# A Geographically Searchable Local News Aggregator

## IPRO 327
## Fall, 2008
## Advisors:  Frieder, Yee

# Goal

- Create a news search engine that understands the contents of articles
  - Allows search by "meaning" versus by "string matching"

  - Example: What happened in Chicago in 1996?

# Current Practice – nytimes.com

# Current Practice – Google

# Why News Articles?

- Contains reasonably factual content
- Articles printed regularly
- Articles have relatively standard format

# Target Audience

- Researchers (e.g., reporters)
- Intelligence agents
- Real estate agents
- Marketers
- ...

# Semester Goals

- Design the system
- Build a working prototype
- Lay a foundation for future work

- Focus – search based on time and location

# Approach

- Consider functional units of system
- Consider IPRO goals

- Make a team for each functional unit
  - Consider strengths/interests of students

# Other Tasks

- Required everyone to do research and present on related topics, such as the semantic web
  - Developed soft skills (reading, writing, presenting)
  - Learned about the problems and benefits of building such a system

# Students

- Nick Bathum (CS, 4th yr)
- Evan Estola (CS, 3rd yr)
- Jaeyeon Kihm (CS, 4th yr)
- Jay Mundrawala (ECE, 2nd yr)
- Yacin Nadji (CS, 4th yr)
- Chris Osswald (BME, 3rd yr)
- Pete Schmitz (CS, 4th yr)
- Cameron Zangenehzadeh (PS, 4th yr)

# Organization

- Developed Plan, Timeline
- Used a free Wiki service for discussion, file-sharing and status reporting
  - Centralized information
  - Avoid mass emails and information loss

# Functional Units

- Design (User Interface)
  - Usability and presentability
- Web Application Architecture
  - Build infrastructure, pull together work of others
- Entity Database
  - Pull "machine readable" information from articles
- Web Crawling
  - Collection of news articles from the Web

# Functional Unit Relationship



13

# Architecture Group

## Evan Estola
## Yacin Nadji

| User Interface |
| Web Application |
| Database |
| Crawler |

# Architecture Group

- Responsible for building Web infrastructure
  - Interconnect other functional units
- Combine work of other groups

# Architecture Approach

- Use leading open source software
  - Sesame semantic database for entities
  - MySQL database for crawled article data
  - Java Server Pages, Java Beans and Tomcat

# Architecture Challenges

- Little to no Web development experience
- Working with other groups
  - Dependency on other groups
  - Communication

# Architecture Outcome

- Success!
  - Working prototype
  - Dynamically builds interface (from Design team) to display article data (from Crawling team) and entities (from Enitity Extraction team)

  - Learned to learn

# Architecture Future Work

- Optimization (speed)
- Build entity display system

# Entity Extraction Team

Jaeyeon Kihm

Jay Mundrawala

User Interface

Web Application

Database

Crawler

# Entity Extraction - Goals

- Given news articles, identify
  - Persons
  - Locations
  - Organizations
  - Dates
- Build database for extracted data

# Entity Extraction – Entity Recognition

Dr. Wai Gen Yee is a professor in Department of Computer Science at Illinois Institute of Technology in Chicago

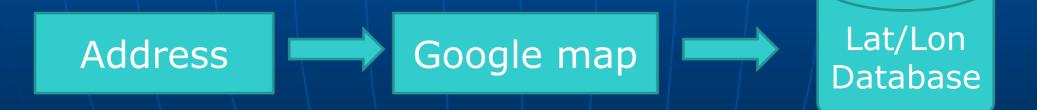# Entity Extraction – Entity Recognition

<Person>Dr. Wai Gen Yee is a professor in <Org>Department of Computer Science at <Org>Illinois Institute of Technology in <Location>Chicago

- Implemented leading open source software (Stanford NER API)
- Person, Location, Organization, Date, Title entities are extracted successfully

# Lat/Lon Coordinate Point

- Generate latitude / longitude coordinate points from location entity.

| Address | → | Google map | → | Lat/Lon Database |

# Entity Extraction - Accomplishments

- Accomplishment:
  - Extract entities from articles
  - Build database to store entities and news articles
  - Generate Lat/Lon coordinate points from address

# Challenges

- Stanford NER does not recognize addresses
  - Wrote custom module
- Poor database performance
  - Added a second custom database

# Entity Extraction – Future works

- Future work
  - Unification and disambiguation of entities
    - Ex: "George Bush" vs. "President Bush"
- Ranking of retrieved results
- Support for different types of entities (e.g., sports teams)
- Support for relationships among entities

# Web Crawling Team

Nick Bathum

Pete Schmitz

| User Interface |
|:---:|
| Web Application |
| Database |
| Crawler |

# Web Crawling Goal

- Retrieve news articles from the Web
- Submit news articles to Entity Extraction system

# Web Crawling Approach

- Select correct software for crawling
  - Nutch selected for its support and scalability
- Determine which news source and how to acquire correct article data

