

# Who is Involved? Semantic Search for E-Discovery

David van Dijk<sup>†‡</sup>  
d.v.vandijk@uva.nl

David Graus<sup>‡</sup>  
D.P.Graus@uva.nl

Zhaochun Ren<sup>‡</sup>  
Z.Ren@uva.nl

Hans Henseler<sup>†</sup>  
j.henseler@hva.nl

Maarten de Rijke<sup>‡</sup>  
derijke@uva.nl

<sup>†</sup>Create-IT, Amsterdam University of Applied Sciences, Amsterdam, The Netherlands

<sup>‡</sup>University of Amsterdam, Amsterdam, The Netherlands

## ABSTRACT

E-discovery projects typically start with an assessment of the collected electronic data in order to estimate the risk to prosecute or defend a legal case. This is not a review task but is appropriately called early case assessment, which is better known as exploratory search in the information retrieval community. This paper first describes text mining methodologies that can be used for enhancing exploratory search. Based on these ideas we present a semantic search dashboard that includes entities that are relevant to investigators such as who knew who, what, where and when. We describe how this dashboard can be powered by results from our ongoing research in the “Semantic Search for E-Discovery” project on topic detection and clustering, semantic enrichment of user profiles, email recipient recommendation, expert finding and identity extraction from digital forensic evidence.

## Categories and Subject Descriptors

H.3 [Information Storage and Retrieval]: H.3.3 Information Search and Retrieval, H.3.1 Content Analysis and Indexing; H.5 [Information Interfaces and Presentation]: H.5.2 User Interfaces

## Keywords

E-discovery; Exploratory search; Entity extraction; Text mining; Semantic search; Technology assisted review; Early case assessment

## 1. INTRODUCTION

Early case assessment (ECA) is a term that is often used in E-discovery and that refers to disclosure of electronically stored information. According to Wikipedia,<sup>1</sup> ECA refers to estimating risk (cost of time and money) to prosecute or defend a legal case. According to [16], ECA is also useful in

<sup>1</sup>[http://en.wikipedia.org/wiki/Early\\_case\\_assessment](http://en.wikipedia.org/wiki/Early_case_assessment)

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

*DESI VI Workshop*, June 8, 2015, San Diego, California.

digital investigations by law enforcement investigations and regulatory bodies.

We argue that ECA is a type of exploratory search [17]. Exploratory search is a form of information retrieval where users start without a clear information need. They do not know beforehand what they are looking for, nor where to find it. In E-discovery this means not only investigators are looking for a needle in the haystack but also that they do not know what the needle looks like.

When looking for technology to support exploratory search, a technology assisted review (TAR) method like predictive coding does not work. At the start of ECA the investigator does not know what exactly to look for and there are no examples that can drive the predictive coding. As an alternative we propose to combine text mining with exploratory search. Text mining provides additional structure to unstructured data that enables interactive filtering by users.

This idea has led to the “Semantic Search for E-Discovery” project on which we have reported earlier at DESI IV [26] and V [10], after an initial paper exploring social network analysis for E-discovery at DESI III [12]. Since DESI V, three Ph.D. students at the Informatics Institute of the University of Amsterdam have been researching text mining approaches that are relevant to semantic search in E-discovery. They are designing new approaches to solve parts of the semantic search problem.

The results of the separate research projects have already been published elsewhere. This paper is intended as a follow up to our previous DESI papers to explain and discuss with the DESI workshop community how the results will be combined in one approach for semantic search in E-discovery.

We first discuss exploratory search and then outline a technical solution that we call the semantic search dashboard. This dashboard supports typical user interaction patterns that enable one or more users to explore an unknown dataset for relevant topics, persons and time frames. We then identify existing approaches from text mining and information retrieval and identify the challenges that need to be solved to realize a usable system.

## 2. EXPLORATORY SEARCH

Investigators in a corporate E-discovery investigation are typically overwhelmed with thousands if not hundreds of thousands of emails that need to be investigated. Often, for these users the E-discovery task does not start as a review task but as an early case assessment task which is a form of exploratory search [17]. Investigators explore the available

information trying to construct a time line that describes who knew what and when, which persons are collaborating, that identifies locations etc. This objective is specifically addressed in time-aware exploratory search and associations of entities over time. These research topics will be discussed in more detail below.

