



Deposited via The University of Sheffield.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/180077/>

Version: Published Version

Article:

Bartliff, Z., Kim, Y. and Hopfgartner, F. (2022) A survey on email visualisation research to address the conflict between privacy and access. *Archival Science*, 22 (3). pp. 345-366. ISSN: 1389-0166

<https://doi.org/10.1007/s10502-022-09387-2>

Reuse

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



A survey on email visualisation research to address the conflict between privacy and access

Zoe Bartliff¹ · Yunhyong Kim¹ · Frank Hopfgartner²

Accepted: 15 January 2022
© The Author(s) 2022

Abstract

Emails, much like communicative genres such as letters that predate them, are a rich source of data for researchers, but they are replete with privacy considerations. This paper explores the resulting friction between privacy concerns and email data access. Studies of email can often be centred on understanding patterns of behaviour and/or relationships between people or groups, and, as such, embody risks of disclosing private information. This is further amplified in humanities research which is concerned with the individual, their work and the circumstances that influence them. Furthermore, previous studies have expounded upon the benefits of visualisations for researching email data, a method which has been reported both as a path to addressing known concerns, as well as, introducing new concerns in privacy. The spectrum of methodologies leave archivists and curators of email data in a quandary, unable to balance accessibility with privacy. The research presented in this paper contributes a systematic approach to examining the relationship between email visualisation research and privacy. It presents a categorisation of email visualisation attributes, and a graded scale of privacy, to be used in conjunction as a framework for interrogating existing research and their associated email collections. The paper aims to instigate the first steps in concretely situating the extent to which research can take advantage of or is challenged by privacy conscious data management.

Keywords Digital archives · Emails · Privacy · Visualisations

✉ Zoe Bartliff
zoe.bartliff@glasgow.ac.uk

¹ School of Humanities, University of Glasgow, Glasgow, UK

² Information School, University of Sheffield, Sheffield, UK

Introduction

As emails increasingly become a feature of archival collections, they promise to become valuable sources for research, particularly for the humanities (Hangal et al. 2012; Schneider et al. 2019; Jaillant 2019; Noonan 2016; Prom 2011; Baker 2015). It has been noted that, so long as there is communication via correspondence, “people will be interested in knowing from whom to whom, at what time, about what topic and for what purpose the message was written” (Zhang 2015). Emails, by their nature, support this kind of enquiry, comprising a log of personal and/or professional activity, relationships and networks. They therefore can be considered a unique window through which to view a person’s life, much in the way that letters were before them (Kennedy 2014).

Accessing email, however, can be problematic because it is usually high in volume, inclusive of a lot of ‘noise’ (e.g. spam, repeated content in complex threads, mark-up, attachments, linguistic variety) and can contain data that may impact upon privacy. The first two of these roadblocks can be managed by finding innovative methods for extracting and exploring the data which are otherwise “[l]ocked up within standard email interfaces” (Perer and Smith 2006). A great range of methods have been developed over the years to support the processing of email data, whether by data cleaning (Tang et al. 2005), collation/classification (Mujtaba et al. 2017), analysis (Borden and Baron 2016) or visualisation such as graphical or other visual representations of data or information (Repke and Krestel 2018). Visualisations, in particular, have regularly been employed as a method for exploring email collections and supporting the identification of high level patterns within a data-set (Louis and Engelbrecht 2011; Kaczmarek and West 2018; Moss et al. 2018; Stadlinger and Dewald 2017). This practice is slowly migrating into the archival setting (Josh Schneider et al. 2017) with the intention of supporting a holistic, creative, perhaps even ‘playful’ (Hendery and Burrell 2019) research.

The third roadblock to access, that of privacy, poses a barrier to these methods of data exploration. The need for privacy management is a pressing concern that needs to be resolved given that breaches can have serious, enduring and harmful consequences for the archive donor, the researcher, the archive itself or, indeed, subjects of the email content who have little agency in its management (Eder et al. 2020). The difficulties arising in response to these types of challenges has led many to err on the side of caution, closing archives rather than risk harm, whereas others advocate “a less risk-adverse attitude to the release of information” in order to facilitate useful research (Moss and Gollins 2017). This paper sits as a counterweight to these dichotomous sentiments and intends to explore a more graded notion of privacy for understanding and managing the risks associated with researching emails. This introduces a nuanced approach to balancing risk whilst still providing a degree of access for researchers that is appropriate to the nature of the material and research to be conducted.

The extensive use of visualisations as a data exploration tool, coupled with the benefits they have shown for handling sensitive information (Avraam et al.

2021), suggests them as a good starting point for developing privacy aware strategies. In fact, visualisations have already been used to address concerns relating to privacy, by allowing for filters to be applied to mediate a researcher's access to data, for example, to support anonymisation or decontextualisation of content (e.g. Chou et al. 2019). These existing approaches, however, are narrowly focused and so do not fully account for the broader and more complex considerations of privacy management, which can encompass legal, ethical, cultural and personal dimensions. Similarly, there has been little discussion with regards to the comparative efficacy of different privacy aware visualisation strategies or the extent to which they might restrict or facilitate access to data for researchers. In fact, in some cases, visual approaches to accessing data have been shown to increase the chances of 'serendipitous discovery' (Hangal et al. 2011) that can be detrimental to privacy management due to their unpredictable nature. The research in this paper assesses how privacy management has been applied to email collections so far with respect to email visualisation research. It contributes a systematic approach to examining:

- The research context and objectives of visualisation-based approaches to email research.
- The levels of protection that can be offered by privacy aware strategies.
- The potential impact of privacy management on the usefulness of the collection for humanities research.

More precisely, this paper contributes a systematic analysis of the key objectives of email research, the relevant categories and research themes that relate to humanities research and beyond. The analysis is focused on areas that reflect the increased interest in the use of visualisations in email research. This categorisation is followed by a proposed scale of privacy that is inspired by disparate approaches to privacy protection evident in wider literature. Finally, these two contributions are brought together to facilitate a better understanding of privacy management as it is currently applied in email visualisation research and frames this within the context of research objectives. The developed privacy scale serves as the first steps towards optimising privacy management to be balanced with researcher needs for access. Whilst the focus of the paper addresses the broader domain of email research, the arising issues are a particular challenge for the humanities given the predilection for identifying and exploring more humanistic, therefore personal, narratives within data.

Methodology

Finding relevant research

The research presented in this paper was conducted in line with the metaanalysis method for "balanced and critical literature reviews" that "serve a crucial function... searching through mountains of potentially contradictory research to uncover nuggets of knowledge that lie buried underneath" (Stanley 2001). The overarching aim

therefore is to aggregate the key themes and findings, as well as to identify recurrent patterns in approach and outcome (Levitt 2018). Through doing so, this paper will offer a consolidated perspective on privacy management for email exploration, using this to identify where gaps in knowledge exist and making recommendations to fill them.

