# Email Mining Toolkit (EMT)

The Email Mining Toolkit (EMT) is a data mining system that computes behavior profiles or models of user email accounts. This toolkit is useful for report generation and summarization of email archives, as well as for detecting email security violations when incorporated with a real-time violation detection system, such as the MET system.

EMT, which includes approximately 13,200 lines of code, is implemented in Java providing a GUI implementing an interface to an underlying relational database application. It provides the means of loading, parsing and analyzing email messages from a wide range of storage formats. It not only demonstrates the statistics of email account behavior, it also computes the volume and velocity of emails exchanged between parties, analyzes specific content and patterns, and explores social relationships between groups of users, and the relative rankings of importance of different individuals in an organization.

Moreover, EMT extends these kinds of analyses to model "user behavior" at a very fine granularity. It models the behavior of individual user email accounts or groups of accounts, and can be used to detect changes in behavior that may be of interest in forensic analyses. These features of EMT provide the means to detect fraudulent misuse and attacks such as viruses and Spam (unwanted) email.

EMT includes 15 different features and models. The statistical models that include stationary and non-stationary user profile are used to generate user behavior models. These models include

• Message Table where individual emails may be automatically classified by built in machine learning subsystems,

• Usage Histogram revealing a user's typical daily email behavior,

• Similar Users which identifies groups of emails users who behave in similar ways ,

• Recipient Frequency providing a detailed analysis of the typical communicants with a user and

• Attachment Statistics detailing attached files serving as a personal file system of a user, as well as the statistical analyses including the birth rate, lifespan, incident rate, prevalence, threat, spread, and death rate useful in identifying interesting attachments and viral attachments.

The analyses built in to EMT concerning groups of accounts and their communication is provided to detect violations of group behavior. These models include

• Enclave Clique groups of users who frequently pairwise exchange messages,
• User Clique the set of accounts a particular user typically emails as a group,
• Email Flow revealing how a single message produces a web of new communication throughout an organization and
• Average Communication Time that views a user's typical response rates to individuals, indicating the relative importance of communicants.
• These models apply algorithms such as Chi Square, Hellinger Distance, Mahalanobis Distance, N-Gram analysis, Naïve Bayes classifier, TF-IDF categorization and graphical cliques analysis. By combining these features, EMT may be applied to a variety of applications and detection tasks.

EMT's graphical user interface provides an easy to use interface to execute these functions and that visualizes results in tabular form with displays of plots and histograms that are easy to understand.

**Related publications:**

Wei-Jen Li, Shlomo Hershkop, Salvotore J. Stolfo, *Email Archive Analysis Through Graphical Visualization*. ACM CCS VizSEC/DMSEC'04[ PDF ]

Salvotore J. Stolfo, Wei-Jen Li, Shlomo Hershkop, Ke Wang, Chia-Wei Hu, Olivier Nimeskern, *Detecting Viral Propagations Using Email Behavior Profiles*, *ACM Transactions on Internet Technology (TOIT)*, May 2004. [ PDF ]
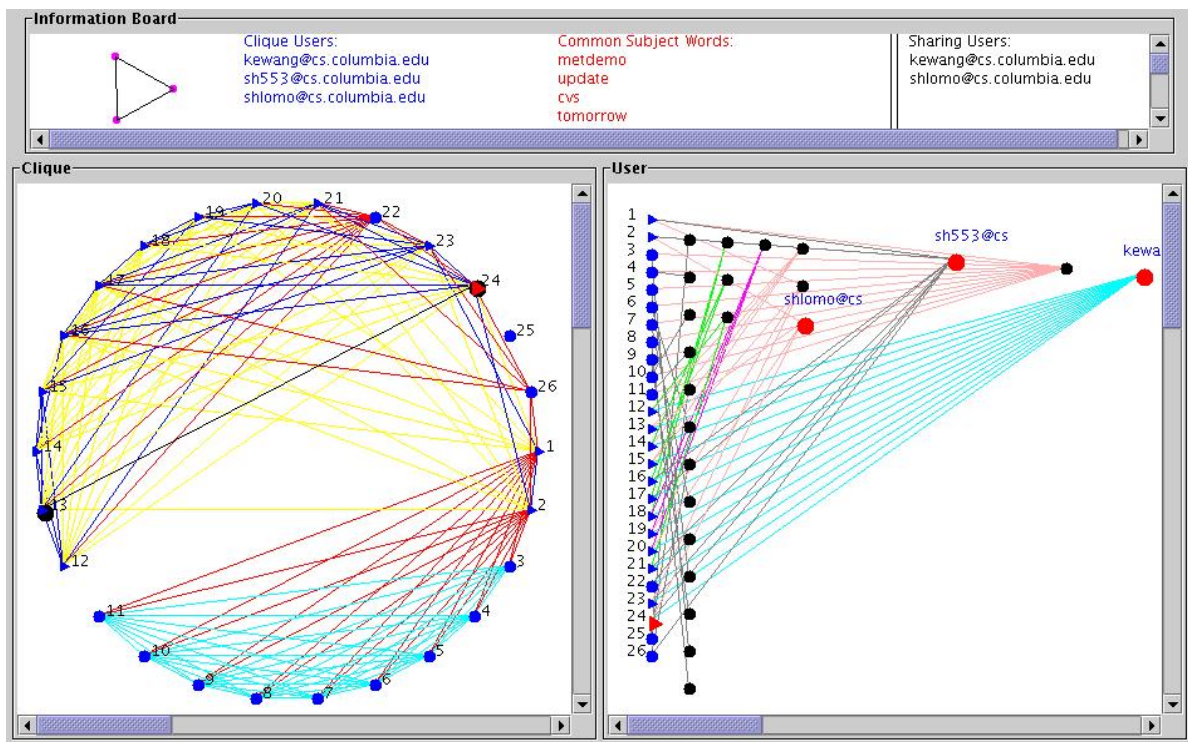
Some EMT screen shots are shown bellow:

