# Email Social Network Extraction and Search

Michal Laclavík, Štefan Dlugolinský, Marcel Kvassay,
Ladislav Hluchý

# IKT Group - Institute of Informatics SAS

## Dept. of Parallel and Distributed Computing

Research and Development Areas:

- Large-scale HPCN and Grid applications
- Intelligent and Knowledge oriented Technologies

Experience from European IST projects:

- **3 project in** FP5: **ANFAS, CrosGRID, Pellucid**
- **6 project in** FP6: **EGEE II, K-Wf Grid, DEGREE (coordinator), EGEE, int.eu.grid, MEDIGRID**
- **4 projects in FP7: Commius, Admire, EGEE III, Secricom**

Several National Projects (SPVV, VEGA, APVT)

## IKT Group Focus:

- Information Processing
- Semantic Web
- Knowledge oriented Technologies
- Parallel and Distributed Information Processing

## Solutions:

- Ontea: Pattern-based Semantic Annotation
- ACoMA: KM tool in Email
- EMBET: Recommendation System

URL: http://ikt.ui.sav.sk

**Director & leader of PDC:**
Dr. Dipl. Ing. Ladislav Hluchý

# Outline

- Social Networks in Emails

- Ontea: Information Extraction

- Business objects in Email Communication

- Building of Email Social Network

- Spread of Activation

- Relation Identification

- Email Social Network Search

- User Interaction with Data

- Evaluation

# Motivation and Approach

**Motivation**

- To exploit information and knowledge included in email communication

**Approach**

- Social Network Extraction
- Entities extraction like People, Organizations, Locations, Contact data
- Forming semantic trees and graphs
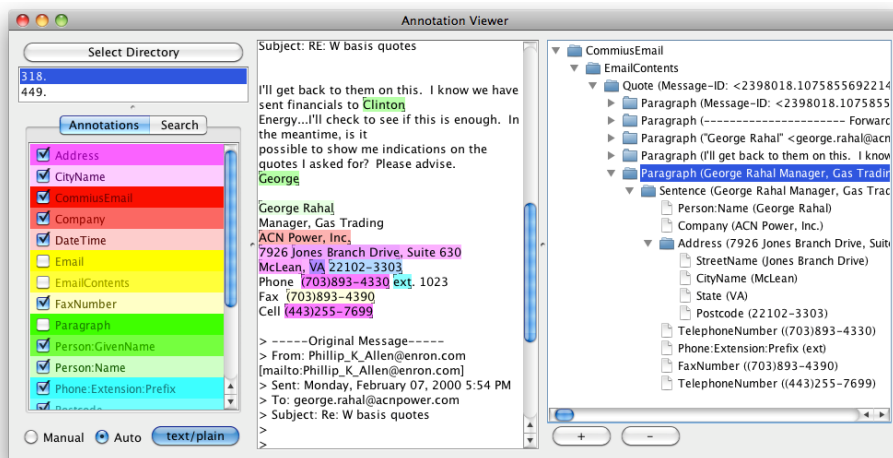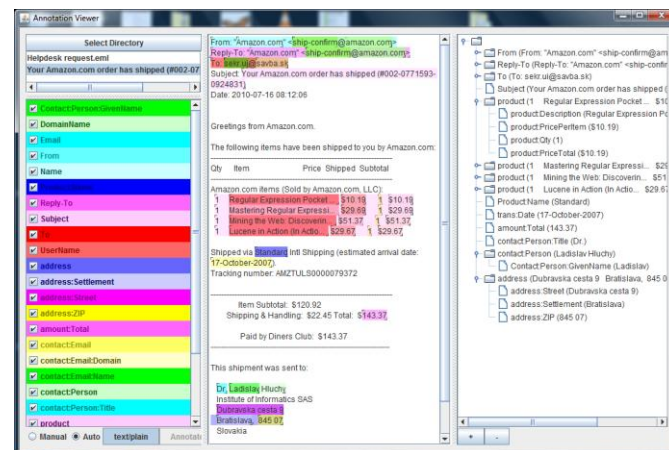- User interaction with graph data

