# FireCite: Lightweight real-time reference string extraction from web pages

Ching Hoi Andy Hong       Jesse Prabawa Gozali       Min-Yen Kan

School of Computing

National University of Singapore

# Outline

- **Introduction**
- **Reference String Recognition**
- **Reference String Parsing**
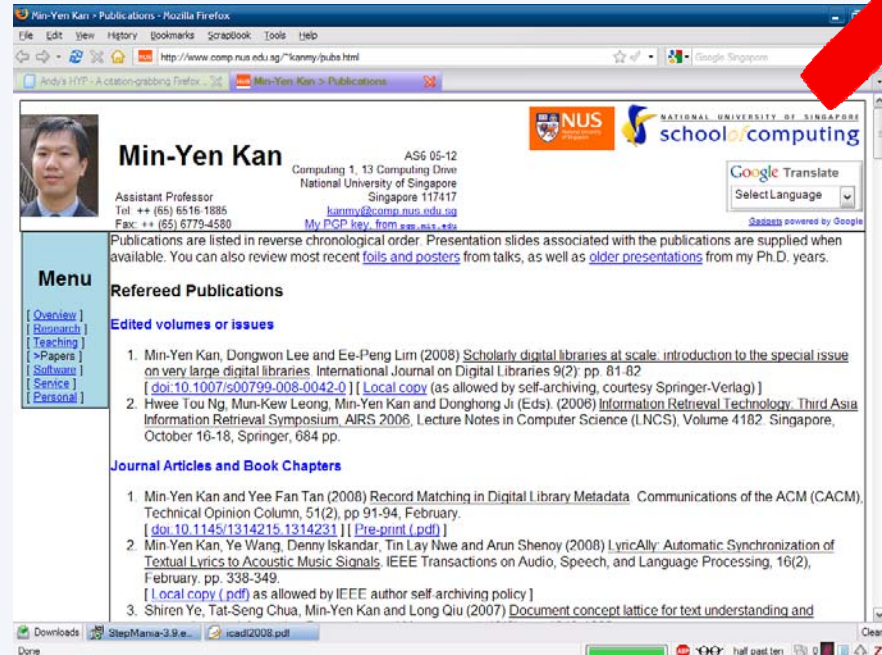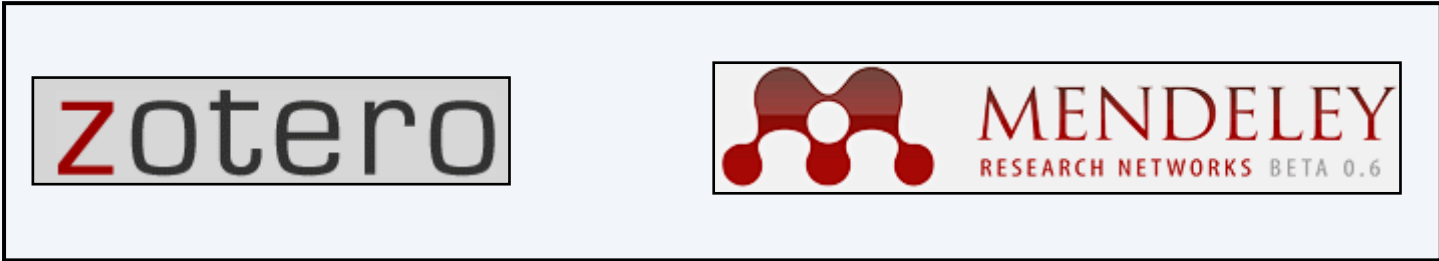- **Firefox Extension**
- **Conclusion**

