

Interactive Filter and Display of Hillary Clinton's Emails: A Cautionary Tale of Metadata

Christopher D. Salahub and R. Wayne Oldford

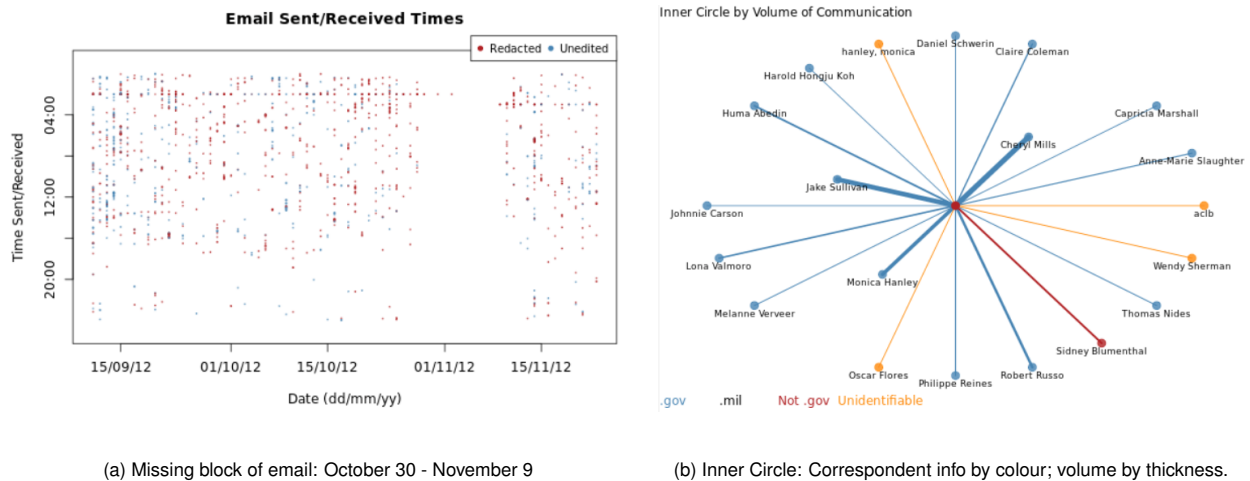


Fig. 1. Metadata: Time, date, redaction, correspondents, volume, source type. Filtered: Sept. 11 to Nov. 23, 2012.

Abstract—We present a web-based visualization that allows the user to interactively filter and display characteristics of 32,795 Hillary Clinton's emails as provided by Wikileaks.

The visualization focuses on the meta-data of each email, including its senders, receivers, and the timestamp the email appeared on the Clinton server. An interactive time range slider filters all email and all displays automatically update to changes in the slider. The main display shows Clinton's most frequent correspondents arranged as nodes of a spoke graph with Clinton at the centre. Volume determines the thickness of each spoke and high volume determines an inner circle whose spokes are shortened. Correspondents and their edges are coloured according to whether that email account could be identified as being an approved Federal government account or not. A second display shows two daily time series: the total number of emails for that day, and the number meeting selection criteria. A third display shows a scatterplot of the time of day versus the day on which that email appeared. Scatterplot points are coloured by whether the email was redacted or not.

Other displays add some information beyond metadata. FOIA exemption codes appear as a selectable list and a barplot shows email counts by FOIA code. The (stemmed) terms having highest frequency in the displayed email, and those having highest tf-idf are listed in separate displays. All displays are interactively filtered by time range and selected FOIA codes.

We illustrate how the filtered displays can be used to generate hypotheses and uncover interesting information. These touch on contentious issues including the handling of classified information, the 2012 attack on the Benghazi U.S. diplomatic compound, and emails apparently missing from those released publicly.

