# ChemXSeer: CyberInfrastructure for Chemical Kinetics

**Prasenjit Mitra**
College of Information Sciences and Technology,
Department of Computer Science and Engineering

In collaboration with: C. Lee Giles, J. Bandstra, K. Mueller, J. Kubicki, S. Brantley, B. Brouwer, S. Nangia, B. Garrison, B. Sun, Y. Liu, Q. Tan, A.R. Jaiswal, Juan Fernandez, L. Bolleli, X. Lu, …

# Motivation

- Data sharing
  - Data repository
  - Querying and finding data efficiently
  - Analysis tools
  - Preservation and archival
- Data extraction
  - From tables
  - From figures

# Vertical Search Engine

- Domain-specific search engine
  - Entity Extraction
  - Indexing
  - Ranking
    - Similarity & Relevance
- ChemXSeer
  - Chemical Formulae and Names
  - Different query semantics
    - Fuzzy search, similarity search
- What are the first-class entities important for your grand challenge?
  - Existing work on diseases, proteins
  - Genes, Enzymes, ...

# Searching Documents

- Efficient search tools
  - Chemical entity search
    - Formula
    - Name
    - Structure

# Formula Search    -Interface



(c) Prasenjit Mitra, PSU                                    5

# Chemical Entity Search

- Search engines do not understand chemical formulae, chemical names
- No fuzzy search capabilities
  - With functional groups, e.g., -OH
- Automatic segmentation of chemical names and indexing
- Structure search algorithms need improvement

# Formula Search -Query Models

- ## Substructure search
  - Search for formulae that may have a substructure
  - E.g. -COOH matches CH3COOH (exact match: high score), HOOCCH3 (reverse match: medium score), and CH3CHO2 (parsed match: low score).

- ## Similarity search
  - Search for formulae with a similar structure of the query formula. Feature-based approach using partial formulae matching.
  - E.g. ~CH3COOH matches CH3COOH, (CH3COO)2Co, CH3COO$^-$, etc.

# Formula Search -Query Models

- Conjunctive search of the four types of formula searches
  - E.g. [*C2H4-6 -COOH] matches CH3COOH, not C2H4O or CH3CH2COOH.

- Document query rewriting
  - E.g. document query water formula:=CH4 is rewritten to water (CH4 OR H4C OR CD4), if formula search of =CH4 matches CH4, H4C and CD4.

# Experiments -Interface

# Data Repository

- Store and publish data
    - Gaussian
    - CHARMM
    - Excel data
        - Soil profiles
        - Dissolution rates