## 2.1 Exploring word meaning through time

The scope of investigations in E-discovery can span several years. Investigations like Enron, Lehman brothers and Mad-off go back from 5 to 10 to sometimes even 30 years. Time is a very important aspect in these investigations and early case assessment is frustrated by the overwhelming amount of information.

The meaning of a word can be inferred from observing its usage, i.e., its distribution in text, and changes to its usage. By visualizing word clouds in an exploratory interface that includes the temporal aspect of the documents users can get a quick insight through summarization [21].

For instance, the context of the word "earnings" may change over time depending on what topics are being discussed in relation to the earnings of a company or organisation. In the course of various months or years, earnings are likely to depend on different projects. Some may even be associated with losses etc.

Once interesting words have been identified in the context of relevant concepts, the presence of these words can be used to discover document clusters in the data set that can assist users to more effectively browse through the available data.

## 2.2 Exploring entity associations over time

Discovering entities in documents and emails can be helpful in exploring an unknown data set. In computer forensics, digital evidence is examined from which person names, email addresses, phone numbers, chat skype id's etc can be extracted automatically from the metadata.

This information can be used to increase the quality of entity extraction when applied to unstructured information in documents and email message texts. Co-occurrence patterns at the document and sentence level can reveal relationships between entities and when they were in existence.

In [22], an exploratory search interface is further extended to retrieve entities that are relevant to a search query and to discover associations between these entities over time. In a historical perspective it is important to anchor political figures in time and to construct a narrative around a certain entity or to investigate the cause and effect of a phenomenon.

Semantic search aims at returning this information directly to a user instead of a list of documents. For instance, in addition to representing the number of responsive documents as a function of time, a histogram representing related entities may be presented that visualises when an entity is associated with responsive documents.

Associations between entities can be computed at the document and sentence level. The authors hypothesize that co-occurrence at the document level results in a more topical association whereas co-occurrence at the sentence level indicates a more relational or functional association. A visual exploratory approach is followed to support users in discovering entities and associations with a temporal filter.

## 2.3 Exploring who was involved and when

In E-discovery we do not only want to find out which entities (or persons) are associated but we would also like to

know who is associated with a particular discussion topic. Discussion topics are typically introduced by a person at a certain moment. Topics may slightly change over time or disappear.

ThemeStreams [7] is a demonstrator focusing on the Dutch political landscape by analyzing political discussions on Twitter. It keeps track of discussions and who is involved in these discussions. One of the aims is to find out who started discussing a particular topic. Who put the topic on the map etc. Four groups are differentiated: persons who have an important position (the politicians), persons who lobby for important issues, the journalists and finally all other people taking part in political discussions (other influencers).

## 3. THE SEMANTIC SEARCH DASHBOARD

Inspired by the work on semantic exploratory search just outlined, we envision a solution to semantic search for E-discovery that enables investigators to quickly scan a large volume of textual documents such as emails, attachments, meeting minutes etc. for topics, who is involved in these topics and when these topics were discussed. A topic would typically be represented to the user as a word cloud that summarizes keywords that are typical for the topic. The solution should offer an interactive interface that enables investigators to explore a case. The interface should be visualizing discussion topics over time and who is involved. Users should be able to drill down on particular time slots, topics and persons. Topics are presented as word clouds. The role of a person is automatically identified by the system (just reading, initiating, answering etc.). Fig. 1 introduces the semantic search dashboard. The dashboard has 8 interactive elements that can be used by the investigator to explore the electronic evidence in a case:

1. Full-text search
2. List of topics
3. Time line
4. Languages
5. List of persons
6. Topic(s) summary
7. Social network analysis
8. Locations

Typical exploratory interaction patterns for investigators could be as follows

- (a) User runs a full-text search to find related topics
- (b) Set time to the relevant start and end period of the investigation scope
- (c) Select a language of interest
- (d) Remove topics from the view that are not interesting
- (e) Remove non-relevant words from the topic summary
- (f) Remove persons that are not relevant
- (g) Merge persons that appear to be the same person

