# Collaborative Social Network Discovery from Online Communications

## *Chris Diehl*

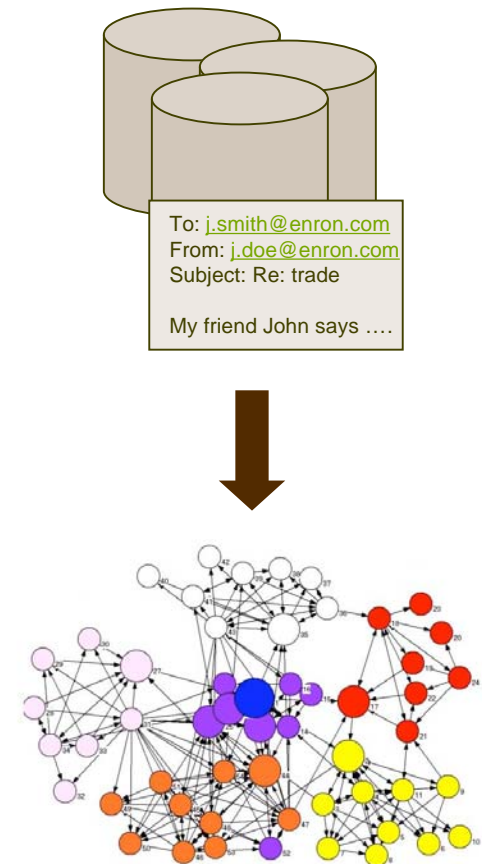### *USMA-ARI Network Science Workshop*

*Collaboration with Lise Getoor and Galileo Namata, University of Maryland – College Park*

**APL**

*The Johns Hopkins University*

**APPLIED PHYSICS LABORATORY**

# The Question

- **Organizations today utilize a number of communication channels**

  - Email, Instant Messaging, Text Messaging, Wikis, Blogs

- **Given access to an organization's online communications, how does one infer relationship and role types within the organization from the data?**

To: j.smith@enron.com
From: j.doe@enron.com
Subject: Re: trade

My friend John says ….

APL

# Data Attributes

- *Structured Data (Metadata)*
  - Sender and recipient(s), datetime
  - Can identify patterns of communication from metadata
  - Metadata provides no relationship context
- *Unstructured Data (Content)*
  - Message subject and body, attachments
  - Content may provide relationship and role information
  - Additional context may be needed to clarify the message
- *Goal is to exploit complimentary cues offered by the metadata and content*

APL

# Identifying Key Actors –
# A Motivating Example

From: Jennifer Fraser

Subject: john arnold bid for 20,000?

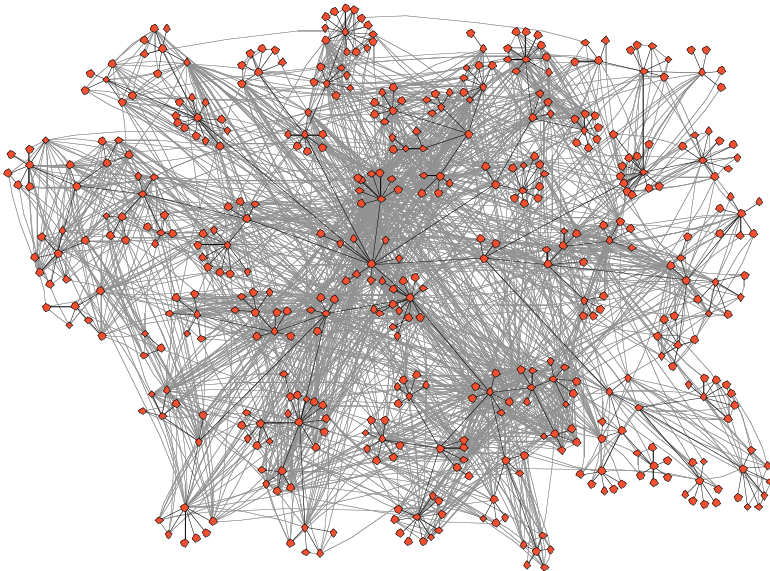**true?  and when do you plan on selling them?**