# Web Crawling Approach

- Get article links
- Use Nutch software to fetch the articles
- Parse relevant content from articles

# Web Crawling Challenges

- Automated extraction of relevant content from Web pages filled with extraneous info
  - Need to identify Web page structure

# Web Crawling Outcome

- Functioning article retrieval system
- Retrieve articles from
  - Reuters
  - LA Times

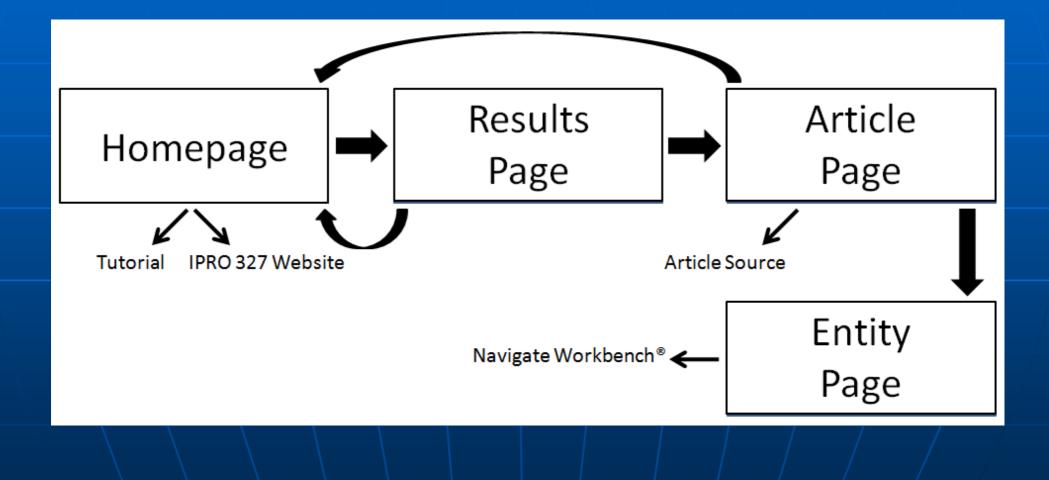# Web Crawling Future Work

- Support for other news sources

# Design Team

## Chris Osswald (BME)
## Cam Zangenehzadeh (PS)

```
User
Interface
   ↕
Web
Application
   ↕
Database
   ↑
Crawler
```

# What do we want to do?

- Show off implemented functionality
- Design a simple, user-friendly interface

# Approach

- Looked at other systems (EveryBlock.com, PowerSet.com, etc.)

- Identify core functionality

- Developed appropriate interface
  - Shows core function of system

# Flow of Information

# The Home Page

- Simple, user-friendly interface for entering query parameters

# The Results Page

- Articles listed in familiar format



SwyneProject

I is munchin' all your news articles...

New Search

41.91862886518302 -87.626953125 100 12/05/1993 12/05/2008

## Results

**USA: CME nearby hogs and pork bellies close limit down.**
Reuters   1996-01-20

Chicago, CHICAGO

**USA: CCC seeks 33,000 T U.S. white wheat for Bangladesh.**
Reuters   1996-01-20

U.S., CHICAGO, Bangladesh

**USA: Widespread aflatoxin found in south Texas corn.**
Reuters   1996-01-20

CHICAGO, Texas, Tx, College Station

**USA: Marcus Theatres expands theater network.**
Reuters

United States, Chicago, Lake Geneva, Detroit, Wisconsin, St. Louis, Minnesota, Colorado, Wis., Illinois, Kansas City, MILWAUKEE

**USA: Corn progress slow, conditions slip in key states.**
Reuters   1996-01-20

U.S., CHICAGO, Illinois, Indiana, Iowa, Ohio, Nebraska

# The Article Page

- Articles shown in familiar format
- Offset with entity list

# Exploring Entities

- Allows user to explore relationships

# Accomplishments

- Created simple interface
  - Learned Web design
- Access to all implemented functionality

- Created IPRO Web site

# Issues Encountered

- Scoping: What is considered "window dressing?"
- Coding of the website itself

- Soft skills
  - Communication
  - Patience
  - Delegating tasks

# Future Work

- Support for additional functionality
  - Categorized results
  - Add search support for additional entity types

# Ethical Issues

- **Professional Issues**
  - Software Copyright
  - Content Copyright
  - Terms of Use
- **Internal Issues**
  - Scheduling Conflicts
  - Interpersonal Conflicts

# Group Accomplishments

- Developed working prototype
- Published deliverables on the Web
- Outlined future work

# Societal Impact

- Inspire future search engines to use semantic data to provide users with specific answers
  - Improve access to information
  - Published our design documents online

# Future Work

- Improve overall speed of system
- Improve accuracy of entity extraction
- Support new the types of queries
- Allow crawling of more types of data sources

# Questions & Answers