The data are extracted from Wikileaks HTML files, cleaned, and stored in a form useful for interactive exploration. A local R shiny server provides the interactive displays as a public service online tool to explore and uncover patterns in the meta-data and summary contents of Clinton's email. Coupled with publicly available sources of information, these interactive tools uncover surprising amounts of information about an individual, especially one holding public office. The ease with which this can be accomplished and shared should serve as a clear warning as to what can be learned about anyone from metadata.

Index Terms—Exploratory data analysis, metadata, text mining, web-scraping, interactive web visualization, R, shiny

1 INTRODUCTION

The 2016 U.S. Presidential election was one of the most contentious in history. The existence and possible content of Hillary Clinton's private email server dogged former Secretary Clinton's bid for the U.S.

presidency and was doubtless a contributing factor to her surprising defeat by Donald Trump.

On March 16, 2016 Wikileaks published a searchable archive [?] containing the contents of more than 30,000 emails (and attachments) that were sent to and from Secretary Clinton on her private server. The documents were provided as pdfs by the U.S. Department of State in response to Freedom of Information Act (FOIA) [?] requests. The State Dept. also provided a searchable web archive of the documents, released in several instalments from May 2015 to March 2017 [?]. Both sites provide a useful tool for anyone searching for particular terms in the documents. The Wikileaks site was put to much use by investigative reporters, and others, to search for topical news items.

What is critically missing from either site is a facility to learn summary, or statistical, features across all, or groups of, emails. A separate

• Christopher D. Salahub, University of Waterloo, Canada
E-mail: csalahub@uwaterloo.ca.

• R. Wayne Oldford, University of Waterloo, Canada
E-mail: rwooldford@uwaterloo.ca.

Manuscript received xx xxx. 201x; accepted xx xxx. 201x. Date of Publication xx xxx. 201x; date of current version xx xxx. 201x. For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org.
Digital Object Identifier: xx.xxx/TVCG.201x.xxxxxxx

site where visual displays of summary features of the emails, which can be filtered by user selection, will go a long way towards filling this need. The site would be *complementary* to the more typical “content search and record display” site as provided by either online archive. Add general internet search in a third browser, and the three combined would provide the means to actively explore the contents of any email archive. One imagines, for example, having three browsers open simultaneously, one for the data base search (e.g. Wikileaks), one for statistical visualization and exploration (proposed here), and one for internet search (to provide important context for the other two services). Each fills a need not met by the other two; the synergy of all three provide the user a powerful set of investigative tools.

In this paper, we describe an implementation of the missing visual analytic service and illustrate it on Secretary Clinton’s e-mails. The implementation is located at `rshiny.math.uwaterloo.ca/clinton` where any visitor may interactively explore summary features of the entire corpus of emails. Different visual displays present different salient features from all selected emails; emails can be selected via a variety of interactive filters. All displays are reactive and update simultaneously in response to user interaction with the data filters. Together, filters and displays provide a visual analysis service that can be used to quickly learn more about, and to discover some (possibly unanticipated) patterns in the emails.

The implementation is described over the next three sections, beginning with discussion of the data in Section ?? which discusses the source of the data, the extraction and cleaning process in Subsection ??, the details of the metadata in Subsection ??, the use of redaction codes in Subsection ?? and content data used in Subsection ???. The displays of the data and their patterns over the entire unfiltered data are then described in Section ???. Within this section, Subsection ?? describes the spoked network plot, Subsection ?? addresses the email volume time series plot, Subsection ?? discusses the scatterplot of daily sending times, Subsection ?? outlines the exemption code barplot, and Subsection ?? explains the content tf-idf and frequency information. Interactive filters which are implemented in the application are discussed in Section ???. Finally, brief analysis of some interesting findings is included in Section ??, alongside a summary of the timeline of the Clinton email imbroglio in Subsection ??.

2 THE DATA

As of March 3, 2017, a total of 32,795 available emails, either to or from Hillary Clinton, have been made publicly available in pdf form as a searchable database [?]. Many of these have been redacted according to the FOIA exemption codes [?].