From: John Arnold

**exaggerations...word travels everywhere doesnt it?
how'd you hear?**

From: Jennifer Fraser

**johnny johhny johnny-- *there is no secrecy when
one is the king of ng* .. your brokers have the
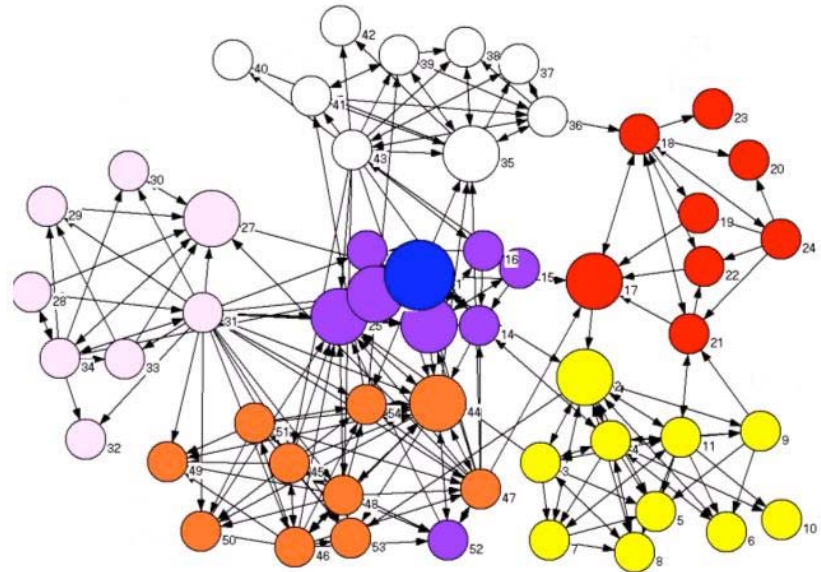biggest moves in the world…**

APL

# Representations: Data and Network

## Communication (Hyper)Graph

## Network (Hyper)Graph



*HP Labs Communication Graph*
*(Adamic and Adar, 2003)*

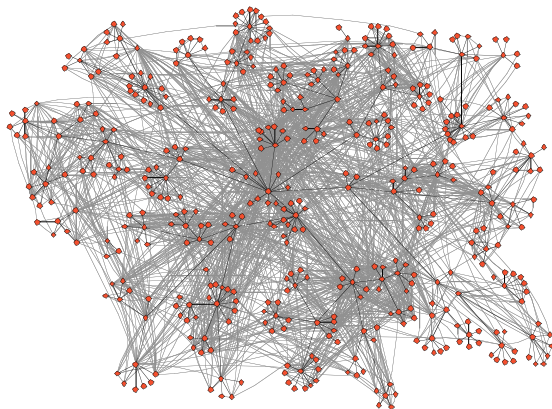**Nodes: Network References**
**Edges: Communication Events**

**Nodes: Entities**
**Edges: Social Relationships**
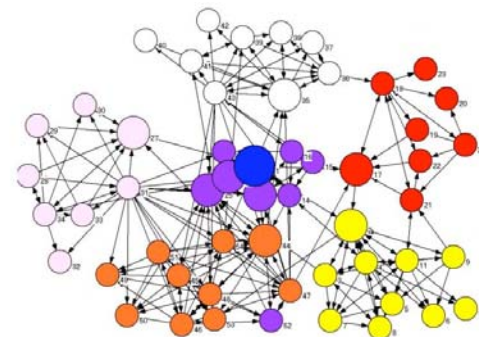
APL

# Collaborative Social Network Discovery

Communication Graph



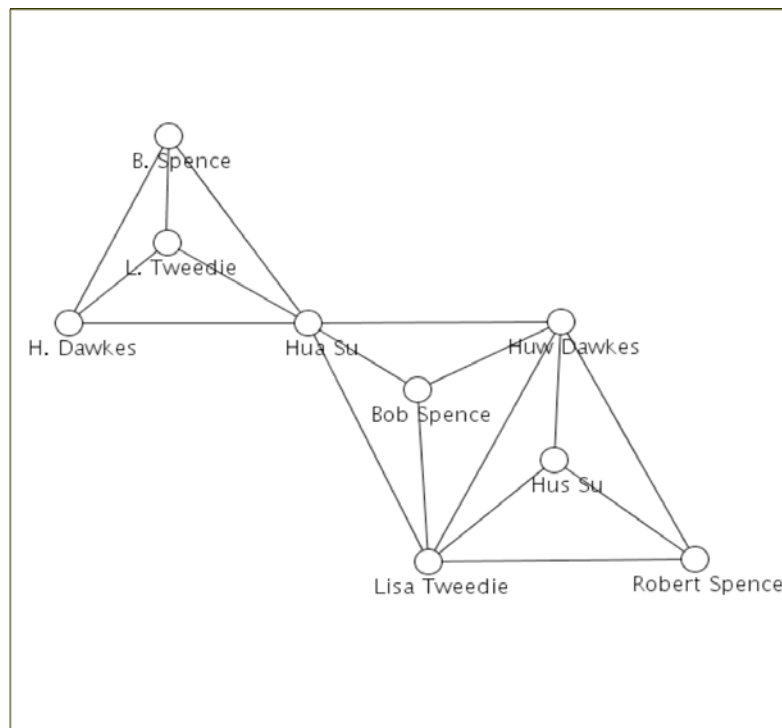Incremental Machine
Learning from Context