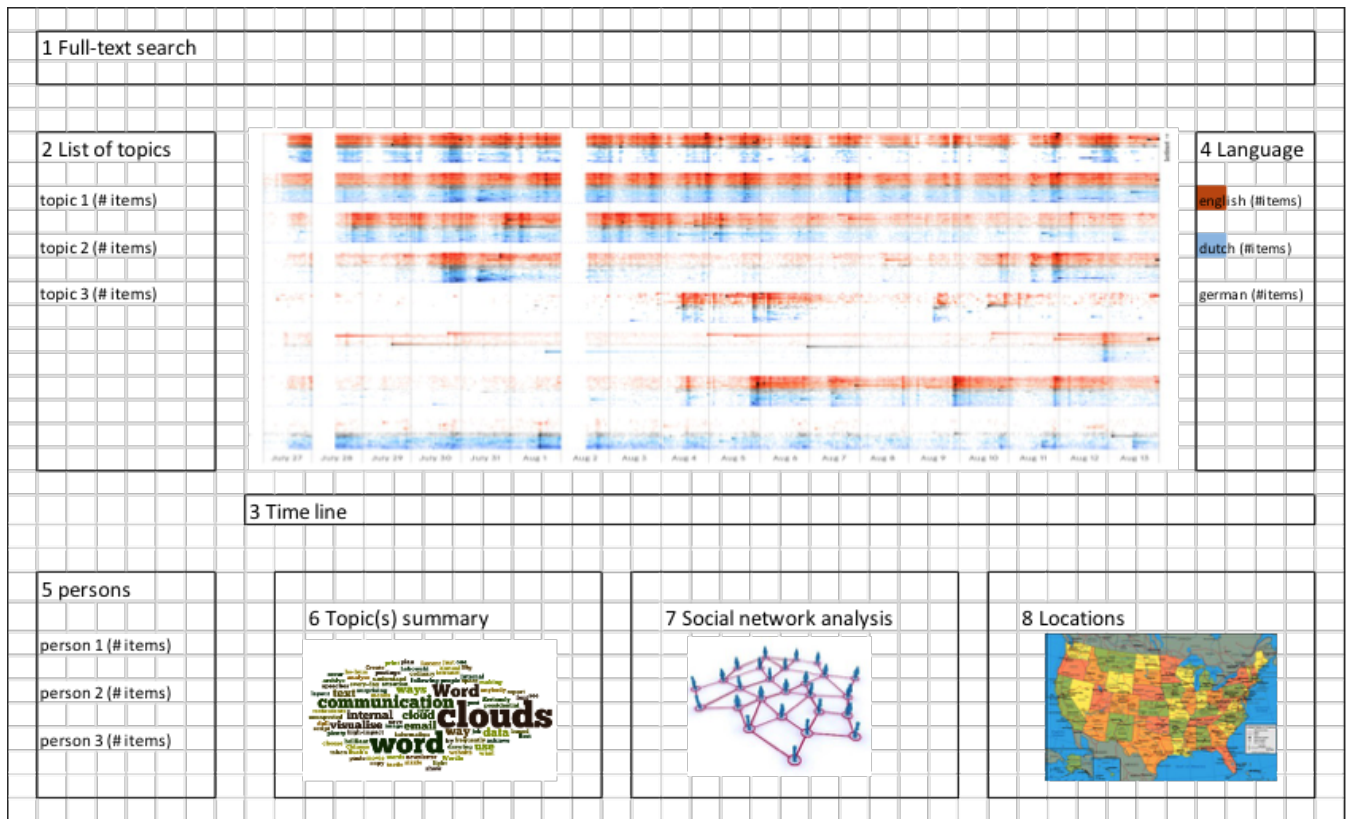


Figure 1: The semantic search dashboard

- (h) Explore the social network of a person
- (i) Restrict search to a geographical location

In an exploratory search approach users combine different interaction patterns by repeatedly eliminating information, which reduces the number of documents while increasing the fraction of relevant documents. This interaction pattern resembles the Scatter/Gather pattern that was introduced in the early 90's as a cluster-based approach to browsing large document collections [5].

#### 4. EXISTING COMMERCIAL SOLUTIONS AND THEIR CHALLENGES

Commercial E-discovery solutions have evolved from full-text search engines to rich user interfaces with support for faceted filtering, document clustering, find similar documents, concept extraction, email thread detection etc. Some tools have integrated named entity extraction and have added filters to the interface that enable users to interactively filter search results.

Named entity extraction still suffers from a relatively high ratio of false positives where meaningless phrases are identified as names. In reality these approaches are highly language dependent and only work when pre-coded dictionaries of names and places are added to the system.