    - Spectroscopy data

# Functionality

- Store in databases
  - Fast access
  - Structured access
    - Query conditions, e.g., 277 < temp < 302
    - Combine information from multiple tables
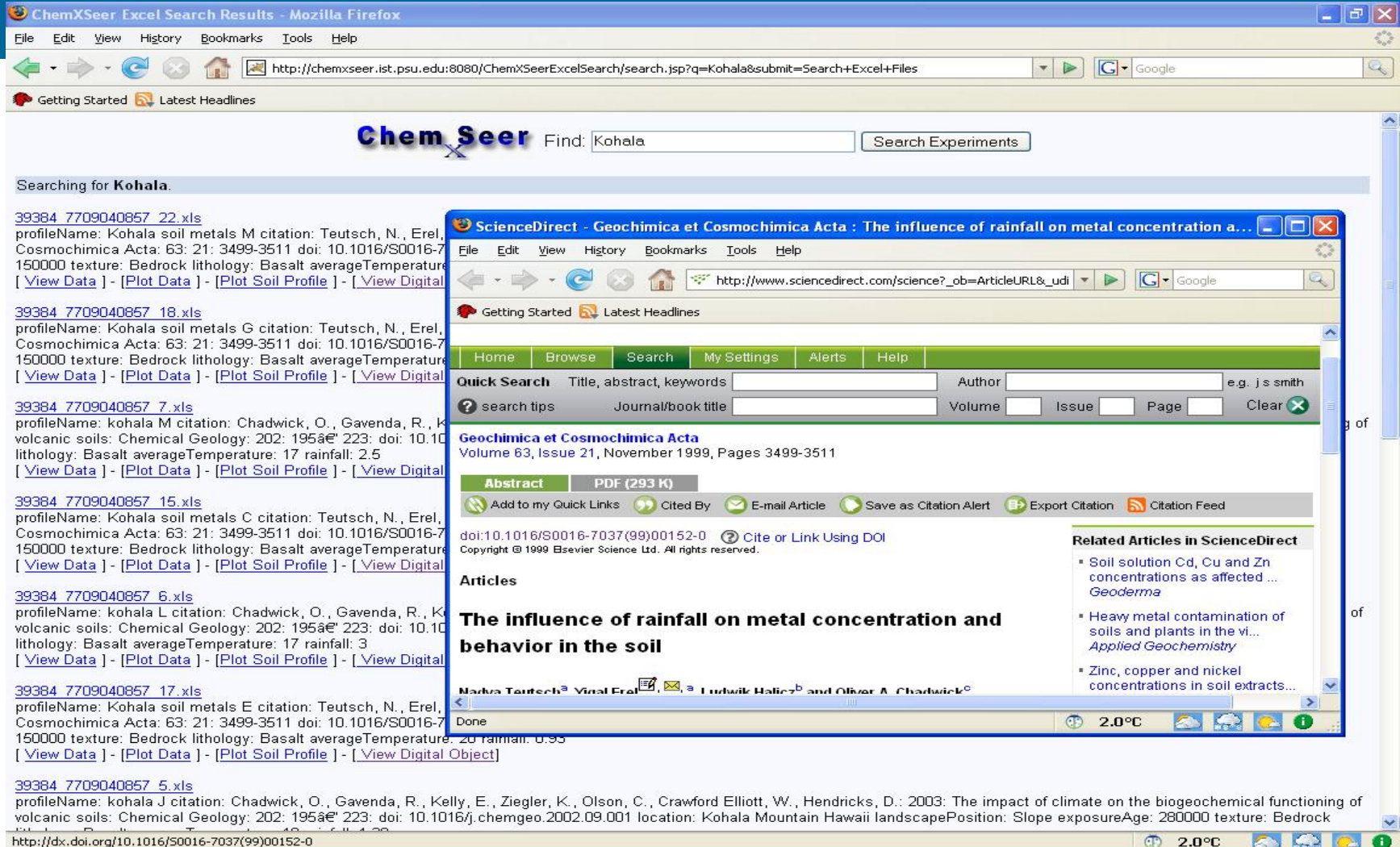  - Query against multiple formats
    - Mediated architecture

# Challenges

- Automatic processing needs fixed data formats
- Ideally one column name
  - Distinct column names across one table
  - Or at least a fixed number of lines for the table header before the data starts
- Fully empty columns and rows are confusing
- Use different datasheets for different tables
- Embedding formulae in the dataset is confusing

# Metadata

- Data describing data
- You need to tell the computer what the columns mean
  - What data is contained in the table?
  - What is the relationship of the different columns and between different sheets of data in Excel?
  - Where was the data obtained?
  - Any other information that may be useful for search and querying later on
- Metadata needs to be in a fixed, specified format
  - Use a form

# Data Search - Interface

# Data Search - Report

# An Example

# Data Search - Plot

# Data Search

- All terms used in the datasheets will be indexed
- Search using those terms will retrieve results
- Then, you can go into an individual sheet and write your queries
- Reports
  – Massage the data & plot
- Command-line interface
  – In SQL

# Data Security

- Login

- Mark data as private or public

- Creation of groups
  - Share data only within your group

# Linking Data and Documents

- Create Digital Object
  - Datasets
  - Documents
  - Supplementary notes
  - Supplementary data
  - Other information
- Data sets use document-identifier
- Document metadata contains pointer to data

