# Chapter 4  Webpage Information Extraction

**Li fang**

**Dept.of Computer Science**

**Shanghai Jiao Tong university**

# Contents

- Overview of Information Extraction tools from Web pages

- Wrapper Induction

- Wrapper Maintenance

# Two kinds of webpages

- Multiple-record page extraction (left)
- One-record page extraction (right)

| Tools | | Degree of Automation | Support for Complex Objects | GUI | XML Output | Support for Non-HTML Sources | Type of Page Contents |
|---|---|---|---|---|---|---|---|
| Languages | Minerva | Manual | Coding | No | Yes | Partial | SD |
| | TSIMMIS | Manual | Coding | No | No | Partial | SD |
| | Web-OQL | Manual | Coding | No | No | None | SD |
| HTML-aware | W4F | Semi-Automatic | Coding | Yes | Yes | None | SD |
| | XWRAP | Automatic | Yes | Yes | Yes | None | SD |
| | RoadRunner | Automatic | Yes | Yes | No | None | SD |
| NLP-based | WHISK | Semi-Automatic | No | Yes | No | Full | ST |
| | RAPIER | Semi-Automatic | No | Yes | No | Full | ST |
| | SRV | Semi-Automatic | No | Yes | No | Full | ST |
| Induction | WIEN | Semi-Automatic | No | Yes | No | Partial | SD |
| | SoftMealy | Semi-Automatic | Partial | Yes | No | Partial | SD |
| | STALKER | Semi-Automatic | Yes | Yes | No | Partial | SD |
| Modeling-based | NoDoSE | Semi-Automatic | Yes | Yes | Yes | Partial | SD |
| | DEByE | Semi-Automatic | Yes | Yes | Yes | Partial | SD |
| Ontology-based | BYU | Manual | Coding | Yes | No | Full | |

Summary of the Qualitative A

Methods of extraction pages

SD:semi-structured data
ST:semi-structured text

# **Language** for wrapper development—for manually constructed IE systems

**Minerva:** combines a declarative grammar-based approach with features typical of <u>procedural programming languages</u>.

**Tsimmis**: includes wrappers that can be configured through <u>specification files </u>written by the user.

**Web-OQL**: originally aimed at performing SQL-like queries over the Web.

# Web-OQL

- **Hypertrees are arc-lableled ordered trees.**



- Tag name
- The piece of HTML code
- Text excluding marku

# Web-OQL (cont.)

- Query: extracts the reviewer names "Jeff" and "Jane" from page *pe2*。

# Overview of Web data extraction tools (cont.)

- **HTML-aware Tools**
- W4F(world wide web wrapper factory): <u>a toolkit</u> for building wrappers.
- XWRAP: <u>a component library</u> that provides basic building blocks for wrapper development.
- **NLP-based tools**

  PAPIER (job posting), SRV, WHISK: suitable for Web pages consisting of <u>grammatical text</u>, such as job listings, apartment rental advertisements, seminar announcements.

# NLP based tools: PAPIER

**BookTitle extraction rule-**

| Pre-filler pattern | Filler pattern | Post-filler pattern |
|---|---|---|
| (1) word: Book | list: len: 2 | word: <b> |
| (2) word: Name | Tag: [nn, nns] | |
| (3) word: </b> | | |

```
01:    <html><body>
02:    <b>
03:        Book Name
04:    </b>
05:    Data mining
06:    <b>
07:        Reviews
08:    </b>
```

- Extraction rule for the book title:

- Preceded by words "Book", "Name", and "</b>"

- Followed by the word "<b>".

- The "Filler pattern" specifies that the <u>title consists of at most two words that were labeled as "nn" or "nns" by the POS tagger</u> (i.e., one or two singular or plural common nouns).

# Overview of Web data extraction tools (cont.)

- **Modeling-based Tools**

- NoDoSE: an interactive tool for semi-automatically determining the structure of Web page.

- DEByE: an interactive tool to extract page contents based on a set of example objects.

- **Ontology-based Tools**

   Ontologies are previously constructed to describe the data of interest, including relationships, lexical appearance and context keywords.