The most advanced solutions try to implement fact finding or even story telling where relations between named entities are automatically extracted. This requires manual coding of fact finding-pattern rules such as company-hires-person. Some tools can produce a topic heat map but they need to be tuned manually and heat maps are not produced in real-

time and therefore difficult to use in an interactive filtering method.

The application of social network analysis, location-based visualization and filtering exists primarily in specialized tools for analysis of call detail records and mobile phone data. With the current advances in social media and the increasing use of smartphones, we expect that these types of filtering will become part of E-discovery investigations as they are already part of law enforcement investigations.

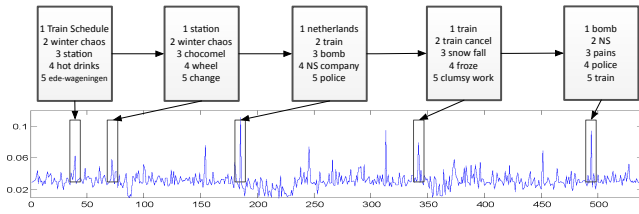
Concluding, we can say that partial implementations of the proposed semantic search dashboard exist but that integrated and scalable solutions are missing. Advanced text mining applications still suffer from poor precision, in particular when dealing with inconsistent (meta)data that is so typical for informal communication data that is often at the heart of E-Discovery data collections.

The outcome of these tools requires manual clean-up, which may be feasible for enterprise search in formal knowledge bases but which is not feasible for large scale E-discovery investigations. This problem is further worsened by the "rolling" nature of data collections in E-Discovery projects that requires an incremental or continuous learning process which is something we intend to address in our research.

#### 5. WORK IN PROGRESS

A number of techniques that are required for the semantic search dashboard presented in Section 3 are readily available from the exploratory search interfaces presented in [7, 21, 22]. They allow users to explore documents over time, searching for topics and finding associated entities.

For early case assessment in E-discovery tasks we want to extend the exploratory search with additional facets and



**Figure 2: An example topic propagation in a public transportation system dataset. The text blocks at the top indicate the top 5 representative terms for the topic being propagated at a specific time period; the bottom side shows the topic distribution over the whole timeline.**

increase the precision of topic detection and entity relation extraction by leveraging meta data that is present in forensic computer evidence.

Since Desi V we have conducted a number of separate research projects in the “Semantic Search for E-Discovery” project that are relevant to the semantic search dashboard.

### 5.1 Topic detection and clustering

Emails, microblogs and other social text streams may contain multiple topics [1, 15, 19], which can be presented as word clouds. Topic detection and clustering over those textual documents is meaningful to explore who is involved in these topics and when these topics were discussed. Work on detecting topics in social text streams can be used but needs to be adapted. In recent years, topic models have proved effective in topic detection and clustering [3, 4, 13, 25]. Among existing approaches to topic models, Latent Dirichlet allocation (LDA) [3] has become an attractive method to classify topics [24]. Topics can be automatically summarized and presented to the user in a visual manner that helps users to quickly identify the relevant topics. The system should support (1) elimination of irrelevant topics (like email footers, standard emails) and (2) identification of anomaly topics, i.e., topics that are not standard. Another challenge is to determine when a new topics starts, how it changes over time and finally when it ends. Topics in text streams may not be stationary. By employing a dynamic extension of the LDA model to track dynamic topics, a recent time-aware topic detection and clustering approach [23] starts to tackle the topic drift and concept drift problems, which have so far been neglected by most previous topic detection and tracking approaches. Using a dataset of tweets related to a major public transportation system in a European country, Fig. 2 shows the propagation process of an example topic. The upper part of Fig. 2 shows the 5 most representative terms for the topic during 5 time periods. The bottom half of the figure plots fluctuating topical distributions over time, which indicates concept drift between two adjacent periods.

