allow us to capture the events for related groups (eg. traders, artists)?
 - person-centric events, personal attribute (e.g., becoming a bishop); relational events, e.g., exchange gift, arguments, book reviews
 - the extent of an event
 - it's not always easy to define what constitutes one event (e.g., a scholar gives 2 copies of a book to a salesman, who sells them to 2 other scholars; how many events?)
- What if these are broken down into incompatible ways.

Mikkel Munthe Jensen

- 600 profiles of nordic professors, 18th century
- event centric = great potential
- good to capture detail, but too much density scares away new users
- web form might address this better than spreadsheets (e.g. hidden options etc.)
- data model developed gradually, with event types and rows added as you go along.
- however this has created some overlaps and other incongruities which should be cleaned up now
- it's important to get this right, because our decisions now will follow us into the future.
- changes in the data model may involve asking people to do something new & different on 'already completed' work
- it needs to be streamlined and the logic behind it needs to be straightened and systematized
- Examples of improvements
- Types of event types:
 - Individual=<1, Institutional=<2, Relational=<4
- Institutions:
 - Difference between institution and geographical location
- Geocoding and mapping
- Professions:
 - Put them together in one category, and related to specific row, not reproduced in different places
- relational event types should focus on the relation, not the profession

TW: Procedural discussion vs. content discussion. Which part should we discuss now? We should seek out for space for more detailed data model discussion; at this meeting we can focus only on procedural issues

HH: How can we go out this meeting to work out which part of the model works well, which part doesn't work well. Need to juxtapose existing models in a very systematic way. Identify complementarity and potential for convergence. ROL excludes women because educational categories don't fit women's education. Need to discuss how current model could be extended to capture this. We need to assemble an inventory of models as well, not just content.

CK: First thing to address, is something crucial missing? Next, what is the core model? In this way, people can start working on it right away.

ML: Can the profession, as a 'über' category, be drawn from the Person category in EMLO?

DO: All cultural heritage institutions have different terminologies. We need to separate terminologies from semantic categories. Concentrate on the semantic framework but don't place restrictions on the terminologies institutions people use.

EM: integration of different datasets can be done when data is needed/used, not only at creation time (difficult to come up with all the relevant datasets to link at the creation time); you can also do it at the creation time when using ontological reference sources (different datasets use shared ontologies and get linked through them); it's hard to reconcile authority records on use time, rather should be done in creation time; the data models are easy to link together in linked data; support for evolution of data models; let's have the EMLO data model as it is; the documentation of the data model is quite sporadic currently; after it's been tidied, it could be published

TW: Could Jetze and Mikkel's STSM reports be completed in dialogue with EM?

EM: temporal attributes of people, relational events; the core is quite solid, but it isn't documented properly; existing scattered docs should be merged

EM: roles of events; event types, list of roles in different event types; EMLO vs. RRL?

TW: would like to see continued in WG2 discussions

Martin Hadley

- IT helping researchers building interactive visualizations (bridging the gap between data and research)
- EMLO case study
 - imported 10 EMLO spreadsheet files into R
 - Shiny by R Studio, with the knowledge of R (well represented in the DH) one can build powerful interactive network/graph visualizations, easily embed interactive content on web pages (paid version, free version limited to 10 hours of interactivity) (editor's note/Jouni: apparently the paid version means the hosting platform; there's a free version available, can be installed on own server)
 - view prosopographical info directly from a network visualization
 - highlight/exclude relationships based on EMLO categories

TW: what is the relation of Palladio and the presented Shiny

EM: Shiny is a different kind of visualization tool than Palladio (e.g., Palladio does not give names for links, only for nodes); the core takeaways:

EM: different needs require different visualizations methods

EM: we are in an era where we can quickly prototype things

HH: the prototype has been built quickly with not many users involved; so think what might be accomplished with properly funded and piloted project with scholars involved

EH: it's important that data is published openly (APIs etc.), so that everyone can freely build new visualizations on top of the data

AF: pilot project on visualizing generic uncertainties, contact us for more info

EM: data quality is important, even simple, widely used SKOS thesauri contain errors that can be found via data validator tools; also the user should understand the data in order to make meaningful visualizations (e.g., what is missing from the data)

CK: some fancy technological solutions (e.g., time sliders) are not actually useful for historians, simple text field might work better; not only tools for visualizing network data, we also need to be sure we have connections between data storage and analysis/visualization tools, the data storage should be able to export dynamic data straight into tools

HH: not only humans, but also institutions were shown in the demo visualization; should be able to toggle what data is visualized

EH: could there be separate databases for the epistolary metadata and its analysis (the interpretation), platform for historical discussion

HH: this kind of platform could be useful across disciplines; enriching with data from other disciplines

EH: we could create a vocabulary for scientific argumentation, data conforming to it and then visualize it

DO: We're building this into our software at ResearchSpace. Due in April. Allow people to atomize their assertions and create a 'provenance of argumentation'.