To achieve a systematic review of literature relating to the analysis of email collections in relation to privacy, a keyword search of the Google Scholar database was conducted. This database has been shown to be particularly representative of interdisciplinary research (Halevi et al. 2017). This was complemented by searches on ACM digital library to validate coverage. The review was conducted based on a varying combination of the keywords: “email~analysis” or “email~visualis(z)ation” or “email collection” AND archives, digital archives, or humanities AND privacy preserving, privacy management, privacy protection or scales of privacy. The phrases involving ‘email’ were restricted to pairings with analysis, visualisation, collection or synonyms of these as alone ‘email’ returned irrelevant results, for example, due to the inclusion of author emails in papers. No limit was imposed on the date range and all sources included in the study were written in English.

This study is focused towards user (archivist and/or end-user) interaction with email data, facilitated by visualisations or visual enhancement of content (highlighting, keyword summation etc.). In the light of this, those articles that do not involve visualisations as a user facing element for facilitating direct analysis of email content were excluded from the data-set, although some have been included in the general discussion. This constitutes a large number of studies that discuss the theory, platform and/or infrastructure behind the research, creation of visualisations, and privacy management of email collections (e.g. Eder et al. 2020) as well as, those papers keyed towards policy and preservation of email.

Following the data gathering and application of exclusion criteria, thirty-nine papers were identified for deeper analysis. Of these, several papers discussed more than one independent email visualisation, each designed with a different research purpose in mind. These were treated individually, resulting in a total of sixty-nine data points. The papers were reviewed and a summary table was created noting the data source used for the paper, the declared purpose of the study, the methods that were employed, whether there was a discussion of privacy, the visualisation techniques that were used and any noted implications of the research. As shown in Fig. 1 the studies that met the criteria for inclusion are well distributed and cover a span of almost 20 years, suggesting that research concerning email analysis using visualisation has been and continues to be a consistent feature of wider scholarship.

Research paper categorisation

The first area of analysis will be to categorise the common themes of interrogation that are applied to email data and the visual methods that supplement them. This will help to identify where likely privacy concerns may arise.

Previous attempts to categorise how visualisations are used to facilitate email research have been roughly cohesive but lacking in consistent vocabulary. For

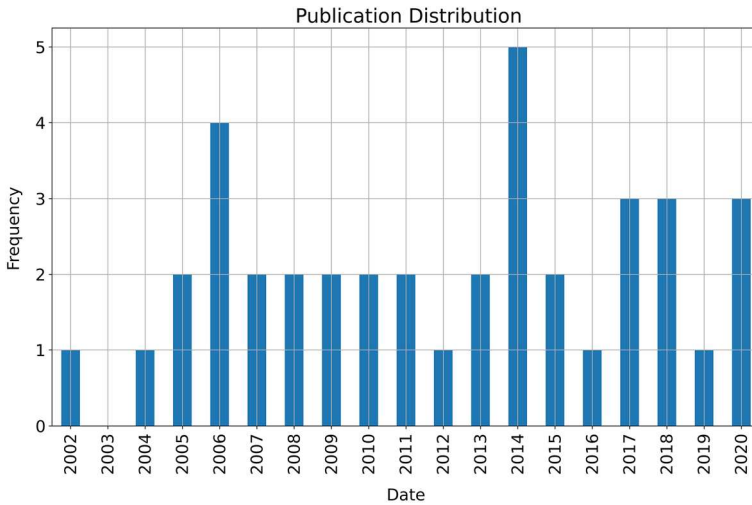


Fig. 1 The number of articles included in the study (y-axis) by year of publication (x-axis)

example, some have identified thread-based visualizations, social-network visualisations, temporal visualisations and contact-based visualisations (e.g. Viégas, Golder and Donath 2006) but some only include the first three of these (e.g. Lu et al. 2019). Others highlight the potential of visualisations to reveal personal information (e.g. Bergstrom 2011) discussing social network visualisation, thread-based visualisation and those that focus on content. One paper (Thanh Tung 2014) finds three categories, namely thread, context and contact, whereas another (Luo et al. 2014) adopts a simple if poignant dyadic categorisation of those that explore people and those that explore events (cf. Luo et al. 2014).

A synthesis of these models revealed that a broad dyadic division could encompass all the identified research focal points. Two branches were organised to reflect an interest in ‘people’ (network structure and patterns that characterise the relationship between people in the network) and, adapting the term slightly from the existing dyadic model (Luo et al. 2014), a focus on email ‘usage’ (understanding content of emails and pattern that characterise the behaviour of people that govern how these are used). These broad categories were then subdivided into five more focused areas of interest to harmonise and homogenise the other pre-existing classifications without sacrificing the nuances found in each. These five categories are defined in the ‘Approaches for gaining insights from emails’ section as a first finding of this meta-analysis but in brief they are social network analysis (SNA), patterns of relationships, patterns of behaviour, content analysis and topic identification (Fig. 2). However, these five categories are not meant to be mutually exclusive; they can overlap and complement each other to attain specific research aims. Instead, they are useful for identifying the type of information being sought and, therefore, the likely methods used and privacy concerns that may arise.

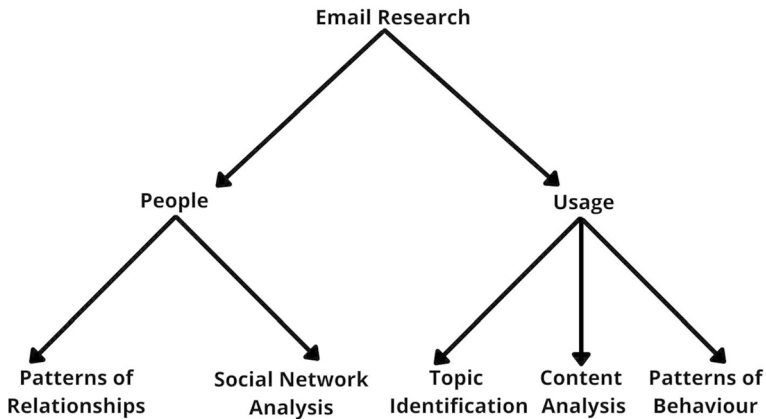


Fig. 2 A visual representation of the dyadic categorisation of email analysis and the more nuanced categories that sit within this

Analysing privacy awareness

The second stage of this paper is a survey of privacy aware strategies, which, based on analysis, have been grouped into a privacy consciousness scale to reflect the extent of privacy management they offer. Throughout the review of literature, including those works that were excluded from the meta-analysis (as described earlier), it became clear that a range of privacy management techniques have been employed and proposed with relation to email analysis. The efficacy of these techniques, however, seems highly dependent on the skill and knowledge of the person applying them and the contextual knowledge possessed by the user (Elliot et al. 2018; Bampoulidis et al. 2020). These methods were collated by their effect on data, to represent five broad types of processes that can be applied to data, each reflecting a different level of privacy consciousness. This privacy consciousness scale (PrivCon) is represented in the ‘Privacy in email analysis’ section as a second finding through the attribution of numerical values 0–4. Such a scale of privacy does not, to the writers’ knowledge exist with regards to text-based material. As such, the scale has been inspired by the privacy protection approaches that have been applied to images (Padilla-López et al. 2015) with established techniques utilised in relation to emails mapped to the equivalent effect of privacy management.