# Figure Search

- Figures are important sources of data and information

- Search for figures by keywords

- Index
  - Captions
  - References to the images
  - Legend/text in images (?)
  - Document-level metadata
    - Title, author, venue, year, publisher

# Data Extraction from Figures

- Axes detection

- Legend detection

- Data point identification

  - Shape detection

  - Resolving overlapping points

- Tic detection

- Label detection

- Identifying the continuation of lines at points of intersection

# Challenges

- Optical Character Recognition (OCR)
  - Small letters
  - Sparse letters
- Overlapping shapes
- Axes Scale
  - Log scale
  - Linear scale
  - Disconnected
- Inversion of dependent and independent variables

# TableSeer

- Search for tables with interesting data

- Extract data into a database

- Plot the data and support analysis

- Currently:
  - Table Search based on keywords

- Future:
  - Data extracted inserted into a database

# TableSeer - Interface

# TableSeer

## Beta online working design of a table search engine

**TableSeer**  | Table Caption ▼ | flow |  | search |  Advanced

Found 25 results for query "TableCaption : flow "

Instituto de Qu mica, Universidade Federal da Bahia, Salvador-BA 40170-290, Brazil d Departamento de Qu mica Analitica, Universidad de Valencia, Dr. Moliner 50, 46100 Burjassot, Valencia, Spain. E-mail: miguel.delaguardia@uv.es ' - ' Analyst ' - ' 2000

*In PAGE 1, LINE 78: ............Table 1 Flow analysis determination of sulfide using the MB method............;*

PDF                                          Preview

**Table 1 Comparative results for the determination of morphine in processliquors with chemiluminescence detection using pulsed flow chemistry(PFC) and conventional flow injection analysis (FIA) methodology**

Pulsed flow chemistry: a new approach to solution handling for flow analysis coupled with chemiluminescence detection

Simon W. Lewis,* a Paul S. Francis, a Kieran F. Lim, a Graeme E. Jenkins b and Xue D. Wang c a Centre for Chiral and Molecular Technologies, School of Biological and Chemical Sciences, Deakin University, Geelong, Victoria 3217, Australia b Precision Devices P/L, 44 Nelson Street, Shoreham, Victoria 3916, Australia c School of Chemical and Biomedical Sciences, Central Queensland University, Rockhampton, Queensland 4702, Australia ' - ' Analyst ' - ' 2000

PDF                                          Preview

TableSeer System Architecture

# Sample Table Metadata Extracted File

**Table 1** Temperature effect on resistance change ($\Delta R$) and response time of tin oxide thin film with 1% $CCl_4$

| Temperature/ °C | $\Delta R^a/\Omega$ | $\dfrac{\Delta R}{(R, O_2)}$ (%) | Response time | Reproducibiliy |
|---|---|---|---|---|
| 100 | 223 | 5 | ~ 22 min | Yes |
| 200 | 270 | 9 | ~ 7-8 min | Yes |
| 300 | 1027 | 21 | < 20 s | Yes |
| 400 | 993 | 31 | ~ 10 s | No |

$^a \Delta R = (R, CCl_4) - (R, O_2)$.

- **<Table>**
- **<DocumentOrigin>Analyst</DocumentOrigin>**
- **<DocumentName>b006011i.pdf</DocumentName>**
- **<Year>2001</Year>**
- **<DocumentTitle>Detection of chlorinated methanes by tin oxide gas sensors </DocumentTitle>**
- **<Author>Sang Hyun Park, a ? Young-Chan Son, a Brenda R . Shaw, a Kenneth E. Creasy,* b and Steven L. Suib* acd a Department of Chemistry, U-60, University of Connecticut, Storrs, C T 06269-3060</Author>**
- **<TheNumOfCiters></TheNumOfCiters>**
- **<Citers></Citers>**
- **<TableCaption>Table 1 Temperature effect o n r esistance change ( D R ) and response timeof tin oxide thin film with 1 % C Cl 4</TableCaption>**
- **<TableColumnHeading>D R Temperature/ ¡ã C D R a / W ( R ,O 2 ) (%) R esponse time Reproducibiliy </TableColumnHeading>**
- **<TableContent>100 223 5 ~ 22 min Yes 200 270 9 ~ 7-8 min Yes 300 1027 21 < 2 0 s Yes 400 993 31 ~ 1 0 s No </TableContent>**
- **<TableFootnote> a D R =( R , CCl 4 ) - ( R ,O 2 ). </TableFootnote>**
- **<ColumnNum>5</ColumnNum>**
- **<TableReferenceText>In page 3, line 11, … Film responses to 1% CCl4 at different temperatures are summarized in Table 1……</TableReferenceText>**
- **<PageNumOfTable>3</PageNumOfTable>**
- **<Snapshot>b006011i/b006011i_t1.jpg</Snapshot>**
- **</Table>**