### 5.2 Semantic enrichment of user profiles

In [9] yourHistory is presented, a Facebook application that aims to generate a tailor-made, personalized timeline of historic events. The application is driven by matching semantically enriched Facebook profiles to historic events from DBpedia. Semantically enriched profiles are constructed by linking textual content extracted from Facebook profiles (e.g., biographic information, likes, posts) to entities and concepts in an open knowledge base (DBpedia). The under-

lying approach of using entity linking to semantically enrich user profiles is a form of document expansion; by including the (linked) entity descriptions from the knowledge base to the Facebook profiles, additional textual information can be leveraged for retrieval. In particular in scenarios where documents (e.g., emails) and queries (e.g., topics) are short, adding textual content to documents or queries can be beneficial, as has been shown in, e.g., web search scenarios [6]. Furthermore, the rich link-structure of knowledge bases can provide additional signals for contextualizing and summarizing information, e.g., by showing how concepts in a single document are related. This latter approach of summarizing textual content through visualizing a network of knowledge base concepts has been explored in an interactive demo in [20].

### 5.3 Email recipient recommendation

In [11] we study email recipient recommendation on enterprise email, where the task is to predict the recipient(s) of an email, given its sender, email content, and all previously seen emails in the network. More specifically, we study and compare the predictive power of communication graph signals (i.e., social network signals), and email content signals.

The communication graph signals consist of different methods of estimating email users’ importance in the network, and different ways of estimating the prior probability of two users communicating. To model email content, we turn to statistical language modeling [28]. More specifically, we compare different ways to represent the content of users’ communication, by computing different language models with different sets of emails. We take, e.g., all of a user’s outgoing emails, their incoming emails, or all emails that have been sent between two users. These different ways to represent email content, and users’ language profiles, has further potential applications in, e.g., identifying anomalous communication. For example, given an email dataset it may be of interest to investigators to retrieve: (1) emails sent between two users that differ significantly from the “average” emails in the dataset, or (2) emails from a single user that differ significantly from the “average” emails of the user, or (3) emails from a single user that differ significantly from the emails in the dataset.

Finally, a model that successfully predicts recipients of emails may be employed to identify unexpected communication patterns, e.g., by computing the model’s confidence of observed emails, and flagging all the emails that receive low confidence scores.

### 5.4 Who is involved and when

The work described in section 2.3 looked at different user groups. It provides a visualisation that shows the activity resp. involvement of each group in a particular discussion topic over time. Characteristics of involvement can be extracted, such as which group started the discussion. We focus on persons instead of groups. Besides who and when, we are interested in how a person is involved in a topic, e.g. we can look at a persons centrality and how it changes over time, or his spread and degree of involvement. Our recent work focusses on the latter. The task of retrieving persons with a high degree of topic involvement strongly relates to the task of expert finding and retrieval [2]. Users and topics are profiled for generating candidate matches. As annotated corpora relevant to the E-Discovery domain are hard to come by, we use a corpus from a domain. Stack

Overflow,<sup>2</sup> a Community Question Answering (CQA) site, provides a rich datasource for user topic involvement.

Using this data, we study the importance of combining a large number of textual, behavioral and time-aware signals for detecting early expertise [27]. We define expertise by the number of a user's best answers rather than time, catering for different user activity behaviors. Our semi-supervised approach leverages textual, behavioral and temporal feature sets. Combining behavioral and temporal features with textual features significantly boosted effectiveness. Our system can accurately predict whether a user will become an expert from a user's first best answer; projected over time, our system makes correct predictions as far ahead as 70 months before a user becomes an expert. In future work, we plan to extend our features to capture more aspects of early expertise, e.g., answer quality, diversity, novelty, and also track how a user's expertise evolves from one topic to another over time, which can yield a strong predictor of early expertise. Although this specific task, data and features might not directly fit to the E-Discovery domain, we expect the lessons learned to be of value for modelling user topic involvement, especially the incorporation of the temporal aspect.

Inspired by work on topic and role detection [18] and on community detection [29] we will further proceed to explore ways of identifying the role of a person in discussions over time. Another aspect we will address is anomaly detection in a user's behavior, e.g., in his involvement in a specific topic.

## 5.5 Extracting identities from digital forensic evidence

Extracting named entities using natural language processing typically suffers from low precision. Most solutions improve quality by using dictionaries of known names. For instance, entity linking in microblogs can be improved by leveraging open source knowledge in Wikipedia, WordNet, DBpedia or other open sources. Unfortunately, in E-discovery open source information does not suffice and dictionaries change for every investigation.

Entity extraction from known metadata in digital forensic evidence (i.e., user accounts, user names in Office documents, names in mobile phone or email address books, addresses in emails etc.) provides a good basis to build a dictionary of entities in a case that can be used for entity extraction [14]. Investigators can be tasked with merging aliases and with the identification of key persons in the investigation (e.g., suspects, witnesses). It then becomes a relatively easy task to discover these entities in unstructured content in the same case and suggest other entities to the user that seem to be related.