Topic III: Towards Linked Open Data

Eetu Mäkelä

- Publish interpretations as data (itself) which feeds into the virtuous cycle.
- Aggregation and counting at scale requires standards for ingestion & integration. Similarly, from the other end, how you publish your interpretation (i.e. as data) affects how others will be able to integrate it later
- Integration at different scales, from LD upwards (common vocabularies, common authorities, for example)
- Common Data models aren't as important, in LD, at the point of use.
- GLAM community has started sharing vocabularies and data using LD technologies
- Which people in EMLO have publications in the French national bibliography (biography?) dataset

- Lots of possibilities, but I don't know what you want to do on my own!
- List of people in EMLO who have letters in the D'Alembert corpus
- The tools are distinct from each other, and require hands-on manual intervention, there is a need to integrate them into production tools

HH: What has to be done to make it properly available for scholarship and usable? But these are the barriers: my serious research is something done over a long period of time (decades). Foundations are laid down at graduate school. Phobia we have is that we park our research and become DH savvy, then we work with you, that takes time, but does tech. move much faster, so our time is 'wasted'? Do we need to 'own' the tool and the infrastructure, as belonging to us?

EM: We want to get there, but that's not possible at present. We need exp. and collab. with CS on these tools. We need to figure out things together, and that involves some rawness

DO: You need to own the knowledge representation, not the software. That changes over time. The knowledge representations system can last for years and years.

HH: Our investment is in curating the data.. the nightmare is that the tech. is outmoded, and your career is over, before the fruits can be reaped. It's a enormous risk for a scholar. Perhaps the answer is to 'own' it as part of a community of like minded scholars. If the interests of the tech. people, change, it's game over, or if the funding changes.

AF: What makes your data accessible is how it is represented in LD. You'd need to lose a great deal to lose your data.

EM: You do stand to lose the tools, but you won't lose your data.

EH: linked data is based on W3C open standards, textual representation, no proprietary database format

Jouni Tuominen

- Different sources for registries, ontologies.. no coordination between them. VIAF, Getty ULAN, CERL.. No interoperability. Semantic and syntactic differences.
- Own databases, own APIs, and editing tools.
- Redundant work being done in different organizations
- Crowdsourced shared ontologies... prototyped as HIPLA, Finnish historical gazetteer service
- ONKI selector widget to query ontologies
- When adding a place, relevant links are useful, could be used for EMLO as well.

Antske Fokkens

- don't throw anything away
- if you can reuse something that exists, do it (EDM, PROV)
- Tool that lets you specify some entities also added, rest mostly reuse.
- Can also draw on news articles.
- Grounded Annotation Framework
- Targeted interpretation of family relations, professions.
- For historical sources, we don't make assertions, but identify entities. Algorithm helps determine subject/object for news stories (contemporary).
- Provenance is important - how the metadata was created is known, but we need to explain provenance by our automated processes
- Detailed level, for programmer/computational linguist. Global level for the historians.
- 20/80 of what? If an tool misses all the positive expressions as part of the 20 then it misses 100% of what is relevant
- Bias you introduce yourself, e.g. default is the netherlands. Fine, but not for all data sources
- Evaluate how well entities, expressions, concepts are recognized, then compared to metadata

- Then see how well we do, compared to manual check by historians

Lessons learned:

- Shared office.. makes collaboration much easier
- Constantly share information
- Keep best extrinsic/intrinsic evaluation in balance
- make the model as compatible as possible with existing data representations
- make it accessible, but accept that it will require some study.. this is a serious tool
- give as much information of the reliability as possible
- start with developing evaluation material, really from day one
- Get a minimally viable product as soon as possible
- This helps you focus on the right things, soon & fast.
- Data-to-document tool, using RDF data. Demonstrator.

TG: Identifiers for people, places. Authority list for events?

EH: Event authority files, not to my knowledge, except sideways, e.g. via LC definitions, for example.