# **Overview** of Web data extraction tools

| Tools | | Degree of Automation | Support for Complex Objects | GUI | XML Output | Support for Non-HTML Sources | Type of Page Contents |
|---|---|---|---|---|---|---|---|
| Languages | Minerva | Manual | Coding | No | Yes | Partial | SD |
| | TSIMMIS | Manual | Coding | No | No | Partial | SD |
| | Web-OQL | Manual | Coding | | | | SD |
| HTML-aware | W4F | Semi-Automatic | | | | | SD |
| | XWRAP | Automatic | | | | | SD |
| | RoadRunner | Automatic | | | | | SD |
| NLP-based | WHISK | Semi-Automatic | | | | | ST |
| | RAPIER | Semi-Automatic | | | | | ST |
| | SRV | Semi-Automatic | No | Yes | No | Full | ST |
| Induction | WIEN | Semi-Automatic | No | Yes | No | Partial | SD |
| | SoftMealy | Semi-Automatic | Partial | Yes | No | Partial | SD |
| | STALKER | Semi-Automatic | Yes | Yes | No | Partial | SD |
| Modeling-based | NoDoSE | Semi-Automatic | Yes | Yes | Yes | Partial | SD |
| | DEByE | Semi-Automatic | Yes | Yes | Yes | Partial | SD |
| Ontology-based | BYU | Manual | Coding | Yes | No | Full | ST/SD |

**WRAPPER induction tools**
WIEN, SoftMealy, Stalker

Table 1: Summary of the Qualitative Analysis

lecture of Internet-based IE technolgies

# Wrapper Technologies

- What is wrapper
- Wrapper induction
- Wrapper maintenance

# What is Wrapper?

- For information integration

  A procedure that is designed for extracting content of a particular information source and delivering the content of interesting in a self-describing representation (eg.XML)
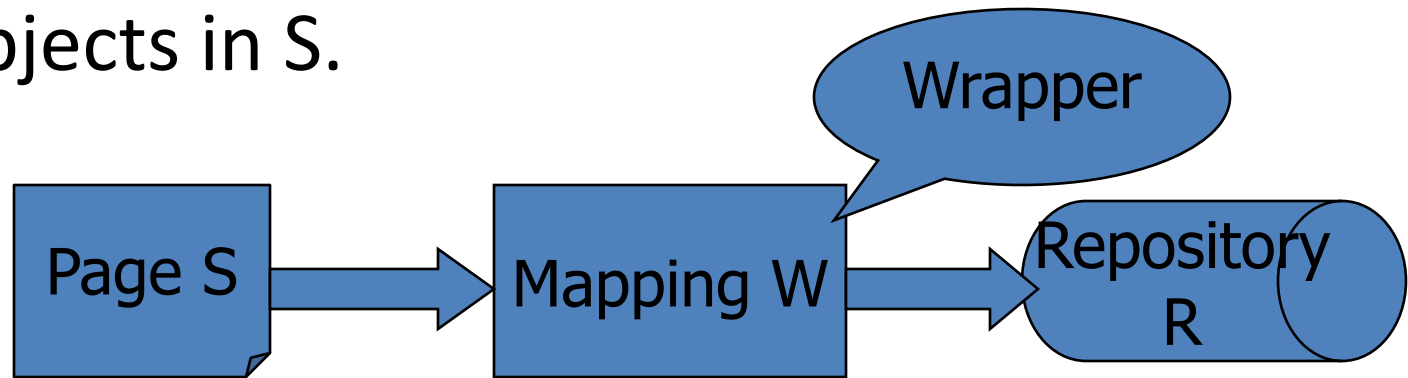
- For Web application

  – An extracting program to extract desired information from Web pages.

  Semi-Structure Doc.– wrapper→ Structure Info.

# For Web Applications:

- Given a Web page S containing a set of implicit objects, determine a mapping W that populates a data repository R with the objects in S.

Wrapper

Page S → Mapping W → Repository R

Similar S pages

# An example for a wrapper



(a)                                          (b)

```
<HTML><TITLE>Some Country Codes</TITLE><BODY>
<B>Congo</B>  <I>242</I><BR>
<B>Egypt</B>  <I>20</I><BR>
<B>Belize</B>  <I>501</I><BR>
<B>Spain</B>  <I>34</I><BR>
```

procedure ccwrap$_{LR}$(page $P$)
  while there are more occurrences in $P$ of '<B>'
    for each $\langle \ell_k, r_k \rangle \in \{('\texttt{<B>}', '\texttt{</B>}'), ('\texttt{<I>}', '\texttt{</I>}')\}$
      scan in $P$ to next occurrence of $\ell_k$; save position as start of $k$th attribute
      scan in $P$ to next occurrence of $r_k$; save position as end of $k$th attribute
  return extracted $\{\ldots, \langle \text{country}, \text{code} \rangle, \ldots\}$ pairs

