

Search Among Secrets

Douglas W. Oard
iSchool and UMIACS
University of Maryland, College Park, USA

Open Government

“Information maintained by the Federal Government is a national asset. My Administration will take appropriate action, consistent with law and policy, to disclose information rapidly in forms that the public can readily find and use.”

From: H <hrod17@clintonemail.com>
Sent: Friday, March 5, 2010 9:56 AM
To: 'vermarr@state.gov'; 'sullivanjj@state.gov'
Subject: Gefilte fish

Where are we on this?

- 62,320 messages on a private server
- Term-list hits yielded 30,490 as “work-related”
 - Provided to the State Department, initially on paper
 - Added emails from State Department archives
 - Of ~1 billion emails per year, ~50,000 are archived (0.005%)
- Complex multi-agency manual review process
 - Target: 1,000 documents per week

- 1** Review a batch (1,000 messages) per week
- 2** Apply proposed redactions to each batch
- 3** Send batch to subject matter experts (SMEs) for consultation
- 4** SMEs return batch to FOIA office
- 5** FOIA office incorporates review recommendations from SMEs
- 6** FOIA office refers to other agencies implicated in emails
- 7** Other agencies review
- 8** Other agencies send feedback to FOIA office
- 9** FOIA office incorporates appropriate recommendations
- 10** FOIA office sends batch to Office of Legal Counsel (OLC) for review
- 11** OLC reviews
- 12** OLC sends feedback to FOIA office
- 13** FOIA office incorporates OLC recommendations
- 14** Divergent recommendations among reviewing entities reviewed
- 15** Final review

Exclusive: U.S. to shift 50 staff to boost office handling Clinton emails

WASHINGTON | BY ARSHAD MOHAMMED



Democratic presidential candidate Hillary Clinton speaks during a press conference at Des Moines Area Community College in Ankeny, Iowa in this August 26, 2015 file photo.

REUTERS/SCOTT MORGAN/FILES

The U.S. State Department plans to move about 50 workers into temporary jobs to bolster the office sifting through Hillary Clinton's emails and grappling with a vast backlog of other requests for information to be declassified, officials said on Tuesday.

The move illustrates the huge administrative burden caused by Clinton's decision to use a private email address for official communications as secretary of state and a judge's ruling in a Freedom of Information Act (FOIA) lawsuit that they be released.

Jeb Bush email dump includes Social Security numbers



38



Getty Images

By Ellee Viebeck - 02/10/15 02:44 PM EST

A trove of 250,000 emails released by prospective 2016 presidential candidate Jeb Bush includes the sensitive personal information of several Florida residents, leaving them vulnerable to identity thieves.

A scan of the **email dump** by technology blogs **The Verge** and **Gizmodo** revealed names, emails and in some cases, Social Security numbers of Bush's correspondents. Many appear to be normal Florida residents unaware their messages to the then-governor would eventually become public.

A Story in Three Parts

- The need for search among secrets
 - Protecting things
 - Protecting parts of things

What's Different Now?



■ Formal ■ Informal



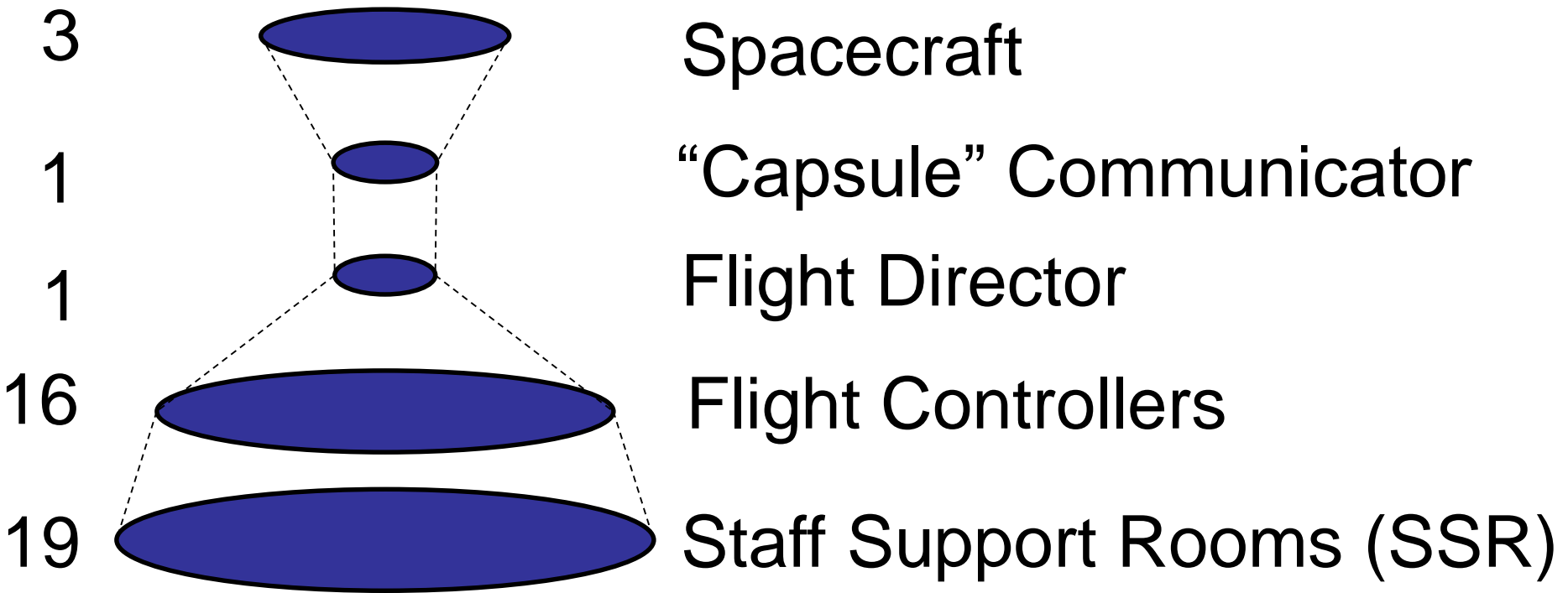
ARCHIVES OF THE UNITED STATES OF AMERICA