Wikileaks [?] has provided the same redacted pdfs and, more usefully for analysis purposes, an HTML version of each pdf. Consequently, have used the Wikileaks database as our data source. All data extraction, cleaning, analysis, and presentation is done using the open source statistical programming language R [?].

2.1 Data extraction and cleaning

The raw HTML of each message was programmatically downloaded from the Wikileaks archive using R packages `RCurl` [?] and `XML` [?]. This took several hours and required the use of manually inserted system pauses to prevent time-outs in the connection, most likely due to Wikileaks DDoS (distributed denial of service) protection software. Besides avoiding such DDoS protection, these pauses are considered web-scraping etiquette and best practice. Once downloaded, the HTML data were processed using the R packages `tm` [?,?], `stringr` [?], and `SnowballC` [?].

After downloading the raw HTML and extracting the data of interest, the resulting data was stored in csv files to provide the flexibility to load the data into any architecture or analysis tool desired. These will be transferred to a relational data base should our server traffic warrant it.

2.2 Metadata

For each email, from the HTML we extract as best we can the identity of the persons sending and receiving the email as well as the date

and time at which the email was processed by the server. The fields used were the address to, address from, contact name to, contact name from, subject line, and time. As well, forwarding chains of email addresses were captured through the identification of any to or from fields followed by emails within the text. In cases where no contact name was present the address was substituted. When the address was missing no imputation was completed. Carbon copy information was not extracted.

The HTML was constructed from email printed out, redacted, and then provided as pdfs by the State Department. Consequently, detailed email header information as would normally be available electronically is mostly missing. All time stamps appear to be the local date and time at which the server sent or received that email (e.g. no time zone or other source time or IP chain information is available). Moreover, because of redaction (typically FOIA exemption B6 [?]), sender and receiver emails may have truncated domains, contain only the person’s name, or be missing altogether. In cases where both From and To are entirely missing, but there is an email chain within the message, we do not impute values (e.g. see <https://wikileaks.org/clinton-emails/emailid/31599>).

Using regular expressions to extract the metadata was occasionally challenging and it is always possible that some fringe cases have been mishandled. On the whole the metadata is fairly consistent and only rarely do some impure and messy addresses and contact names arise due to irregular spacing or placement of text within the HTML code. One such fringe case is Huma Abedin’s `clintonemail.com` email account which will show in the displays as an overly long string. We have chosen not to special case this but leave it as is.

Where possible, the email addresses have also been parsed so that they may be classified into one of four categories: those which are .gov, those which are .mil, those which are identifiable as coming from a domain that is neither .gov nor .mil, and those whose domain was not identifiable from the data.

2.3 Redaction information

Emails that are redacted are marked as such by the string “RELEASE IN PART” and by the presence of one or more of nine FOIA exemption codes B1, B2, ..., B9 marking the place where text is missing (redacted). This provides two further pieces of quasi-metadata (since it is not actual email content) that can be used in analysis.

2.4 Content information

The entire (redacted) content of each email is available on the State Department and Wikileaks sites and the user is encouraged to view it there (the pdf forms are more informative, especially in appreciating the extent of the redactions). To provide some coarse statistical summaries of the content, all word tokens are extracted and partially processed to reduce their number. For example, “stopwords”, as identified by the `stopwords` function under either the “en” and “SMART” settings, or a set of custom stopwords, were removed. The remaining words were stemmed by the `stemDocument` function from `tm` and `SnowballC`. While neither the stopword lists nor the stemming tools are particularly well tuned to this corpus of emails, they nevertheless provide some insight as to the topics covered.

3 DISPLAYS

TODO Each of these four displays is described and commented on in context of whole time line Four displays of the metadata and quasi-metadata

3.1 Inner circle

This display takes up to the 20 most frequent correspondents in the selected emails and arranges their email addresses at the ends of equiangular spokes around the hub that is Secretary Clinton’s email. The hub is coloured red to show that this account is known not to be a government sponsored account (either .gov or .mil). Every account that is identifiably not a government sponsored account is coloured red. Those that are identifiably government are coloured either blue

(for .gov) or black (for .mil, if any appear). Those which can not be determined as either government or not, are coloured orange.

