

# **Optimizing Information Extraction over Evolving Text**



**Fei Chen**  
**Database Group**  
**Computer Science Department**  
**University of Wisconsin – Madison**

# Lots of Text



**amazon.com.**

WELCOME YOUR STORE BOOKS ELECTRONICS

SEARCH BROWSE SUBJECTS BESTSELLERS

Get 5% off select bouncers

**Machine Learning**  
by Tom M. Mitchell

LOOK INSIDE!

List Price: \$422.15  
Price: \$122.15

This item ships for

Availability: Usually ship

Used & new from \$76.99

Edition: Hardcover

See more product det

Buy Now Pay Later: Make no payment

Great Buy

Buy this book with



**PubMed.gov**  
U.S. National Library of Medicine  
National Institutes of Health

Search: PubMed

Display Settings:  Abstract

Genome Res. 2010 Jan 14. [Epub ahead of print]

**Genome-wide discovery of human heart enhancers.**

Narlikar L, Sakabe NJ, Blanski AA, Arimura FE, Westlund JM, Nobrega MA, Ovcharenko

NCBI

velopment rely on th  
ifying the cis-regul  
: expression is an  
velopment. We des  
ents that share

**80% of the world's data, and growing!**

<Operator> Thank you

<Caller> is John am

<Operator> Yeah you ar

<Caller> Morning ho

<Operator> Oh I'm I a

<Caller> um ah of co

ways... pror

one in and

I put in thi

what is the p

you know...

doesn't mat

in and I hav

ing can you

**MEDICAL RECORDS INFORMATION**

<b>PATIENT NAME: ADKINS, PAUL J</b>		<b>PATIENT NO: 000001</b>
<b>ADDRESS: 1121 E MADISON EDWARDSVILLE, IL 62025</b>		
<b>HM TEL: 618-692-5545</b>	<b>WK TEL: 618-251-4784</b>	<b>SSN: 654-31-5818</b>
<b>OCCUPATION:</b>	<b>DOB: 01/16/1965(43 yrs)</b>	<b>CHARTNO: PA1234</b>
<b>GUARANTOR: ADKINS, PAUL J</b>		<b>SEX: MALE</b>

**HISTORY & PHYSICAL ON 03/03/2007**

**PROVIDER: MELMAN, IRVING G** **ENTERED BY: IGM**

**SUBJECTIVE**

Mr. Adkins is a 41 yo male who noticed a recent bulge at his tailbone, prior to being seen he noted purulent drainage from this area. There has been no previous drainage from this area and no history of a known pilonidal abscess or cyst. He is not a diabetic. He was seen by his primary care MD and started on antibiotics. He states that the area has improved significantly since first developing.

# Exploiting Text by Information Extraction (IE)

## Disease News

Jan 31: Turkey confirmed an incident of foot&mouth disease.

Jan 30: H1N1 identified in California  
Turkey Flock.

IE program

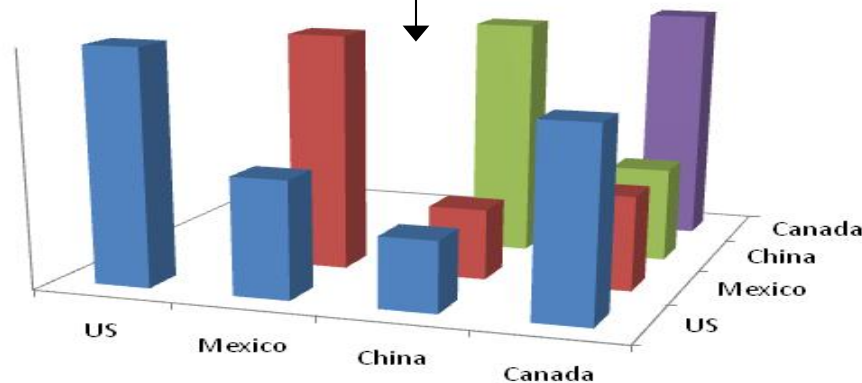
## Incidents

disease	location
foot&mouth	Turkey
H1N1	California
...	...

```
SELECT location
FROM Incidents
WHERE disease = 'H1N1'
```

California

Data Mining



Web Search

Advertising

Monitoring

Visualizing

⋮

# Current State of Art

- **Many players**
  - **AI/DB/DM/IR/NLP/Web communities:** CMU, Columbia, Cornell, IIT, JHU, Max-Panck, PSU, Stanford, UCLA, UIUC, UMass, UMichigan, USC, UT Austin, UToronto, UWashington, UWisconsin...
  - **industry:** AT&T, Google, HP Labs, IBM, Microsoft, Yahoo!...
- **Mainly centered on improving extraction accuracy**
- **Recent work starts to consider improving runtime...**

# Current Work on Improving Runtime

## Pruning documents

[Agichtein *et al.* ICDE-03]  
[Etzioni *et al.* WWW-04]  
[Ipeirotis *et al.* SIGMOD-06]

## Efficient pattern match

[Gravano *et al.* VLDB-01]  
[Cho *et al.* ICDE-02]  
[Chandel *et al.* ICDE-06]

## Parallel processing

[Gruhl *et al.* IBMSJ-04]  
[Lin SIGIR-09]

## Relational-style optimization

[Shen *et al.* VLDB-07]  
[Reiss *et al.* ICDE-08]

## IE over evolving text

[Chen *et al.* ICDE-08]  
[Chen *et al.* SIGMOD-09]  
[Chen *et al.* TechReport-10]

# IE over Evolving Text : An Example

## Disease News

Jan 31: Turkey confirmed an incident of foot&mouth disease.

Jan 30: H1N1 identified in California  
Turkey Flock.

IE

## Incidents

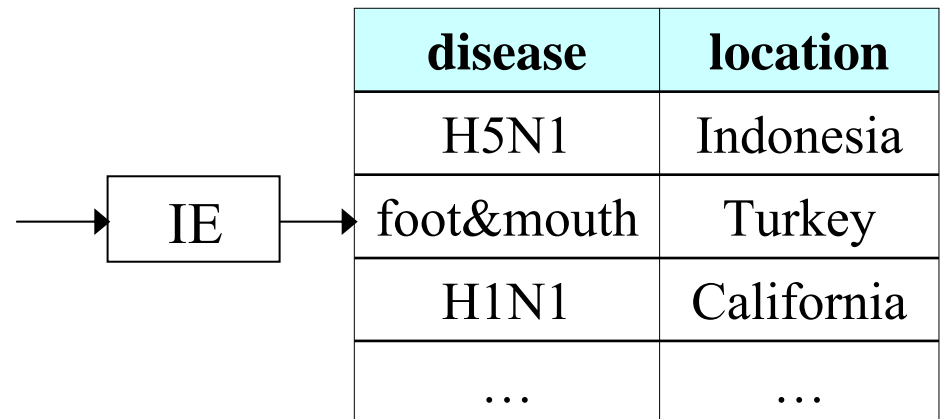
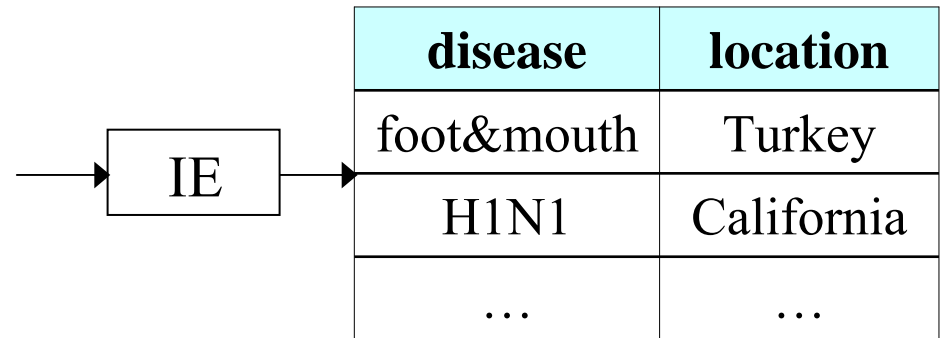
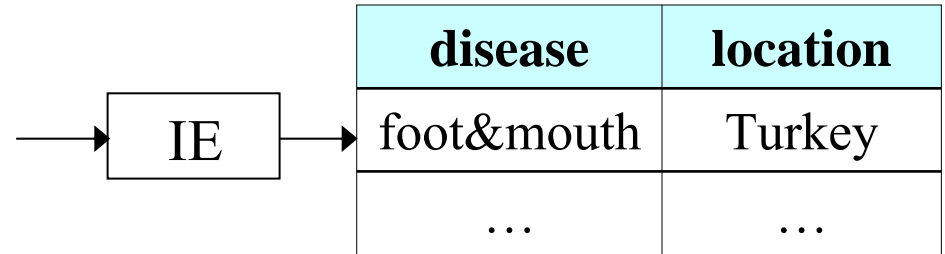
disease	location
foot&mouth	Turkey
H1N1	California
...	...



Alert me if there is an H1N1 incident

day 1

<b>Recent Incidents</b>
<b>Flu Situation Update</b>
Influenza activity remained at the same levels.
Reduction of Inventory at the
<b>Disease News</b>
Feb 1: A new H5N1 case confirmed in Indonesia.
Jan 31: Turkey confirmed an incident of foot&mouth disease.
Jan 30: H1N1 identified in California Turkey Flock.