THIS BUILDING WAS DESIGNED BY THE ARCHIVES OF OUR GOVERNMENT
AND DEDICATED TO THE PRESERVATION OF OUR HISTORY





Hourglass Model of Mission Control





IMPROVING DECLASSIFICATION

A REPORT TO THE PRESIDENT
FROM THE PUBLIC INTEREST DECLASSIFICATION BOARD



"A popular Government, without popular information or the means of acquiring it, is but a Prologue to a Farce or a Tragedy; or perhaps both. Knowledge will forever govern ignorance; And a people who mean to be their own Governors, must arm themselves with the power which knowledge gives."

James Madison to W.T. Barry
AUGUST 4, 1822

DECEMBER 2007



ISSUE NO. 2: Prioritizing the Declassification Review of Historically Significant Information.

There is no satisfactory means at present of identifying historically significant information within the vast body of information that is being reviewed and declassified. Accordingly, no priority is given to the declassification and release to the public of such information.

“Value-Sensitive” Prioritization

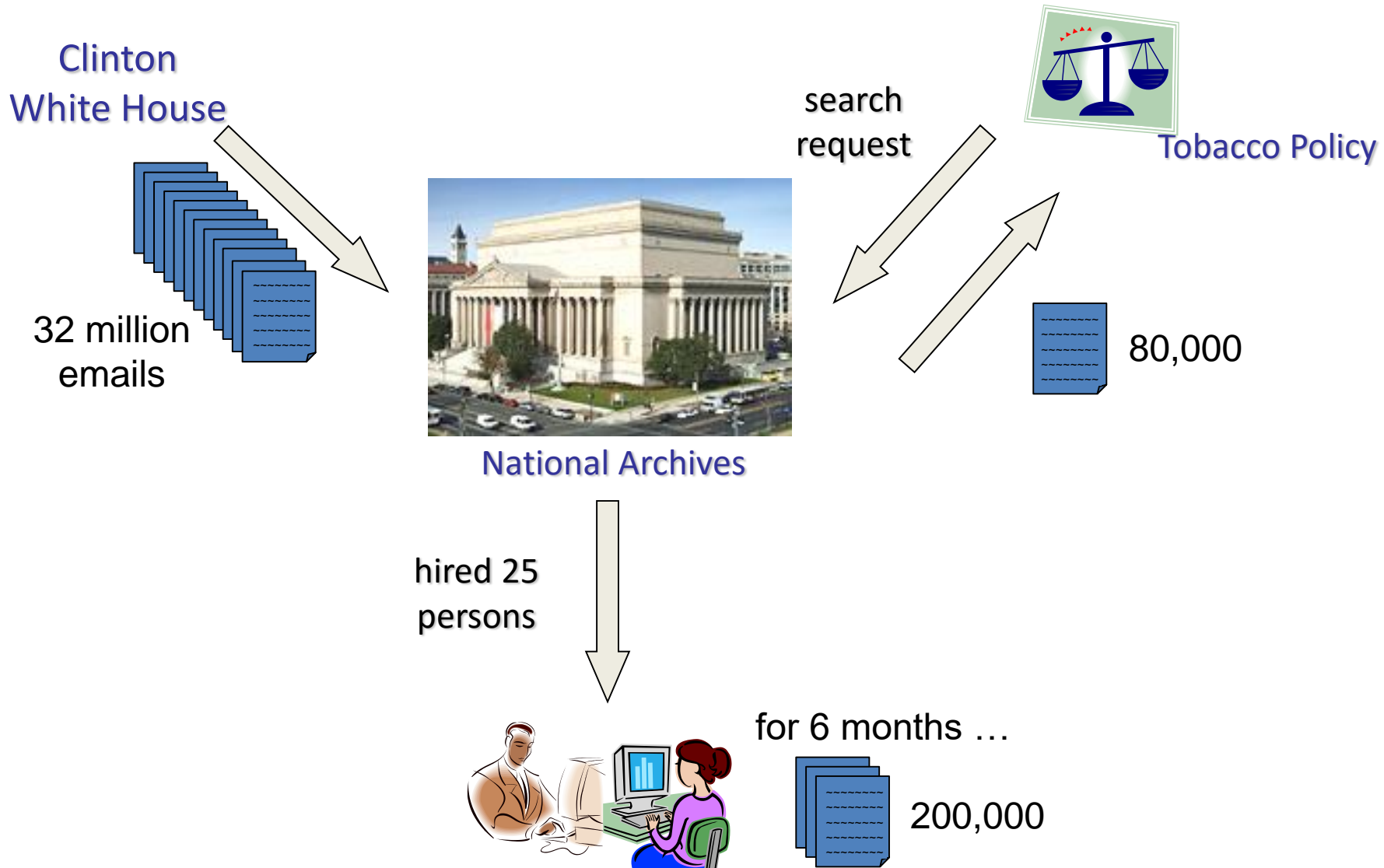
Such a system might operate as follows:

