# InVEST: Intelligent Visual Email Search and Triage

**Jay Koven, Enrico Bertini, Luke Dubois and Nasir Memon**

**NYU TANDON SCHOOL Of ENGINEERING CSE Department**

# Presentation Outline

- Motivation
- Related Work
- New Methodology
- Initial Results
- Future Research

# Motivation

Email account have grown drastically

- Free Gmail accounts are now 15GB

- Average Gmail account has >10,000 emails

Investigative datasets are even bigger

- Millions of emails

Users and Investigators are becoming overwhelmed

- Finding specific information requires persistence

- Finding a "Evidence" much worse

# Motivation

Email search methods haven't changed in 30 Years

- Grep / REGEX search still the core

- Long lists of results or worse no results

Machine learning is not the complete answer

- Emails are short and can be cryptic

- Techniques tend to work only on limited problems

- No clues to what was not found

# Investigation Not Search!

- Starts with some knowledge but incomplete
- Must find all emails related to investigation
  - Need to know more than just content
    - Fill in the "blanks"
      - Relationships between corresponders
      - Who or what were they talking about
      - When were they talking about it
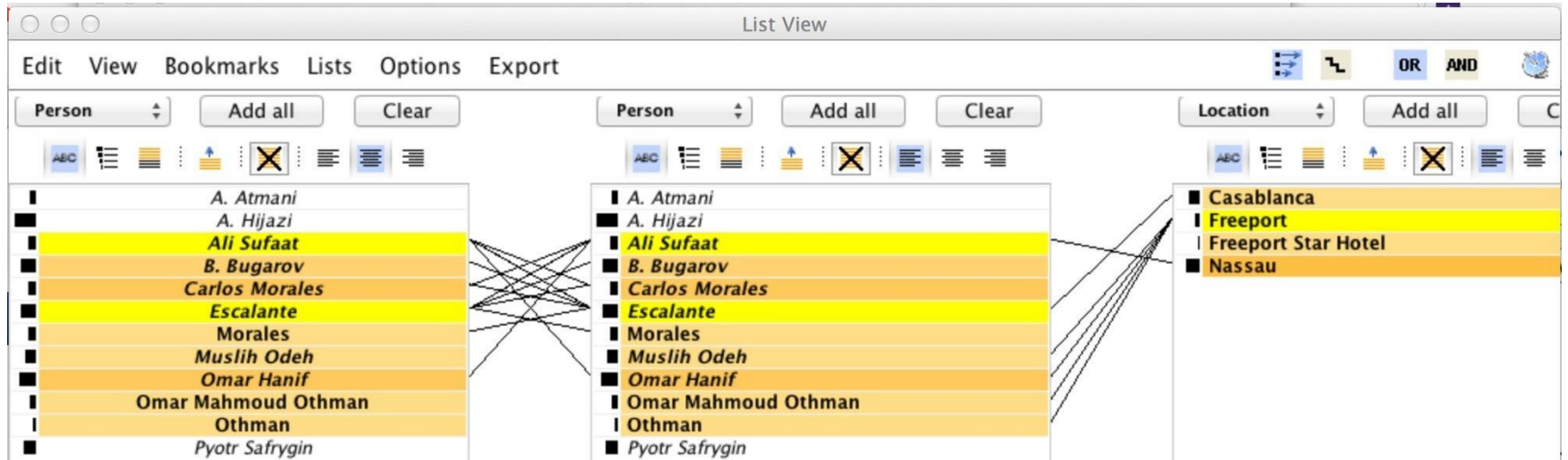      - Need to find hidden connections

# Related Research

- Machine learning analysis of text and email data
    - Unsupervised (LDA)
    - Supervised (sLDA)
- Forensic Visualizations
    - Mostly work on static data (Tableau)
- Visual Analytics
    - Overview, Zoom, Filter (Wrong direction)

