# E-discovery viewed as integrated human-computer sensemaking: The challenge of 'frames'

**Simon Attfield**
UCL Interaction Centre*
s.attfield@cs.ucl.ac.uk

**Ann Blandford**
UCL Interaction Centre*
a.blandford@cs.ucl.ac.uk

University College London Interaction Centre, MPEB 8th floor, University College London, Malet Place, London, WC1E 7JE, UK

## ABSTRACT

In addressing the question of the design of technologies for e-discovery it is essential to recognise that such work takes place through a system in which both people and technology interact as a complex whole. Technology can promote discovery and insight and support human sensemaking in this context, but the question hangs on the extent to which it naturally extends the way legal practitioners think and work. We describe research at UCL which uses this as a starting point for empirical studies to inform the design of supporting technologies. We report aspects of an interview field study with lawyers who worked on a large regulatory investigation. Using data from the study we describe document review and analysis in terms of a sequence of transitions between different kinds of representation. We then focus on one particular transition: the creation of chronology records from documents. We develop the idea that investigators make sense of evidence by the application of conceptual 'frames' (Klein et al's, 2006), but whilst the investigator 'sees' the situation in terms of these frames, the system 'sees' the situation in terms of documents, textual tokens and metadata. We conclude that design leverage can be obtained through the development of technologies that aggregate content around investigators' frames. We outline further research to explore this further.

## INTRODUCTION

Electronic Data Discovery (EDD, or e-discovery) has been defined as a process (or series of processes) in which electronic data is sought, located, secured, and searched with the intent of using it as evidence in a civil or criminal legal case, or as part of a court-ordered or government sanctioned inspection [Conrad, 2007].

The rapid increase in the volume of electronically stored information within modern enterprises has led to a situation in which preparing for and executing e-discovery represents a considerable challenge for modern organizations and legal companies, and it is one that is set to increase. It is likely that we will increasingly see companies falling prey to legislation if they cannot uncover all electronically stored information (ESI) relevant to a legal or regulatory matter within a specified timeframe [Baron, 2008].

Advances in digital technologies which have brought about this challenge, however, also offer part of the means for addressing it. The e-discovery technology industry is seeing year-on-year increases in turnover. Software revenues in 2006 were estimated at around $150 million, with further vigorous growth predicted [Socha & Gelbmann 2007]. Technologies attracting particular interest in this arena include media restoration tools, dedicated document management systems, advanced information retrieval systems (such as concept search and information extraction), information visualization and case analysis tools.

In addressing the question of how to design technology for e-discovery, however, it is essential to recognise that e-discovery work takes place through the operation of a system in which both people and technology interact as a complex whole. In this context, the role of technology is to provide tools and resources that can be usefully appropriated by legal professionals, often working in teams, in constructing strategies and processes that address their goals more effectively. Understanding how technologies can offer additional leverage depends on how those technologies impact on and reshape such systems for the better.

In considering technology developments a significant research object is the e-discovery process viewed as a complex worksystem. Such a perspective becomes particularly pertinent where people are required to engage in intense cognitive activities such as information assimilation, theorising and reasoning, as occurs during the review and analysis of large document collections. As in all branches of knowledge work, technology can promote discovery and insight and support human sensemaking, but the question hangs on the extent to which it integrates within and naturally extends the way that legal practitioners think and work.

We argue that the design of systems to support this kind of work needs to be predicated upon an understanding of the cognitive and social aspects of e-discovery in practice. This

mandates a detailed understanding of the task as it unfolds, including associated processes of sensemaking, teamwork, how people currently coordinate different tools and resources to meet their aims, and what barriers and difficulties arise in doing so. In essence, the need is to examine how work is done in order to speculate how it might be done better [Rassmussen et al 1994].

With this in mind, we are conducting research in field and laboratory settings with the aim of better understanding evidence review and analysis in e-discovery in order to support reasoning about the design of supporting technologies. In this paper we provide an example of that work relating to an interview field study we performed with lawyers who worked on a large regulatory investigation.

In analysing the data from this study two complementary perspectives emerged. The first, which we have reported elsewhere [Attfield et al, 2008a], focuses on how the investigation work was structured. This concerns how the investigators made the investigation tractable by decomposing it into multiple, emerging lines of enquiry or 'issues' distributed across a team. Significant issues relate to how the investigation was decomposed along emerging lines of enquiry resulting from ongoing discoveries, and on the challenge of integrating outcomes from multiple investigation threads to form an integrated perspective.