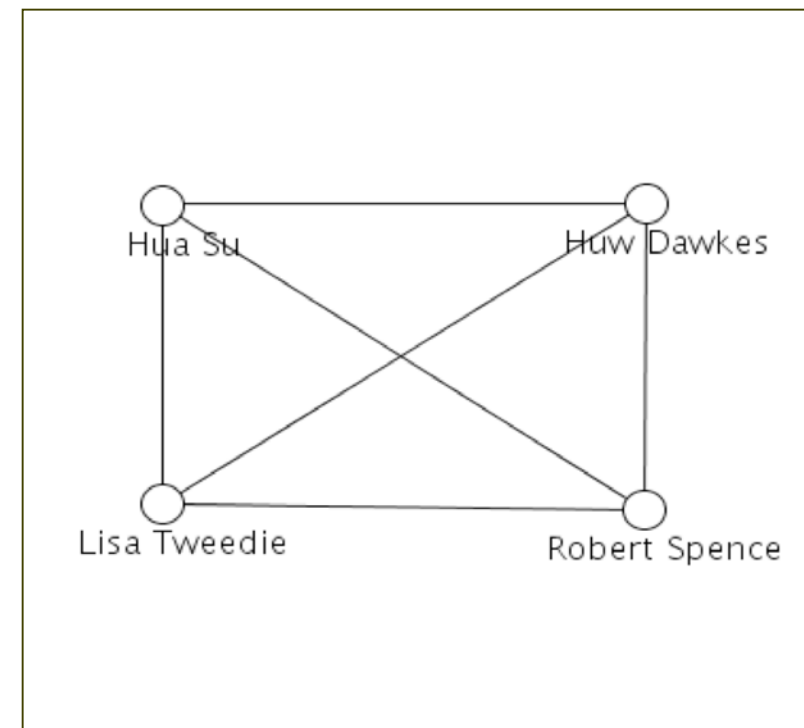Entity Resolution
Relationship Identification