# TableRank

- Rank tables by rating the <u>&lt;query, table&gt;</u> pairs, instead of the &lt;query, document&gt; pairs: preventing a lot of false positive hits for table search, which frequently occur in current web search engines

- The similarity between a &lt;table, query&gt; pair: the cosine of the angle between vectors

$$sim(tb_j, Q) = cos(tb_j, Q) = \frac{\sum_{i=1}^{s} w_{i,j,k} w_{i,q,k}}{|tb_j||Q|}$$

- Tailored term vector space => table vectors:
  - Query vectors and table vectors, instead of document vectors

29

# An Example of the Citation Network

# Challenges

- Identification of cells
  - Cells getting fused
- Irregularly shaped tables
- Horizontal letters
- Identify columns
  - Sometimes columns and rows are fused
  - Define heuristics to detect such cases
- Identifying what the columns signify
- units of the columns

# Proposed Architecture and Framework

# Digital Libraries

- Connecting data and digital documents
  - Mining legacy data from tables in digital documents
  - Data submission system where creator can link data and documents
- Use ORE model
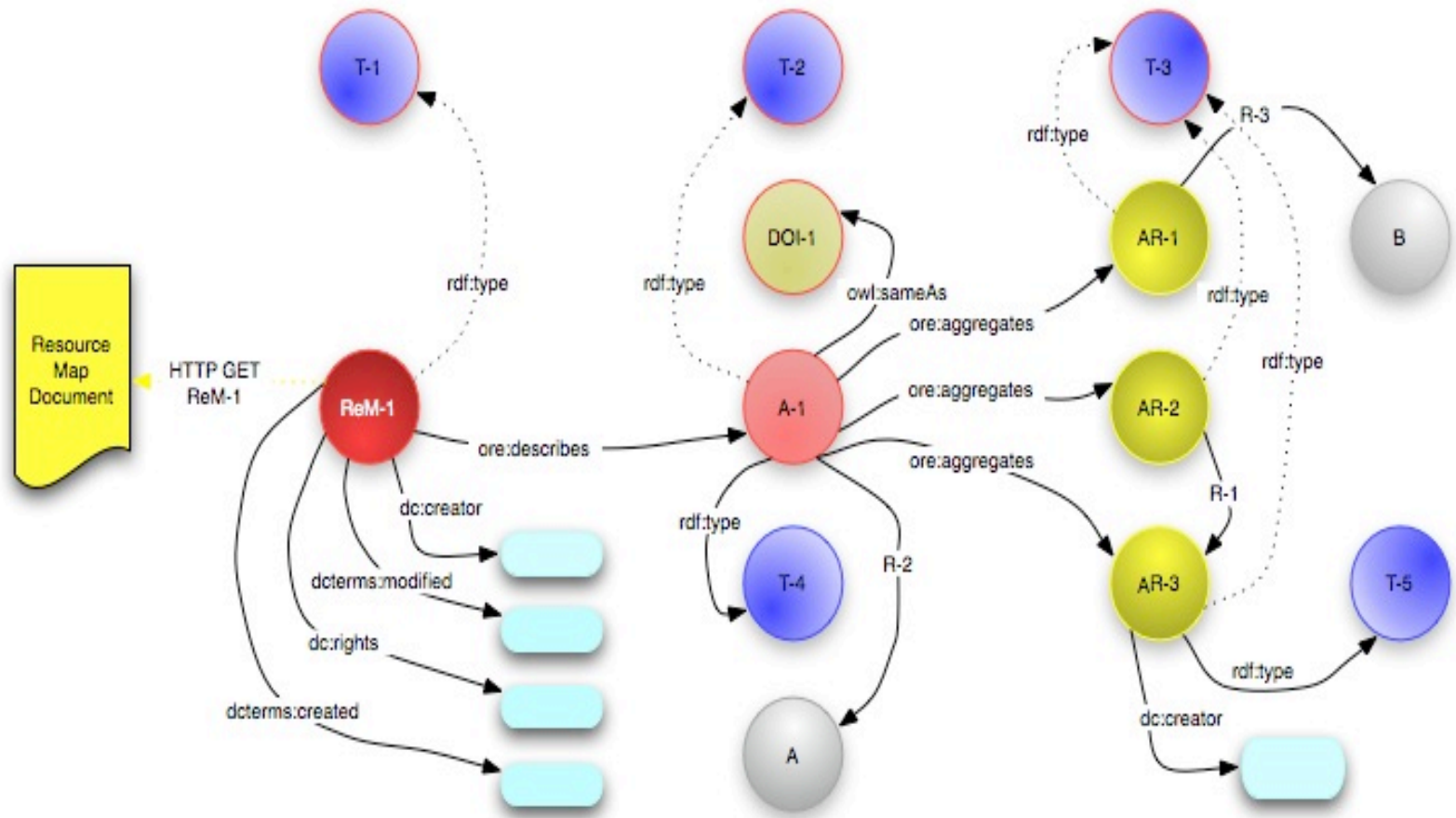
# Data-Document IntegrationNeeds