The second perspective is that of process. Complex knowledge work often occurs in stages which form an iterative sequence of transformations between different kinds of intermediate representation [Attfield, 2008b]. These representations can embody task objectives (such as questions), constructed sub-sets of data with particular meaning (such as search results) or recorded findings and interpretations (such as notes, narratives and structured knowledge representations). As a representation is created

or changed, so it provides raw material for further work, creating new representations and so on. In this way resources act as stepping-stones on an iterative path of sensemaking and a key part of that is the information processing that is performed in order to transition from one representation to another.

In this paper we focus on this second perspective. We first describe the process of document review and analysis in overview in terms of a sequence of transitions between different kinds of representational resource. Next we focus in on one transition in detail, describe how it was done, and use this discussion to reflect on alternative technologies that might offer additional leverage.

## INVESTIGATION PROCESS

Based on the interviews, we developed a description of the document review process as shown in figure 1, in the form of a 'process-resource' model. In this figure, boxes represent resources and arrows represent transitions between them. For example, given a set of investigation issues and a document universe, keyword searching (t2) resulted in sets of search results. Given a set of search results, initial manual review (t3) produced a set of documents coded as relevant.

These resources changed constantly and were interdependent but only through transformations (i.e. t1 to t7) which were created by the investigators. The transformations were achieved through some form of information processing, whether this be the investigators reviewing a resource and recording the outcome of their thinking, or by their additional use of automated processing, such as information retrieval.

Each transformation, then, has the effect of using one or more resources in order to shape another, with each
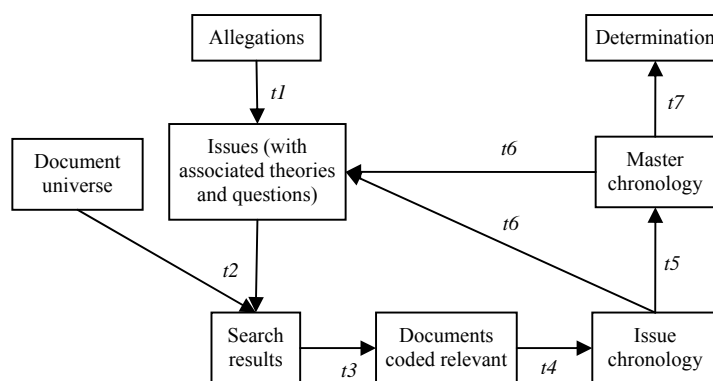


Figure 1. An overview of the investigation process.

representing an intermediate step that ultimately links allegation with a determination.

In overview, the transitions were:

t1 - Given the allegations, the investigators defined and recorded a set of issues that they wanted to investigate and associated questions they wanted to ask.

t2 – Given these questions, queries were submitted to the document universe to return documents relevant to each of the issues.

t3 - Returned documents were individually read and coded for relevance to the issues (within a document management system).

t4 – Relevant documents were then used to infer entries within issue-specific event chronologies.

t5 – Selected entries within separate issue chronologies were 'escalated' into a single master chronology designed to record the most significant aspects of the developing narrative.

t6 – By reflecting upon the narratives within the chronologies as they evolved, the investigators were able to identify apparent gaps, inconsistencies and periods of potential interest. This helped them to develop theories which guided the refinement of the investigation issues and associated questions.

t7 – given the knowledge acquired, the investigators formed a view concerning the allegations.

In fact, the structure of the investigation process evolved somewhat over time. What we present here is the process in its mature form. We have also restricted ourselves to a description of the investigation as it applied to electronic documents, omitting reference to witness interviews which were nevertheless an important, if non-technological, source of information.

A number of things are happening in this process, but broadly we see it as a process of information reduction achieved by different kinds of filtering and abstraction, directed by reflective interpretation on the part of the investigators.

Two things are important to note. First the investigators constructed each step for a reason—this being in general terms to help them move in a direction in which they wanted to go. Hence we can learn about their needs from what they did. The second point is that although they had discretion to design the process as they saw fit, they did so within the constraints of the tools available to them at the time and whatever costs there were associated with their appropriation and use. Hence we can use the process to consider other tools which may have supported their needs better.

**Focussing on transition t4**

In considering where new technologies might offer leverage we might focus attention in detail on any part of the process we have described to consider how things might be changed (or even change the process as a whole). Our interviewees consistently cited manual document review (stages t3 and t4 in figure 1) as imposing the major overhead in terms of time and effort. Over the course of the investigation 130,000 documents were reviewed in all. This represents a significant reduction on the document universe, but is nevertheless a very significant number of documents. Here we will consider transition t4 in more detail.

T4 involved the creation of semi-standardized event records based on the review of documents which themselves had been coded as relevant to one or other issue. The investigators constructed chronologies as a table using Microsoft Excel according to a preformed schema. An example of an anonymised record which reflects this schema is shown in figure 2.