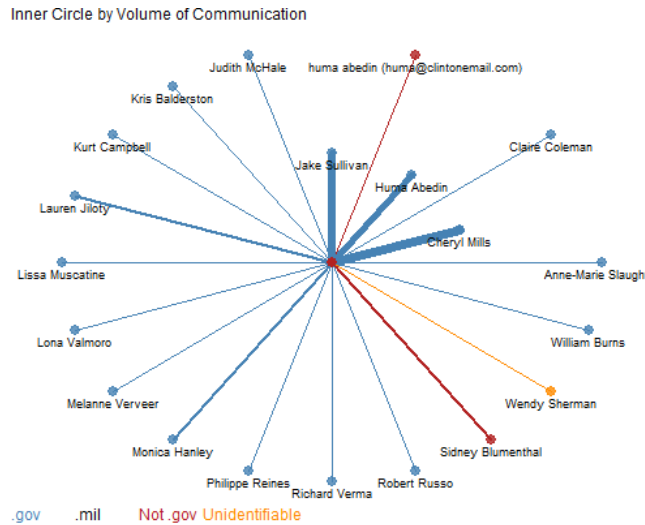


Fig. 2. Inner Circle over all correspondence

The figure shows the top 20 (or “Inner Circle”) hereover all emails in the collection. Note that Clinton aide Huma Abedin is seen to be using both a .gov account for some email and a clintonemail account for others and so appears as two different colour spokes. The other identifiably non-government account is that of Sidney Blumenthal. This is of particular interest, as his closeness to Clinton in providing advice on all kinds of matter has been the source of controversy [?] given that he had been refused a State department position and so lacked State Department clearance. Wendy Sherman, though appointed by Clinton to Under Secretary of State for Political Affairs in 2011, her email address could not be identified as either definitively government or definitively non-government and so is shown in orange. It is important to note, however, that the most likely reason an email is not identifiable is that it has been redacted (likely B6). While no guarantee that the email is not a government email, it is more likely that it was viewed as personal or private information and so probably not government.

The greater the number of emails between Clinton and the correspondent, the wider is the spoke. Whenever the difference in volume is great enough the correspondents will also be separated into two groups. Those with the greatest correspondence will have shorter spokes and hence be visually “closer” to Clinton. As the display shows, Clinton’s closest inner circle are her closest aides Huma Abedin, Cheryl Mills, and Jake Sullivan. **TODO Chris, can you be more precise about the algorithm?**

3.2 Email volume

Figure ?? has two lines, the first is a line showing the total volume of emails for each day, and the second shows the number of emails satisfying the classification filter criteria for a given day. This allows users to view which days have an exceptionally high or low number of classified conversations. Returning to the motivation of hypothesis generation, referencing world events on any days of interest provides a number of interesting insights on how public figures like Clinton react and respond to world events.

This plot shows highly variable nature of the daily email volumes, with a number of peaks and troughs suggesting regions in which further investigation might be insightful.

3.3 Email times

As shown in Figure ??, this display shows the time stamp of each email separated into its calendar date on the horizontal axis and its (24hr) time of day on the vertical. Each email is a point in the plot and is

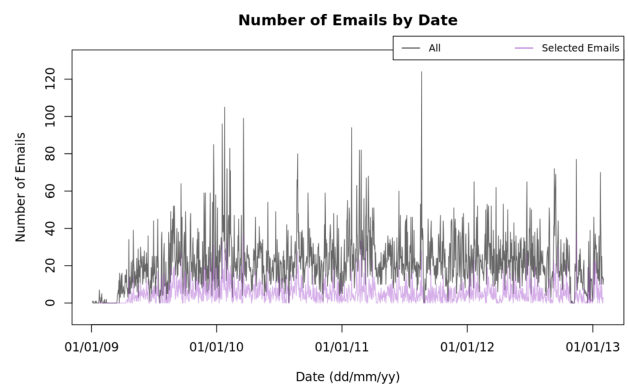


Fig. 3. Email Volume: All emails (top) ; Sent by Clinton (bottom)