The final stage of this paper results from this application of the scale and is a discussion addressing the extent of privacy management commonly used in relation to visualisation-based email analysis. The research in the metaanalysis has been mapped to the PrivCon scale to demonstrate how it can be used to reveal characteristics of privacy awareness in the research area. The intention of this is to highlight any gaps that exist and note opportunities for the development of a

privacy aware strategy that maximises researcher access and minimises risk to privacy (Emam and Arbuckle 2013).

Approaches for gaining insights from emails

People

The first category of this branch, **social network analysis (SNA)**, is an established and prolific area of research often applied to email that is consistently mediated with visualisations (e.g. Chen and Yang 2010; Creamer et al. 2007; Golbeck et al. 2018; Heer 2005; Heibi 2017; Humphreys et al. 2008; MacLean et al. 2011; Magalingam, Rao, and Davis 2014; Mesarina et al. 2009; Mishra 2008; Park et al. 2015; Smilkov 2014; Stadlinger and Dewald 2017; Straub 2016; Tsetini 2015; Štorga et al. 2013; Wen et al. 2020; Wise 2014; Zhang 2018). This kind of analysis ‘studies the connections between particular units or groups’ and an individual’s role within that network (Boryczka et al. 2016, cf. Tabassum et al. 2018). The next category is **patterns of relationships (PoR)** which addresses broader attempts to study interaction between people through, for example, the ebb and flow of contact or how communication varies between communities (e.g. Lu et al. 2019; Luo et al. 2014; Perer et al. 2006; Save 2020).

All eighteen of the items identified as social network analysis SNA in the meta-analysis employ a variation of a network graph (Fig. 3A) as a way to visualise the data. These can produce different layouts (e.g. random, force directed or tree) sometimes including hierarchical (Fig. 3E; Creamer et al. 2007; MacLean et al. 2011; Straub 2016) and/or arc structures (Mesarina et al. 2009). Some are combined with other visualisation methods like timelines or bubble plots (Fig. 3G; Zhang 2018; Štorga et al. 2013). The visualisations for studying PoR are much broader in range, with comparable numbers belonging to multiple types in Fig. 3 and beyond. For example, bar graphs and histograms (Fig. 3D; Bulkley 2006; Mondal et al. 2017; Perer et al. 2006; Stadlinger and Dewald 2017; Straub 2016) and scatter or bubble plots (Fig. 3f, g; Bulkley 2006; Perer and Smith 2006; Viégas et al. 2006). There are also iconographic representations of people or types of conversation (Mesarina et al. 2009; Mandic and Kerne 2004; Viégas 2005; Save 2020), as well as, mountains and rose petals visualisations (Viégas 2005; Lu et al. 2019).

Usage

The first category in this branch, **topic identification (Topic ID)**, focuses on the identification of spam or phishing with a smaller proportion of work aimed towards ‘Multi Folder Categorization’. There are only very few attempting other classifications of content (Mujtaba et al. 2017). There are also studies aimed at management of email inboxes (e.g. Schreck 2010) rather than the discovery of information by external stakeholders, although can offer integral insights into the collection as a

whole (e.g. Kalyan and Chandrasekaran 2007). Topic identification can support archive management and research by offering a rapid method of identifying superfluous content (e.g. Schreck 2010), and differentiating potentially risky content, as well as, give an overview of the nature and underlying themes of the collection. Only two papers were in this category with one (Schreck 2010) using heat-maps to represent the likelihood of content falling into the spam or ham, and, the other (e.g. Kalyan and Chandrasekaran 2007) using line graphs (Fig. 3B) to track user behaviour and identify anomalous content.

Second category of the branch, **email content analysis (Contents)** (e.g. Weisgerber and Butler 2009; Luo et al. 2014; Butavicius et al. 2012; Thanh Tung 2014; Lu et al. 2019) aims to help users navigate a collection and withdraw meaningful data whether as a search or summary mechanism. This area draws heavily from the field of text analysis, with visualisations aimed at text summarisation, clustering of content and identification of themes or sentiments and is therefore the most varied in approach. Visualisations include wordclouds (Fig. 3H; Lu et al. 2019; Wen et al. 2020; Weisgerber and Butler 2009) and/or plotting keyword, sentiment or topic on timelines (Luo et al. 2014; Hangal et al. 2011; King 2013; Lu et al. 2019; Thanh Tung 2014; Viégas et al., 2006). Other examples include an analysis of content clustering (Fig. 3F; Butavicius et al. 2012), a pie chart (Fig. 3C; Cadman, MacDonald and Soomai 2020) and a combined SNA/content analysis network graph (Fig. 3A; Golbeck et al. 2018).

The third category of this branch, **patterns of behaviour (PoB)**, focuses on exploring and explaining how an individual engages with their email collection (e.g. Kalyan and Chandrasekaran 2007; Luo et al. 2014; Lu et al. 2019; Bulkley 2006; Heibi 2017). As a 'distant' approach, PoB analysis involves the exploration of contextual data that is embedded into emails. This can, amongst other things, be timeline, location, or flag (e.g. importance, spam, attachment, read/unread) based and often involves frequency and dispersion analysis of these areas. Like the PoR category, PoB involves a range of approaches to visualisation. There are a higher number of scatter/bubble plots (Fig. 3F; Bulkley 2006; Lu et al. 2019; Magalingam et al. 2014; Perer and Smith 2006) and line plots (Fig. 3B; Lu et al. 2019; Kalyan and Chandrasekaran 2007; Bulkley 2006; Mondal, Shukla, and Lodha 2017; Perer and Smith 2006; Stadlinger and Dewald 2017; Straub 2016) compared to PoR and also two studies that combine analysis of behaviour with SNA (Heibi 2017; Smilkov 2014).