## 6. FUTURE WORK

Future work in the "Semantic Search for E-Discovery" project will focus on the validation of the presented semantic search techniques for exploratory search in an E-discovery context. This requires, first of all, the development of an integrated prototype that can be shared with user groups in the project. To support the development of this prototype, the open source UFORIA<sup>3</sup> framework [8] will be used that has been developed at the Amsterdam University of Applied

Sciences to support indexing, interactive filtering and visualization of multi-faceted data.

In addition to testing semantic search techniques for early case assessment, we also intend to investigate the application of semantic search for enhancing predictive coding in E-discovery review tasks. We expect that the application of semantic search for early case assessment and feedback from users will provide us with useful insights. We intend to validate our approach by participating in the new Total Recall Track which is one of the eight tracks in the 2015 edition of TREC.<sup>4</sup>

The Total Recall task resembles the TREC 2011 legal track in that the objective is to find as nearly as possible all relevant documents while examining as few as possible. However, the level of automation will be higher by providing a web interface that simulates a human reviewer. Alternatively, participants submit an encapsulated version of their system in a virtual machine which will be run by the TREC assessors on non-public data.

**Acknowledgments** This research was supported by the European Community's Seventh Framework Programme (FP7/2007-2013) under grant agreement nr 312827 (VOX-Pol), the Netherlands Organisation for Scientific Research (NWO) under project nrs 727.011.005, 612.001.116, HOR-11-10, 640.006.013, 612.066.-930, CI-14-25, SH-322-15, Amsterdam Data Science, the Dutch national program COMMIT, the ESF Research Network Program ELIAS, the Elite Network Shifts project funded by the Royal Dutch Academy of Sciences (KNAW), the Netherlands eScience Center under project nr 027.012.105, the Yahoo! Faculty Research and Engagement Program, the Microsoft Research PhD program, and the HPC Fund.

## 7. REFERENCES

- [1] J. Allan. Introduction to topic detection and tracking. In *Topic detection and tracking*, pages 1–16. Springer, 2002.
- [2] K. Balog, Y. Fang, M. de Rijke, P. Serdyukov, and L. Si. Expertise retrieval. *Foundations and Trends in Information Retrieval*, 6(2–3):127–256, August 2012.
- [3] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022, 2003.
- [4] J. Chang, S. Gerrish, C. Wang, J. L. Boyd-graber, and D. M. Blei. Reading tea leaves: How humans interpret topic models. In *NIPS*, 2009.
- [5] D. R. Cutting, D. R. Karger, J. O. Pedersen, and J. W. Tukey. Scatter/gather: a cluster-based approach to browsing large document collections. In *SIGIR '92 Proceedings of the 15th annual international ACM SIGIR conference on Research and development in information retrieval*, SIGIR '92, pages 318–329, New York, NY, USA, 1992. ACM.
- [6] J. Dalton, L. Dietz, and J. Allan. Entity query feature expansion using knowledge base links. In *Proceedings of the 37th International ACM SIGIR Conference on Research & Development in Information Retrieval*, SIGIR '14, pages 365–374, New York, NY, USA, 2014. ACM.
- [7] O. de Rooij, D. Odiijk, and M. de Rijke. Themestreams: Visualizing the stream of themes discussed in politics. In *SIGIR'13: 36th international ACM SIGIR conference on Research and development in information retrieval*. ACM, 07/2013 2013.
- [8] A. Eijkhoudt and T. Suerink. Uforia: Universal forensic indexer and analyzer. *Journal of Computer Virology and Hacking Techniques*, 9:59–63, 2013.
- [9] D. Graus, M.-H. Peetz, D. Odiijk, O. de Rooij, and M. de Rijke. yourhistory – semantic linking for a personalized timeline of historic events. In *Proceedings of the LinkedUp Veni Competition on Linked and Open Data*