A board consisting of prominent historians, academicians, and former Government officials would be appointed by the Archivist to **determine which events or activities of the U.S. Government should be considered historically significant** from a national security and foreign policy standpoint for a particular year.

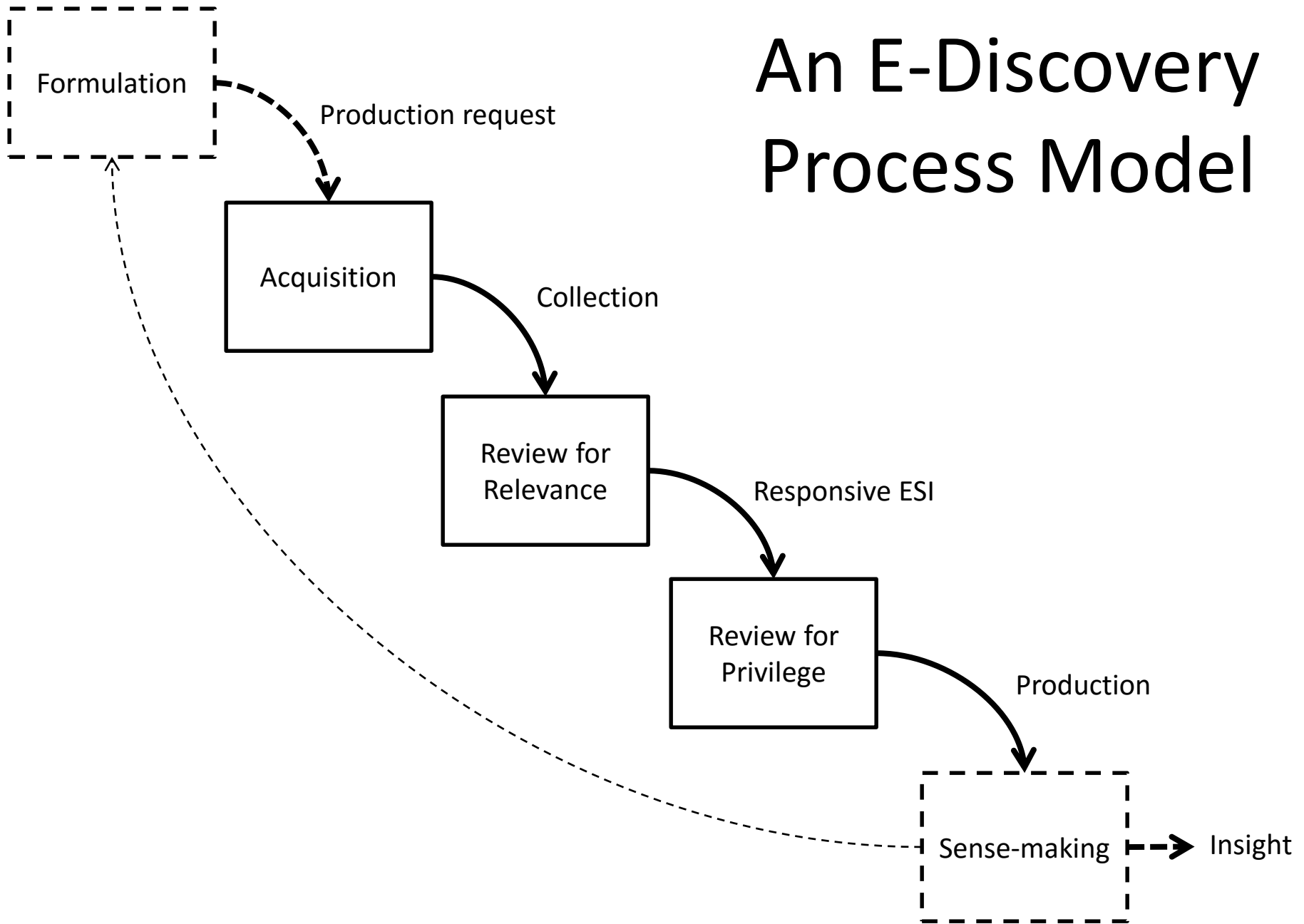
U.S. Freedom of Information Act

“identify the subject(s) or record(s) as clearly and specifically as possible -- for example, all previously released National Intelligence Estimates (NIEs) on the former Soviet Union's space program.”