# Wrapper Induction

- Web wrappers wrap…
  - "Query-able" or "Search-able" Web sites
  - Web pages with large itemized lists
- The primary issues are:
  - How to build the extractor quickly?
  - Wrapper induction algorithms search a hypothesis space of possible wrapper programs for a wrapper that has high extraction accuracy on a set of training pages.
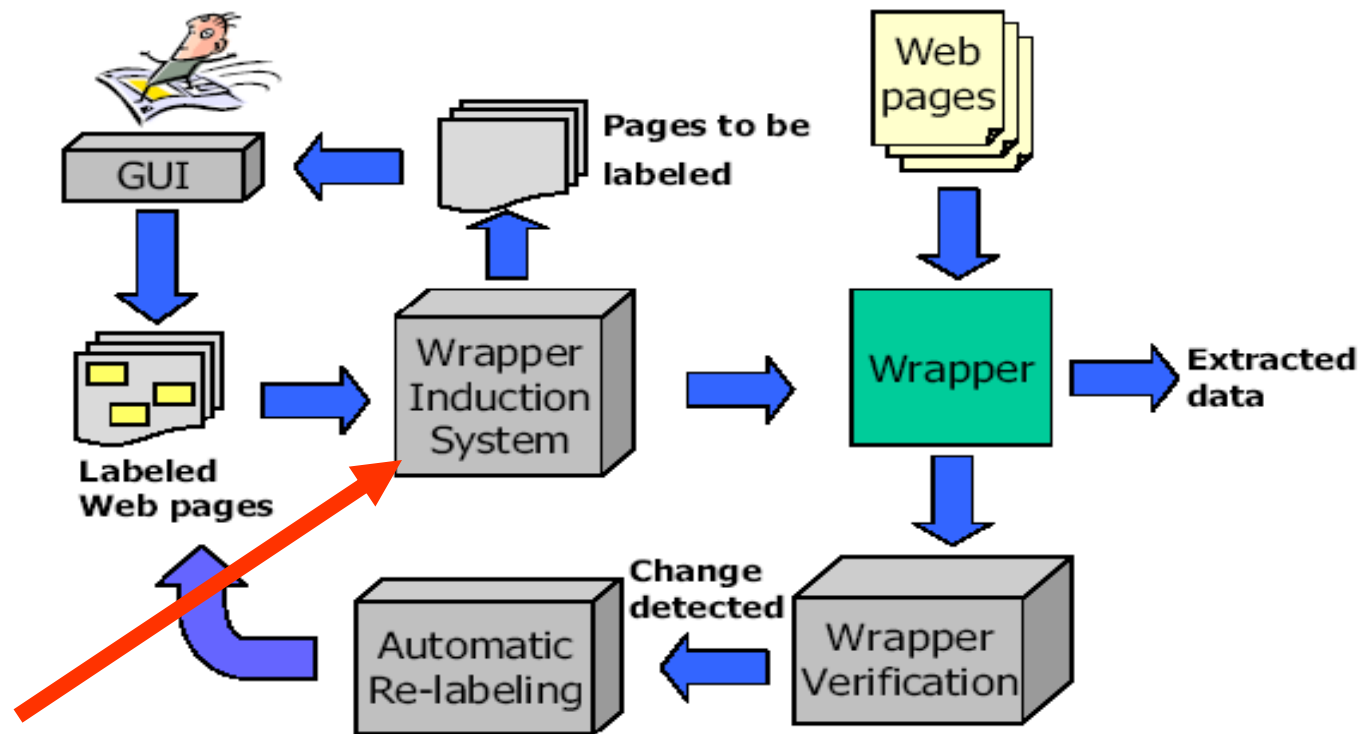
# Wrapper Induction: Methods

- **Manually writing wrappers**

- Tedious, time consuming task, eg. TSIMMIS, Minerva, ...

- **Wrapper  programming languages**
- Florid (a logic-programming formalism), pillow (an HTML/XML programming library for logic programming systems) ...

- **Machine learning methods**
- Stalker, Softmealy, WIEN ...

-  **Supervised interactive wrapper**
- W4F (uses an SQL-like query called HEL), Xwrap (uses a procedural rule system), ...

# Wrapper Induction Tools

- ## WIEN:

- Input: a set of pages where data of interest is labeled to serve as examples

- Output:a wrapper that is consistent with each labeled page.

- ## SoftMealy

- Using finite-state transducers (FST) which takes a sequence of tokens as input and matches the context separators with contextual rules to determine state transitions

- ## Stalker

- The wrapper induction techniques used in WIEN and SoftMealy are further developed in Stalker

# Wrapper Induction: machine learning methods (Stalker)



Figure 1: The Lifecycle of a Wrapper

**Our focus here**

lecture of Internet-based IE technolgies

# Learning Extraction Rules
## ---from pages

- Aim:

Defining a set of extraction rules that precisely define how to locate the information on the page.

→ How to describe the content of a page?

# Describing the content of a page: Embedded Catalog Tree

- Embedded catalog (EC): a tree-like structure to represent a Web page.