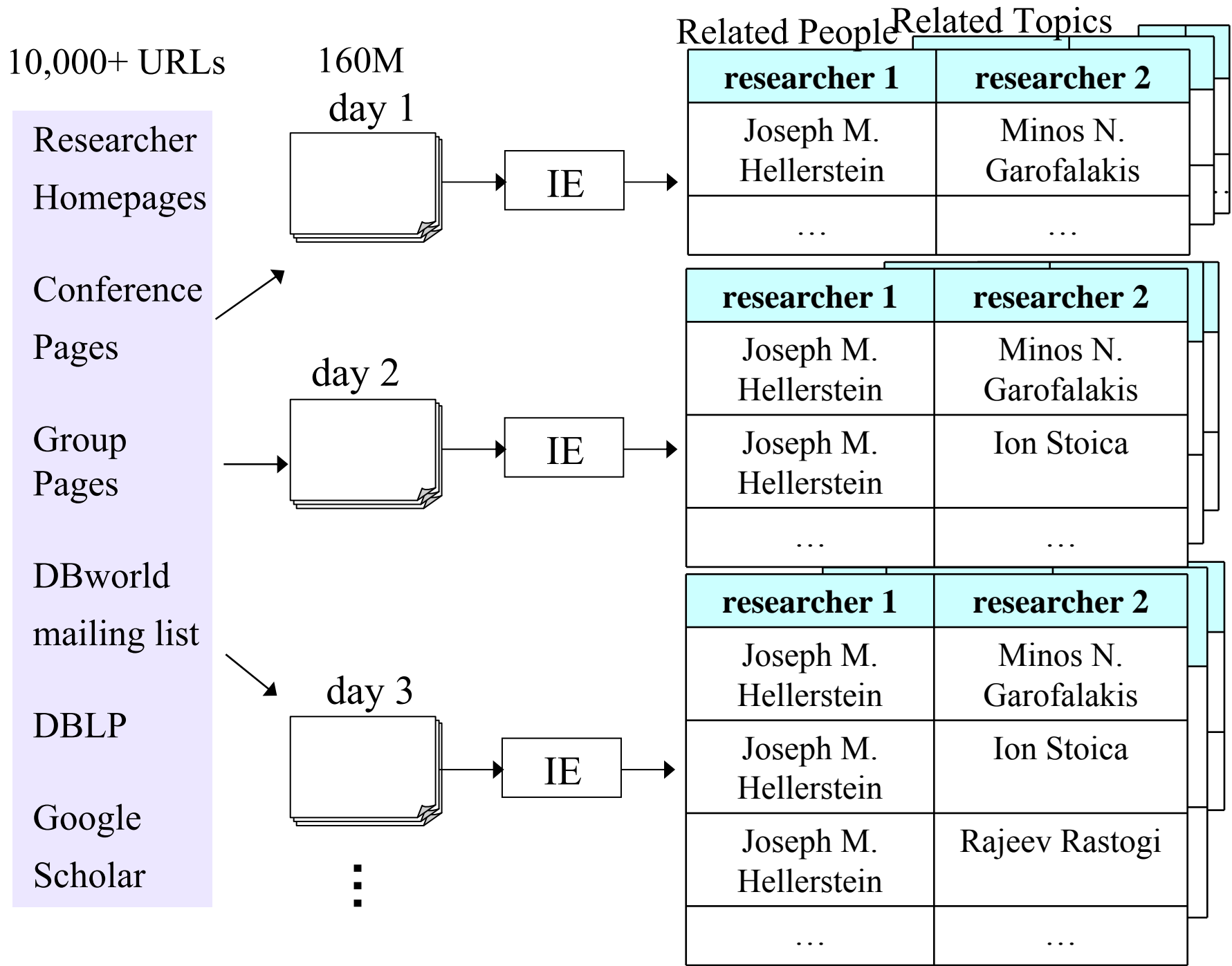
# IE over Evolving Text: Another Example

The screenshot shows a Mozilla Firefox browser window with the following details:

- Address Bar:** [http://dblife.cs.wisc.edu/person/Joseph\\_M.\\_Hellerstein](http://dblife.cs.wisc.edu/person/Joseph_M._Hellerstein)
- Page Title:** DBLife: Joseph M. Hellerstein
- Page Content:**
  - Header:** DBLife logo, Search bar, and navigation links (Help, The...).
  - Profile:** Joseph M. Hellerstein, with links to Bing, Citeseer, DBLP, Google, Google Scholar, Kosmix, Wikipedia, and Yahoo!.
  - Media:** A row of images including a portrait, a bar chart, another portrait, and book covers. A "more" link is present below the images.
  - Recent News:**
    - [Scalable, Distributed Data Structures for Internet Service Construction](#) cited 7 times - [details](#)
    - [Brief announcement: prefix hash tree](#) cited 5 times - [details](#)
    - [Practical Predicate Placement](#) cited 2 times - [details](#)
    - [High-Performance Sorting on Networks of Workstations](#) cited 5 times - [details](#)
    - [A Case for Intelligent Disks \(IDISKS\)](#) cited 4 times - [details](#)
    - [Blobworld: A System for Region-Based Image Indexing and Retrieval](#) cited 15 times - [details](#)
    - [Concurrency and Recovery in Generalized Search Trees](#) cited 3 times - [details](#)
    - [Data gathering tours in sensor networks](#) cited 5 times - [details](#)
    - [Beyond Average: Toward Sophisticated Sensing with Queries](#) cited 9 times - [details](#)
  - Right Sidebar:**
    - Professor:** <http://db.cs.berkeley.edu/jmh/>, University of California-Berkeley USA, Papers cited 14,646 times, [H-Index](#) of 53.
    - Related People:**
      - [Minos N. Garofalakis](#)
      - [Ion Stoica](#)
      - [Rajeev Rastogi](#)
      - [Vijayshankar Raman](#)
    - Related Topics:**
      - [query processing](#)
      - [database systems](#)
      - [streams](#)
      - [querving](#)
    - Services:**
      - [CIDR 2009 \(PC\)](#)

The browser's status bar at the bottom shows "Done" and a McAfee security icon.

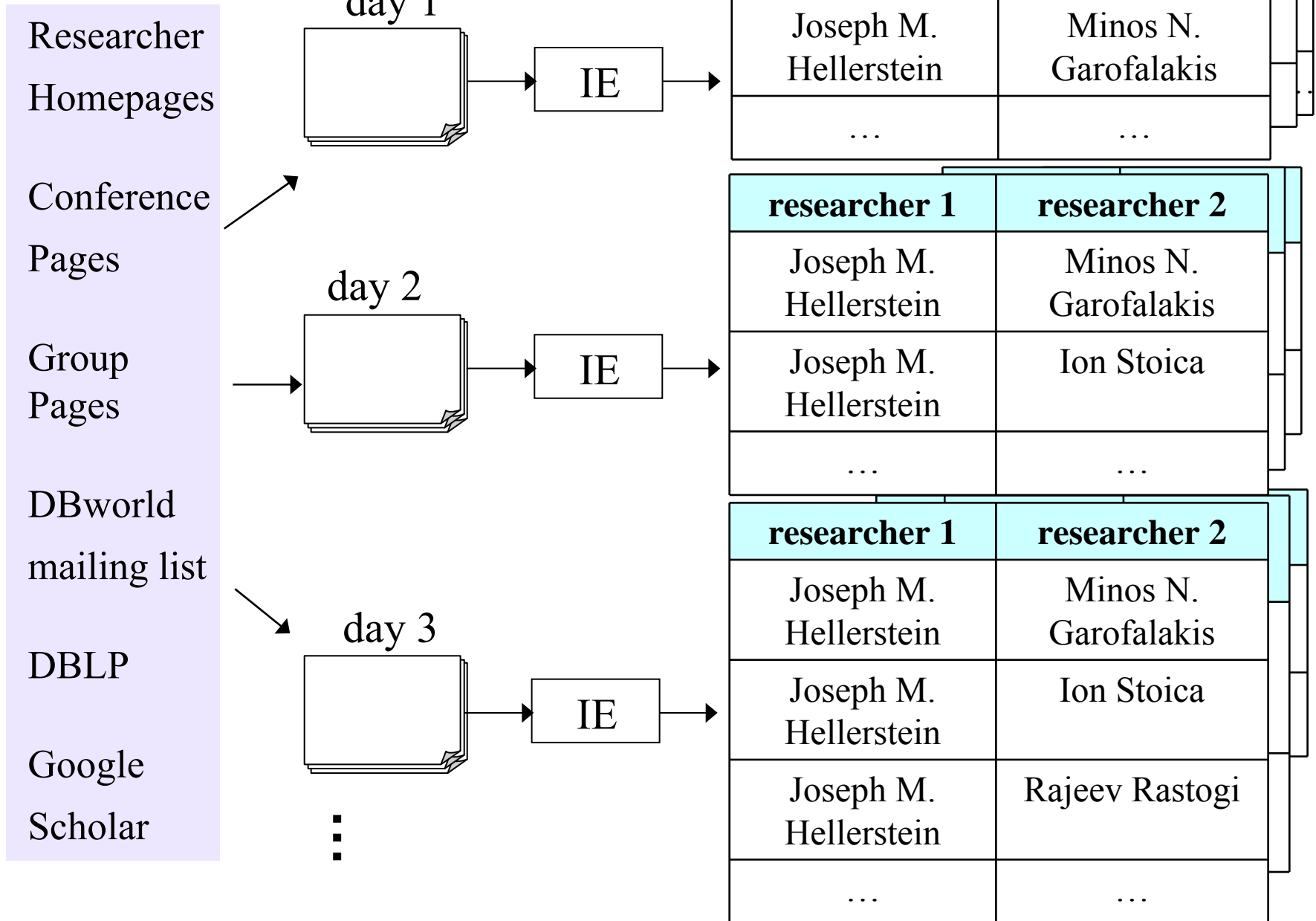




# Many Other Applications of IE over Evolving Text

- **Impliance @ IBM**
  - manages information on enterprise intranets
  - finds the latest information
- **IWP @ Univ. of Washington and YAGO @ MPI**
  - extract structures from Wikipedia
  - keep extracted structures up-to-date
- **Blogscope @ Univ. of Toronto**
  - monitors the blogosphere
- ...