# Email Social Networks



- Email Social Networks are less explored
  - **Several scientific publications: Apache mailing list, Enron, …**
  - **Commercial: Xobni (contacts and attachments)**
- Benefit
  - **Web Social Network Sites: owned by third parties**
  - **Email SN: owned by organization, individual or community**
  - **Additional level of interaction and context is present in emails**
- Information and Knowledge
  - **People, locations, contacts, product, services, attachments or links**
  - **Interactions**
  - **Time**
  - **Discovering relations can bring significant benefits**
  - **Spread of Activation – simple way to discover relations**
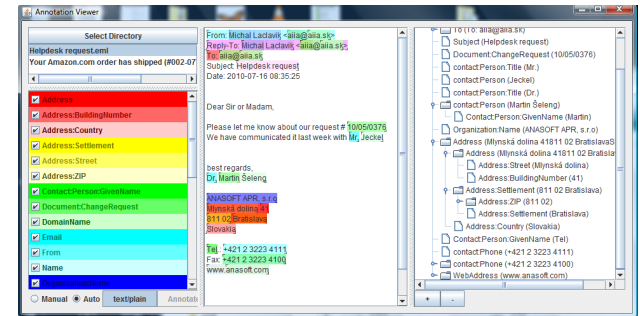
# Ontea: Information Extraction (Features)

- ❖ Regex patterns
- ❖ Visual Annotation Tool
- ❖ Integration with external tools
  - ❖ **GATE, Stemers, Hadoop …**
- ❖ Gazetteers
- ❖ IE System configuration
- ❖ Automatic loading of extractors
- ❖ Patterns
- ❖ Multilingual tests
  - ● **Spanish**
  - ● **Slovak**
  - ● **English**
  - ● **Italian**

# Business objects in Emails

- Study on 6 organizations show:
  - **Objects can be identified by patterns and gazeteers**
  - **It is possible to define set of common objects**
- Objects identified:
  - **Organization:**
    - org:Name, org:RegNo, org:TaxNo
  - **Person:**
    - person:Name, person:Function
  - **Contact:**
    - contact:Phone, contact:Email, contact:Webpage
  - **Address:**
    - address:ZIP, address:Street, address:Settlement
  - **Product:**
    - product:Name, product:Module, product:Component, product:BOID
  - **Document:**
    - doc:Invoice, doc:Order, doc:Contract, doc:ChangeRequest
  - **Inventory:**
    - inventory:ResID, inventory:ResType
  - **Other business object**
    - ID: BOID

**Acoma is not part of Paper but related to NextMail**

- **Useful hints with links are included in enriched email**
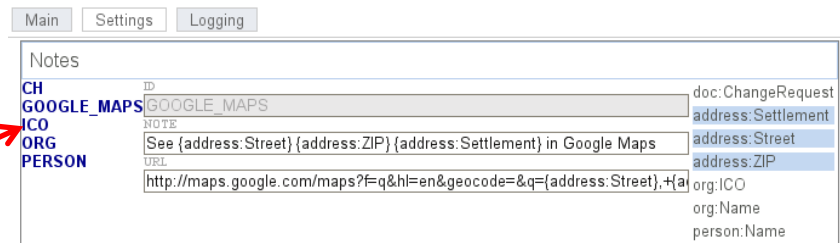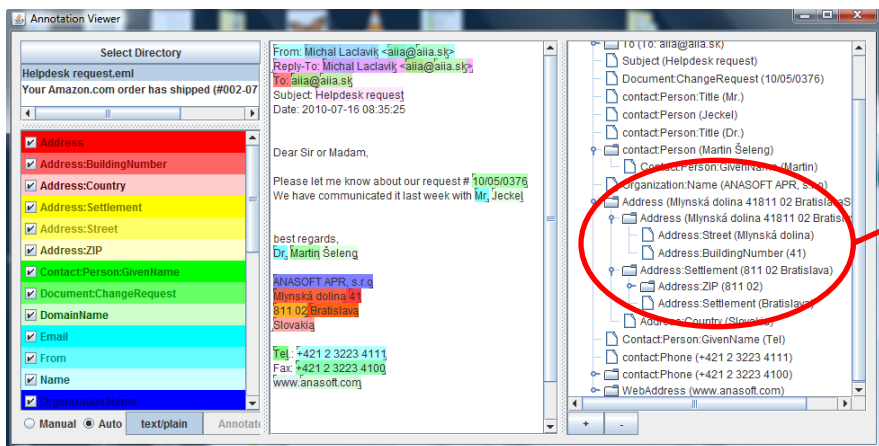
- **Links lead to internal or external systems (Internet, Intranet)**