- Leaves: items of interest for the user

- Internal nodes: lists of k-tuples where each item in the k-tuple can be either a *leaf* or another *list L.*

# Embedded Catalog Tree (for example)



**LA WEEKLY**

## LA Restaurants

Search Criteria: Name: killer shrimp Location: Any Cuisine: Any

**KILLER SHRIMP**
523 Washington Blvd., Marina del Rey
(310) 576-2293

Food for the gods — fresh, sweet, tender, succulent, big Louisiana shrimp floating in a heavenly spicy sauce. You want it, Killer's got it, you deserve it. Around for eight years, Killer Shrimp is a popular hot spot and has become one of L.A.'s landmark dining experiences — tourists and natives all seem to know that this is the place to satisfy cravings for the real thing. Indoor and patio dining. Lunch and dinner seven days. Beer and wine, takeout, parking. MC, V.
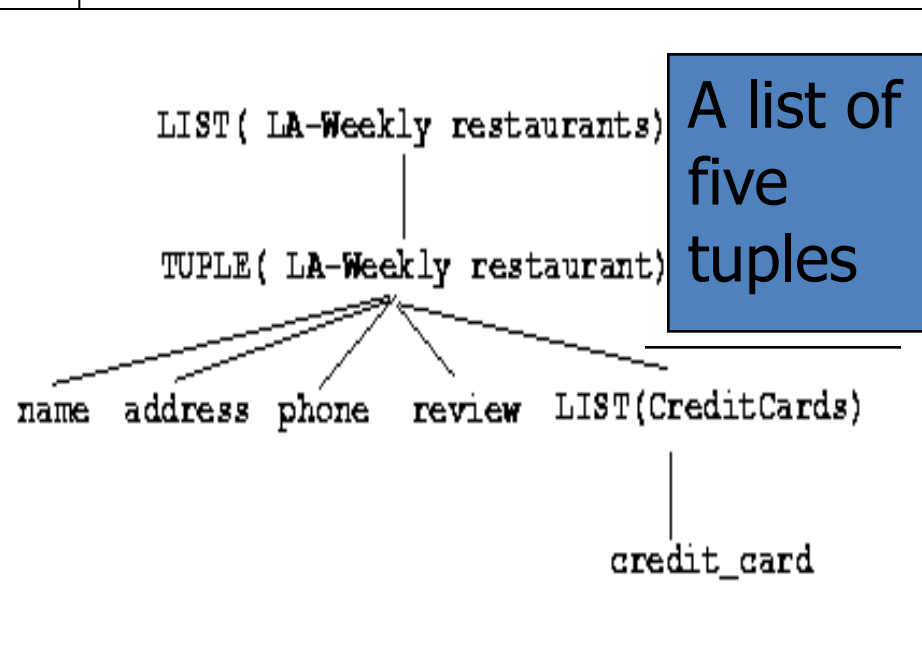
**KILLER SHRIMP**
4113 N. Pacific Coast Hwy., Redondo Beach
(310) 798-0008

Food for the gods — fresh, sweet, tender, succulent, big Louisiana shrimp floating in a heavenly spicy sauce. You want it, Killer's got it, you deserve it. Around for eight years, Killer Shrimp is a popular hot spot and has become one of L.A.'s landmark dining experiences — tourists and natives all seem to know that this is the place to satisfy cravings for the real thing. Indoor and patio dining. Lunch and dinner seven days. Beer and wine, takeout, parking. MC, V.

**KILLER SHRIMP**
4000 Colfax Ave., Studio City
(818) 508-1570

Food for the gods — fresh, sweet, tender, succulent, big Louisiana shrimp floating in a heavenly spicy sauce. You want it, Killer's got it, you deserve it. Around for eight years, Killer Shrimp is a popular hot spot and has become one of L.A.'s landmark dining experiences — tourists and natives all seem to know that this is the place to satisfy cravings for the real thing. Indoor and patio dining. Lunch and dinner seven days. Beer and wine, takeout, parking. MC, V.

LIST( LA-Weekly restaurants)

**A list of five tuples**

TUPLE( LA-Weekly restaurant)

name  address  phone  review  LIST(CreditCards)

credit_card

Figure 2: $\mathcal{EC}$ description of LA-Weekly pages.

et-based IE technolgies

# Extracting Rule based on EC

- **A rule:** for each node $x$ in the EC Tree, the wrapper needs a rule $r$ that extracts that particular node from its root, $p$ is a path from the root to the leaf.

- **A list iteration rule:** decomposes $p$ into individual tuples, and then apply $r$ to each extracted tuple for each list node.