The reason for creating chronologies was so that the investigators could have a compact representation of events they considered to be of significance to their investigation. This then provided a resource for considering what they knew, for developing theories from this and in other ways establishing what it was they wanted to find out (transitions t6).

The resource for creating event records (transition t4) was a list of documents (predominantly emails) displayed in overview as a chronological file listing within a document management system. The task of the reviewer was to review each document in turn and, where appropriate, create a record of any event of potential significance to the investigation. For example, this might be a meeting proposed by email between one protagonist and another.

An appreciation of which events held significant to the

| Date | Time | Event/Document | People Involved/ Author/Recipient | Evidence / File Reference |
|---|---|---|---|---|
| 8th Nov | 7.45 | {company A} Meeting in {country A} (time is {person B} flight departure from {location A} to {location B}) with return to {location A} for 12.55 on 9th Nov. {person I} to pick up {person B} at Airport | {person I}, {person B} and {person H} in {location C} | Email between {person I}, {person I}, {person H} and {person F}/Doc ID 169246 |

Figure 2. An anonymised event entry from
one of the chronologies

investigation (and hence what to record) evolved over time as the investigators' understanding developed and they reviewed their interests. What we focus on here, however, is what happens when an investigator first discovers information about a potentially significant event. The information contained in the message acts as a cue to the investigator about something that should be recorded. However, they are also aware that the information they have found is not the complete picture. For a meeting, the investigators we interviewed described a number of things they might like to know such as where and when it took place, who attended, what was discussed and what were the outcomes? Some or all of this information might be missing from the initial cue, and may be found distributed across a number of other messages. In addition, the initial lead may have been misleading: there may have been a change in plans or the meeting may not have actually taken place at all.

Following Klein et al's [2006] model of the process of sensemaking, we think of the investigator's concept of an event as an instance of a 'frame'. Frames are structures that we impose on the world in the process of understanding it. They are triggered by cues and act as plausible interpretations of those cues. A significant property of a frame which is important here is that they extend beyond the data from which they were cued. The ability to interpret things in this way is a fundamental human capacity. But as a consequence of this they can be wrong, perhaps as a result of a misleading cue.

Returning to the investigation, following the initial discovery or cue, a question arises about how to proceed. The initial document provides an important lead, prior to which the investigator knew nothing of the event. We may say that at this point they have a theory that a significant meeting took place. But this theory gives rise to a need for further information, specifically in order to address the need of elaborating and validating the interpretation.

In this situation the investigators we interviewed described two strategies. Given potential difficulties in locating other documents about a given event, one strategy was simply to record the event as a conjecture and move on. Investigators would raise an event record in a chronology (marked as a conjecture) and continue reviewing documents as before in the hope that they, or someone else, would come across further relevant information later. The second strategy was to construct further keyword and/or date delimited queries designed to re-filter the collection in a way that might bring relevant documents to the surface.

Whilst the second strategy offers continuity to the investigator in terms of focus by supporting a single chain of thought, it is also non-trivial. The investigator sees the situation under investigation in terms of events, whilst the system they are using sees the situation in terms of documents, textual tokens and metadata. Consequently, the investigator must translate their question (of all documents relevant to a particular event) into something understood by the system—referred to more generally as a 'compromised need' [Taylor, 68]. This can require some cognitive effort and result in what is at best an approximation.

**Reflections on design**
This example suggests a general principle which we can apply to such problems. That is—where a user is making sense of information through the application of a particular type of frame (or frames), leverage can be obtained by linking information around possible instances of that frame (or frames) in the data. Of course, there may be a number of types of frame that are important to an investigator. Other frames we identified from our interviews included entire business activities (such as contracts), particular time periods surrounding major events within those activities, and protagonists or potential protagonists under investigation. In many of these cases, information discovered within the collection cued the investigators to their existence, raising them as foci for further investigation (the investigators started from an almost entirely blank slate). And in all cases further information was distributed across the document collection.

We note here that frames that were of significance to the investigators were reflected in the way the investigators structured the knowledge representations they generated i.e. chronologies structured in terms of individual events and individual chronologies dedicated to information about specific business activities and people. Hence, in understanding how the investigator seeks to translate information at transition t4, we may need to look no further than the representational form they seek to create.

Construing the investigation from the perspective of the investigator, then, it is a question of how frames are cued and how this leads to the need for their elaboration and/or validation. From this, we can ask the question of the extent to which information retrieval technologies support this thought process.

We have addressed this question to some extent in relation to the need for the investigator to translate their needs into terms understood by the system—terms which are characteristically low-level in their characterization of document content.