# Acoma: Hint Recommendation

## Acoma is not part of Paper but related to NextMail

# Email Social Graph/Network

# Email Social Network Search: Features

- Social network of communicating people with relation to other entities

- Discovering relation in the graph/network using spread of activation

- Showing relations restricted to concrete type, e.g. telephone numbers related to a person

- User interaction with data (merging, deleting entities) with immediate impact on discovered relations

- Navigation over related entities

- Full-text search of the entities

- User interface for search

# GUI Features

# Algorithm and Evaluation

- All described in last year WI-IAT 2010 publication:
  - **Laclavik et al. Use of Email Social Networks for Enterprise Benefit, IWCSN 2010**

- Algorithm
  - **breadth-first**
  - **Node fires only once**

- Information Extraction Evaluation
  - **Evaluation on set of 50 Spanish emails**
  - **Strict match 50-90%**
  - **Intersect match 80-94%**

- Spread of Activation (relevance identification) Evaluation
  - **50 Spanish emails (phone/name):**
  - **Precision 60% (due to lower recall in IE)**
  - **Precision 85% (achievable with better IE)**
  - **self-healing (with new incoming emails)**
  - **28 English emails: precision 77%**

# Performance evaluation

- Focus of this paper
- Experiments with 5 different sizes of dataset:
  - **Number of visited nodes grows too high**
  - **Number of fired nodes grows acceptable**
  - **Search time ~ visited number of nodes**
  - **Scalability not possible with current implementation but achievable**

| | | | | | |
|---|---|---|---|---|---|
| **Number of Mailboxes** | 1 | 5 | 7 | 10 | 15 |
| **Number of Emails** | 3 033 | 9 939 | 20 521 | 36 532 | 50 845 |
| **Number of Verticles** | 41812 | 159 776 | 369 932 | 608 146 | 835 025 |
| **Number of Edges** | 98566 | 380 254 | 971 929 | 1 796 403 | 2 514 031 |
| **Processing time (ms)** | 81 672 | 430 025 | 1 199 463 | 1 948 847 | 2 680 171 |
| **Processing time (minutes)** | 1 | 7 | 20 | 32 | 45 |
| **One Email processing time** | 27 | 43 | 58 | 53 | 53 |
| | | | | | |
| **Person:Name=>Mike Grigsby** | | | | | |
| **Search Response Time** | **144** | **446** | **758** | **1 396** | **1 696** |
| Results | 344 | 463 | 494 | 781 | 761 |
| Fired | 6 363 | 20 732 | 19 045 | 23 466 | 23 839 |
| Visited | 112 280 | 281 060 | 476 324 | 939 642 | 1 174 400 |
| Visited Unique | 18 382 | 53 772 | 82 219 | 145 192 | 178 829 |
| **Search Slowed down x Times** | **1** | **3,1** | **5,3** | **9,7** | **11,8** |
| **Fired x Times** | **1** | **3,3** | **3,0** | **3,7** | **3,7** |
| Number of messages x Times | 1 | 3,3 | 6,8 | 12,0 | 16,8 |
| Number of verticles x Times | 1 | 3,8 | 8,8 | 14,5 | 20,0 |
| Number of edges x Times | 1 | 3,9 | 9,9 | 18,2 | 25,5 |
| | | | | | |
| **TelephoneNumber=>713 780-1022** | | | | | |
| **Search Response Time** | **5** | **8** | **8** | **12** | **13** |
| Results | 4 | 4 | 4 | 4 | 4 |
| Fired | 116 | 150 | 157 | 181 | 183 |
| Visited | 6 318 | 8 776 | 9 550 | 13 424 | 14 710 |
| Visited Unique | 698 | 954 | 1 059 | 1 424 | 1 513 |
| **Search Slowed down x Times** | **1** | **1,5** | **1,6** | **2,3** | **2,5** |
| **Fired x Times** | **1** | **1,3** | **1,4** | **1,6** | **1,6** |
| Number of messages x Times | 1 | 3,3 | 6,8 | 12,0 | 16,8 |
| Number of verticles x Times | 1 | 3,8 | 8,8 | 14,5 | 20,0 |
| Number of edges x Times | 1 | 3,9 | 9,9 | 18,2 | 25,5 |
| | | | | | |
| **Address=>6201 Meadow Lake, Houston, TX 77057** | | | | | |
| **Search Response Time** | **7** | **14** | **28** | **40** | **59** |
| Results | 23 | 38 | 71 | 91 | 170 |
| Fired | 236 | 515 | 701 | 896 | 1 546 |
| Visited | 8 134 | 15 571 | 32 336 | 40 563 | 58 571 |
| Visited Unique | 1 097 | 1 952 | 6 526 | 8 029 | 11 295 |
| **Search Slowed down x Times** | **1** | **2,1** | **4,3** | **6,0** | **8,9** |
| **Fired x Times** | **1** | **2,2** | **3,0** | **3,8** | **6,6** |
| Number of messages x Times | 1 | 3,3 | 6,8 | 12,0 | 16,8 |
| Number of verticles x Times | 1 | 3,8 | 8,8 | 14,5 | 20,0 |
| Number of edges x Times | 1 | 3,9 | 9,9 | 18,2 | 25,5 |
| | | | | | |
| **Email=>ina.rangel@enron.com** | | | | | |
| **Search Response Time** | **106** | **552** | **1 162** | **2 156** | **3 017** |
| Results | 732 | 1 764 | 2 668 | 2 809 | 2 952 |
| Fired | 5 165 | 16 062 | 17 629 | 19 716 | 20 997 |
| Visited | 91 199 | 369 584 | 865 300 | 1 694 065 | 2 326 867 |
| Visited Unique | 13 355 | 54 987 | 81 757 | 134 876 | 168 955 |
| **Search Slowed down x Times** | **1** | **5,2** | **11,0** | **20,3** | **28,5** |
| **Fired x Times** | **1** | **3,1** | **3,4** | **3,8** | **4,1** |
| Number of messages x Times | 1 | 3,3 | 6,8 | 12,0 | 16,8 |
| Number of verticles x Times | 1 | 3,8 | 8,8 | 14,5 | 20,0 |
| Number of edges x Times | 1 | 3,9 | 9,9 | 18,2 | 25,5 |