Interfaces

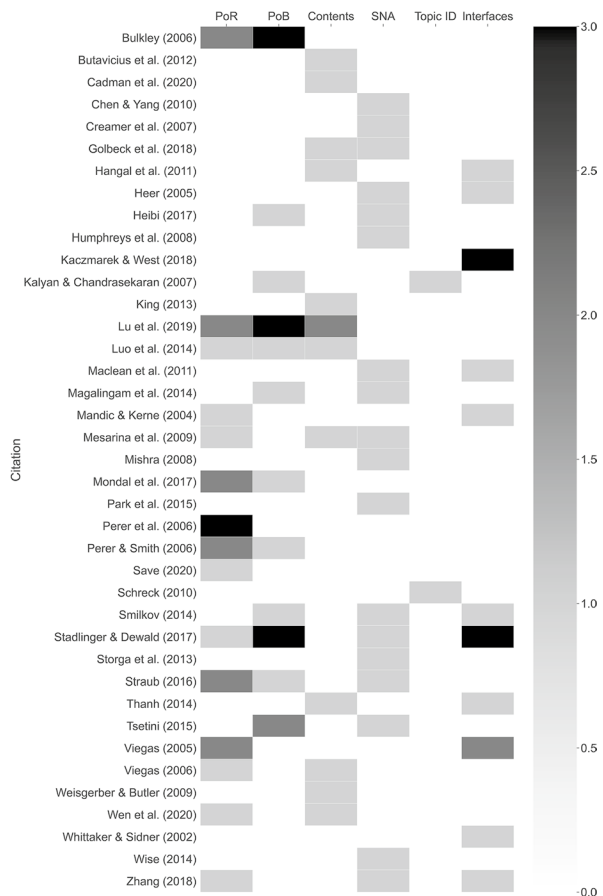
In addition to the two branches of categories identified above, there are numerous instances where visualisation methods are combined into an interface. Interfaces are not a category of interest, per se, but rather a tool to overcome the technical requirements involved in the employment of data analysis and visualisation methodologies (Hutchinson 2020) that bridges across the categories (e.g. Kaczmarek and West 2018; Heer 2005; Thanh Tung 2014; Zhang 2018; MacLean et al. 2011; Whitaker et al. 2002; Hangal, Lam, and Heer 2011; Viégas 2005; Mandic and Kerne

2004; Smilkov 2014). *Memories Using Email* (Hangal et al., 2011) and *Email: Process, Appraise, Discover, Deliver* (ePADD) are two particularly popular interfaces that have been designed for archival use and adapted for user exploration of email data. In fact, ePADD has become one of the more widely used interfaces for archival institutions to manage their notable collections such as that of Wendy Cope at The British Library or Salman Rushdie at Emory University (Schneider et al. 2019).

Analysis

As seen in Fig. 4, these categories of email visualisation are not evenly represented within the literature. The focus is heavily weighed towards understanding relationships (PoR then SNA), with PoB third in popularity. Word frequency analysis, conducted with stopwords removed, shows that words such as ‘social’ (62 tokens), ‘network(s)’ (44), ‘communication’ (30), ‘people(s)’ (23) and ‘relationship(s)’ (27) are the most prolific in the abstracts of the data-set, with ‘social’ surpassed only by

Fig. 4 A heat-map demonstrating the categorisation of the meta-analysis papers (cf. the ‘Approaches for gaining insights from emails’ section). The colour intensity of the block indicates the number of visualisations in that paper that fall within each category, namely patterns of relationships (PoR), patterns of behaviour (PoB), contents, social network analysis (SNA), topic identification (Topic ID) and interfaces



'email' (140). In contrast, those words that would appear to be concerned with usage of emails (e.g. content (16), time (14), history (12), organisational (14)) are mid-dling within the frequency list.

The relationship between visualisations and categories is not clear cut with certain visualisations stretching across more than one category. For example, the Email-Map project clusters data points both by similarity of content (PoB) and closeness of relationships (PoR) to determine 'how things evolved and progressed in the past' (Luo et al. 2014). These multi-functional visualisations have been advocated for SNA as a means to enrich the data and gain deeper insights into the collection (e.g. Van den Elzen and Van Wijk 2014; Wen et al. 2020). These combined approaches offer multiple perspectives on the data, broadening the information that might be explored and extracted.

Privacy in email analysis

The scales presented in this section are not intended to be mutually exclusive, but rather meant to approximate increasing levels of protection intended by the techniques associated with them. In terms of the privacy protection they offer, there may be some conceptual overlap perceived between them depending on the skill with which the methods have been executed. When conducted with a high level of skill and a discerning eye, anonymisation, a technique that falls within PrivCon1, may well be a highly effective method of privacy management. The principle applied in the privacy scales developed here is to increasingly minimise the risk of identity, attribute, inference disclosure (Arbuckle and El Emam 2020) by introducing techniques to obscure explicit knowledge regarding individual data points. Equally, however, the task can be poorly achieved, making the reconstruction of these aspects not only possible but likely. Another factor, one which cannot be accounted for in a standardised way is the skill and knowledge base of the user (Elliot et al. 2018; Bampoulidis et al. 2020). Even with the more effective applications of privacy protection, there is a risk that a user may be able to re-identify individuals and/or make inferences about groups from surrounding data points (Arbuckle and El Emam 2020). What the more advanced techniques on the PrivCon scales below offer is a higher baseline of privacy management as you go from 0 to 3.

PrivCon0

PrivCon0 is reflective of data that is either accessed in its native environment or has been supplied to researchers in a state that mirrors this, with minimal intervention. Of the papers included in this study, 13 of the 39 papers sit in this level. This approach offers complete access to the email data and so offers no protection with regards to privacy. However, the papers that are characterised by PrivCon0 are, on the whole, either designed for forensic investigation (e.g. Stadlinger and Dewald 2017) in nature or self-reflective, intended to allow an individual to explore their

own email data (e.g. Perer and Smith 2006; Hangal et al. 2011). As such, whilst such a data-set would be the ideal for a humanities researcher, the likelihood of archives providing complete and unfettered access to an email collection is unlikely.

PrivCon1

The **PrivCon1** group includes situations whereby the data have been altered or removed in order to obscure the identity of individuals contained within. When compared to a scale of image protection (Padilla-López et al. 2015), PrivCon1 is reminiscent of the practice of blurring or pixelating an image because the overall content is comprehensible and recognisable but certain details or features are obscured and/or removed. This group is the most popular of the approaches comprising 30 out of 39 papers. Specific methods that fall into this approach include, anonymisation, pseudonymisation and the redaction or removal of data. Anonymisation or pseudonymisation, in particular, appears to be a popular method for protecting privacy at the point of publication (e.g. Hangal et al. 2011; Humphreys et al. 2008; Smilkov 2014; Zhang 2018). Despite the use of this privacy strategy, these papers have been classified as PrivCon0 because that is a truer reflection of the levels of privacy experienced by the researcher, who is the focus of the research aims. Many of the studies included in the meta-analysis engage with PrivCon1 data-sets because the organisation who donated the material implemented their own privacy policies to alter data prior to donation. A prime example of this is the Enron data-set. Whilst it was made open access by the Federal Energy Regulatory Commission, making it ideal for developing email exploration tools, the Enron data-set has had a number of personal emails removed ahead of publication (Shetty and Adibi 2004). However, within an archival context and, in particular with reference to personal archives, it is likely that the archive will need to maintain the full collection but restrict access to certain aspects of the data that contain personal or private data that cannot be disclosed. In this situation, it is a matter of judgement on the part of the archivist to interpret the content and identify potential risks to privacy, perhaps without prior or sufficient knowledge of the individuals whose privacy may be at risk or the context in which they are involved in the collection.