E-Discovery



An E-Discovery Process Model



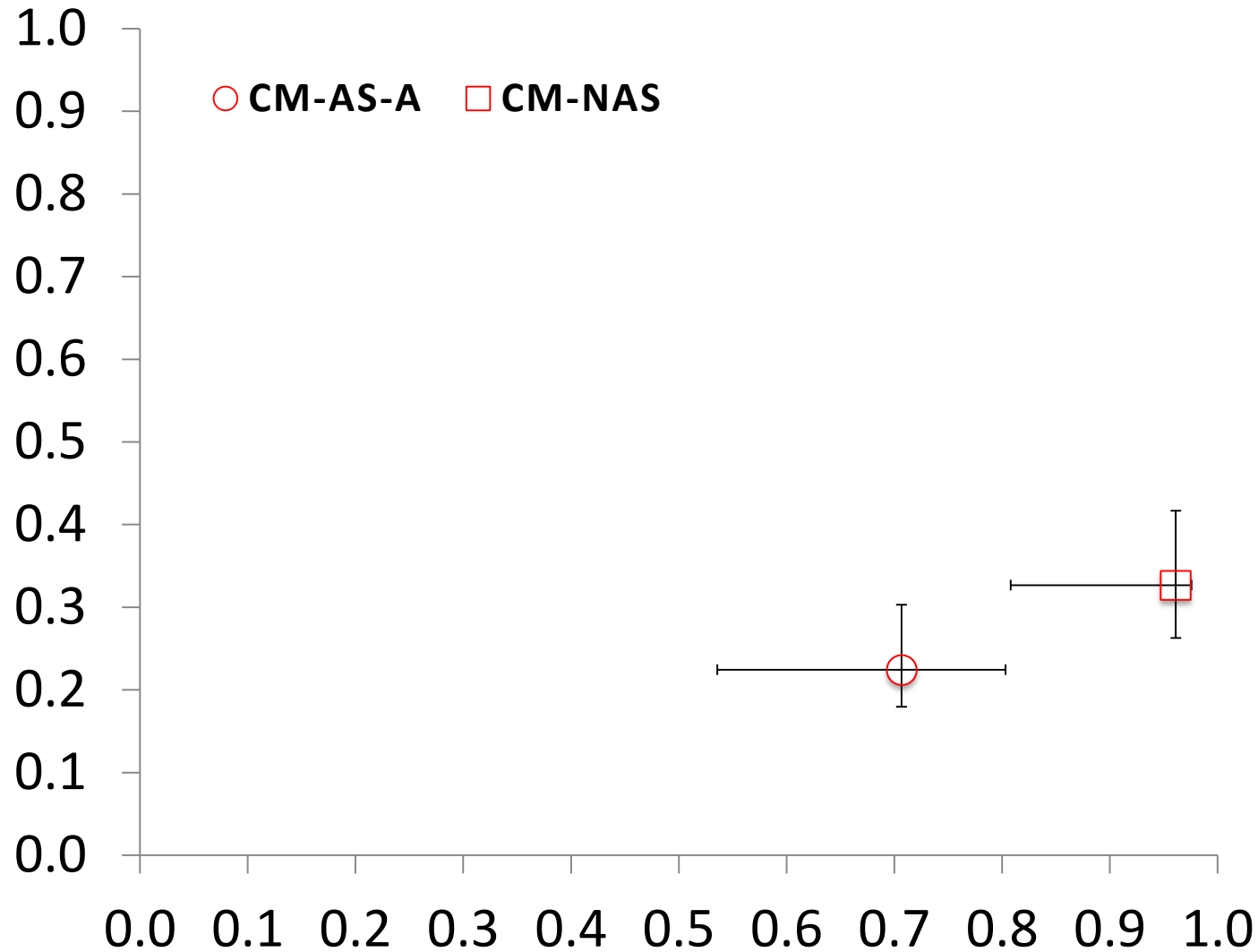
A Story in Three Parts

- The need for search among secrets



➤ Protecting things

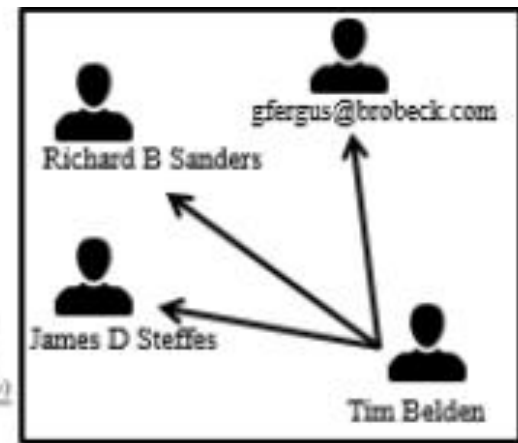
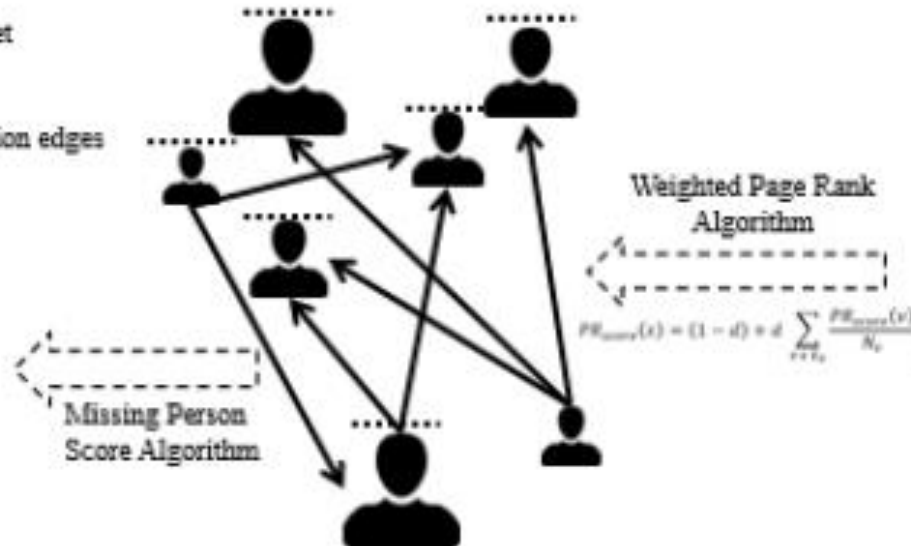
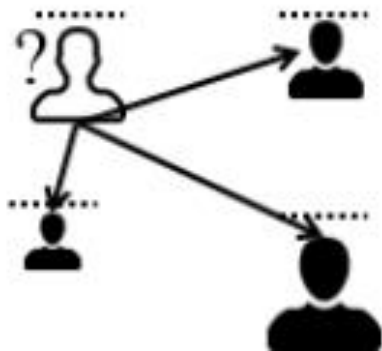
- Protecting parts of things

Automated Privilege Classification



Estimating Privilege Propensity

-  → Missing Person in Training-set
-  → Person Seen in Training-set
- → Implies Missing communication edges