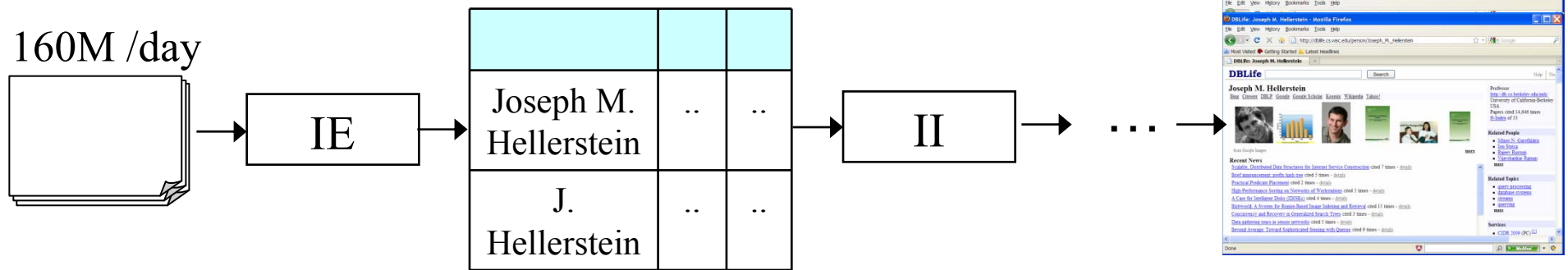
# No Good Current Solution



# No Good Current Solution

- **Inefficient**

- e.g., takes 8+ hours in DBLife everyday



- **Unsuitable for time-sensitive applications**

- e.g., stock and auction
- cannot finish extracting in time

- **Unsuitable for interactive debugging**

- long debug loop

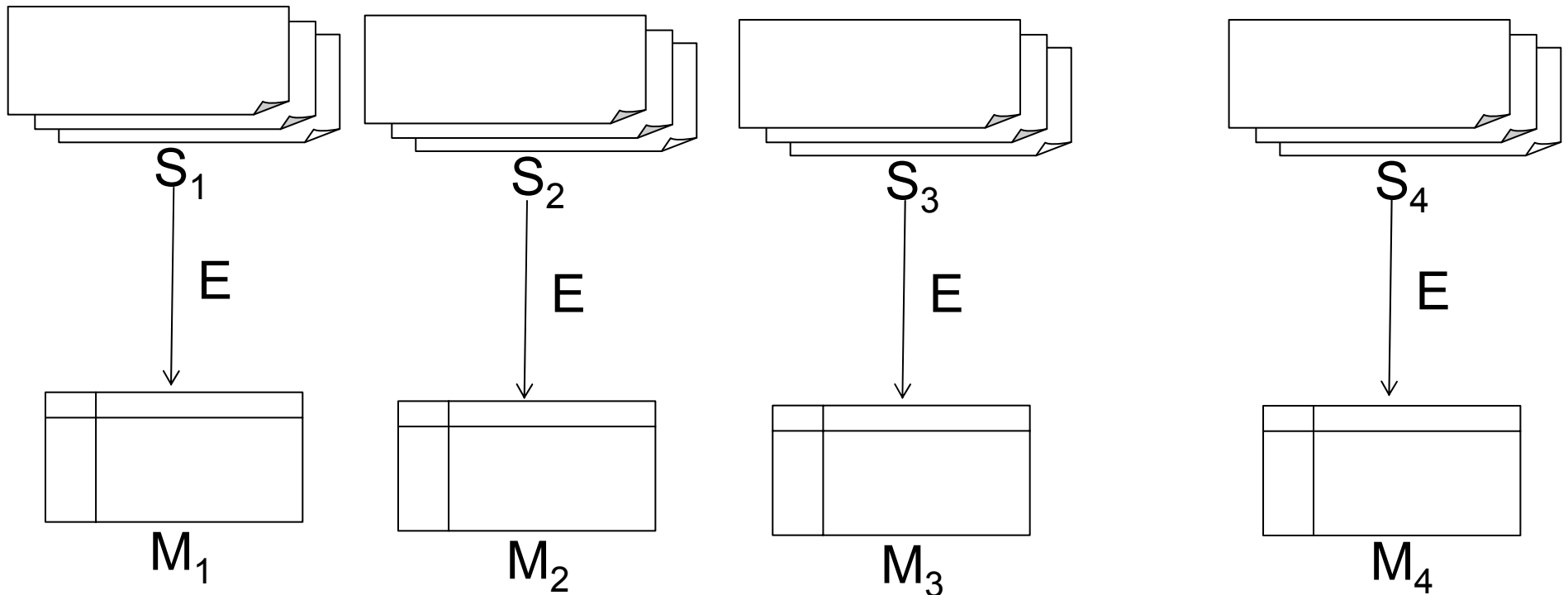
# **My Dissertation Contributions**

**Developed efficient IE solutions over evolving text that**

- recycle previous IE efforts
- guarantee correctness
- deal with large text corpora
- deal with large IE programs
- deal with complex learning-based IE programs

**Results in [ICDE-08], [SIGMOD-09] and [TechReportA-10]**

# Problem Definition



Can we get  $M_4$  faster than applying  $E$  to  $S_4$  from scratch?

# Cyclex [ICDE08]: Recycling Extraction

slickdeals.net (day 1)

Hot Deals
Dec 10: Sony TV \$850 at Amazon.
Dec 9: iPhone \$199 at BestBuy.

p

match

slickdeals.net (day 2)

Hot Deals
Dec 11: All Dell LCDs \$299.
Dec 10: Sony TV \$850 at Amazon.
Dec 9: iPhone \$199 at BestBuy.
Expired!

q

Deals

product	seller
TV	Amazon
iPhone	BestBuy

E

Deals

product	seller
LCDs	Dell
TV	Amazon
iPhone	BestBuy

copy

E

copy

copy

E

# Not Always Correct to Copy Everything!

E extracts a deal if (a) product & seller are within a window of at most 20 chars  
(b) no “Expired!” is within 10 chars around the window.

slickdeals.net (day 1)

**Hot Deals**

Dec 10: Sony TV \$850 at Amazon.

Dec 9: iPhone \$199 at BestBuy.

p

Deals

product	seller
TV	Amazon
iPhone	BestBuy

slickdeals.net (day 2)

**Hot Deals**

Dec 11: All Dell LCDs \$299.

Dec 10: Sony TV \$850 at Amazon.

Dec 9: iPhone \$199 at BestBuy.

Expired!

q

Deals

wrong!

product	seller
LCDs	Dell
TV	Amazon
iPhone	BestBuy



# Extractor Properties: Context

- Many extractors only examine small “context windows” surrounding a mention to extract the mention.

E extracts a deal if (a) product & seller are within a window of at most 20 chars  
(b) no “Expired!” is within 10 chars around the window.

## Hot Deals

Dec 11: All Dell LCDs \$299.

Dec 10: Sony TV \$850 at Amazon.

Dec 9: iPhone \$199 at BestBuy.

Expired!



E's context = 10 chars

- The text outside the context of a mention m is irrelevant for E to extract m.

# Revisit The Example

E extracts a deal if (a) product & seller are within a window of at most 20 chars  
 (b) no “Expired!” is within 10 chars around the window.

slickdeals.net (day 1)

**Hot Deals**

Dec 10: Sony TV \$850 at Amazon.

Dec 9: iPhone \$199 at BestBuy.

p

E's context =  
 10 chars

Deals

product	seller
TV	Amazon
iPhone	BestBuy

slickdeals.net (day 2)

**Hot Deals**

Dec 11: All Dell LCDs \$299.

Dec 10: Sony TV \$850 at Amazon.

Dec 9: iPhone \$199 at BestBuy.