On the whole, PrivCon1 maintains a high level of detail, allowing researchers to access the majority of the data and to apply familiar, often manual methodologies. This high level of access is likely to have led to the prevalence of meta-analysis studies that fall into this level. However, the issue is that, because much of the data remains unaltered, privacy management of this kind can be bypassed to reveal private information. Pseudonymisation and anonymisation in particular have been shown to be an ineffective, 'naive' (Task 2015) method of privacy preservation, particularly with reference to social network data (Ying et al., 2009; Backstrom et al. 2007; Ma et al. 2017; Wang et al. 2018; Fu et al. 2020; Ying and Wu 2009). Additionally, the reliance upon the skill and knowledge of the person applying the privacy management strategy means that 'human judgement becomes a strong factor' (Kaczmarek and West 2018) which is another key reason

why the level of protection offered by this approach is so low. Not only is a consistent application hard to achieve, but an archivist or curator is unlikely to be a subject specialist and so may not be able to identify what data could become sensitive information.

PrivCon2

PrivCon2 involves the grouping or amalgamation of data to the point that individuals become 'lost in the crowd', minimising the risk that details might be identified. When compared to the image-based scale of privacy (cf. Padilla-López et al. 2015), PrivCon2 equates to the embossing or silhouetting—supporting privacy by unifying certain features. Only six papers within the dataset utilise this approach, four of which employ data which also display characteristics of PrivCon1. Of these four, two use the Enron data-set which has had private messages removed at the request of the email creator (Tsetini 2015; Heibi 2017) and another an undisclosed commercial data-set (Kalyan and Chandrasekaran 2007). The final is a borderline case as they avoided analysis of 'sensitive' content such as attachments, links or images, although it should be noted that researchers were still allowed full access to the email body and metadata (Mondal et al. 2017).

On the whole, PrivCon2 provides a holistic perspective on the data, revealing overarching patterns at the expense of the detail. Where this scale of management has been applied in research, it has been to derive profiles of activity, for instance the average rate of activity or the flow and/or type of contact across the data-set. Whilst lacking in the specifics to be found in PrivCon0 & 1, PrivCon2 approaches might provide to the humanities researcher unique insights that might otherwise have been obscured by more focused methods (Bartliff et al. 2020).

PrivCon3

The **PrivCon3** form of privacy management involves shifting the data through the use of an algorithm, statistical model or encryption, in a way that maintains the statistical characteristics of the data-set but the detail does not consistently reflect the original. To refer once more to the image-based scale of privacy (cf. Padilla-López et al. 2015), this approach would equate to the skeleton or avatar methods, where the data are only a representation of its original form. Whilst none of the studies represented in this study employ this technique to manage the privacy of their data studies into differential privacy (one method within PrivCon3) have been particularly common in SNA studies. This approach is typified by providing detailed and/or holistic perspectives. Whilst the researcher can rest assured that overall the representation respects the patterns found within the original data-set, they cannot draw reliable observations about specific data-points. It should be noted that there is one example (Thanh Tung 2014) where researchers encrypt their data, but this refers only to data storage in their interface. From the perspective of the user, the email data can be considered open.

PrivCon4

PrivCon4, the antithesis of PrivCon0, refers to the practice of keeping an archive ‘dark’, inaccessible to researchers without special, often on site, permissions. This approach is used because ‘[t]he most intuitive way to preserve privacy is to stop the propagation’ (Li et al. 2020). It has been widely used to protect sensitive archival content when full review is impractical, as is often the case with digital material (Moss and Gollins 2017), especially emails. The issue with this approach is that it restricts the value of the data and limits the extent to which researchers can engage with the collection. This is especially the case with digital data whereby the relevance of the data and the ability to access the data could be a relatively small window (Moss and Gollins 2017). This level of privacy is therefore the least desirable, but unfortunately has become the norm for email data (Jaillant 2019).

Analysis

The range of approaches to email analysis and visualisation evident within the study can be seen in Fig. 5. Although wider scholarship displays privacy management strategies across the PrivCon scale, the studies that were included in the meta-analysis show a clear preference for the more open ends of the scale. It is partly for this reason that the scale was determined through a wider exploration of email analysis work that included more theoretical papers and those that addressed matters of infrastructure.

In Fig. 6, the PrivCon scales are represented across the categories of research defined above. The figure demonstrates that the two categories concerned with people (PoR and SNA) are quite similar in distribution, prioritising PrivCon1, with a good proportion of PrivCon0 and two papers each that utilise PrivCon2 approaches. The Topic ID category would appear to be the most privacy aware of the set, only demonstrating usage of PrivCon1 and 2 approaches. Those studies that are interested in the content of emails have a higher proportion of PrivCon0 papers, suggesting the perceived need for increased levels of access for this kind of research. All six of the papers that are at PrivCon2 fall within the PoB category as well as one or more of the other groupings. This is perhaps due to the more holistic nature of this kind of research.

Discussion

Despite the readily acknowledged importance of privacy, both legal and ethical (Buchanan et al. 2007; Weisband and Reinig 1995; Schneider et al. 2019; Noonan 2016; Baker 2015; Biber and Luker 2014), only half of the papers included in this review (19 out of 39) explicitly discuss privacy management. Where privacy is discussed, it is rarely the focus of the study, but rather presented as an obstacle preventing useful research and leading to email archives being withheld from the public (Biber and Luker 2014; Jaillant 2019; Moss and Gollins 2017; Schneider et al.

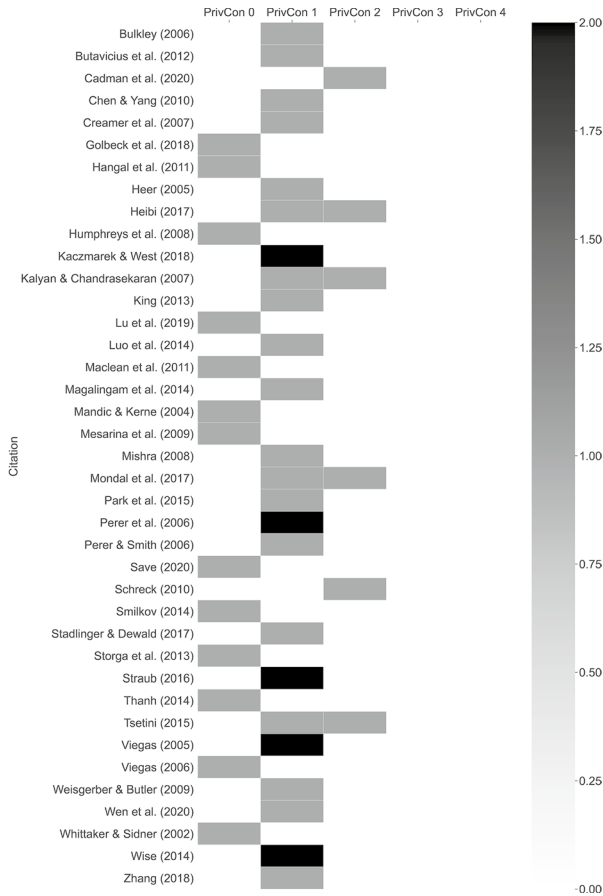


Fig. 5 A heat-map demonstrating the distribution of the reviewed papers across the PrivCon scale. The colour intensity of the block indicates the number of papers that fall within each category, namely patterns of relationships (PoR), patterns of behaviour (PoB), contents, social network analysis (SNA), topic identification (Topic ID) and interfaces

2019). Indeed, when considering email data from the perspective of humanities researchers, whose standard methodologies involve the close and usually manual examination of data, the scale of privacy may well be considered inversely related to the degree of useful access. The responsibility of privacy management strategies is therefore to balance the need for privacy against that for making the data useful, a concept that has elsewhere been referred to as the “Goldilocks principle” (Emam and Arbuckle 2013).