# Introduction: The Problem

- **Recognition and parsing of references found on the Internet**

- **Criteria:**
  - Accurate
  - Fast

**Journal Articles**

| | | | |
|---|---|---|---|
| IR | NL | ML | Fuchun Peng and Xiangji Huang; **Machine Learning Approaches to Automatic Text Classification for Asian Languages**, *Journal of Documentation,* Volume 63, Issue 3, pages 378-397, (2007) |
| IR | ML | | Fuchun Peng and Andrew McCallum; **Information Extraction from Research Papers using Conditional Random Fields,** *Information Processing and Management,* 42(4), pages 963-979, (2006) |
| NL | ML | | Shaojun Wang, Dale Schuurmans, Fuchun Peng and Yunxin Zhao; **Combining Statistical Language Models via the Latent Maximum Entropy Principle**, *Machine Learning Journal,* Vol. 60, No. 1-3, pages 229-250, (2005)  Special Issue on Learning in Speech and Language Technologies. |
| ML | | | Shaojun Wang, Dale Schuurmans, Fuchun Peng and Yunxin Zhao; **Learning Mixture Models with the Regularized Latent Maximum Entropy Principle**, *IEEE Transactions on Neural Networks ,* Vol. 15, No. 4, pages 903 - 916, (2004). Special Issue on Information Theoretic Learning |

# Introduction: Related Work

# Outline

- **Introduction**
- **Reference String Recognition**
- **Reference String Parsing**
- **Firefox Extension**
- **Conclusion**

# Reference String Recognition: Definition

- **Are there reference strings?**

- **Where are the reference strings?**

# Reference String Recognition: Methodology

1. **Web page exclusion**

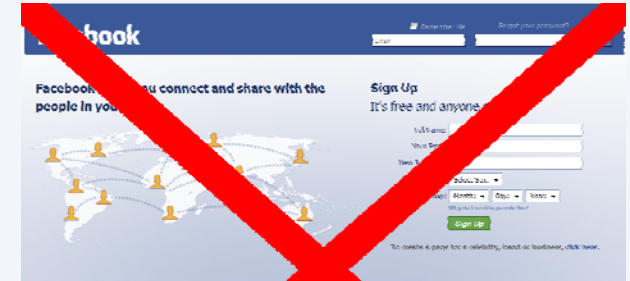   - Keyword + URL Matching



**Journal Articles**

| IR | NL | ML | Fuchun Peng and Xiangji Huang; **Machine Learning Approaches to Automatic Text Classification for Asian Languages**, *Journal of Documentation,* Volume 63, Issue 3, pages 378-397, (2007) |
| IR | ML | | Fuchun Peng and Andrew McCallum; **Information Extraction from Research Papers using Conditional Random Fields,** *Information Processing and Management,* 42(4), pages 963-979, (2006) |
| NL | ML | | Shaojun Wang, Dale Schuurmans, Fuchun Peng and Yunxin Zhao; **Combining Statistical Language Models via the Latent Maximum Entropy Principle**, *Machine Learning Journal,* Vol. 60, No. 1-3, pages 229-250, (2005)   Special Issue on Learning in Speech and Language Technologies. |
| ML | | | Shaojun Wang, Dale Schuurmans, Fuchun Peng and Yunxin Zhao; **Learning Mixture Models with the Regularized Latent Maximum Entropy Principle**, |

# Reference String Recognition: Methodology

## 2. Splitting

- Split web page text into segments
- GOAL: Each segment to contain at most 1 reference string, and each reference string to exist in only 1 segment.