Machine-Assisted Privilege Review



Attachments (1)

Privileged

Not Privileged

No Decision

Date: Wed, 25 Apr 2001 03:43:00

From: JBennett <JBennett@GMSSR.com> Government Affairs and Relations Executive

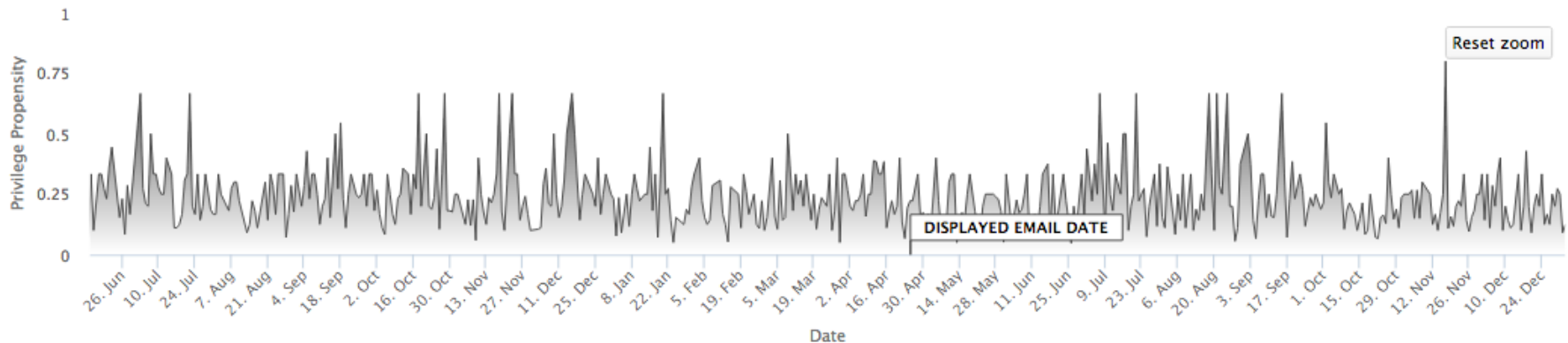
To: Harry.Kingerski@enron.com, **Jeff.Dasovich@enron.com**, Leslie.Lawner@enron.com, tjohnso8@enron.com, James D Steffes <**James_D_Steffes@enron.com**>, Scott Stoness (E-mail) <sstoness@enron.com>, Sue Mara (E-mail) <**smara@enron.com**>