AS: At what different scales is it useful for historians to model events.

AF: From sneeze to second world war..

TW: (Coming back to Howard's earlier concern) Encyclopaedias have been around for a long time, the community carries this kind of knowledge by itself

DO: Generalizations at too high level, like DC, do lose data. 'What kind of compromises are you willing to make to represent your data in a structured format'. Some projects make too many compromises, e.g. convert to EDM you lose knowledge you can't get back (events). What level do we need to capture to get reasonable results. Howard wants to ask difficult questions, which sound like modeling data.. (then?)

HH: is the BiographyNet tool/method able combine different perspectives on same person

AF: yes, not sure about how they are shown in visualizations

EH: how many event types was found in BiographyNet?

AF: Thousands (based on WordNet), but generic.

Topic IV: Future Prospects, Topics, and Partners

Sandra Toffolo

- ProsoNet
- Rather than taking a large, ideal approach (best model, best method, best technology..), taking a pragmatic approach.
 - small organizations with limited resources participating in prosopographical data creation
 - lots of independent heterogeneous databases in CESR, overlaps between them
 - no interoperability between them, no easy access to them simultaneously
- We can't ask people to make major changes to their databases
- But some silver linings as well, because the DBs are all 'best suited' for their research needs
 - different languages
 - technical database solutions
 - different sets of reference vocabularies (CERL, VIAF, etc.) used
- construction of ProsoNet repository, definition of metadata fields, metadata format, authority lists
- transform metadata
- export metadata from individual databases into ProsoNet repository
- Goal is to avoid adding further duplicates of the same names. Just the first and last name is being mapped to DC

- manual disambiguation
- creation of ProsoNet authority list with ID numbers
- add ProsoNet ID number to original CSER database
- new metadata export
- regular updates to ProsoNet repository (automatic updates, but manual disambiguation)

AB: different fields of different databases are mapped into Dublin Core?

ST: yes, common denominator of metadata fields were mapped

TW: is lastname + firstname (+ specifiers included in those fields) enough for disambiguating people?
Very difficult, in particular if one looks at medieval period.

ST: True. We'll have to take that into consideration.

CK: complexity vs. ... historians like complex data models for their own data; core model for prosopographical data, basic information to identify people across sources and corpora, the usage of linked data

DO: semantic web and linked data communities encourage people to get data fast out there; this may lead to oversimple data models; vocabularies presented today are not designed for open collaboration in research environment, but for internal cataloging; if you want to do historical research, you need complex data models

HH: reliability, robustness, and durability are the key requirements

DO: what's the core: as we've heard today, there exists so different datasets

AF: linked data technology allows for generic model on top of more detailed ones

DO: the current systems on the web usually do not do that though

CK: huge disconnect between CS and library information science, authority records (disambiguate and reconcile data between sources)

DO: library science guys think they can apply their techniques into DH, but it does not work

HH: can we begin consolidating multiple datasets by creating rules how persons can be matched (similar name, flourished at the same time, etc.); core set of categories(?)

HH e.g., BiographyNet's data model could be prime candidate if there is such

HH: disciplinary matrix, we need to support multiple perspectives: current, historical

DO: i am worried about fixed core model; three simple fields is not good, if their usage isn't clear, need to have context; implicit semantics made explicit in the database; lots of cultural heritage institutions have forgot meanings of some of their database fields; semantic web tries to answer to this (add context to data)

HH: if your criticism on current systems is that they're not rigorous and semantic, we can take that as input to RRL

DO: in linked data and research context, you need to have provenance and be able to separate own data and open world (linked data cloud)

AF: provenance is important

TW: asked Antske to send her 21 research questions mentioned in the presentation

ML: we could ask users what kind of ... would they like to have

CK: different needs: 1) ... (simpler questions), 2) questions like: how several networks intersect beyond epistolary contact - for these one needs to have complex prosopographical data

EH: data models should be polished, documented (e.g., description of the fields) and published to start the discussion with the community, get feedback and corrections; e.g. developing CIDOC CRM is hard work; Getty open data documentation has, e.g., duality (the place/person/etc. in reality is separated from the conceptual view of it in a vocabulary)

HH: asked Eero to crystallize that into paragraph or two

Tanya presented a demonstrator of a facet search on prosopographical data, based on Blacklight open source software

Discussion: Future Steps for WG2 and Horizon 2020

Eero closed the meeting