**Journal Articles**

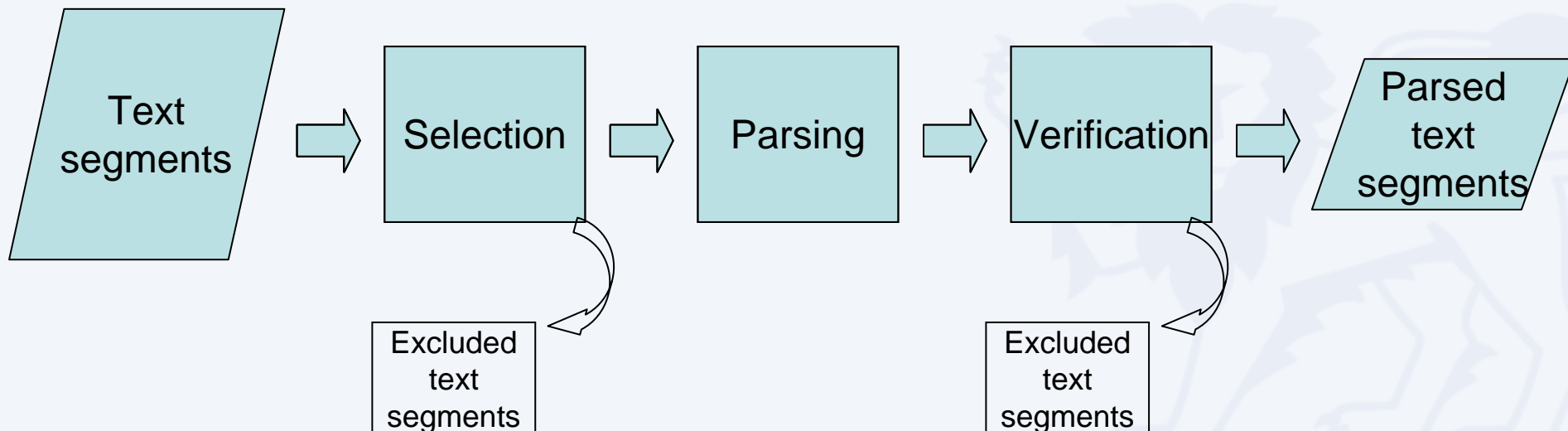| | | | |
|---|---|---|---|
| IR | NL | ML | • Fuchun Peng and Xiangji Huang; **Machine Learning Approaches to Automatic Text Classification for Asian Languages**, *Journal of Documentation,* Volume 63, Issue 3, pages 378-397, (2007) |
| IR | ML | | • Fuchun Peng and Andrew McCallum; **Information Extraction from Research Papers using Conditional Random Fields,** *Information Processing and Management,* 42(4), pages 963-979, (2006) |
| NL | ML | | • Shaojun Wang, Dale Schuurmans, Fuchun Peng and Yunxin Zhao; **Combining Statistical Language Models via the Latent Maximum Entropy Principle**, *Machine Learning Journal,* Vol. 60, No. 1-3, pages 229-250, (2005)   Special Issue on Learning in Speech and Language Technologies. |
| ML | | | • Shaojun Wang, Dale Schuurmans, Fuchun Peng and Yunxin Zhao; **Learning Mixture Models with the Regularized Latent Maximum Entropy Principle**, |

# Reference String Recognition: Methodology

3.  **Selection**

    - Token length and word length of segment

4.  **Verification**

    - Reject segments that do not have a title or authors

```
┌──────────┐      ┌──────────┐      ┌──────────┐      ┌──────────┐      ┌──────────┐
│   Text   │  =>  │Selection │  =>  │ Parsing  │  =>  │Verification│ =>  │  Parsed  │
│ segments │      │          │      │          │      │          │      │   text   │
└──────────┘      └──────────┘      └──────────┘      └──────────┘      │ segments │
                        │                                    │          └──────────┘
                        v                                    v
                  ┌──────────┐                         ┌──────────┐
                  │ Excluded │                         │ Excluded │
                  │   text   │                         │   text   │
                  │ segments │                         │ segments │
                  └──────────┘                         └──────────┘
```

# Reference String Recognition: Evaluation

- **Test set: 40 staff homepages from 4 universities**
- **Reference strings found: 364/379 (96.0%)**
- **False positives: 269 (42.5%)**

| System | Precision | Recall | F1-measure |
|---|---|---|---|
| FireCite (All 40 pages) | 0.575 | 0.960 | 0.719 |
| FireCite (Only 20 pages with reference strings) | 0.655 | 0.960 | 0.779 |

# Outline

- **Introduction**
- **Related Work**
- **Reference String Recognition**
- **Reference String Parsing**
- **Firefox Extension**
- **Conclusion**

# Reference String Parsing: Definition

- Purpose
  - Store reference according to metadata fields
  - Assist reference string recognition

- Only identify authors, title, date

Jesse Prabawa Gozali and Min-Yen Kan (2007) A Rich OPAC User Interface with AJAX, In Proceedings of the Joint Conference on Digital Libraries (JCDL '07). Vancouver, Canada, June, pp. 329-330. Short paper.