Cc: MDay <MDay@GMSSR.com>

Subject: Advice **Letter** Protest on SCE PX **Credit**

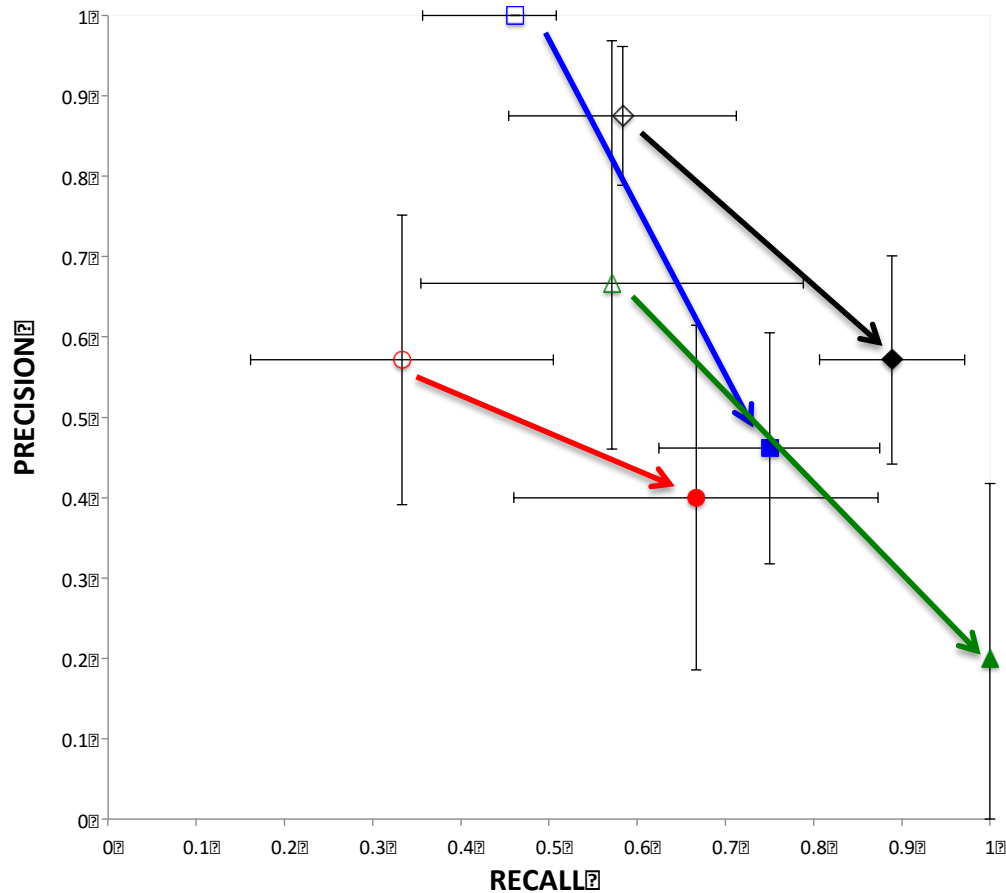
Attached is the revised version of the protest of SCE's advice **letter** on the PX **credit**. Please **provide** all **comments** ASAP. If you cannot get them to me by 12:30 PDT, then send to Mike Day for incorporation.