Expired!

q

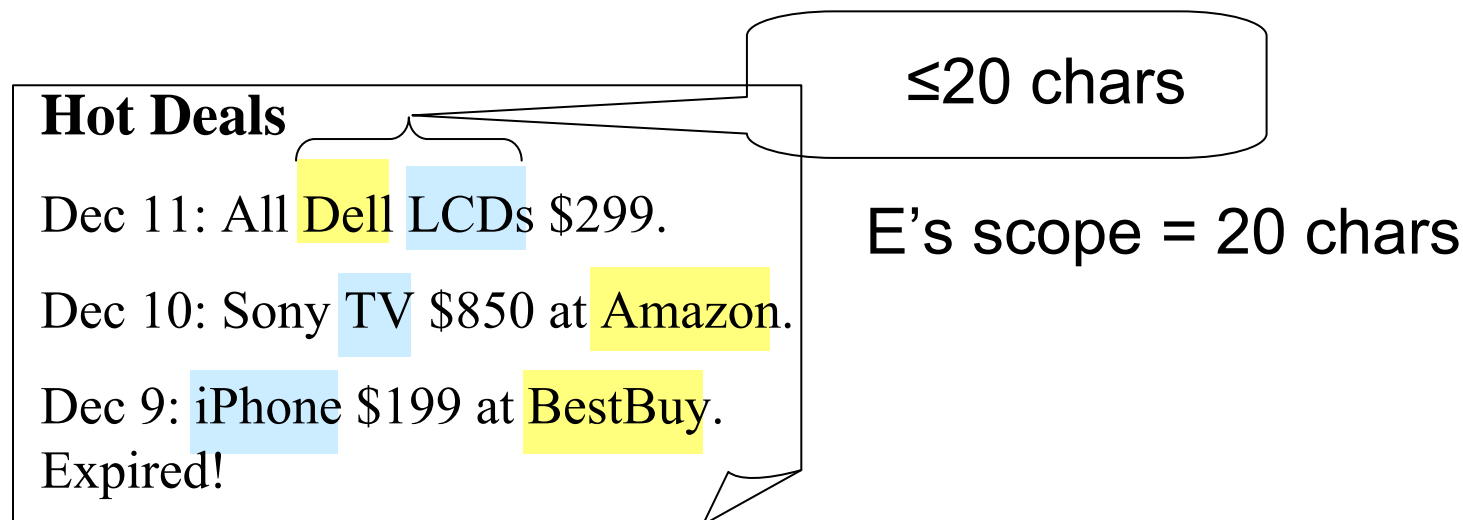
Deals

product	seller
LCDs	Dell
TV	Amazon

# Extractor Properties: Scope

- **Mention attributes often appear in close proximity.**
  - an extractor E has scope  $\alpha$  iff any mention produced by E at most spans  $\alpha$ .

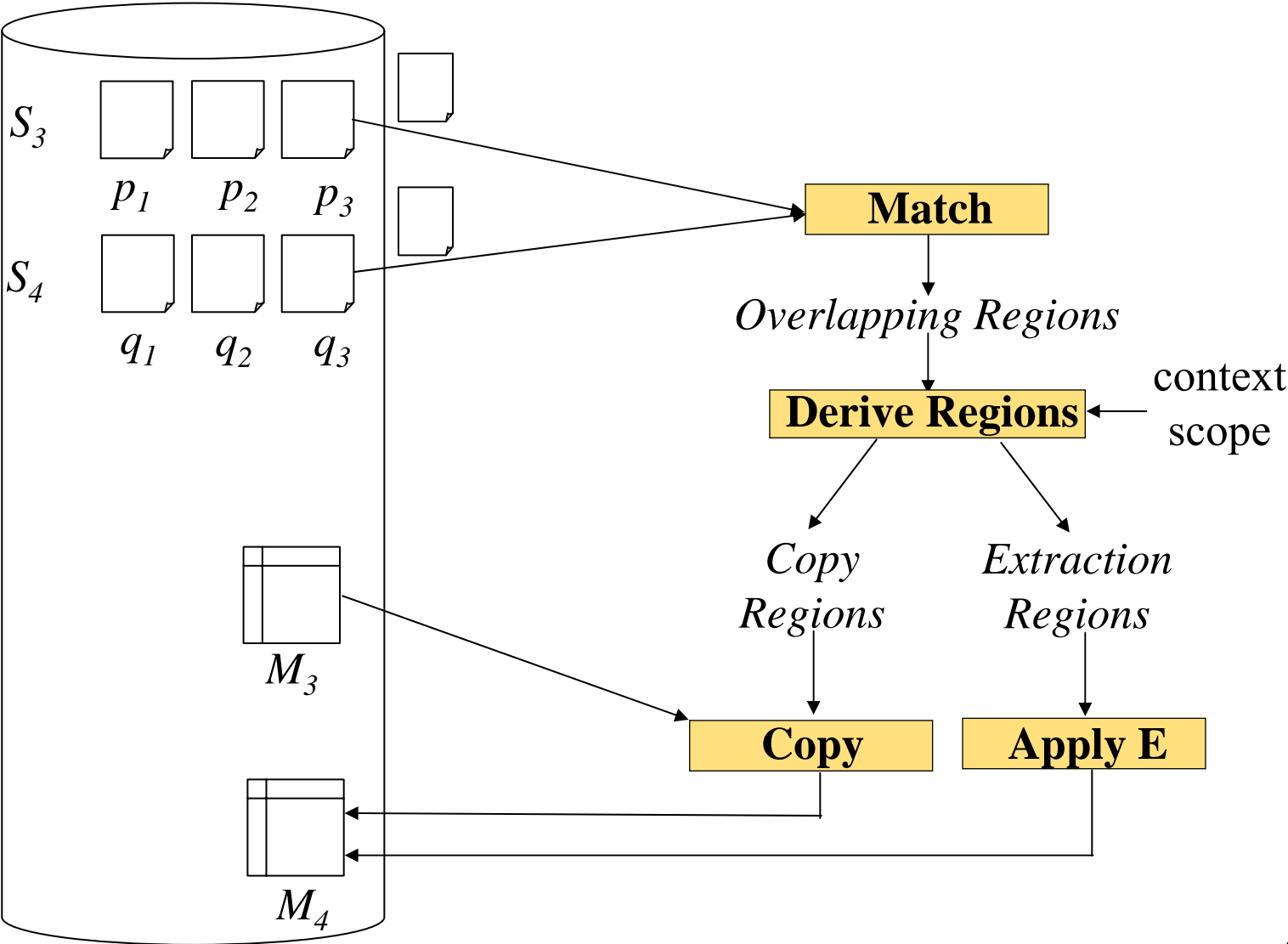
E extracts a deal if (a) product & seller are within a window of at most 20 chars  
(b) no “Expired!” is within 10 chars around the window.



# Obtain Scope and Context

- **Cyclex takes E, its scope & context**
- **Scope & context are provided by**
  - who writes E
  - who knows how E works
- **Even with loose scope and context**
  - can still guarantee correctness
  - can still reduce runtime
  - tighter scope/context → more recycling