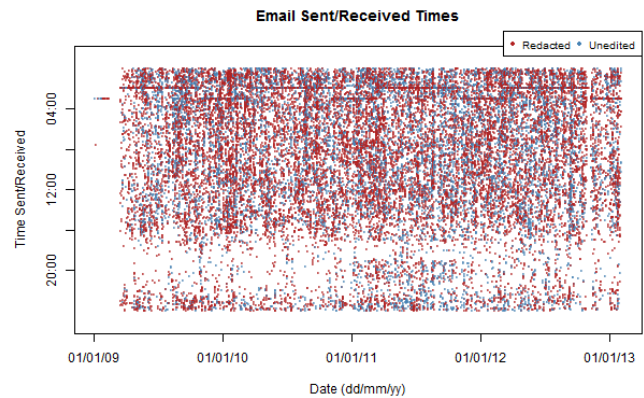


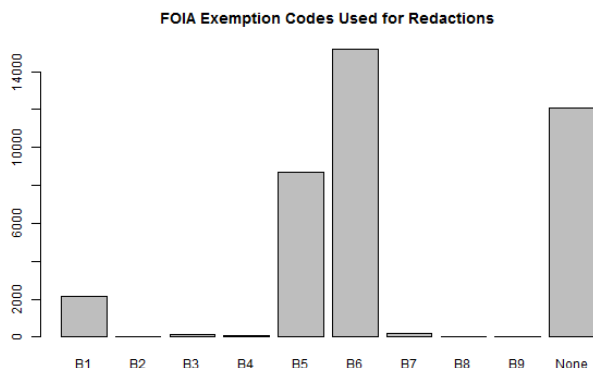
Fig. 4. Email time patterns: all emails

coloured red if any part of that email was redacted, and blue if it was released in full. Alpha blending prevents over-plotting from obscuring patterns.

A few patterns are easily discerned. For example, the least amount of email appears for a few hours around 20:00 hours, or 8 PM. Evenings appear to be when email traffic is lightest. Otherwise, for the most part, it seems fairly uniformly distributed throughout. Midnight is at the top (and bottom) of the plot, so it would seem that email will occur mainly from midnight until about 4 or 5 PM. Note also the surprisingly regular horizontal lines that appear across the top of the plot. This regularity is suggestive of some automated schedule for the server which causes the email to stack up before being recorded at either 2 or 3 AM. The location of these lines switches exactly whenever daylight savings time switches in North America.

3.4 FOIA exemptions

The United States Freedom of Information Act (FOIA) exemption codes [?] are used by the United States Federal Government to remove sensitive information from released documents accessed using a Freedom of Information Act request. The sensitive matters being addressed by the codes are national security and foreign policy matters for B1, personnel practices for B2, exempted by statute for B3, trade secrets or financial information obtained in confidence for B4, inter- or intra-agency memorandums for B5, personal privacy for B6, records compiled for law enforcement for B7, prepared in relation to financial monitoring institutions for B8, and geophysical and geological information concerning wells for B9 exemptions. The frequencies of each of these codes within filtered emails are shown in the final display. Note that as emails can have numerous redaction codes, this barplot can be useful to identify which codes co-occur over a selection.



It is clear that redactions made for national security and foreign policy reasons account for only a small portion of redactions made. The most frequent redactions are those made to protect personal privacy, which occur in almost half of all emails, and those made with the vague B5 exemption code, which occur in roughly a quarter of all emails. The other codes appear in very few emails and just over a third of the emails are not redacted at all.