# Example for extraction rules

R1=SkipTo(Address) SkipTo(I)

R2=SkipTo(Address : <I>)

R3=SkipTo(Cuisine : <I> )
SkipTo(Address : <I>)

Figure 2: Two Sample Restaurant Documents Fr... at Guide.

```
...Cuisine:<i>Seafood</i><p>Address:<i>  12 Pico St. </i><p>Phone:<i>...
...Cuisine:<i>Thai    </i><p>Address:<i> 512 Oak Blvd.</i><p>Phone:<i>...
...Cuisine:<i>Burgers</i><p>Address:<i> 416 Main St. </i><p>Phone:<i>...
...Cuisine:<i>Pizza</i><p>Address:<b> 97 Adams Blvd. </b><p>Phone:<i>...
```

# Example for extraction rules (cont.)

R4=SkipTo(Cuisine : <I> _Capitalized_</I> <p> Address : <I>)

R4 is defined based on a 9-token landmark that uses the wildcard_Capitalized_, which is a placeholder for any capitalized alphabetic string.

Disjunctive rules: either R1 or R2

To deal with variations in the format of the documents, disjunctions are allowed to use.

# Extraction Rules as Finite Automata

- ## Landmarks: each argument of a skipTo()

  - ### A sequence of tokens and wildcards

- ## Landmark automata

  - ### A non-deterministic finite automata

  $l_{i,j}$

  $S_i \rightarrow S_j$

  The transition takes place if the automaton is in the state $S_i$ and the landmark $l_{i,j}$ matches the sequence of tokens at the input.

# Landmark Automata (linear LA)

- A linear LA has one accepting state.

- from each non-accepting state, there are exactly two possible transitions: a loop to itself, and a transition to the next state;

- each non-looping transition is labeled by a landmarks;

- all looping transitions have the meaning "consume all tokens until you encounter the landmark that leads to the next state".

# Rules and its automaton

R1::=skipTo(()),
R2::=skipTo(phone)
skipTo(<b>).

*Disjunctive rule
either R1 or R2*



: An $\mathcal{SLG}$ for the start of the area code.

- The initial state $S_0$ has a branching-factor of k.

- It has exactly k accepting states. (one per branch)

- All k branches that leave the $S_0$ are sequential LA.

# Learning Extraction Rules

User marking

E1: 513 Pico, <b>Venice</b>, Phone: 1-<b> 800 b>-555-1515

E2: 90 Colfax, <b> Palms </b>, Phone: ( 818 ) 508-1570

E3: 523 1st St., <b> LA </b>, Phone: 1-<b> 888 </b>-578-2293

03 La Tijera, <b> Watts </b>, Phone: ( 310 ) 798-0008

r examples of restaurant addresses.

Marked samples

-it accepts the positive examples in E2 and E4
-it rejects both E1 and E3 because R1 can not be matched on them. R2 can do.

STALKER algorithm

Extraction Rules

R1::=skipTo(()), R2::=skipTo(-<b>)

lecture of Internet-based IE technolgies

# Process of the example (STALKER)

1. <span style="color:red">Generating a linear LA</span> that covers as many as possible of the four positive examples.

2. Create another linear LA for the remaining examples, and so on.

3. Once STALKER covers all examples. It returns the disjunction of all the induced LAs.

# STALKER Algorithm

**LearnRule(** *Examples* **)**
- let *RelVal* be an empty rule
- WHILE *Examples* $\neq \emptyset$
  - *aDisjunct* =**LearnDisjunct(***Examples***)**
  - remove all examples covered by *aDisjunct*
  - add *aDisjunct* to *RelVal*
- return **OrderDisjuncts(***RelVal***)**

**LearnDisjunct(** *Examples* **)**
- let *Seed* $\in$ *Examples* be the shortest example
- *Candidates* = **GetInitialCandidates(** *Seed* **)**
- DO
  - *BestRefiner* = **GetBestRefiner(** *Candidates* **)**
  - *BestSolution* = **GetBestSolution(** *Candidates* $\cup$ {*BestSolution*} **)**
  - *Candidates* = **Refine(***BestRefiner, Seed***)**
  WHILE **IsNotPerfect(***BestSolution***)** AND *BestRefiner* $\neq \emptyset$
- return **PostProcess(***BestSolution***)**

See example

lecture of Internet-based IE technolgies

# STALKER Algorithm (cont.)

**BestRefiner()**
Prefer candidates that have:
- larger coverage
- more early matches
- more failed matches
- fewer wildcards
- shorter unconsumed prefixes
- fewer tokens in *SkipUntil*()
- longer end-landmarks

**BestSolution()**
Prefer candidates that have:
- more correct matches
- more failures to match
- fewer tokens in *SkipUntil*()
- fewer wildcards
- longer end-landmarks
- shorter unconsumed prefixes

*Figure 6.* The STALKER heuristics.