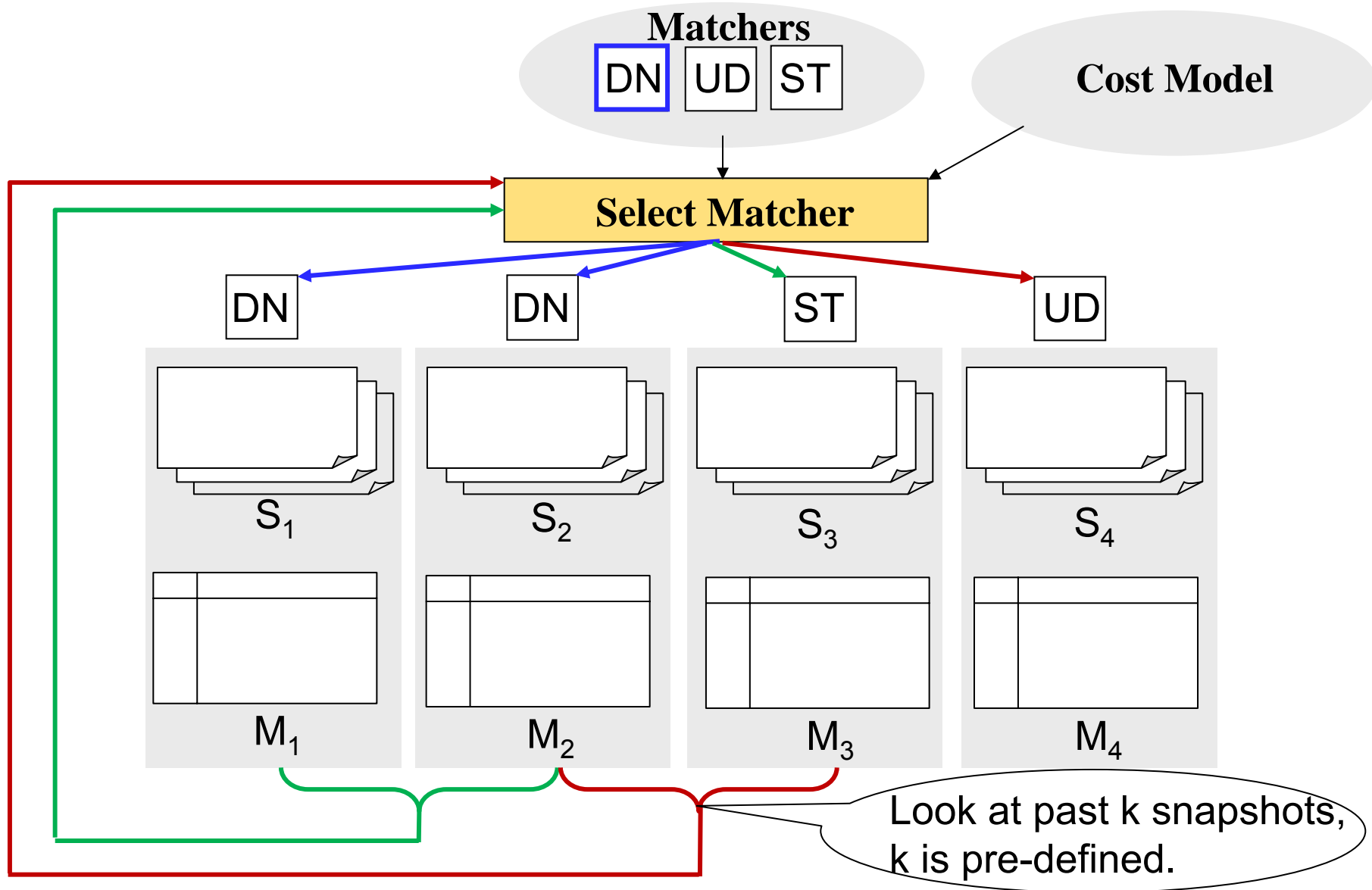
# Baseline Approach



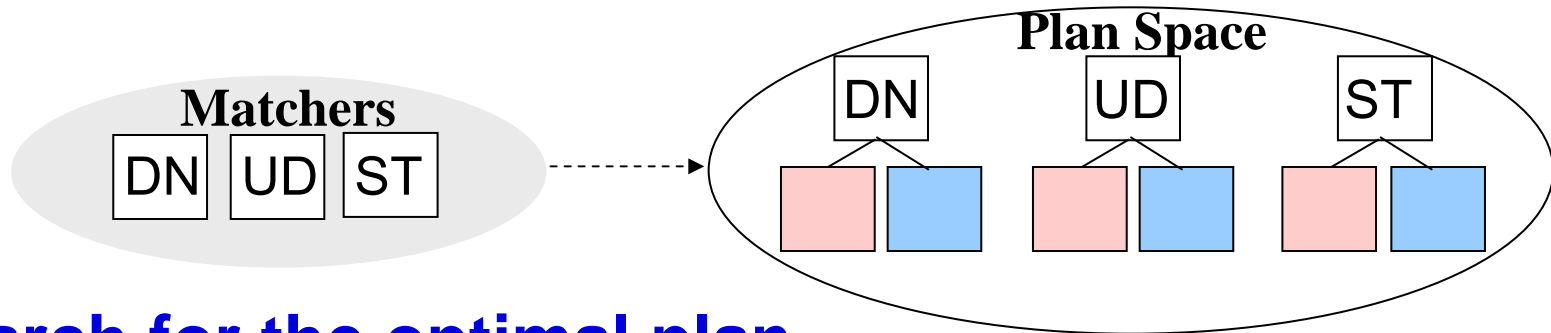
# Matchers

- **Consider 3 matchers (more can be added)**
  - DN (Doing Nothing): **0** runtime, **no** overlapping regions
  - UD (Unix Diff): fast, find some overlapping regions
  - ST (Suffix Tree): relatively slow, **all** overlapping regions
- **Trade off runtime vs. size of overlapping regions**
- **No matcher is always optimal**

# Select An Optimal Matcher: Adaptive and Cost-based Optimization



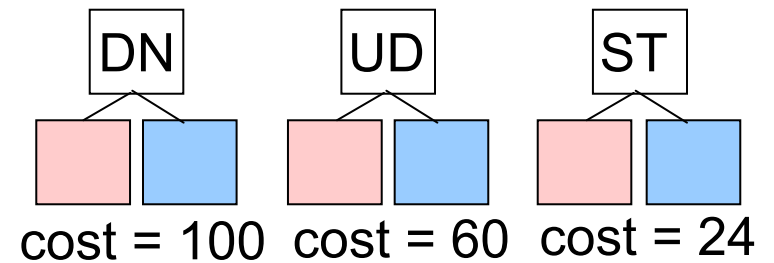
# Select An Optimal Matcher: Cost-based Selection