Validated Network
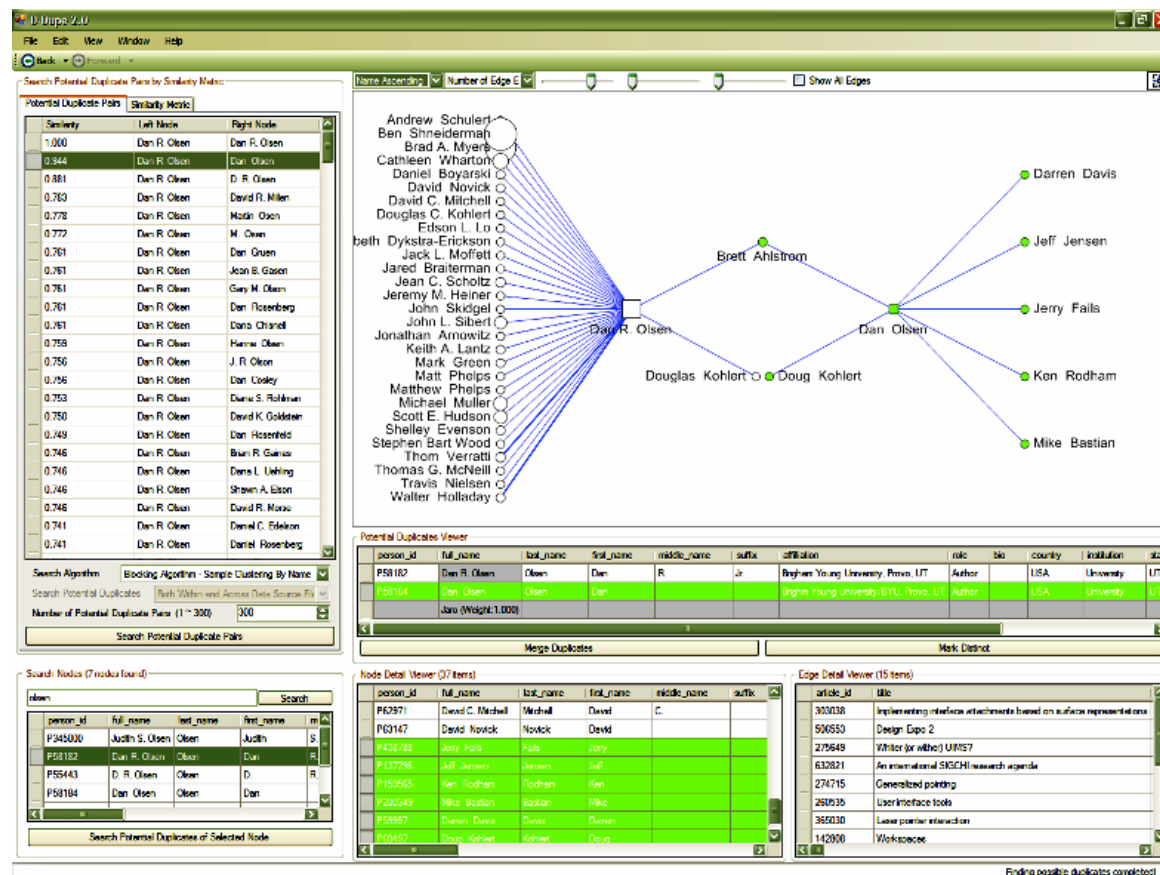
APL

# Entity Resolution: InfoVis Co-Author Network Fragment



Before

After

APL

# D-Dupe: An Interactive Tool for Entity Resolution



http://www.cs.umd.edu/projects/linqs/ddupe

APL

# Entity Resolution: Name and Network References

**Network References**

Datetime: 2001-01-23 09:45:00

Sender: sara.shackleton@enron.com

Recipients: tana.jones@enron.com

Subject: Hedge Funds

**Name References**

**Tana**: Other than your email attached, have you had other discussions with **Mark** or credit about hedge funds?  **Sara**

- Every individual has two classes of references

- To define an individual's identity and draw broader connections across emails, we need to first associate name and network references

Reference: C. P. Diehl, L. Getoor, G. Namata, "Name Reference Resolution in Organizational Email Archives," SIAM Data Mining 2006

APL

# Context Challenges

Datetime: 2000-06-19 09:52:00

Sender: tana.jones@enron.com

Recipients: marie.heard@enron.com

Subject: Just a tease!!!

Wouldn t you like to know which of the two Susan s gave her notice today
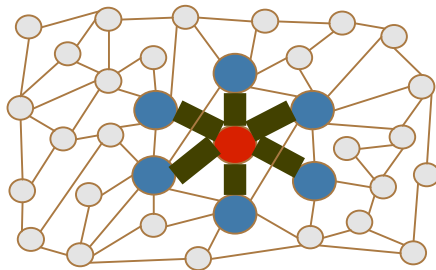
Datetime: 2001-02-28 09:32:00

Sender: liz.taylor@enron.com

Recipients: john.arnold@enron.com

Subject: Greg s Bill

Johnny, What does Greg owe you for the champagne? Is it $896.00? Liz

APL

# Relationship Identification - Incremental Ego Network Exploration



**Relationship Ranking**

| Rank | Relationship with Ego (Christian Yoder) |
|------|------------------------------------------|
| 1 | Elizabeth Sager |
| 2 | Richard Sanders |
| 3 | Steve Hall |
| 4 | Mark Haedicke |
| 5 | Dave Fuller |
| 6 | Tracy Ngo |

**Message Ranking**

| Rank | Message Subject |
|------|-----------------|
| 1 | Happiness |
| 2 | System Outage Risk |
| 3 | Mark Taylor Visit |
| 4 | Question about a deal we did |

