# Data Sets for Large-Scale Social Network Analysis

## Christopher DuBois

Ph.D. Student

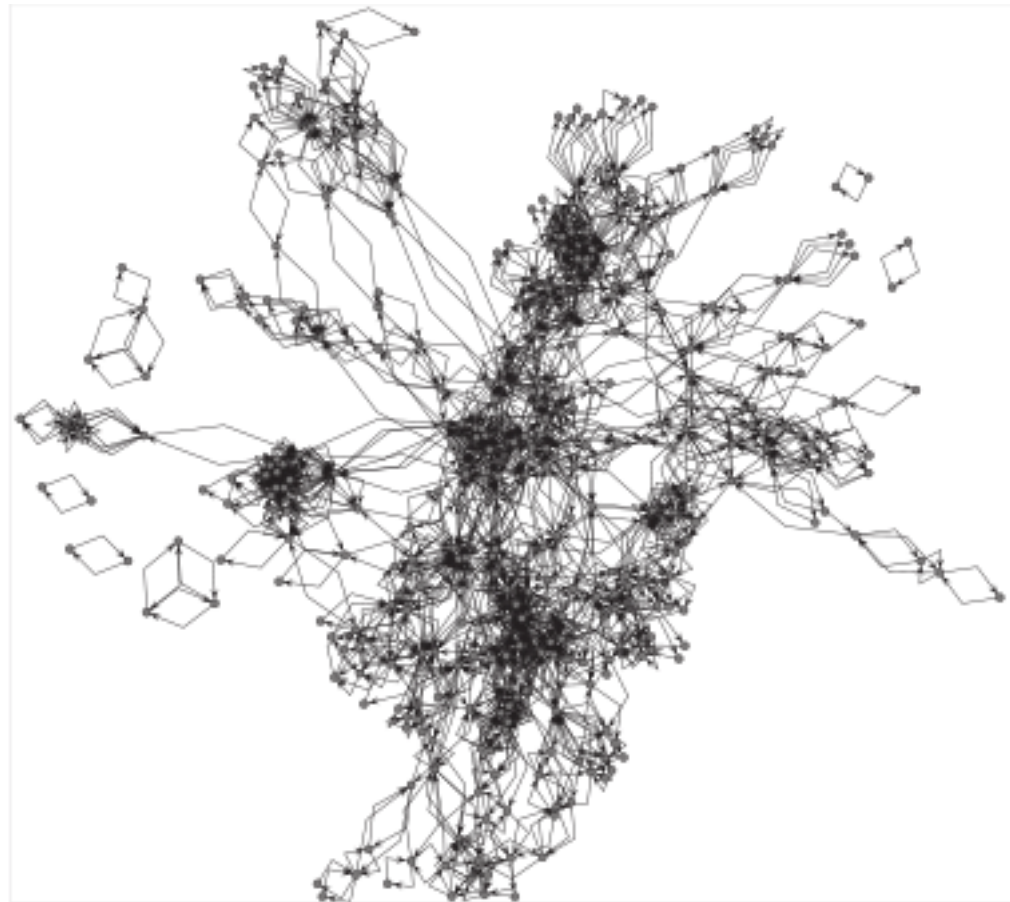Department of Statistics, UCI

# Overview

Data sets from online interaction

- Email, wikis, social networking sites, etc.

- Motivations for these data sets

- Problems applying current methods

- Opportunities

# Enron Emails

- Collection of anonymized emails
- Over 1000 employees
- 2 years

- Modeling challenges
- Baseline for comparison



O'Madadhain, Hutchins, Smyth. SIGKDD, 2004.

# Blogosphere Data

- 12 million blogs
- 40,000 posts per day


- Modeling challenges:
  - exact temporal info

Spinn3r.com

# Wikipedia

- 10 million articles
- 5.8 million users
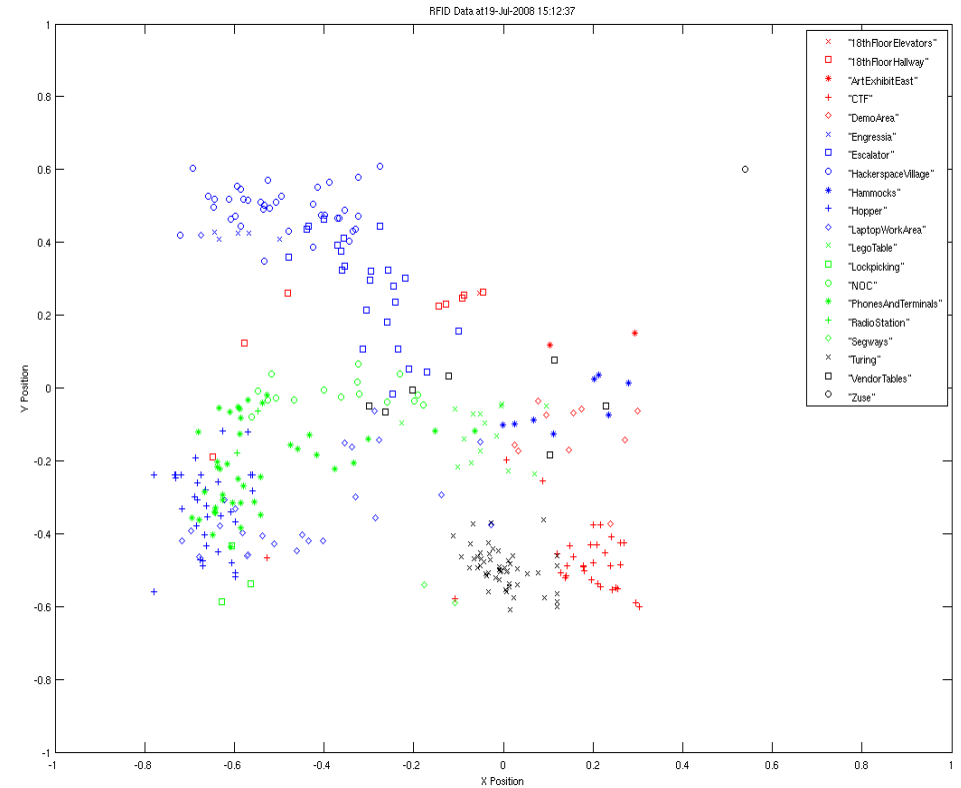- (~500,000 highly active)
- Since 2001