- Queries:
  - Get me the data on the dissolution rate of kaolinite
  - Get me the data & experimental results dissolution rate of kaolinite reported in papers referred to by paper x.

- Link information from tables in databases to documents
- Table-level Links
  - All data from a table from a single article
- Row-level Links
  - Data in table linking to multiple articles

# Use ORE

- Dataset – resource
- Aggregates of linked resources
  - Database
    - Published
    - Raw
  - Article
    - Preprints
  - Supplementary notes
- Resource Map

# Digital Objects

# Linking to Ontology

# Data Annotation

# Ontology-based Mediated Search

- Ontologies used to express vocabularies
  - Mediate federated databases
  - Mediated search
    - Multiple data formats
      - Gaussian, CHARMM, Excel, XML, ASCII …

# Data Search Example in Chem<sub>X</sub>Seer



Search document content and metadata

# Utilizing the Social Network

- Automatic discovery of people who have similar interests (optional feature)
  - Who reads the same people/papers?
  - Who cites the same people/papers?
  - Who produces similar work?
  - Who has similar browsing habits?
  - Zhou, Manavoglu, et al, 2006; Zhou, Ji, et al, 2006
- Data traditionally used for collaborative filtering
  - Instead: connect people, build community awareness
  - Discover new people/trends of interest
- Shy users may cloak their profiles, if desired

# Personalization

- MyCiteSeerX

# Conclusion

- Repository for data and documents
  - Linking data to documents
- Intelligent tools for data extraction
  - Data analysis capabilities can be built in
- Advanced search capabilities
  - Chemically aware search

# Thanks!

- http://chemxseer.ist.psu.edu
- http://citeseerx.ist.psu.edu


- pmitra@ist.psu.edu
- giles@ist.psu.edu

# Conclusion

- Repository for data and documents
  - Linking data to documents
- Intelligent tools for data extraction
  - Data analysis capabilities can be built in
- Advanced search capabilities
  - Chemically aware search

# Appeal

- Please submit your data and documents!!!

- The portal will not succeed without data

# Needs

- What do you need?


- pmitra@ist.psu.edu