# STALKER Algorithm (cont.)

- Refine() function: obtain better disjuncts either by making its landmarks more specific (landmark refinements), or by adding new states in the automaton (topology refinements).

- Landmark refinements

- Topology refinements

# Landmark Refinement

E1: 513 Pico, \<b\>Venice\</b\>, Phone: 1-\<b\> 800 \</b\>-555-1515

E2: 90 Colfax, \<b\> Palms \</b\>, Phone: ( 818 ) 508-1570

E3: 523 1st St., \<b\> LA \</b\>, Phone: 1-\<b\> 888 \</b\>-578-2293

E4: 403 La Tijera, \<b\> Watts \</b\>, Phone: ( 310 ) 798-0008

*Figure 4.* **Four examples of restaurant addresses.**

- R4 = SkipTo\<b\>

Refine as :

$$\mathbf{R7} = SkipTo( \text{ - } \texttt{<b>})$$

$$\mathbf{R8} = SkipTo( \; Punctuation \; \texttt{<b>})$$

$$\mathbf{R9} = SkipTo( \; Anything \; \texttt{<b>})$$

# Topology Refinements

E1:     513 Pico, **Venice**, Phone: 1-**$\boxed{800}$**-555-1515

E2:     90 Colfax, **Palms**, Phone: ($\boxed{818}$) 508-1570

E3:     523 1st St., **LA**, Phone: 1-**$\boxed{888}$**-578-2293

E4:     403 La Tijera, **Watts**, Phone: ($\boxed{310}$) 798-0008

*Figure 4.* **Four examples of restaurant addresses.**

- R4 = skipTo&lt;b&gt;

Refine as :

R10: $SkipTo(Venice)\ SkipTo(\texttt{<b>})$

R11: $SkipTo(\texttt{</b>})\ SkipTo(\texttt{<b>})$

R12: $SkipTo(:)\ SkipTo(\texttt{<b>})$

R13: $SkipTo(-)\ SkipTo(\texttt{<b>})$

R14: $SkipTo(,)\ SkipTo(\texttt{<b>})$

R15: $SkipTo(Phone)\ SkipTo(\texttt{<b>})$

R16: $SkipTo(1)\ SkipTo(\texttt{<b>})$

R17: $SkipTo(Numeric)\ SkipTo(\texttt{<b>})$

R18: $SkipTo(Punctuation)SkipTo(\texttt{<b>})$

R19: $SkipTo(HtmlTag)\ SkipTo(\texttt{<b>})$

R20: $SkipTo(AlphaNum)\ SkipTo(\texttt{<b>})$

R21: $SkipTo(Alphabetic)\ SkipTo(\texttt{<b>})$

R22: $SkipTo(Capitalized)\ SkipTo(\texttt{<b>})$

R23: $SkipTo(NonHtml)\ SkipTo(\texttt{<b>})$

R24: $SkipTo(Anything)\ SkipTo(\texttt{<b>})$

Each initial candidate is a 2-state landmark automaton that is either a token t that ends one prefix(p) or a wildcard that matches such a t

**Example of rule induction**

E1 ...1-<b> 800 </b>-555-1515

... 818 ) 508-1570

...b> 888 </b>-578-2293

one: ( 310 ) 798-0008

...ddresses.

R1 (0) →[1]    R2 (0) —<b>→[1]    R3 (0) —Symbol→[1]    R4 (0) —HtmlTag→[1]

Figure 8: Initial candidate-rules generated in the first DO...WHILE iteration.

R1 is a perfect disjunct as a result for first iteration.

...net-based IE technolgies

A perfect rule which matches examples

During the second iteration with E1 and E3 example, the initial candidate rules R5 and R6

Refinement:tries to obtain better disjuncts either by marking its landmarks more specific (Landmark refinement) or by adding new states in the automaton (Topology refinements)

Initial Candidates:

R5  (0) —<b>→ [1]

R6  (0) —HtmlTag

Topology Refinements

R7  (0) —Phone→ (1) —<b>→ [2]

R8  (0) —:→ (1) —<b>→ [2]

R9  (0) —</b>→ (1) —<b>→ [2]

...