When contrasting the scale of privacy for image data (Padilla-López et al. 2015) to that provided here for email data, the adage that ‘a picture is worth a thousand words’ seems particularly apt. Image data, whilst complex in its own way, is actually quite uniform in relation to privacy—the areas that require protection are often quite clear and, with the right algorithm, privacy management may be uniformly



Fig. 6 A bar chart showing the application of different privacy aware strategies (y-axis) and how these are distributed across the categorisations of research interests (x-axis)

applied (e.g. a person's physicality, license plates). Textual data are typified by its great variety, requiring a thousand words where there exists only one image. Even data that may be considered uniform (e.g. names, addresses, telephone number) can be presented in multiple ways and in a range of languages, scripts and/or codes. This complexity makes the appraisal of textual data a more time consuming and resource intensive process (Schneider et al. 2019).

Email, whilst primarily textual in nature, can be multimedia in content, providing further challenges. This complexity makes the usage of privacy management strategies quite difficult and/or technically challenging, leading to the limited range of data seen in this study. When testing the efficacy of a visualisation technique, a popular approach (9 out of 39 studies) is to allow participants to apply the technique to their own email collection, therefore eliminating the need for privacy considerations. In other situations, the data are drawn from openly available, reviewed, and partially redacted collections such as the Enron data-set (8 out of 39). Some studies have engaged with specific organisations and received partial access or access through a proxy (e.g. Wen et al. 2020) to limit researcher access to potentially sensitive email data. Each of these data-sets lend themselves to PrivCon levels 0 and 1 because the data controller has extensive privileges over what (if anything) is made available to researchers as well as an elevated knowledge base for determining where private emails might be found. For example, in the case of the Enron data-set, the employees were given the opportunity to identify sensitive data within their own subset of the corpus. This is not to say that such data-sets are without risk, as can be seen from the enormous range of personal insights revealed by analysis of the Enron

corpus. Rather, there is a degree of confidence that the risks associated with these data-sets is low and reasonably well understood.

Personal email collections (other than one's own) are much less common within the literature for the exact opposite reason. The level of risk is often unknown meaning that the custodians must err on the side of caution. In fact, only two studies contained in this meta-analysis engage with personal archives (e.g. Perer et al. 2006; Hangal et al. 2011). In the former situation the owner of the archive is a co-author in the paper which is acknowledged to be 'not a luxury we expect most historians and social scientists to have'. This lack of variety in the data is a blind-spot that may be present in part due to the unavailability of open archives containing personal email data-sets and in part due to the higher privacy risks associated with such a collection (Schneider et al. 2019; Moss and Gollins 2017; Borden and Baron 2016).

It has been stated that '[h]istoric emails can only be valuable if you are actually able to access them. However, while information retrieval capabilities are significantly advanced, they typically require a user to have at least a rough idea of the information they are trying to retrieve' (Carpenter et al. 2012). For instance, the importance of searching, whether for key terms or themes, has been highlighted as a method to promote 'faster, more selective reading' for business historians (Nix and Decker 2021). Similarly, there has been a call for archivists to 'consider very different types of access' in line with researcher needs (Langdon 2016). Visual analytics are designed to promote exploratory search and browsing behaviour from the user (Hangal et al. 2015; Borden and Baron 2016) supporting the extraction of knowledge on a variety of levels both holistic and specific and bypassing the limitations of search centred methods of exploration (Winters and Prescott 2019). Such exploratory approaches 'capitalise on the characteristics of digital sources' (Nix and Decker 2021) and embrace the flexibility of the digital format. However, in terms of privacy, this flexibility and efficacy of presentation is actually more concerning than if a user were required to follow a single stream of inquiry. Visualisations are often lauded for their ability to reveal previously unknown patterns and relationships within data, but it is this feature that could jeopardise privacy if incorrectly managed. The issue may then not be one of custodians of data adapting and compromising their "risk-adverse attitude" (Moss and Gollins 2017), but instead of humanities researchers supplementing their practice to 'embrace mixed approaches common in social science' (Jaillant 2019) in order to capitalise on the information available to them at the appropriate scale of privacy.

Conclusions

This paper presents a systematic and comprehensive study of literature relating to the use of visualisations for the study of archival email data and the privacy implications of such work. Through this research, it has been possible to identify the key approaches utilised in the study of email that uses visualisation. These have then been examined in line with a scale of privacy that was separately derived from a wider literature search in order to theorise on usage priorities and the extent to which these marry with requirements for privacy. The meta-analysis has revealed that most

studies bypass issues of privacy through the selective choice of data. However, this does not help to mitigate the issues facing humanities researchers with regards to email access. The email collections of authors, artists, businesses etc. that represent a rich source of data for humanities researchers continue to be held in the dark due to privacy concerns that cannot be avoided. The question therefore arises as to whether the more privacy aware methods of managing and presenting data might be consistently and strategically applied to archival collections, whilst still accommodating humanities methodologies. The relatively scarce body of humanities focused email research presently in circulation makes this question quite difficult to address. For the relatively low measures of privacy management, those that fall into the Priv-Con0 or 1 categories, there is evidence to suggest that the answer is yes, but with the caveat that certain aspects of the sensitive content may be reconstructed and/or uncovered, therefore rendering the measures somewhat moot. What needs to be investigated more fully is the extent to which the more privacy aware approaches to presenting email data might find some cohesion with humanities methodologies, therefore opening the possibility for previously dark archives to be brought into the light.

Acknowledgements This work was supported by the Arts and Humanities Research Council [AH/R007012/1].

Declarations

Conflict of interest The authors declare that they have no conflict of interest.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Arbuckle L, El Emam K (2020) Building an anonymization pipeline: creating safe data. O'Reilly Media
- Avraam D, Wilson R, Butters O, Burton T, Nicolaidis C, Jones E, Boyd A, Burton P (2021) Privacy preserving data visualizations. *EPJ Data Sci* 10(1):2
- Backstrom L, Dwork C, Kleinberg J (2007) Wherefore art thou r3579x? anonymized social networks, hidden patterns, and structural steganography. In: Proceedings of the 16th international Conference on World Wide Web, pp 181–190
- Baker F (2015) E-mails to an editor: safeguarding the literary correspondence of the twenty-first century at the University of Manchester library. *New Rev Acad Librariansh* 21(2):216–224
- Bampoulidis A, Bruni A, Markopoulos I, Lupu M (2020), Practice and challenges of (de-) anonymisation for data sharing. In: International Conference on Research Challenges in Information Science', Springer, pp 515–521