Not only relationships, but also text and time



Visualization for edits to the Wiki-article "Evolution" (IBM Research)

# Conference Tracking

- 3 day conference
- 800 attendees
- position data every 30 seconds via RFID tags
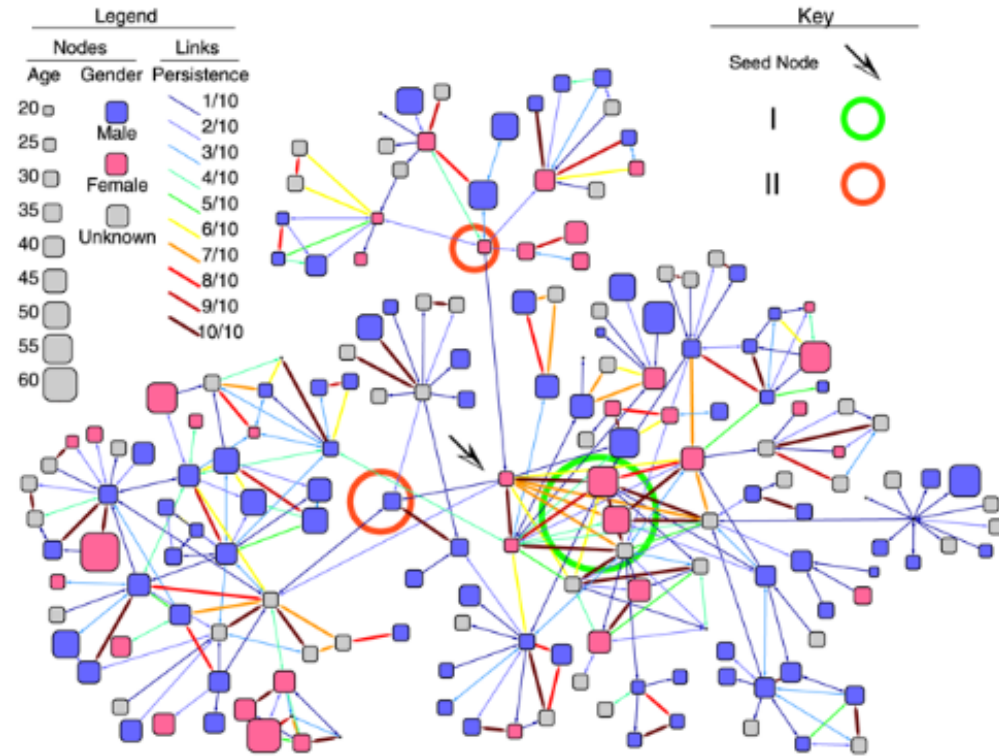- empirical face-to-face network



Additional attributes to consider

AMD Hope Team, 2008
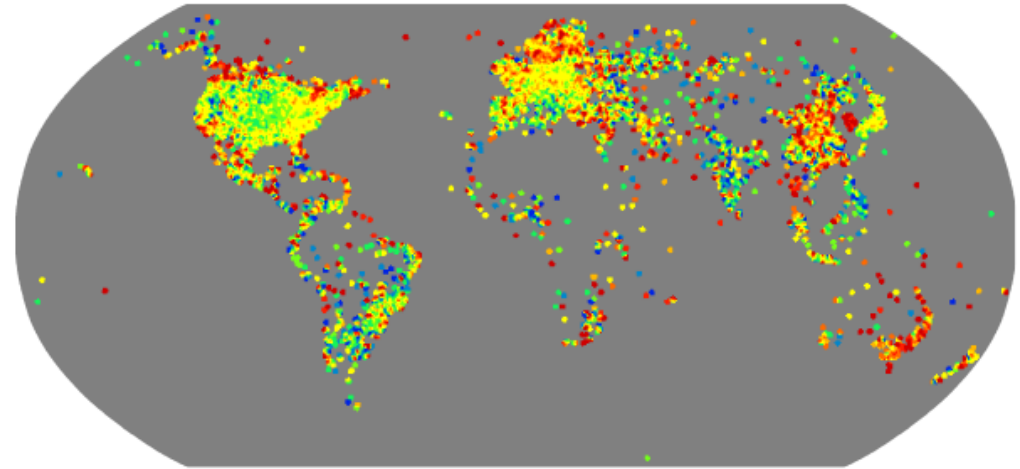
# Mobile phone calls

8,000,000 calls among more than 1,000,000 people.