R16 (0) —HtmlTag→ (1) —<b>→ [2]

Landmark−Refinements

R17 (0) —- <b>→ [1]

R18 (0) —Symbol <b>→ [1]

Figure 7: Rule induction (second iteration).

lecture of Internet-based technolgies

# Seed examples →

## Identifying highly informative examples

- The most informative examples illustrate exceptional cases
- Active learning :analyzes the set of unlabeled example to automatically select examples for the user to label
- forward and backward rules:

Fwd R1=SkipTo(Address)SkipTo($<I>$)

Bwd R1=BackTo(Phone) BackTo(_Number_)

**If two rules disagree on the sample, which is selected for user to label –highly informative training example.**

# Results reported from STALKER

- From 28 sources, 206 extraction rules: 182 rules (100% correct),18 rules (>90%),3% rules are<90%.
- Active learning:

Average accuracy from 85.7% → 94.2%

# STALKER features

- the ability to wrap <span style="color:red">a larger variety</span> of sources.

- capable of learning most of the extraction rules based on just a couple of examples.

- Using single-slot rules, keep high accuracy.

- improving the efficiency based on active learning for hardest items.

# Other Wrappers

- WIEN: learns the landmarks by searching *common prefixes* at the *character level*, needs more training data.

- SoftMealy: its extraction rules are less expressive than STALKER, complex to deal with missing items and various orderings of items

# Test page

## Quote Server: Tabular style document

Ticker Symbols : *(Up to 5 tickers may be entered separated by spaces)*

| | Get Quotes | ✔No Fractions |

| TICK | LAST | CHG | % | VOL #TRDS | BID ASK | LOW HIGH | PREV OPEN | 52LOW HIGH | EPS DIV | DATE TIME |
|------|------|-----|---|-----------|---------|----------|-----------|------------|---------|-----------|
| DJ   | 47.2500 | +1.0000 | +2.16 | 140,800 181 | na na | 46.0000 47.3125 | 46.2500 46.0000 | 41.5625 59.0000 | 0.09 0.96 | 02/22 16:02 |
| DJM  | 11.1875 | +0.1250 | +1.13 | 28,600 12 | na na | 11.1250 11.3125 | 11.0625 11.1250 | 9.5000 11.5000 | 0.01 na | 02/22 15:19 |
| DJT  | 4.3125 | −0.0625 | −1.42 | 167,500 112 | na na | 4.1250 4.4375 | 4.3750 4.3125 | 2.7500 10.8750 | −1.79 na | 02/22 15:57 |
| DK   | 6.0000 | +0.1250 | +2.13 | 21,500 25 | na na | 5.8125 6.0625 | 5.8750 5.8125 | 5.0625 16.7500 | −3.03 na | 02/22 15:54 |
| DL   | 29.5000 | +0.1875 | +0.64 | 382,100 206 | na na | 29.2500 29.8125 | 29.3125 29.5000 | 19.5000 32.0000 | 1.02 0.32 | 02/22 16:01 |

**Market Watch.** A detailed look at Market activity.

APL −Ticker Search −Advertise on QS −Internet Stock Report −Questionaire

lecture of Internet-based IE technolgies

# Test Pages

## Internet Address Finder: Tagged-list style document

1. Name: 'Lithium' J Smith
   E–Mail: aulmer@u.washington.edu
   Last Update: 08/01/95

   Organization: **University of Washington**

2. Name: 'Sir Brand' Gregrobin Smith
   Alt. Name: Smith Gregrobin
   E–Mail: sirbrand@u.washington.edu
   Organization: university of washington
   Last Update: 06/21/96

   Organization: **University of Washington**

3. Name: (raig Smith
   E–Mail: chs@maxwell.cs.uoregon.edu
   Last Update: 08/01/94

   Organization: **University of Oregon**

4. Name: – Richard Smith
   Alt. Name: Richard
   E–Mail: GBORDERS@SFASU.EDU
   Last Update: 11/12/95

   Organization: **Stephen F. Austin State University**

5. Name: – David S Smith
   Alt. Name: David S
   E–Mail: dssmith@INDIANA.EDU
   Last Update: 11/16/95

   Service Provider: **Indiana University**

# Result Comparison

◈ Quote Server

- **Stalker: 10 example tuples, 79%, 500 test**
- WIEN: the collection beyond learn's capability
- SoftMealy: multi-pass **85%**, **single-pass 97%**

◈ Internet Address Finder

- **Stalker: 85% ~ 100%, 500 test**
- WIEN: the collection beyond learn's capablity
- SoftMealy: multi-pass **68%**, single-pass **41%**,

# Result Comparison (cont.)

- ◆ Okra  (tabular pages)
    - **Stalker: 97%, 1 example tuple**
    - WIEN: **100%** , 13 example tuples, 30 test
    - SoftMealy: single-pass **100%**, 1 example tuple, 30 test
- ◆ Big-book (tagged-list pages)
    - **Stalker: 97%, 8 example tuples**
    - WIEN: **perfect**, 18 example tuples, 30 test
    - SoftMealy: single-pass **97%**, 4 examples, 30 test
        multi-pass **100%**, 6 examples, 30 test

# A General View of Wrapper (as Summarization)

- Machine learning method for Wrapper Induction

DeLa, RoadRunner,…

Iepad, Olera, …



Stalker, SoftMealy,Wien,Whisk

# Overall Comparison

Three dimensions: the difficulty of IE task, the techniques used, the effort made by the user for the training process and necessity to port IE across different domains.