We may ask what other kinds of technologies might be developed which may be more helpful. The question here is one of raising the bar in terms of system intelligence in order to achieve the potential for aggregating documents in terms which more closely approximate to the concepts of the investigator. In this way, transformations performed earlier in the process (i.e. search) would organize the data in a way better adapted for subsequent work.

A number of possibilities exist here. First, systems that offer representations of email documents in terms of subject threads may offer some advantage. Analysis of the Enron collection, however, has suggested that the average length

of an email thread (in organizations at least) is typically quite short. Also systems that are capable of semantically clustering documents (e.g. Attenex) may be of value, depending on the extent to which emergent document clusters map to investigators' conceptual frames.

Another alternative is to use systems that perform information extraction (IE). IE system process free text and use techniques in computational linguistics in order to identify pre-defined elements of meaning [Gaisauskus & Wilks, 1998]. Jigsaw [Stasko et al], for example is an investigators tool specifically designed to graphically represent the results of information extraction over a free text collection. Elsewhere, capabilities for identifying temporal and event references in text have been demonstrated at 83% accuracy against hand-annotated data (Mani and Wilson, 2000).

## DISCUSSION AND FUTURE WORK

We believe that a promising approach to the design of more appropriate systems for e-discovery work is to structure them around the terms or concepts in which the investigators understand the subject-matter of the investigation. The significance of these concepts is that they provide the vocabulary through which investigators see the world. The sensemaking process in e-discovery can be seen as one of translating large amounts of unstructured data into representations structured in these terms. The transitions represented in figure 1 can be seen as a process of filtering and abstracting information into these terms.

The approach we have illustrated involves the identification of the sensemaker's typical frames and the operations that they want to perform on them. By providing an analysis of related information needs and the way these develop, frames provide a foundation for reasoning about the design in terms of the typical cognitive paths users follow during sensemaking.

If systems can indeed be configured around the kinds of concepts that e-discovery investigators themselves apply to data, then they are likely to provide a higher platform on which investigators can apply their own expertise in sensemaking and allow them to work to a higher conceptual level [Rasmussen et al., 1994]. The ideal is that investigators can pursue investigations with fewer interruptions imposed by constraints of the systems that they use. And by identifying documents that are relevant to emerging concepts of the investigation, there is an opportunity to reduce the very high overhead of document review.

We intend to explore these ideas further in future work. We are about to embark on a further investigation case study. Of key interest will be the way in which investigators conceptualised their problem as expressed through the way that they talk about them and the ways in which concepts are rarified in the process of generating useful representations of knowledge.

We are also planning a laboratory study in which non-lawyer participants will perform a mock investigation using a subset of the Enron email collection. Manipulations in this study will involve the presentation of a document collection according to visual indexes based around different kinds of document aggregation, including email threads, semantic clusters and event references. Our aim will be to understand the value that these provide in the process of cueing, elaborating and validating users' conceptual frames.

## REFERENCES

Attfield, S., Blandford, A. & De Gabrielle, S (2008a) Investigations within investigations a recursive framework for scalable sensemaking support. Sensemaking Workshop, ACM SIGCHI Conference 2008.

Attfield S., Fegan S. & Blandford A. (2008b) Idea generation and material consolidation: Tool use and intermediate artefacts in journalistic writing. Cognition, Technology and Work (online first).

Baron, D. (2008) UK firms report jump in spend on e-discovery systems. Computer Weekly, March 2008

Conrad, J.G. (2007) E-Discovery revisited: A broader perspective for IR researchers. DESI: Workshop on Supporting Search and Sensemaking for Electronically Stored Information in Discovery Proceedings

Gaizauskas, R. & Wilks, Y. (1998) Information extraction: Beyond document retrieval. Journal of Documentation, 54(1), 70-105.

Klein, G., Phillips, J. K., Rall, E. L. and Peluso, D. A. (2006) A Data-frame theory of sensemaking. In Expertise Out of Context: Proceedings of the Sixth International Conference on Naturalistic Decision Making (Pensacola Beach, Florida, May 15-17, 2003). Lawrence Erlbaum Associates Inc, US, 2007, 113-155.

Mani, I. and Wilson. G. (2000) Robust processing of news. In Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics (ACL2000, Hong Kong), 69-76.

Rasmussen, J. Pejtersen, A.M. & Goodstein, L.P. (1994) Cognitive Systems Engineering. New York, Wiley.

Socha, G. & Gelbmann, T. (2007) The 2006 Socha-Gelbmann electronic discovery survey report. Socha Consulting, Saint Paul, MN.

Stasko, J., Gorg, C. & Liu, Z (2008) Sensemaking across text documents: Jigsaw. Sensemaking Workshop, ACM SIGCHI Conference 2008.