General email client window / Machine Learning analysis

Graphical clique analysis

## Information Board

| | Clique Users: | Common Subject Words: | Sharing Users: |
|---|---|---|---|
| | kewang@cs.columbia.edu | metdemo | kewang@cs.columbia.edu |
| | sh553@cs.columbia.edu | update | shlomo@cs.columbia.edu |
| | shlomo@cs.columbia.edu | cvs | |
| | | tomorrow | |

Email flow analysis



## User & Email Table

Copy Email Table

| Index | msch... | aghosh | eeskin | cmich... | mgs28 | sal | shl |
|---|---|---|---|---|---|---|---|
| 0 | sender | rcpt | rcpt | rcpt | | | |
| 1 | rcpt | rcpt | rcpt | rcpt | sender | | |
| 2 | rcpt | rcpt | rcpt | rcpt | rcpt | sender | rcpt |
| 3 | rcpt | rcpt | rcpt | rcpt | rcpt | sender | |
| 4 | rcpt | rcpt | rcpt | rcpt | rcpt | rcpt | rcp |
| 5 | rcpt | rcpt | rcpt | rcpt | sender | rcpt | rcp |
| 6 | rcpt | rcpt | s & r | rcpt | rcpt | rcpt | rcp |
| 7 | rcpt | rcpt | rcpt | rcpt | rcpt | rcpt | rcp |
| 8 | rcpt | rcpt | rcpt | rcpt | sender | rcpt | rcp |
| 9 | rcpt | rcpt | rcpt | rcpt | rcpt | rcpt | rcp |
| 10 | rcpt | rcpt | rcpt | rcpt | sender | rcpt | rcp |
| 11 | rcpt | rcpt | rcpt | rcpt | rcpt | rcpt | rcp |
| 12 | rcpt | rcpt | rcpt | rcpt | rcpt | sender | rcp |
| 13 | rcpt | rcpt | rcpt | rcpt | rcpt | rcpt | rcp |
| 14 | rcpt | rcpt | rcpt | rcpt | rcpt | sender | rcp |
| 15 | rcpt | sender | rcpt | rcpt | rcpt | rcpt | rcp |
| 16 | rcpt | rcpt | rcpt | rcpt | rcpt | sender | rcp |
| 17 | s & r | rcpt | rcpt | rcpt | rcpt | rcpt | rcp |
| 18 | rcpt | rcpt | s & r | | rcpt | | rcp |
| 19 | rcpt | rcpt | rcpt | rcpt | rcpt | sender | rcp |

Subject: algorithm question

Content:
Hi guys,

## Email Flow

lick here for email exchange info

Similar Users

Usage Histogram

Usage Frequency Analysis

Virus simulation and detection

Virus detection