Click and Drag to Zoom in



* The above graph plots Privilege propensity

Annotations Increase Recall



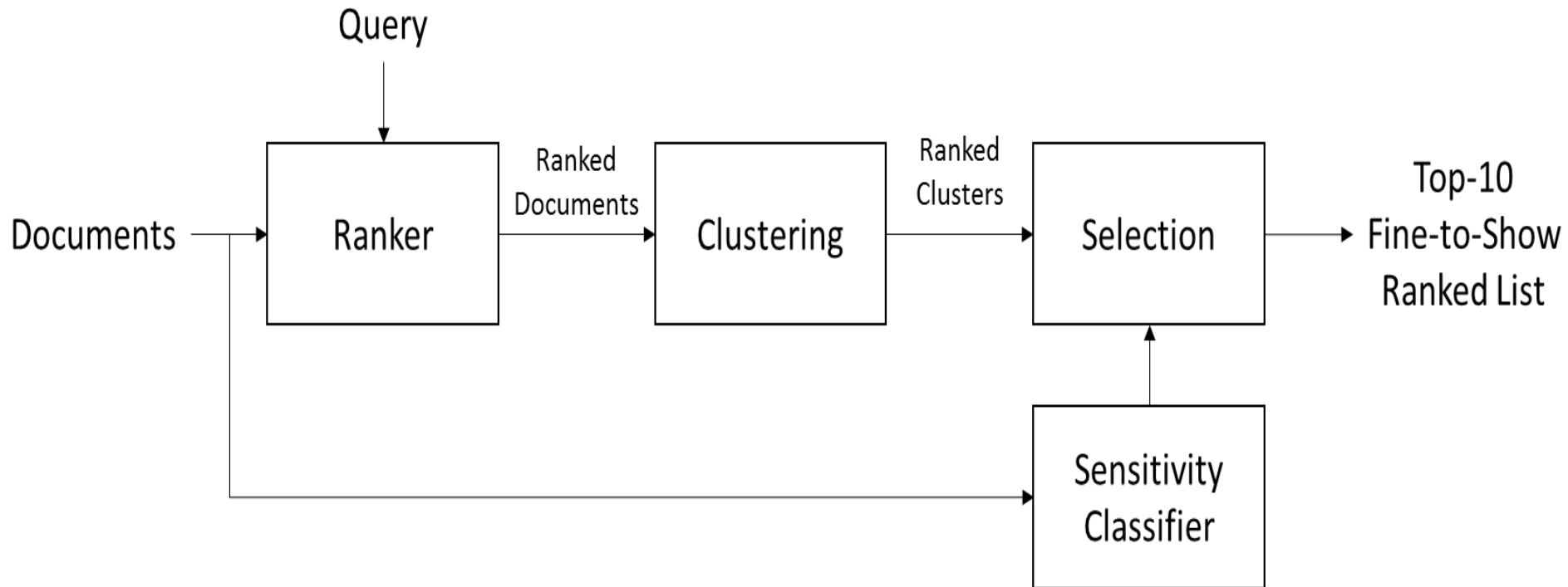
S_1 Judgments as Ground Truth

RECALL ↑

PRECISION ↓

- E1-Baseline
- E2-Baseline
- △ I1-Baseline
- ◇ S2-Baseline
- E1-AID
- E2-AID
- ▲ I1-AID
- ◆ S2-AID

Balancing Relevance and Sensitivity



Evaluation

- Our goal is to show some relevant documents
 - To do that, we must risk showing some sensitive docs
- Precision-oriented search could limit the risk
- Loss function from operational risk management
 - Severity, mitigation, continuous improvement

Discounted Cumulative Gain

| | Highly Relevant (h) | Moderately Relevant (m) | Not Relevant |
|---------------|---------------------|-------------------------|--------------|
| RETRIEVED | $+G_h$ | $+G_m$ | 0 |
| NOT RETRIEVED | 0 | 0 | 0 |

$$DCG_k = \sum_{i=1}^k \frac{g_i}{d_i}$$

Cost-Sensitive DCG

| | Highly Relevant (h) | Moderately Relevant (m) | Not Relevant |
|---------------|---------------------|-------------------------|--------------|
| RETRIEVED | $+G_h$ | $+G_m$ | 0 |
| NOT RETRIEVED | 0 | 0 | 0 |

$$DCG_k = \sum_{i=1}^k \frac{g_i}{d_i}$$

| RETRIEVED | Highly Relevant (h) | Moderately Relevant (m) | Not Relevant |
|------------------------|---------------------|-------------------------|--------------|
| Fine to Show | $+G_h$ | $+G_m$ | 0 |
| Somewhat Sensitive (s) | $-C_s$ | $-C_s$ | $-C_s$ |
| Very Sensitive (v) | $-C_v$ | $-C_v$ | $-C_v$ |

$$CS - DCG_k = \sum_{i=1}^k \left(\frac{g_i}{d_i} + c_i \right)$$

| NOT RETRIEVED | Highly Relevant (h) | Moderately Relevant (m) | Not Relevant |
|------------------------|---------------------|-------------------------|--------------|
| Fine to Show | 0 | 0 | 0 |
| Somewhat Sensitive (s) | 0 | 0 | 0 |
| Very Sensitive (v) | 0 | 0 | 0 |

A Story in Three Parts

- The need for search among secrets
- Protecting things
- Protecting parts of things

Stanford ePADD: Collection Overview

[DEMO](#)

[Home](#) > [Robert Creeley Papers, Email Series](#)

Collection Title Robert Creeley Papers, Email Series

Collection Number M0662.E

Extent October 23, 1994 to February 20, 2013
49,646 messages (19,491 outgoing and 30,155 incoming).
10,608 attachments 4,239 images.
24,010 people.