3.5 Term frequency and TFIDF

The 20 terms (stemmed words excluding stopwords) which appear in the greatest number of the emails selected are displayed as those having “Highest Frequency”. Those terms which appear frequently

20 Highest Frequency Terms

will, state, pm, said, secretari, can, call, depart, time, presid, govern, offic, work, usa, one, also, meet, new, us,

20 Highest TFIDF Terms

msg, pr, pager, folder, fyi, nternet, tx, folderid, rim, ticker, pls, cheryl, send, dev, sorri, fw, soon, jake, delet, true

within some emails but less frequently across emails will have a high tf-idf (term frequency - inverse document frequency) score. The top 20 scoring of these in the selected emails are shown in the “TFIDF” display. **TODO Chris, check that I have described these two measures correctly** The top 20 shown here are for all emails in the corpus; these will change depending on the filtering.

The terms give some limited insight into the contents of the email. As seen here, for example, the first names of two of Clinton’s closest aides, Jake Sullivan and Cheryl Mills, appear under TFIDF but not as high frequency terms. One problem with tf-idf for this e-mail corpus is that many emails contain email chains which grow as each person replies which could magnify the within email frequency for some terms.

As mentioned in Section ??, the stemming and stopword removal provided by the package `tm` in R is also challenged by this messy and non-standard data. The list of stopwords had to be supplemented so as to avoid uninformative action verbs such as ‘will’ and ‘can’. The stemming algorithm was also challenged by typographical errors in emails and the proliferation of acronyms (e.g. for individuals and government abbreviations). Even so, the terms nevertheless occasionally turn up something of interest as in the case of “windrush”, which turned out to refer to “Windrush Ventures” (e.g. see [?, ?]), a company owned by former U.K. Prime Minister Tony Blair. “Windrush” appears at the end of emails as part of Mr. Blair’s electronic signature for his email id “aclb” (e.g. see the “Inner Circle” of Figure ??, b).

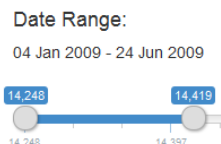
4 FILTERS

There are only three filters but each one affects all data displays. As each filter is changed, every display redraws itself on the filtered data,

immediately in reaction to the change. In this way, the filters can be used together to focus on particular subsets of the emails or simply to observe how the display patterns change with the filter being applied.

4.1 Time filtering

A sliding time window filters the emails displayed to those whose date

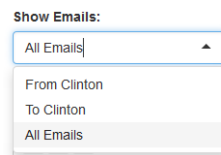


lies between the two end points. The size of the time window is changed by moving either end, either by mouse click and drag or by selecting an end to move with the arrow keys. The whole time window is moved by selecting the middle bar and either dragging it or moving it with the arrow keys.

This filter is simple but powerful. All emails are displayed when the range covers all dates. Moving the end points towards one another allows the user to focus the displays on any particular range, down to as fine as all emails on a single day. With a fixed range of days, for example a two week period, dragging the middle bar from left to right will have each display smoothly update over time.

4.2 Sent or received

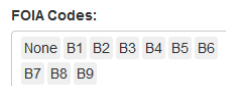
A drop down menu is used to select those emails which Clinton sent, or



received, or either sent or received. The default is “All Emails”, which is the same as no filter being applied for sent or received. Choosing “From Clinton” selects for display only those mails sent by Secretary Clinton, choosing “To Clinton” selects only those emails which she received. This might be used, for example, to explore whether the inner circle of correspondents changes depending on whether Secretary Clinton is emailing them, or they are emailing her.

4.3 FOIA exemptions

Multiple filtering by FOIA exemption codes is supported in the app with a box selection area. Users can add FOIA exemption codes currently



not included to the data by clicking the box and selecting a code from the menu which is then displayed below the box. To remove codes, the user simply clicks the code within the display box and presses delete, which moves the code to the selection menu. The box therefore displays all codes currently included in the filter.