**Evidence Discovery**

*From: Christian Yoder [christian.yoder@enron.com]*

*To: Elizabeth Sager [elizabeth.sager@enron.com],*

*Genia Fitzgerald [genia.fitzgerald@enron.com]*

*Subject: Happiness*

**Happiness is looking at the new legal org chart (which Jan just now dropped on my desk). I always approach these dry documents as though they were trigrams resulting from throwing the coins and consulting the I-Ching. At the top of the trigram which I find myself listed in I see a single name: Elizabeth Sager, and at the bottom I see the name Genia FitzGerald. ... cgy**
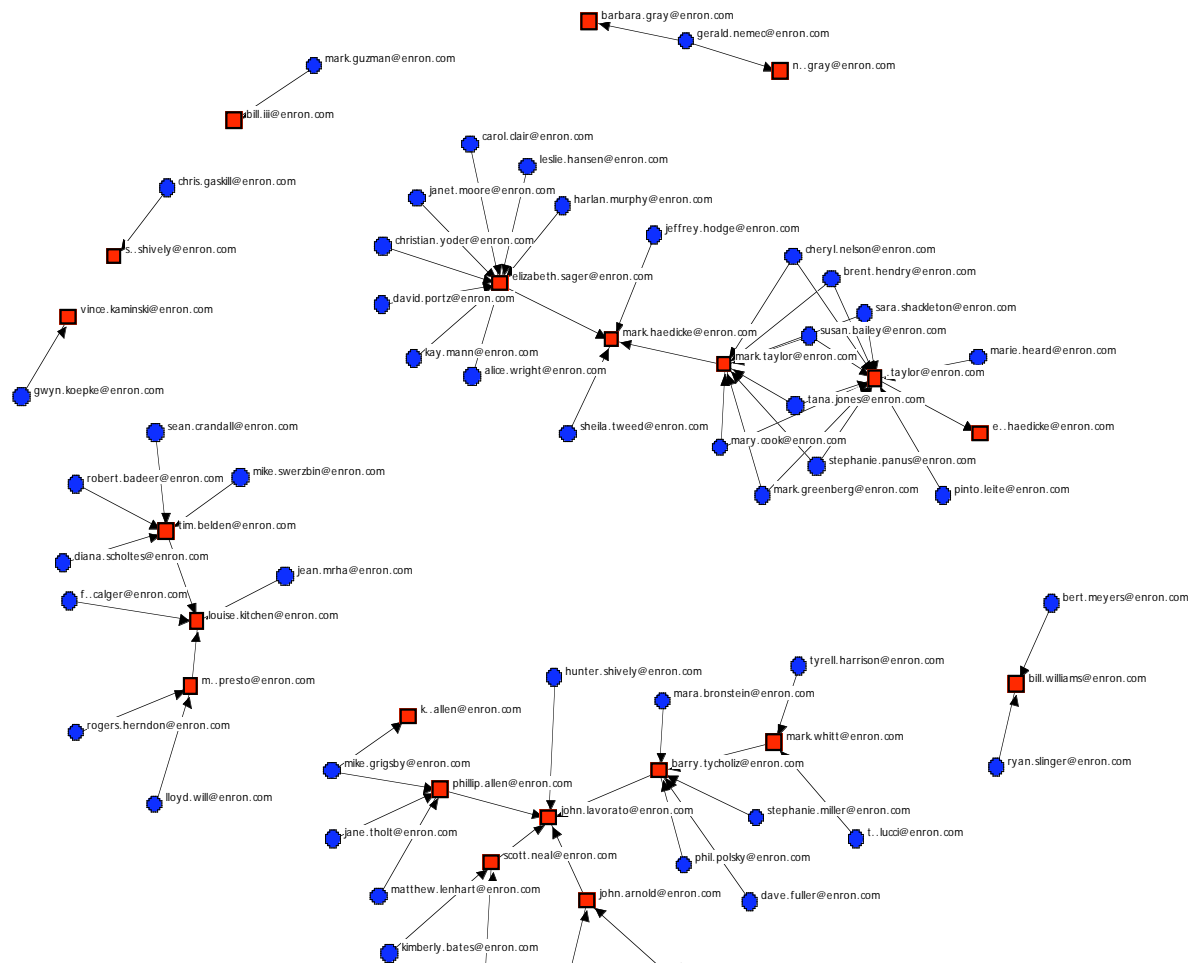
Reference: C. P. Diehl, G. Namata, L. Getoor, "Relationship Identification for Social Network Discovery," AAAI 2007

APL

# Enron Manager-Subordinate Communications Relationships

# Relationship Identification - Manager-Subordinate Relations

- **Preference Learning**
  - Supervised learning of relationship ranker
  - Given initial set of labeled ego networks
  - Ranking dyadic relationships
- **Traffic-Based Approach**
  - Message frequency
  - Number of recipients
  - Exchanges between relationship participants and common recipients
- **Content-Based Approach**
  - Term frequency vector for set of messages corresponding to the relationship
  - Exploits text from sender to recipient

| Approach | Mean Reciprocal Rank |
|---|---|
| Content-Based with Attribute Selection | 0.719 |
| Content-Based | 0.660 |
| Traffic-Based | 0.518 |
| Random Selection | 0.211 |
| Worst Case | 0.141 |

APL

# Future Directions

- **_Incremental, Active Learning_**
  - Relationship-Level and Message-Level Annotations
  - Automated Model Selection
  - Automated Feature Selection
- **_Visualization_**
  - Communications Graph Exploration
  - Network Graph Construction
- **_Interaction Paradigms_**
  - Unified Workflow for Entity Resolution and Relationship Identification

APL