Links [Finding Aids](#)
[SearchWorks](#)

Browse

Categories

[Entities](#)

[Correspondents](#)

Word Cloud of Entities

Afghanistan **America** American **Amherst Street** Baghdad Berkeley
Best **Bill** Black Mountain **Bob** Bob P.S. **Boston** Brooklyn **Brown**
Brown University **Buffalo** **Bush** CA California Canada **Charles**
Charles Bernstein Charles Olson Chicago **Congress** Creeley Dad **David** England
Europe Florida **France** George George Bush George W. Bush Germany **Hannah**
Harvard House **Iraq** Israel Italy Jim **John** LA London Los Angeles MSN
Maine Manhattan Marfa Mark **Michael** Middle East Mike **Mr. Creeley** N.Y.
NYC NYU **New York** New York City New York Times Olson Paris Pen
Penelope Petaluma **Peter** Philadelphia Portland Providence **Robert**
Robert Creeley Russia Saddam Hussein **San Francisco** Sarah Seattle
Senate Susan TX Texas Tom Toronto Tricia **U.S.** UB UK UN **US** USA
United States University University of California Press Vietnam **Waldoboro**
Washington White House Williams Yahoo

Bulk Search Entities

[Search](#)

Date: March 2, 2001 11:24am

From: [Katbarc@...](#)

To: [creeley@...](#)

Subject:

.... Robert,

.....
.....
..... Central Michigan (...'
.....)
.....
.....

..... **Portland**
..... " "
..... \$
..... \$
.....

.....
Kathryn

Kathryn Barcos
The Steven Barclay Agency

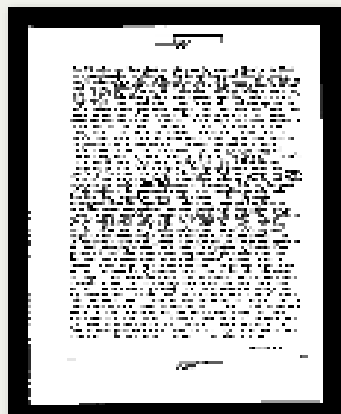
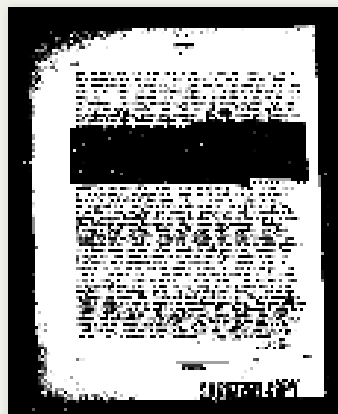
.....
Petaluma,

.....
.....

Redaction 1 / 30

1990

1993



danger. The USSR does not foment revolution but the United States always looks for outside forces whenever certain upheavals occur. One example of USSR's determination not to interfere in internal affairs of other countries is Iran, an ally of the United States. The Soviet Union does not want a revolution there and does not do anything in that country to promote such a development. However, the people of that country are so poor that the country has become a volcano and changes are bound to occur sooner or later. The Shah will certainly be overthrown. By supporting the Shah, the United States generates adverse feelings toward the United States among the people of Iran and, conversely, favorable feelings toward the USSR. This, of course, is to the US's own disadvantage. The Soviet Union does not sympathize with dictators or tyranny. This is the crux of the matter. No agreement seems to be possible on this

Document Info

View Side-by-Side

Report Invalid Match

JAN 1961 APR JUL OCT JAN 1962 APR JUL OCT JAN 1963 APR JUL OCT

Evaluating Redaction

- Re-identification attacks
- K-anonymity
 - Goal: Set a lower bound on uncertainty
 - Method: Suppress identifying attributes
 - Challenges: weak ties, future data
- Differential privacy
 - Goal: public data for preliminary analysis
 - Method: preserve approximate statistics (e.g., IDF)
 - Challenges: low counts, high dimensionality

Consequences for IR Evaluation

- 20th Century IR Evaluation
 - Perfect retrieval is our ideal
 - All and only the relevant documents
 - We measure how close to the ideal we come
 - Precision, Recall, MAP, NDCG, ...
- 21st Century IR Evaluation
 - Maximizing access to relevant content is our goal
 - Protecting sensitive content expands search space

Conclusion

- IR systems are too good
 - They could find what they shouldn't
 - To prevent this, we do not index some things
 - What we don't index, we can't find
- Two lines of research are needed
 - Techniques for search among secrets
 - Techniques for evaluation search among secrets

SHHHH!



ShhhCLEF

2017?

For More Information

- FnTIR survey on E-Discovery
 - <http://ediscovery.umiacs.umd.edu>
- Stanford ePADD (Email Redaction)
 - Demo: <http://epadd.stanford.edu/>
 - Code: <https://library.stanford.edu/projects/epadd>
- The Redaction Engine
 - <http://www.history-lab.org/>