# Reference String Parsing: Methodology

- **Tokenising**

Atlas , L . , and S . Shamma ,

" Joint Acoustic and Modulation Frequency , "

EURASIP JASP , 2003 .

- Advantages
  - Reduce number of computations
  - Allow information-richer learning features

# Reference String Parsing: Methodology

- **Labelling**
  - J48 decision tree classifier
  - CORA corpus (500 reference strings) as training corpus

- **Repairs**
  - "Title" and "Authors" fields are contiguous

Thuy Dung Nguyen and Min-Yen Kan /author ( 2007 /date ) Keyphrase Extraction in Scientific Publications/title .
In Proc/misc . of International Conference on Asian Digital Libraries /~~title~~ misc ( ICADL '07/misc ).
Hanoi/misc , Vietnam/misc , December/misc . pp/misc . 317-326/misc .

# Reference String Parsing: Evaluation

## 6 Faculty Staff Publication Pages

| Page (No. of ref. strings) | Token-level F-measure | | | |
|---|---|---|---|---|
| | Title | Authors | Date | All Tokens |
| A (72) | 0.902 | 0.893 | 0.988 | 0.708 |
| B (52) | 0.953 | 0.957 | 0.990 | 0.960 |
| C (29) | 0.684 | 0.304 | 0.774 | 0.651 |
| D (68) | 0.753 | 0.968 | 0.889 | 0.917 |
| E (8) | 0.692 | 0.875 | 1.000 | 0.889 |
| F (45) | 0.847 | 1.000 | 0.989 | 0.966 |
| Overall (274) | 0.836 | 0.916 | 0.948 | 0.878 |

# Reference String Parsing: Evaluation

FLUX-CiM Computer Science Dataset (300 citations)

| System Name | Field-level F-measure | | | |
|---|---|---|---|---|
| | Title | Authors | Date | Overall |
| **FireCite** | **0.92** | **0.96** | **0.97** | **0.94** |
| ParsCit | 0.96 | 0.99 | 0.97 | 0.94 |
| FLUX-CiM | 0.93 | 0.95 | 0.98 | 0.97 |

# Reference String Parsing: Evaluation

| Parser | Classifier Type | Size of classifier model (KB) | Size of dictionary (KB) |
|--------|-----------------|-------------------------------|-------------------------|
| FireCite | Decision Tree | 12.6 | 0.0 |
| FLUX-CiM | Knowledge-Based | >786 (estimated) | 0.0 |
| ParsCit | Conditional Random Fields | 7339 | 1722 |

# Reference String Parsing: Evaluation

| Page Type | Time taken (milliseconds) | | |
|---|---|---|---|
| | Minimum | Maximum | Average |
| Pages with reference strings | 90 | 544 | 192 |
| Pages without reference strings | 6 | 222 | 74 |
| All pages | 6 | 544 | 133 |

# Outline

- **Introduction**
- **Reference String Recognition**
- **Reference String Parsing**
- **Firefox Extension**
- **Conclusion**

# Extension: Demo

# Conclusion

- **Results**

    - Fast and lightweight reference string parser

    - Reference string recogniser with good recall

    - Basic, expandable Firefox extension

# Conclusion

- **Future work**

  - Reference String Recognition

    - More rules to improve precision

  - Reference String Parser

    - Use web page reference strings as training data

    - Recognise implicit/common metadata

  - Firefox Extension

    - Add more features to the extension

Questions?