# Related Work

## Jigsaw - Stasko et al

# What is Intelligent Visual Email Search and Triage
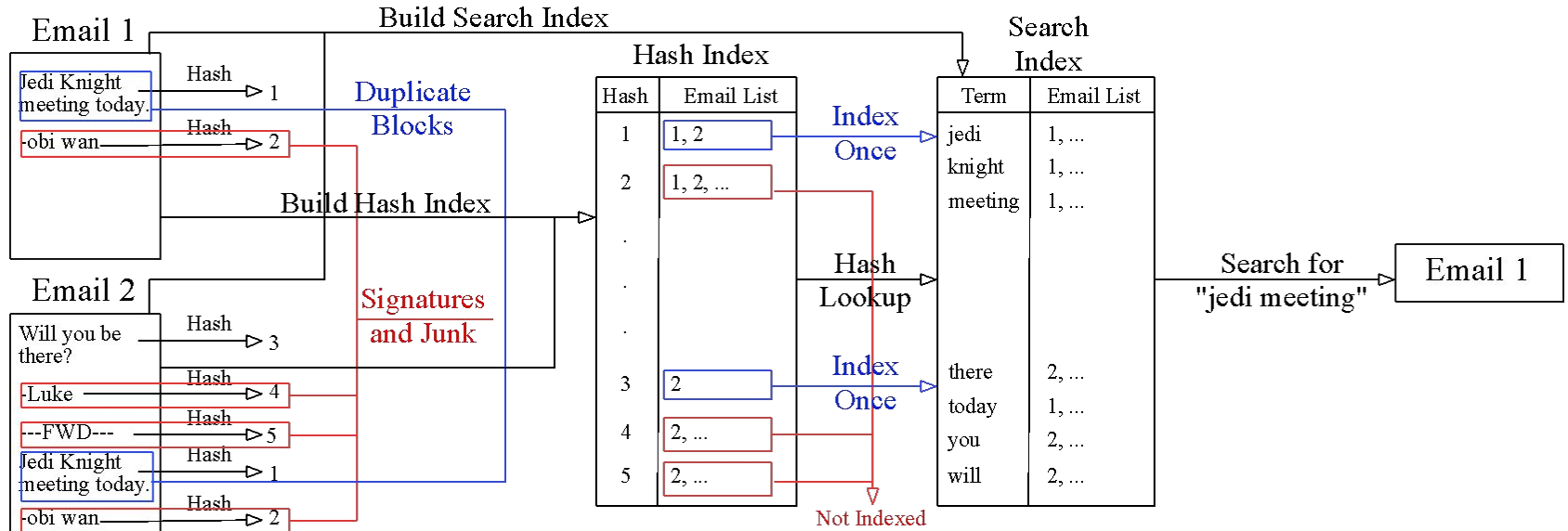
A Methodology to investigate email data sets
- Combine different types of information into clear results
  - Content (Subjects, Body Content and Entities)
  - Social Network (Corresponders)
  - Relationships between all of the above
- Find important entities and keywords
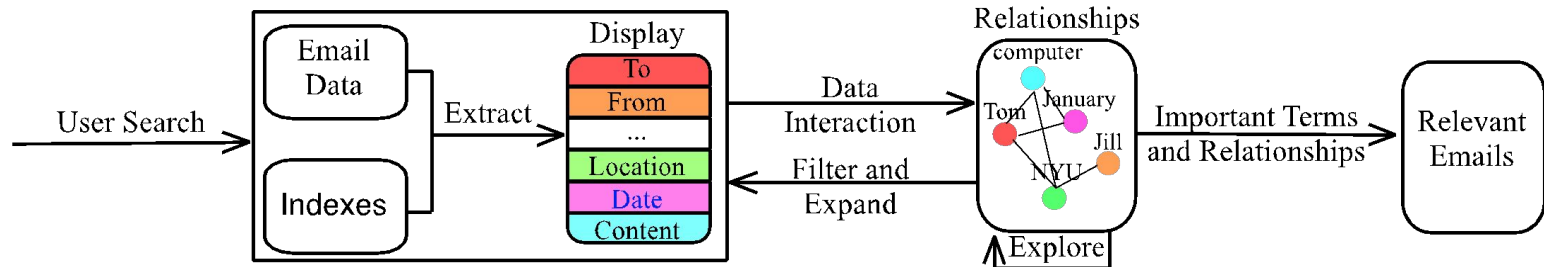  - Ranking
  - Guidance