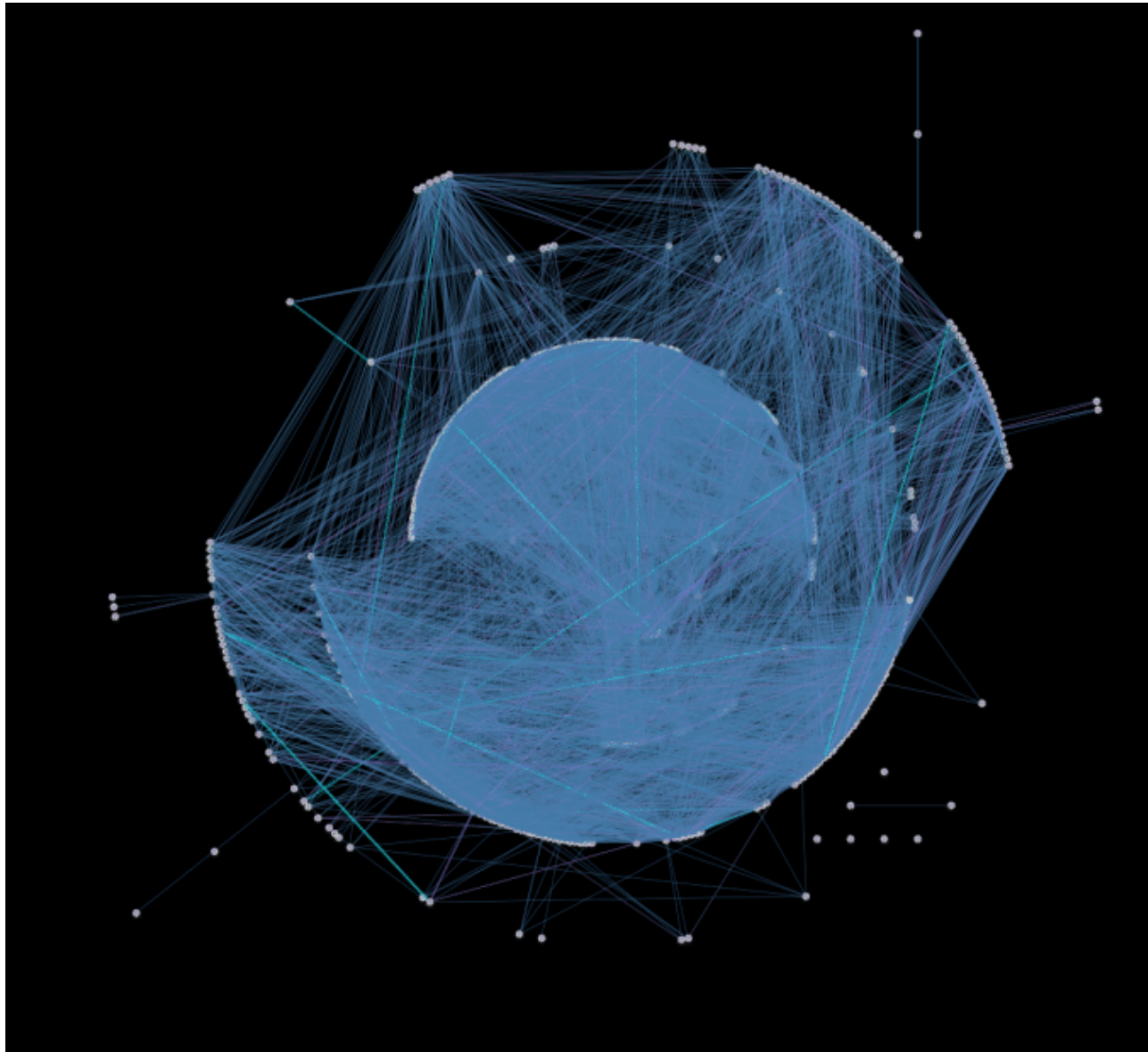Hidalgo, Rodriguez-Sicker. The dynamics of a mobile phone network. Physica A 2007.

# MSN Messenger

- Data from the month of June, 2006
- 180 million users
- 1.3 billion conversations



Leskovec, Horvitz. Planetary-Scale Views on a Large Instant-Messenging Network.  WWW 2008.

# Social Networking Sites



My Facebook network. (Nexus)

# Wrapping up...

- Prevalence of data regarding online interaction
- Large populations, small timescales, many attributes

- Need for models and methods that can leverage these massive data sets for studying theories of large social systems

# Open Access Network Data

- Browse, download network data
- Publicly available
- Facilitate network research