# New Developments not included in the paper

- Faster algorithm
- Takes graph topology into account
- Breadth First
- Ends after it visit certain number of nodes (set to 10,000 experimentally)

- Gives similar results as original algorithm
- Possibility for improvements:
  - **It should take edge and vertex weight into account**
  - **Ignores multiple edges between nodes**

```java
private void computeRelatedBreadthFirst(Result start) {
    LinkedList<Result> rLL = new LinkedList<Result>();
    rLL.addLast(start);
    int count = visitNodeCount;
    rM.put(start, (double) count);
    vNodes++;
    while (!rLL.isEmpty() && count >= 0) {
        Result r = rLL.removeFirst();
        visited.add(r);
        int nCount =  g.g.getNeighborCount(r);
        double v = rM.get(r)/(double)nCount;
        if (v < threshold) //if value is to low we do not activate more
            continue;

        if (nCount<=count) {
            Collection<Result> rC = g.g.getNeighbors(r);
            for (Result result : rC) {
                if (!visited.contains(result)) {
                    rLL.addLast(result);
                }
                visited.add(result);
                double val = v;
                if (rM.containsKey(result))
                    val += rM.get(result);
                rM.put(result, val);
                vNodes++;
            }
            count -=nCount;
        }
    }
}
```

# Conclusion

- Email Archives
  - **Valuable source of knowledge**
  - **Hidden Social Networks owned by Enterprise or Individual**
  - **Information Extraction and Social Network Analysis can help**

- Experiment
  - **Pattern-based Information Extraction**
  - **Social Network Extractor**
  - **Spread of Activation**
  - **Scalable Relation identification with acceptable success rate**

- Applications
  - **Recommendation and Search in Emails**
  - **Population of Databases (Cold start problem)**
  - **Possibility to extend social network graph with processed document repositories and other business data**
  - **Business Intelligence and Knowledge Management**