<sup>2</sup><http://stackoverflow>

<sup>3</sup><http://www.uforia.nl/#/search>

<sup>4</sup><http://trec.nist.gov/pubs/call2015.html>

- for Education held at the Open Knowledge Conference, volume 1124, Geneva, Switzerland, 09/2013 2013. CEUR Workshop Proceedings.
- [10] D. Graus, Z. Ren, M. de Rijke, D. van Dijk, H. Henseler, and N. van der Knaap. Semantic search in e-discovery: An interdisciplinary approach. In *ICAIL 2013 Workshop on Standards for Using Predictive Coding, Machine Learning, and Other Advanced Search and Review Methods in E-Discovery (DESI V Workshop)*, 06/2013 2013.
- [11] D. Graus, D. van Dijk, M. Tsagkias, W. Weerkamp, and M. de Rijke. Recipient recommendation in enterprises using communication graphs and email content. In *Proceedings of the 37th International ACM SIGIR Conference on Research & Development in Information Retrieval, SIGIR '14*, pages 1079–1082, New York, NY, USA, 2014. ACM.
- [12] H. Henseler. Network-based filtering for large email collections in e-discovery. *Artificial Intelligence and Law*, 18(4):413–430, 2010.
- [13] M. Hoffman, F. R. Bach, and D. M. Blei. Online learning for latent dirichlet allocation. In *NIPS*, pages 856–864, 2010.
- [14] J. Hofste, H. Henseler, and M. van Keulen. Computer assisted extraction, merging and correlation of identities with tracks inspector. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Law*, pages 247–248. ACM, 2013.
- [15] P. Krafft, J. Moore, B. Desmarais, and H. M. Wallach. Topic-partitioned multinet network embeddings. In *NIPS*, 2012.
- [16] D. Lawton, R. Stacey, and G. Dodd. E-discovery in digital forensic investigations. Technical Report CAST publication number 32/14, UK Home Office, Centre for Applied Science and Technology, 2015.
- [17] G. Marchinini. Exploratory search: from finding to understanding. *Communications of the ACM*, 49(4):41–46, 2006.
- [18] A. McCallum, X. Wang, and A. Corrada-Emmanuel. Topic and role discovery in social networks with experiments on enron and academic email. *Journal of Artificial Intelligence Research*, pages 249–272, 2007.
- [19] Q. Mei, D. Cai, D. Zhang, and C. Zhai. Topic modeling with network regularization. In *WWW*, 2008.
- [20] D. Odijk, E. Meij, D. Graus, and T. Kenter. Multilingual semantic linking for video streams: Making “ideas worth sharing” more accessible. In *Proceedings of the 2nd International Workshop on Web of Linked Entities (WoLE 2013)*, 2013.
- [21] D. Odijk, G. Santucci, M. de Rijke, M. Angelini, and G. Granato. Time-aware exploratory search: Exploring word meaning through time. In *SIGIR 2012 Workshop on Time-aware Information Access*, Portland, OR, USA, 08/2012 2012.
- [22] R. Reinanda, D. Odijk, and M. de Rijke. Exploring entity associations over time. In *SIGIR 2013 Workshop on Time-aware Information Access*, 08/2013 2013.
- [23] Z. Ren, M.-H. Peetz, S. Liang, W. van Dolen, and M. de Rijke. Hierarchical multi-label classification of social text streams. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, pages 213–222. ACM, 2014.
- [24] Z. Ren, D. van Dijk, D. Graus, N. van der Knaap, H. Henseler, and M. de Rijke. Semantic linking and contextualization for social forensic text analysis. In *Intelligence and Security Informatics Conference (EISIC), 2013 European*, pages 96–99. IEEE, 2013.
- [25] M. Steyvers, P. Smyth, M. Rosen-Zvi, and T. Griffiths. Probabilistic author-topic models for information discovery. In *KDD*, 2004.
- [26] D. van Dijk, H. Henseler, and M. de Rijke. Semantic search in e-discovery. In *DESI IV Workshop on Setting Standards for Searching Electronically Stored Information In Discovery Proceedings*, 2011.
- [27] D. van Dijk, M. Tsagkias, and M. de Rijke. Early detection of topical expertise in community question answering. In *Submitted*, 2015.
- [28] C. Zhai. Statistical language models for information retrieval a critical review. *Foundations and Trends in Information Retrieval*, 2(3):137–213, 2008.
- [29] Z. Zhao, S. Feng, Q. Wang, J. Z. Huang, G. J. Williams, and J. Fan. Topic oriented community detection through social objects and link analysis in social networks. *Knowledge-Based Systems*, 26:164–173, 2012.