- Bartliff Z, Kim Y, Hopfgartner F, Baxter, G (2020) Leveraging digital forensics and data exploration to understand the creative work of a filmmaker: a case study of Stephen Dwoskin's digital archive. *Inf Process & Manag* 57(6)
- Bergstrom AD (2011) Social mirrors: visualization as conversation feedback. PhD thesis, University of Illinois at Urbana-Champaign
- Biber K, Luker T (2014) Evidence and the archive: ethics, aesthetics, and emotion. *Aust Fem Law J* 40(1):1–14
- Borden BB, Baron JR (2016) Opening up dark digital archives through the use of analytics to identify sensitive content. In: 2016 IEEE Int Conf on Big Data (Big Data), IEEE, pp 3224–3229
- Boryczka U, Probiez B, Kozak J (2016) Automatic categorization of email into folders by ant colony decision tree and social networks. In: *Intelligent Decision Technologies 2016*, Springer, pp 71–81
- Buchanan T, Paine C, Joinson AN, Reips U-D (2007) Development of measures of online privacy concern and protection for use on the internet. *J Am Soc for Inf Sci and Technol* 58(2):157–165
- Bulkley N (2006) Email and output: communication effects on productivity. PhD thesis, University of Michigan
- Butavicius MA, Lee MD, Pincombe BM, Mullen LG, Navarro DJ, Parsons KM, McCormac A (2012) An assessment of email and spontaneous dialog visualizations. *Int J Hum-Comput Stud* 70(6):432–449
- Cadman R, MacDonald BH, Soomai SS (2020) Sharing victories: characteristics of collaborative strategies of environmental nongovernmental organizations in Canadian marine conservation. *Marine Policy* 115:103862
- Carpenter L, Jackson TW, Matthews G, Thomas D, Spencer A (2012) The role of IT in email preservation and archiving. in '18th International Conference on Automation and Computing (ICAC)', IEEE, pp 1–6
- Chen X, Yang C-Z (2010), Visualization of social networks. In: *Handbook of social network technologies and applications*, Springer, pp 585–610
- Chou J-K, Wang Y, Ma K-L (2019) Privacy preserving visualization: a study on event sequence data. in 'Comput Gr Forum', Vol 38, Wiley Online Library, pp 340–355
- Creamer G, Rowe R, Hershkop S, Stolfo SJ (2007) Segmentation and automated social hierarchy detection through email network analysis. In: *International Workshop on Social Network Mining and Analysis*, Springer, pp 40–58
- Eder E, Krieg-Holz U, Hahn U (2020) Code alltag 2.0—a pseudonymized German-language email corpus. In: *Proceedings of the 12th Language Resources and Evaluation Conference*, pp 4466–4477
- Elliot M, O'hara K, Raab C, O'Keefe C M, Mackey E, Dibben C, Gowans H, Purdam K & McCullagh K, (2018) Functional anonymisation: personal data and the data environment. *Comp Law & Secur Rev* 34(2):204–221
- Emam K, Arbuckle L (2013) Anonymizing health data: case studies and methods to get you started. O'Reilly Media
- Fu L, Zhang J, Wang S, Wu X, Wang X, Chen G (2020) De-anonymizing social networks with overlapping community structure. *IEEE/ACM Trans Netw* 28(1):360–375
- Golbeck J, Gerhard J, O'Colman F, O'Colman R (2018) Scaling up integrated structural and content-based network analysis. *Inf Syst Front* 20(6):1191–1202
- Halevi G, Moed H, Bar-Ilan J (2017) Suitability of Google scholar as a source of scientific information and as a source of data for scientific evaluation—review of the literature. *J of Informetrics* 11(3):823–834
- Hangal S, Chan P, Lam MS, Heer J (2012) Processing email archives in special collections. In: 'DH', pp 208–211
- Hangal S, Lam MS, Heer J (2011) Muse: reviving memories using email archives. In: *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology*, pp 75–84
- Hangal S, Piratla V, Manovit C, Chan P, Edwards G, Lam MS (2015) Historical research using email archives. In: *Proceedings of the 33rd Annual ACM Conference Extended Abstracts on Human Factors in Computing Systems*, pp 735–742
- Heer J (2005) Exploring Enron: visual data mining of e-mail. <https://homes.cs.washington.edu/~jheer/projects/enron/> Accessed 07 Jan 2022
- Heibi I (2017) A visual framework for graph and text analytics in email investigation. Master's thesis, University of Bologna
- Hendery R, Burrell A (2019) Playful interfaces to the archive and the embodied experience of data. *J Doc*
- Humphreys T, Leung L, Weakley A (2008) Embedding expert users in the interaction design process: a case study. *Des Stud* 29(6):603–622