# How?

- Create and effective visual analytic pipeline
  - Allow the analyst to explore the data interactively
  - Immediate feedback
  - Separate corresponders from Entities
  - Show relationships between keywords, entities and corresponders
  - Give intelligent guidance through ranking

# Preprocessing: Index and De-Junk

# Visual Analytic Pipeline

- The starting point
- Filter and Expand
- Interacting with the results
- Cross Links between displays

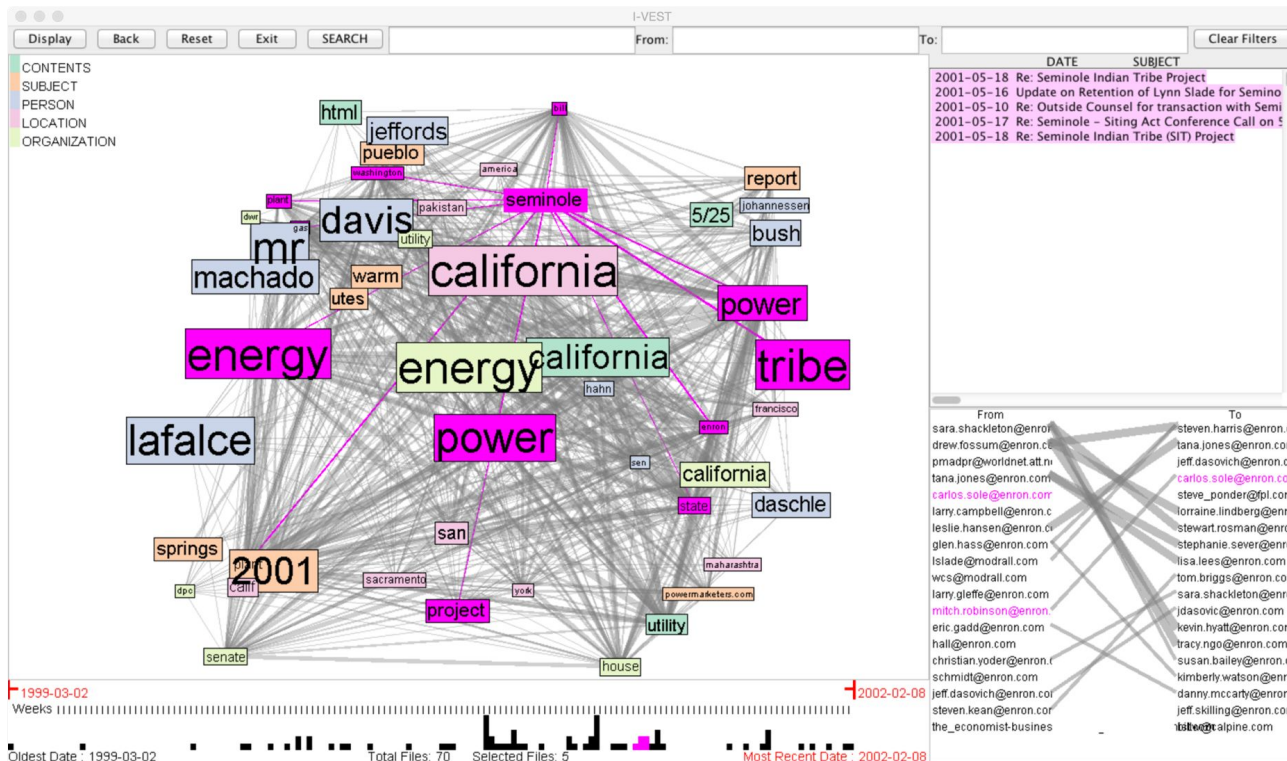# Entity Extraction and Separation

- Entity Integration
  - Why separate entities from corresponders?
  - Entity Extraction Methodology (Stanford NER)
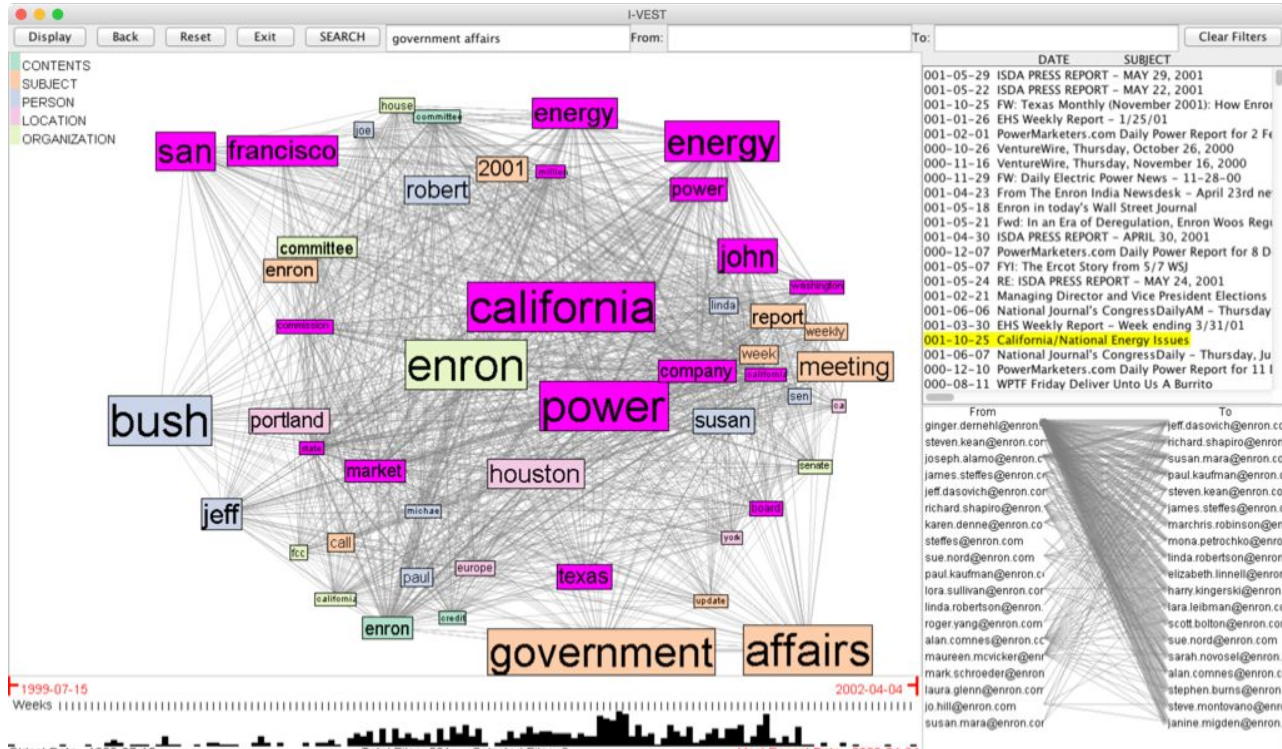  - Entity Identification (When is Tom also Thomas?)

# Intelligent Guidance

- Ranking of terms (Finding "Important")
  - For guiding analysts with list
  - Graph expansion
  - TF-IDF
    - Works (Sort of)
    - Problems
      - Should term frequency be measured by email or by returned set?
      - What makes a term important in Investigation?

# InVEST Views

# InVEST Views

# InVEST Future Research

- User experiments
  - HCI usability testing
  - Understand the users mental models
- Improve Ranking algorithms
  - Probabilistic approach using priors?
  - Interactively Driven Clustered (Faceted) Search
- User Defined Entities

# Acknowledgements

- My advisors and co-authors
- Anonymous Reviewers
- My Shepard Timothy Leschke
- NSF