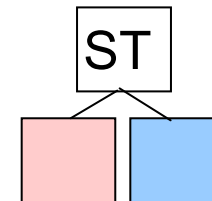
- **Search for the optimal plan**

1. enumerate all possible plans

2. use cost model to estimate cost of each plan



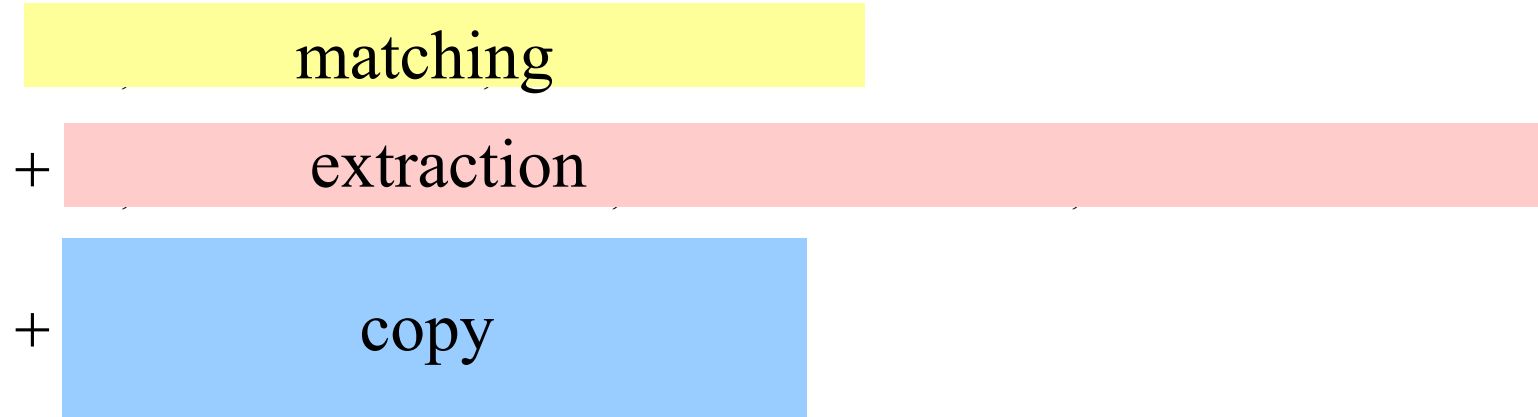
3. select plan with minimal cost





# Cost Model

- Captures matching/extraction/copy times

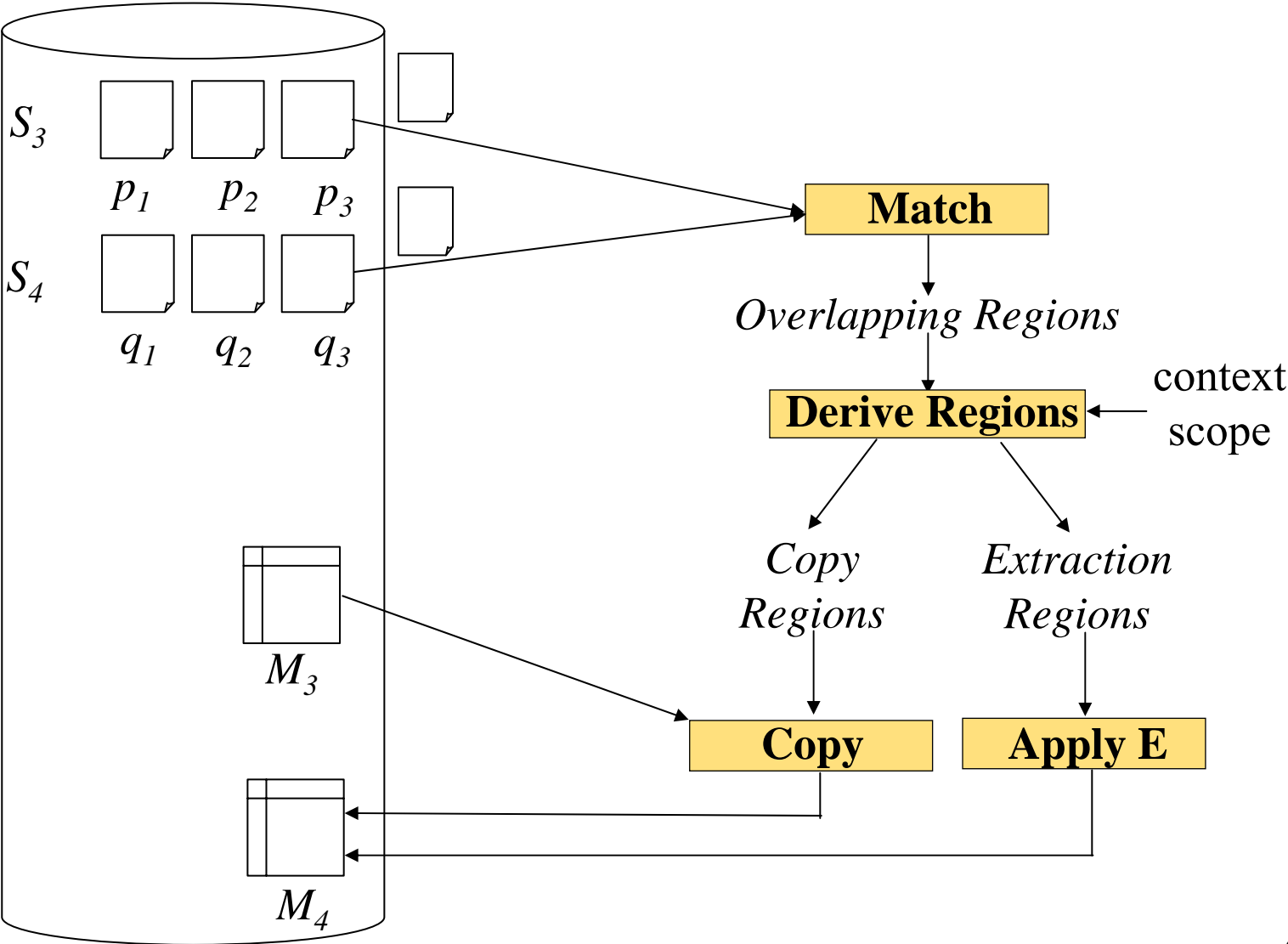


- Captures text properties

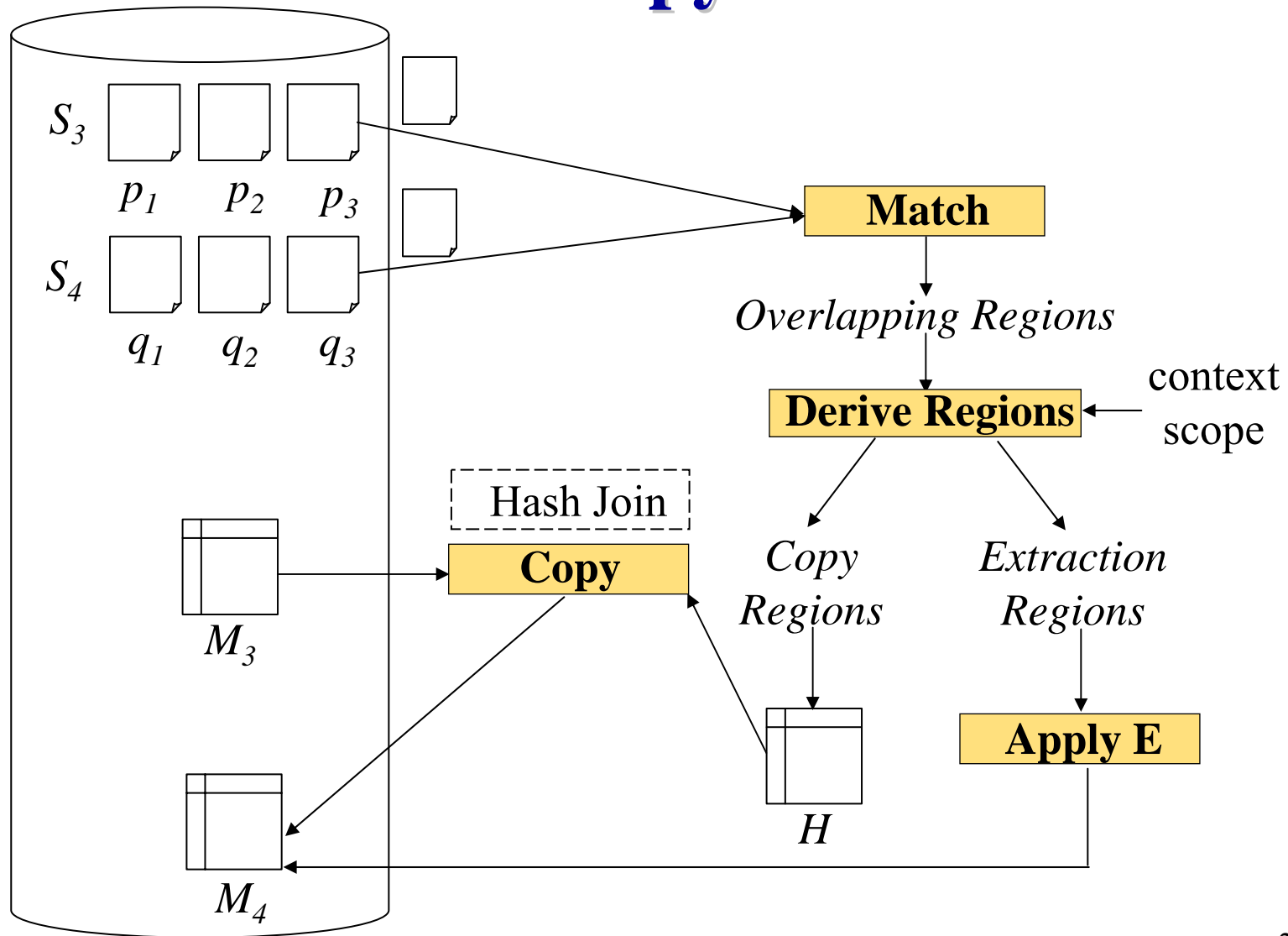
$$\hat{w}_{1,ex} \cdot \underbrace{m \cdot f \cdot l \cdot \hat{g}}_{\text{length of new regions on matched pages}} \cdot \text{length (in byte)/page}$$

# of pages in  $S_4$       fraction of old pages      fraction of new regions

# Baseline Approach



# Interleave Matching, Extraction and Copy



# Experimental Setup

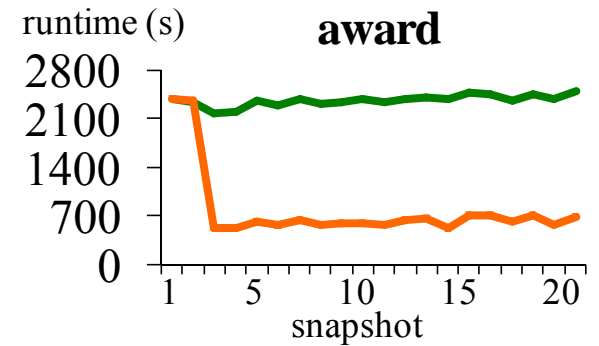
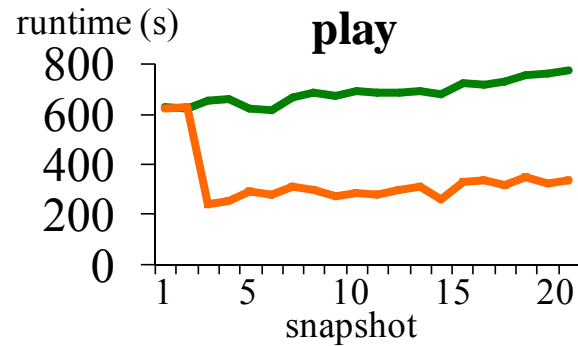
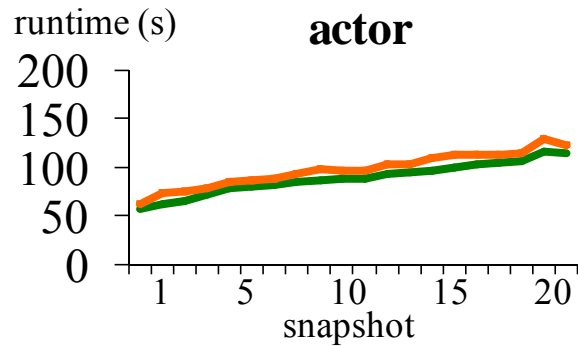
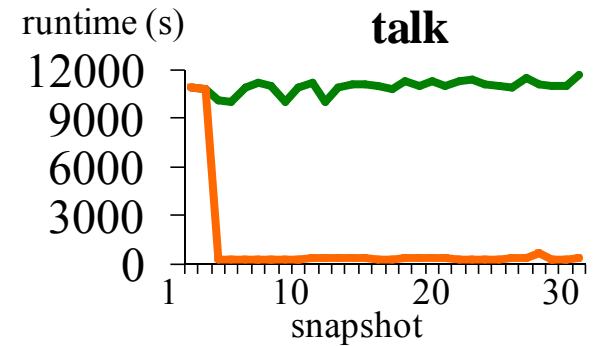
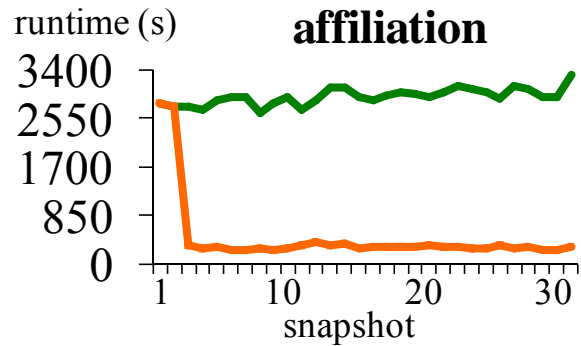
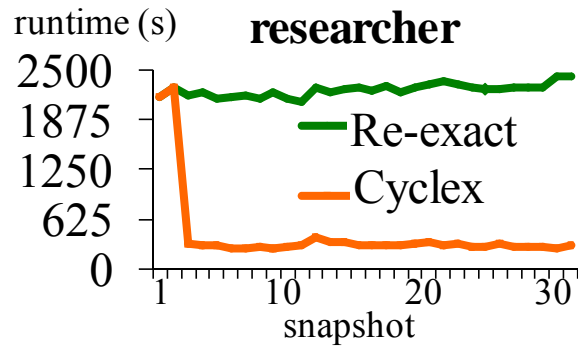
- Datasets**

Data Sets	DBLife	Wikipedia
# Data Sources	980	925
# Snapshots	30	20
Time between snapshots	1 day	21 days
Avg # Page per Snapshot	10155	3038
Avg Size per Snapshot	180M	35M