4.4 Auxiliary information

TODO e.g. Clinton’s foreign schedule A number of cosmetic adjustments using auxiliary information can also be made. Users can choose to adjust the scatterplot of time by day, which re-defines 6 pm as the time at which the day changes. This simply groups the emails best, placing all emails into one mostly uninterrupted region by redefining

the day about the midsection of the region of lowest email traffic. Another cosmetic adjustment is the inclusion of Clinton's foreign trip schedule as blue regions on the email volume plot. This schedule was taken from the official State Department website [?] and allows users to distinguish periods spent out of the United States (coloured in blue) from those spent within the United States (uncoloured).

5 ANALYSIS

5.1 Context: A brief timeline on the private email server

While the following timeline is somewhat abbreviated, it should serve to raise some of the major issues and concerns related to the private email server and its contents. It also introduces some of the principal characters in the whole affair. More complete and in depth timelines are readily available on the internet (e.g. [?, ?, ?, ?, ?]).

On November 21, 2008, the New York Times reported that Hillary Clinton had decided to accept the position of U.S. Secretary of State. On January 13, 2009 the internet domain name `clintonemail.com` was registered [?]; eight days later Senator Clinton was confirmed as Secretary of State.

Public knowledge that a private email server being used by Secretary Clinton and others for State Department and personal communications did not surface until March 2015 [?] during the course of a U.S. Congressional investigation [?] of the September 11, 2012 attack by militants on U.S. compounds in Benghazi Libya.

The State Department had difficulty fulfilling public FOIA and House Benghazi Committee requests [?] for Secretary Clinton's government emails because she had exclusively used the private server for all her email. On March 10, 2015, Clinton told reporters that she turned over 30,490 emails to the State Department and deleted 31,830 emails deemed to be personal [?]. Clinton had tasked three lawyers Cheryl Mills (Clinton's former chief of staff), David Kendall (Clinton's personal lawyer), and Heather Samuelson (a State Department staffer during Clinton's tenure) to make the determinations as to which emails were work related and which were not [?, ?].

On March 10, 2015 the House Benghazi Committee requested that the private email server be turned over to a neutral third party to determine which emails are personal and which are government records [?], but was informed March 27 by David Kendall that no emails remained on the private server for any kind of review [?]. Between March 25 and 31, 2015, Paul Combetta (then the server's system administrator), erased all backup copies using BleachBit (see www.bleachbit.org).

Combetta, Mills, and Samuelson will later be granted partial immunity by the Justice Department during the FBI investigations into the private email server, as were two others: Bryan Pagliano (original server manager) and John Bentele (former director of Information Resources Management for the State Department's Executive Secretariat) [?, ?, ?].

On April 12, 2015 Clinton announces that she is running for the U.S. Presidency. On July 24, 2015, inspectors general for the State Department and the national intelligence agencies announce finding classified information in the emails and that the information they found was classified at the time sent [?], though her campaign declared that they must have been classified after the fact. On August 19, 2015, Clinton calls the allegation of mishandling classified information a "disagreement between agencies" [?].

Nearly one year later, July 5, 2016, FBI Director James Comey recommended that no charges be laid against Clinton on use of private email server [?]. On October 28, 2016, Comey revealed that in a separate investigation into former Congressman Anthony Weiner, that emails belonging to his wife Huma Abedin have been found on his laptop. Since Abedin was a close aid to former Secretary Clinton, FBI investigations were reopened into the private server usage but closed again by Comey on November 6, 2016 without charges being laid [?]. In both cases, Comey and the FBI are criticized by pundits from different political parties.

5.2 Highest peak

Within the plot of email volume over time, one very large peak stands out. Historically important, possibly had an impact on the 2016 US

presidential election.