- Hutchinson T (2020) Natural language processing and machine learning as practical toolsets for archival processing. *Rec Manag J* 30(2):155–174. <https://doi.org/10.1108/RMJ-09-2019-0055>
- Jaillant L (2019) After the digital revolution: working with emails and born digital records in literary and publishers' archives. *Arch Manuscr* 47(3):285–304
- Kaczmarek J, West B (2018) Email preservation at scale: preliminary findings supporting the use of predictive coding. In: 'iPRES'
- Kalyan C, Chandrasekaran K (2007) Information leak detection in financial e-mails using mail pattern analysis under partial information. In: 'AIC'07: Proceedings of the 7th Conference on 7th WSEAS International Conference on Applied Informatics and Communications', Citeseer, pp 104–109
- Kennedy R (2014) Affecting evidence: Edith Thompson's epistolary archive. *Au Fem Law J* 40(1):15–34
- King V (2013) Self-portrait with mortar board: a study of academic identity using the map, the novel and the grid. *High Edu Res Dev* 32(1):96–108
- Langdon J (2016) Describing the digital: the archival cataloguing of born digital personal papers. *Arch Rec* 37(1):37–52
- Levitt HM (2018) How to conduct a qualitative meta-analysis: tailoring methods to enhance methodological integrity. *Psychother Res* 28(3):367–378
- Li J, Hu X, Xiong P, Zhou W (2020) The dynamic privacy-preserving mechanisms for online dynamic social networks. *IEEE Trans Knowl Data Eng*
- Louis A, Engelbrecht AP (2011) Unsupervised discovery of relations for analysis of textual data. *Digit Investig* 7(3–4):154–171
- Lu Q, Zhang Q, Luo X, Fang F (2019) An email visualization system based on event analysis. In: CCF conference on comp support co-op work and soc comput'. Springer, pp 658–669
- Luo SJ, Huang LT, Chen BY, Shen HW (2014) Emailmap: visualizing event evolution and contact interaction within email archives. In: 2014 IEEE Pacific Visualization Symposium. IEEE, pp 320–324
- Ma J, Qiao Y, Hu G, Huang Y, Sangaiah AK, Zhang C, Wang Y, Zhang R (2017) De-anonymizing social networks with random forest classifier. *IEEE Access* 6:10139–10150
- MacLean D, Hangal S, Teh SK, Lam MS, Heer J (2011) Groups without tears: mining social topologies from email. In: Proceedings of the 16th International Conference on Intelligent User Interfaces, pp 83–92
- Magalingam P, Rao A, Davis S (2014) Identifying a criminal's network of trust. In: 2014 Tenth International Conference on Signal-Image Technology and Internet-Based Systems, pp 309–316
- Mandic M, Kerne A (2004) Visualizing rhythms of intimacy in email communication. *Interface Ecol Lab, Center for Digi Libr*
- Mesarina M, Jain J, Sayers C, Close T, Recker J (2009) Evaluating a personal communication tool: sidebar. In: International Conference on Human-Computer Interaction' Springer, pp 490–499
- Mishra A (2008) Workflow pattern mining using e-mail communications. State University of New York at Buffalo
- Mondal S, Shukla M, Lodha S (2017) Privacy aware temporal profiling of emails in distributed setup. In: Proceedings of the 2017 ACM conference on information and knowledge management, pp 1229–1238
- Moss M, Gollins T (2017) Our digital legacy: an archival perspective. *J Contemp Arch Stud* 4:1
- Moss M, Thomas D, Gollins T (2018) Artificial fibers—the implications of the digital for archival access. *Front Digit Humanit* 5:20
- Mujtaba G, Shuib L, Raj RG, Majeed N, Al-Garadi MA (2017) Email classification research trends: review and open issues. *IEEE Access* 5:9044–9064
- Nix A, Decker S (2021) Using digital sources: the future of business history? *Bus Hist* 1–24
- Noonan DW (2016) Email: an appraisal approach. *J Arch Organ* 13(3–4):146–151
- Padilla-López JR, Chaaoui AA, Gu F, Flórez-Revuelta F (2015) Visual privacy by context: proposal and evaluation of a level-based visualisation scheme. *Sensors* 15(6):12959–12982
- Park P, Weber I, Macy M (2015) The mesh of civilizations in the global network of digital communication. *PloS One* 10(5), e01122543
- Perer A, Shneiderman B, Oard DW (2006) Using rhythms of relationships to understand e-mail archives. *J Am Soc for Inf Sci Tech* 57(14):1936–1948
- Perer A, Smith M A (2006) Contrasting portraits of email practices: visual approaches to reflection and analysis. In: Proceedings of the working conference on Advanced visual interfaces, pp 389–395
- Prom CJ (2011) Preserving email. Digital Preservation Coalition Heslington
- Repke T, Krestel R (2018) Topic-aware network visualisation to explore large email corpora. In: EDBT/ICDT workshops, pp 104–107

- Save MV (2020) People oriented email: a social approach to email interfaces. Master's thesis, North Carolina State University
- Schneider J, Adams C, DeBauche S, Echols R, McKean C, Waugh MJ, D. (2019) Appraising, processing, and providing access to email in contemporary literary archives. *Arch Manuscr* 47(3):305–326
- Schneider J, Chan P, Edwards G, Hangal S (2017) Eppard: computational analysis software enabling screening, browsing, and access for email collections. In: *iPRES*
- Schreck T (2010) Visual-interactive analysis with self-organizing maps advances and research challenges. In: *Self-Organizing Maps, IntechOpen*
- Shetty J, Adibi J (2004) The Enron email dataset database schema and brief statistical report. *Inf Sci Int Tech Rep Univ Southern California* 4(1):120–128
- Smilkov D (2014) Understanding email communication patterns. Master's thesis, Massachusetts Institute of Technology
- Stadlinger J, Dewald A (2017) A forensic email analysis tool using dynamic visualization. *J Digi Forens Secur Law* 12(1):6
- Stanley TD (2001) Wheat from chaff: meta-analysis as quantitative literature review. *J Econ Perspect* 15(3):131–150
- Štorga M, Mostashari A, Stanković T (2013) Visualisation of the organisation knowledge structure evolution. *J Knowl Manag*
- Straub KM (2016) Data mining academic emails to model employee behaviors and analyze organizational structure. Master's thesis, Virginia Tech
- Tabassum S, Pereira FS, Fernandes S, Gama J (2018) Social network analysis: an overview. *Wiley Interdiscip Rev Data Min Knowl Discov* 8(5):e1256
- Tang J, Li H, Cao Y, Tang Z (2005) Email data cleaning. In: *Proceedings of the 11th ACM SIGKDD international conference on knowledge discovery in data mining*, pp 489–498
- Task CM (2015) Privacy-preserving social network analysis. PhD thesis, Purdue University
- Thanh Tung V (2014) Email search visualization: an efficient way for searching email. Master's thesis, Helsingfors universitet
- Tsetini M (2015) Computer forensics on financial crimes. Master's thesis, Thessaloniki, Greece
- Van den Elzen S, Van Wijk JJ (2014) Multivariate network exploration and presentation: from detail to overview via selections and aggregations. *IEEE Trans vis Comp Gr* 20(12):2310–2319
- Viégas FB (2005) Revealing individual and collective pasts: visualizations of online social archives. PhD thesis, Massachusetts Institute of Technology
- Viégas FB, Golder S, Donath J (2006) Visualizing email content: portraying relationships from conversational histories. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pp 979–988
- Wang M, Tan Q, Wang X, Shi J (2018) De-anonymizing social networks user via profile similarity. In: *2018 IEEE Third International Conference on Data Science in Cyberspace (DSC)*, IEEE, pp 889–895
- Weisband SP, Reinig BA (1995) Managing user perceptions of email privacy. *Commun ACM* 38(12):40–47
- Weisergerber C, Butler S (2009) Visualizing the future of interaction studies: Data visualization applications as a research, pedagogical, and presentational tool for interaction scholars. *Electronic J of Commun* 19(1–2)
- Wen Q, Gloor PA, Fronzetti Colladon A, Tickoo P, Joshi T (2020) Finding top performers through email patterns analysis. *J Inf Sci* 46(4):508–527
- Whittaker S, Jones Q, Nardi BA, Terveen LG, Creech M, Isaacs E, Hainsworth J (2002), Contactmap: using personal social networks to organize communication in a social desktop. In: *CSCW Videos*, p 7
- Winters J, Prescott A (2019) Negotiating the born-digital: a problem of search. *Arch Manuscr* 47(3):391–403
- Wise S (2014) Can a team have too much cohesion? The dark side to network density. *Eur Manag J* 32(5):703–711
- Ying X, Pan K, Wu X, Guo L (2009) Comparisons of randomization and k-degree anonymization schemes for privacy preserving social network publishing. In: *Proceedings of the 3rd workshop on social network mining and analysis*, pp 1–10
- Ying X, Wu X (2009) On link privacy in randomizing social networks. In: *Pacific-Asia conference on knowledge discovery and data mining*. Springer, pp 28–39

- Zhang J (2018) Miteams: quick organizational mapping by combining email and survey data. Master's thesis, Massachusetts Institute of Technology
- Zhang J (2015) Correspondence as a documentary form, its persistent representation, and email management, preservation, and access. *Rec Manag J* 25(1):78–95. <https://doi.org/10.1108/RMJ-03-2014-0015>

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.