- Extractors**

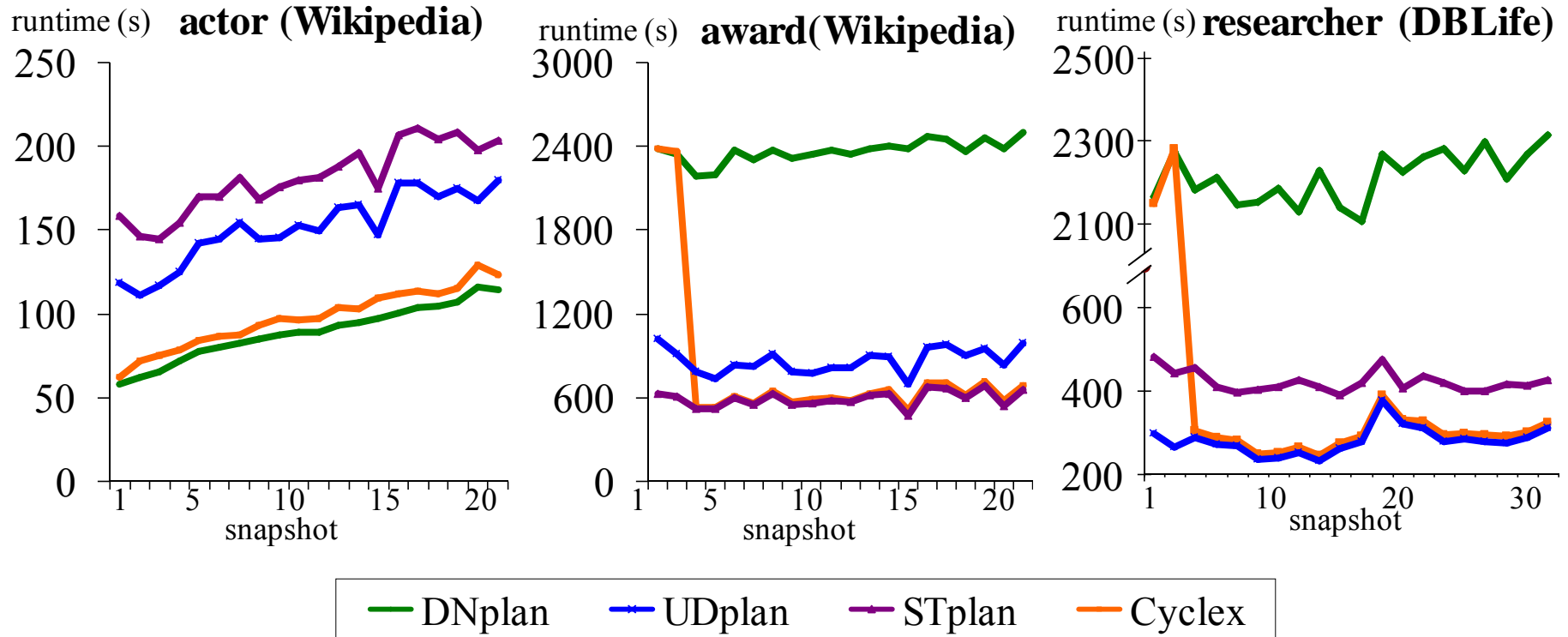
	DBLife			Wikipedia		
	researcher	affiliation	talk	actor	play	award
Scope $\alpha$	32	93	400	35	96	250
Context $\beta$	3	7	10	3	4	10

# Benefit of Recycling IE Results



- In 5 out of 6, outperformed extract-from-scratch by 50-90%

# Importance of Optimization



- No matcher is uniformly optimal

# Summary of Cyclex

- **Guarantee correctness**
  - model extractors using **scope and context**
- **Choose a good way to match pages**
  - adaptive optimization
  - cost-based decisions using **text specific** cost model
- **Efficiently interleave matching, copy and extraction**
  - a way to **scan** data **once**

# My Dissertation Contributions

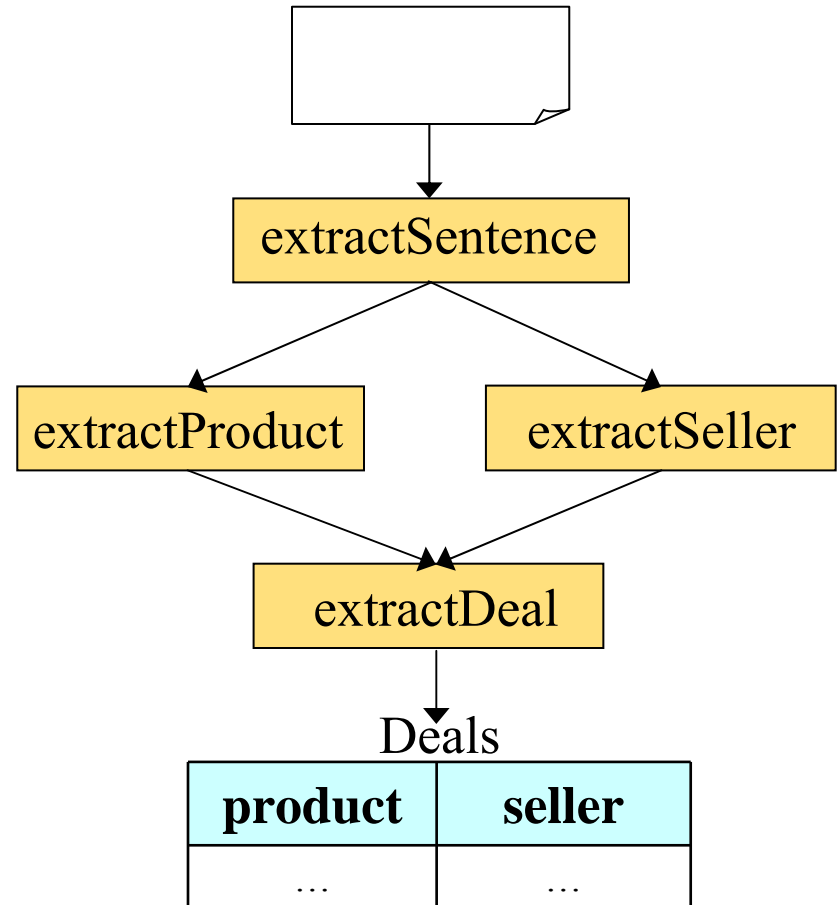
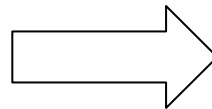
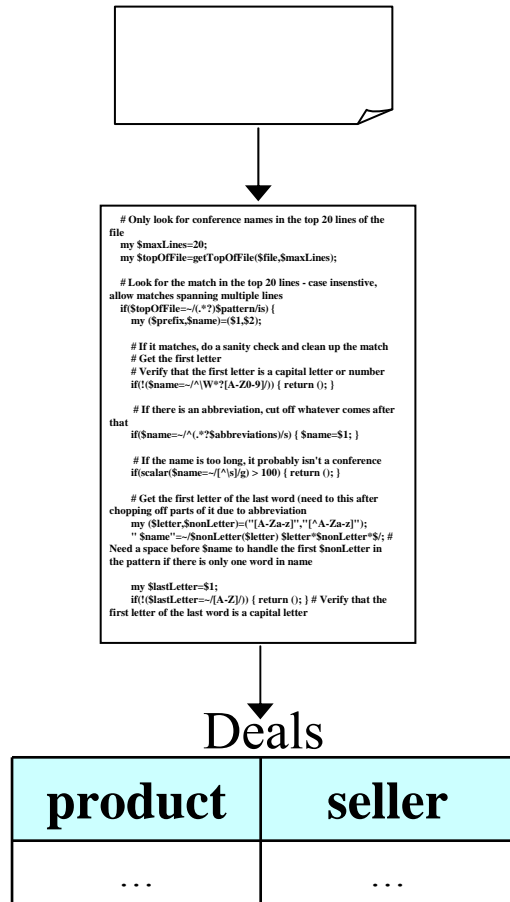
**Developed efficient IE solutions over evolving text that**

- recycle previous IE efforts
- guarantee correctness
- deal with large text corpora
- **deal with large IE programs**
- deal with complex learning-based IE programs