5.3 No email sent by Clinton

Using 1 month window slider and looking for flat spots at zero Dates:

- None until mid April 2009
- June 8 or 9, June 15 to about the 20th, 2009
- July 17-20, around 27, 2009
- Day or 2 around Oct 13, 14, 2009
- April 10-15 or so? 2011
- one or two days late August, 2011. (no mail)
- end of march first half of April 2012
- August 27-30 2012
- low first week of September, 2012
- 2 weeks end of october, up to Nov 10 approx 2012
- Dec 7- 17? 2012, Again about Dec 25 2012 to Jan 1, 2013

Contentious issues:

- private email server
- classified documents (outside of state)
- Sidney Blumenthal
- Benghazi spin handling of media (Susan Rice)
- scrubbing of her server (bleachbit)
- missing email from online email State department

Learn:

- inner circle over time (state or not)
- spike in email around Libyan revolution
- gap in the email
- daily email patterns, server behaviour (daylight savings time)

Filter:

- time
- redacted or not
- FOIA exemption codes

Content:

- Term frequency, TFIDF

Tool:

- Web-based, interactive filter and display tool

Discovery from visualization & connecting discovery sources

- Nothing on this ... Preparation for Benghazi (security considerations)
- House Oversight and Government Reform Committee (standing committee) Darrell Issa, Chair Jason Chaffetz, Chair (Gowdy a member) (discovered email server) ... - House Select Committee on Benghazi (Summer 2014) Trey Gowdy, Chair

6 OLD STUFF

TODO This material is to be reused wherever An interactive web application in R shiny was constructed utilizing these visualizations and the extracted data. Incorporated in this app was a short article explaining the use of the app, providing links to background information, and outlining some of the findings of the more interesting trends observed. This framework provides any individual with interest the ability to generate hypotheses by exploring the data without necessarily searching for any particular conclusion.

In generating and exploring this data, a considerable amount of respect is gained for the power of the often overlooked metadata of emails, especially with public figures. Simply having access to the network of communications, simple indicators of identifiable features of the individuals communicating, and times sent is enough to generate interesting questions which simple Google searches of important events and days can shed a great deal of light on. Leveraging data beyond this metadata, including the Freedom of Information Act (FOIA) exemption codes (see <https://vault.fbi.gov/explanation-of-exemptions>) used to redact the emails, only bolsters this simple metadata driven investigation. Metadata is often viewed as less significant or somehow less intrusive data, but the exploration performed in this paper demonstrated, at the very least to the authors, what a powerful tool metadata is for generating hypotheses about data. Given the lack of privacy present in the modern internet age and the amount of information most individuals make publicly available, there is a clear warning here about the utility of metadata for bad or for good. For a public official like Hillary Clinton, the data become even more interesting due to a number of contentious issues which arose during her tenure as the Secretary of State, independent of and as a result of her use of a private server.

7 CONCLUDING REMARKS

In providing the service, we purposely focused the displays (for the most part) on simple characteristics of the emails, the so-called metadata. There are several related reasons for this focus.

First, metadata is often more reliable, regular, and, being easily collected, more generally available. For Secretary Clinton's emails, some metadata is lost, either because their source is a printed form or because it has been redacted. Examining metadata also arguably intrudes least upon the privacy of the individual correspondents, compared at least to the email's content.

Second, since at least the Snowden revelations (e.g. [?]), public discourse has grown on the potential value of metadata to those who have it. By providing a web service where anyone can see for themselves what might be learned from metadata, we hope to contribute positively to this important discussion. Moreover, the work-related only metadata of a public figure as senior as Secretary Clinton will hopefully resonate more strongly with other public figures engaged in the debate (e.g. [?, ?, ?, ?, ?]) than might that of an ordinary citizen.

Finally, it is important that users realize that the value of metadata is not just in itself, where one might expect to easily discern general day-to-day habits such as one's circle of correspondents. Rather it is that when coupled with other sources, which are abundantly and publicly available for Senator Clinton, much more can be learned than from any one source alone. Those who have knowledge of, or access to, other sources may filter metadata to test previously held hypothesis; conversely, exploration of metadata could uncover patterns that generated hypotheses to be tested elsewhere.

ACKNOWLEDGMENTS

This work was supported in part by a Discovery grant from the Natural Sciences and Engineering Research Council of Canada..