**Results in [ICDE-08], [SIGMOD-09] and [TechReportA-10]**



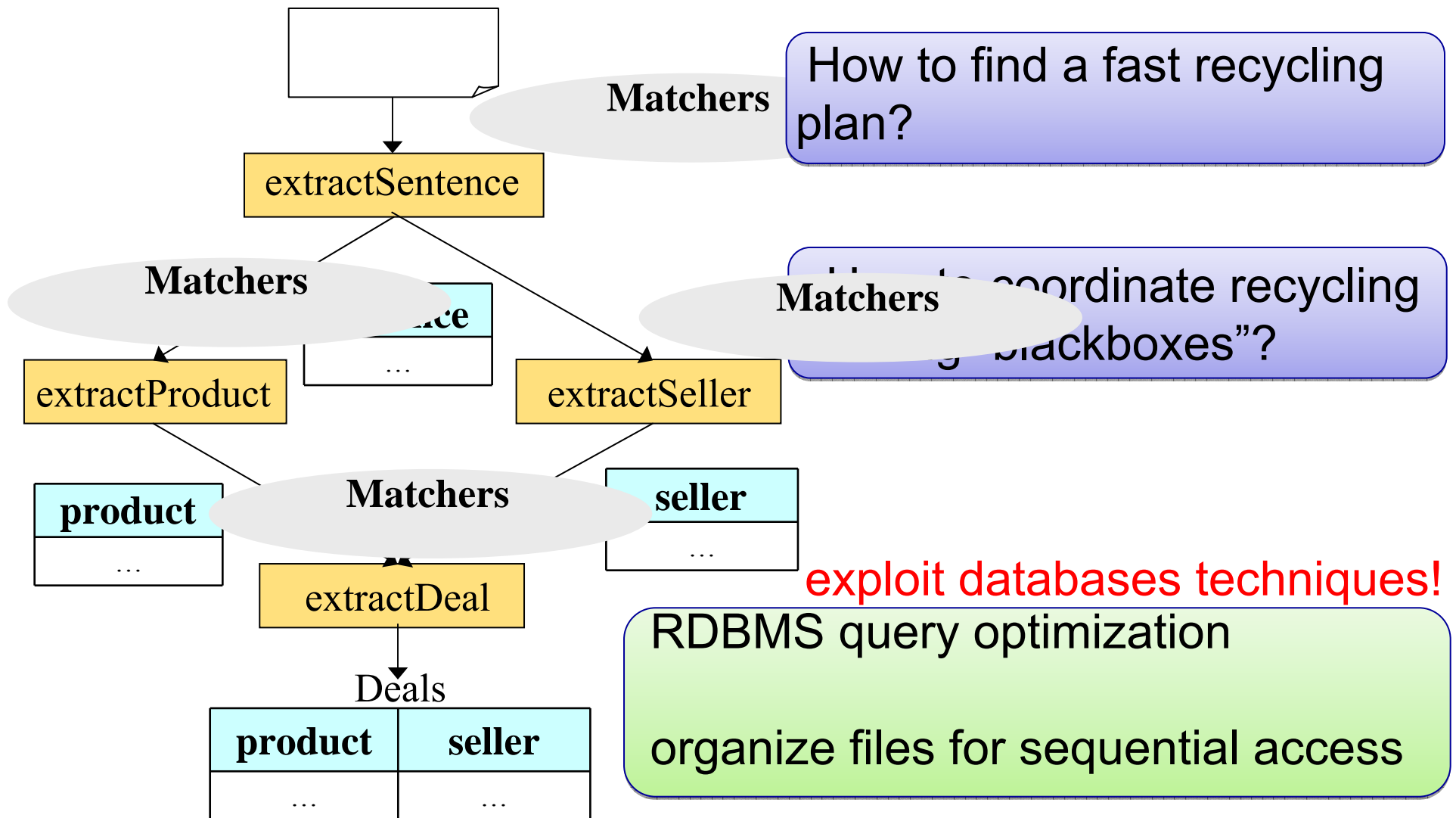
# Large and Complex IE Programs



- **Real-world IE programs are complex**

- Avatar@IBM: 25+ blackboxes
- DBLife@WISC: 45+ blackboxes

# Delex [SIGMOD09]: Decompose and Recycle



# Experiment Setup

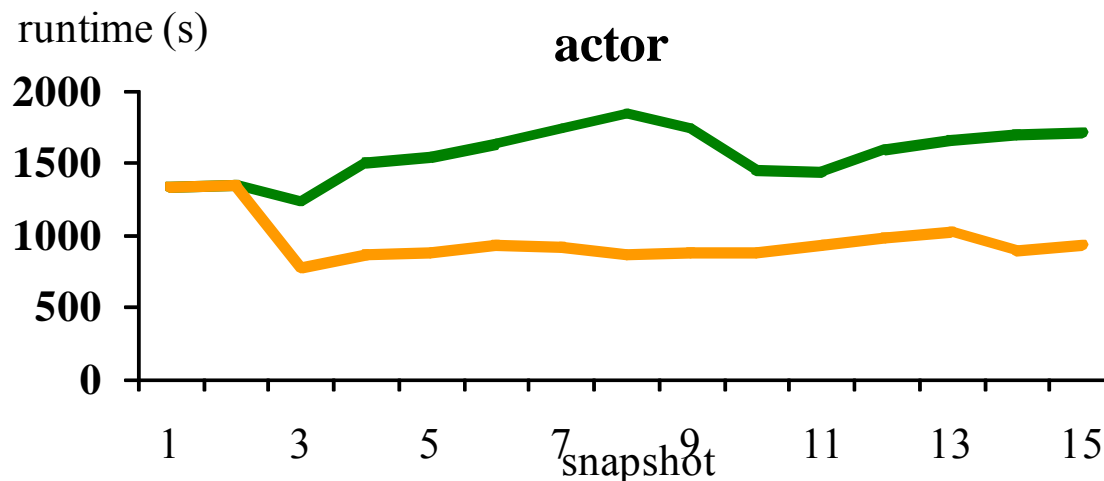
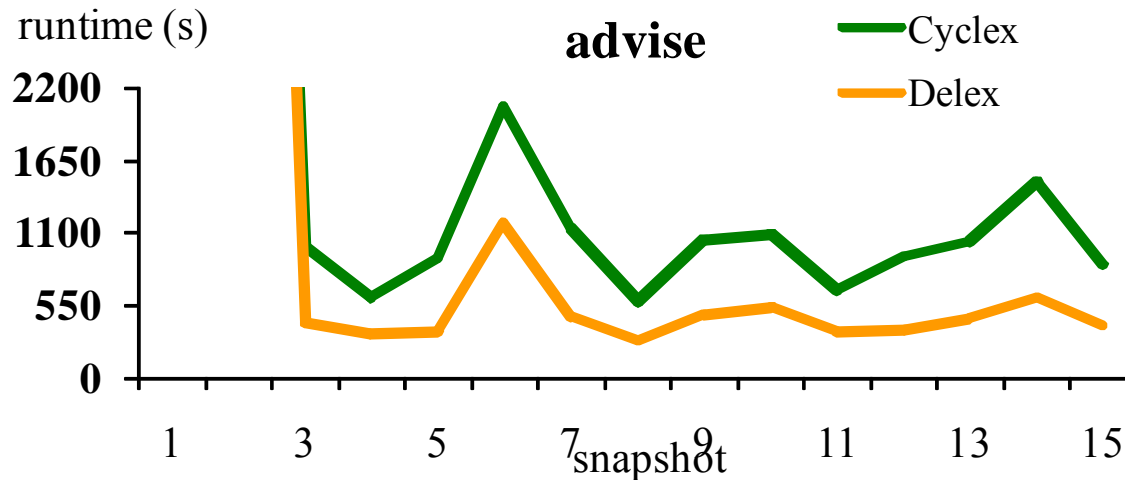
- Datasets**

Data Sets	DBLife	Wikipedia
# Snapshots	15	15
Time between snapshots	2 days	21 days
Avg # Page per Snapshot	10155	3038
Avg Size per Snapshot	180M	35M

- IE Programs : Rule-based and Learning-based IE Programs**

	DBLife (Rule-based)			Wikipedia (Rule-based)			Wikipedia (Learning -based)
	talk	chair	advise	blockbuster	play	award	actor
# of IE “Blackboxes”	1	3	5	2	4	6	5

# Runtime Comparison



- **Delex drastically cuts runtime of Cyclex by 45-71%**  
**(See paper for more experiments)**

# Related Work

- **Early works on IE**
  - [Bikel UAI-97] [McCallum KDD-00]...
  - focus on improving accuracy
- **Recent works on IE**
  - [tutorial in KDD-06, SIGMOD-06] [Gravano et al, SIGMOD-06]...
  - consider developing efficient IE
  - do not consider evolving text corpora
- **Evolving text data**
  - [McCann VLDB05] [Lim WWW03]
  - consider problems other than IE
- **Incremental view maintenance**
  - [Gupta&Mumick][Widom et al, SIGMOD95]...
  - only consider relational operators
  - assume changes to the inputs are available

# My Other Contributions

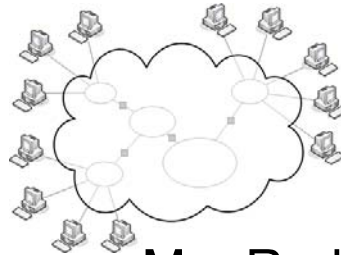
- **Web community systems**

- DBLife [[DE Bulletin-06](#), [VLDB-07a](#), [CIDR-09](#)]



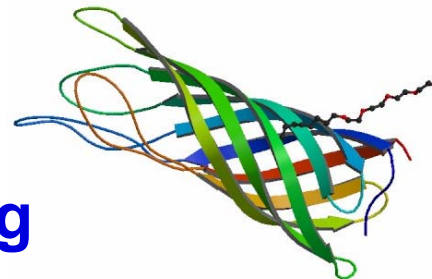
- **Text mining in cloud**

- scalable mining algorithms on MapReduce clusters [[TechReportB-10](#)]



- **Biology sequences mining**

- incorporate domain knowledge into mining algorithms [[KDD-03](#), [Bioinformatics-04](#)]
- interpret complex mining models with human-understandable rules [[SFUThesis](#)]



# Future Work

- **Maintain IE and information integration (II)**
- **Develop/Maintain IE/II over large clusters**
  - build efficient, scalable and robust IE/II
  - **declaratively** develop and **automatically** optimize IE/II
- **Develop user-friendly and efficient data analysis tools**
  - analyze non-structured data, e.g., text and biology sequences
  - let non-database users such as biologists ask query **easily** and obtain answers **efficiently**

# Conclusions

- **IE over evolving text is increasingly important!**
- **Developed solutions that**
  - recycle previous IE to reduce runtime
  - guarantee correctness
  - deal with large text corpora / large IE programs
  - deal with complex learning-based IE programs
- **Database techniques are increasingly critical for developing efficient IE solutions**
  - in collaboration with other communities: NLP, AI, Web, IR, KDD, HCI