Conclusion:
- Template-based pages have high automation degree.
- IE cross-site pages and free texts, semantic features are required.
- Manual IE systems can be applied to all kinds of inputs
- Semi-supervised and unsupervised IE systems can be applied only to template-based pages
- Unsupervised systems usually apply superficial features.

lecture of Internet-based IE technolgies

# Problem?

- The Web are very dynamic: contents, page structures
- Original wrappers can stop working: rely on Web page structures
- Re-generating wrappers is not easy: heavy workload to system developers

Changed Documents → Original Wrapper → **Extract nothing …**

→ Original Wrapper → **Incomplete results**

......... .........

→ Original Wrapper → **Incorrect results**

lecture of Internet-based IE technolgies

# Example

May Morning (1972) directed by
Featuring : Jane Birkin; John
- DVD – $ 15.38–23.26
- VHS – $ 14.98–18.99

**The original wrapper fails**

May Morning (1972)
Directed by: Ugo Liberatore
Featuring: Jane Birkin, John Steiner, Ros...

DVD from $8.99
VHS from $9.19

- Monitoring a set of generic features

- Machine learning techniques to learn a set of patterns that describe the information that is being extracted from each of the relevant fields.

- …

# How to solve it? (discussion)

# Wrapper Mainte...

- DataProg algorithm, which ... information (**patterns**) about ... field from a set of examples of the field →

Street address: 12 Pico St.,512 Oak Blvd, 416 Main st. and 97 Adams Blvd.→ (_Number_ _ capitalized_) (Blvd.) or (St.)

- **wrapper verification**: Is a wra... correctly?

detecting when a wrapper stops extracting data correctly from a Web source?

- **Wrapper maintenance**: how to a... a wrapper when the pages have changed?

identify new examples of the data field in order to rebuild the wrapper if it stops working.

# Example for Wrapper Maintenance Strategy



an example of the original site, the extracting rule for a book title and the extracted results from the example page.

The source and incorrectly extracted result after the titles's font and color were changed.

Rule changed.

# Wrapper Maintenance Methods (Kushmerick's method)

- Each data field was described by a collection of global features, such as word count, average word length, and density of types.

- Calculated the mean and variance of each feature's distribution over the training examples.

- Individual feature probabilities are then combined to produce a value.

- If the value exceeds a threshold, the wrapper is correct, otherwise, it is failed.

# A prototype for tracking changes to webpages – *Microsoft Research*

Diff-IE is a prototype Internet Explorer add-on that:

- Highlights the changes to a webpage since the last time you visited it.

- Enables you to view and compare previously cached version of a page.

→**Tracking changes to webpages.**

lecture of Internet-based IE technolgies

# Download DIFF-IE

From:  Microsoft research

http://research.microsoft.com/en-us/projects/diffie/default.aspx

- How it was implemented?

- Cache: stores the previous versions of the page, in order to highlight how a page has changed.
- Comparison component: is responsible for detecting and highlighting the changes.
- Toolbar component: is the portion of the application with which the user interacts.

# Comparison Component (1)

## Web page representation

- DiffIE identifies changes to text-based Web content at the Document Object Model (DOM) level.  Pages are represented internally as a tree of hash values to support this DOM-level comparison of text across pages.

- The text nodes of a Web page:  the leaves

- The content of these nodes are hash

  algorithm.

MD5:
A message of arbitrary length → 128-bit fingerprint

Two messages will not
Produce the same
Fingerprint.

# Comparison Component (2)

**Detecting Differences:**

- Starting at the root node, DiffIE compares the pre-computed subtree hash of the live version and the cached version.

- If same, DiffIE terminates comparison of the corresponding subtree, since identical hashes implies the content must not have changed.

# Comparison Component (3)

- 4 Types of Differences: only addition and changes are highlighted.



Figure 3. An illustration of the types of changes that can occur at the DOM level of a Web page.

# Application 1

- Monitoring a page for change, to keep track of the latest stock prices, or latest updates on the page.



lecture of Internet-based IE technolgies

# Application 2

- See new or different search results.

# Application 3

## Find changes in long lists of text

It can be hard to see changes in long lists of text, but Diff-IE identifies these automatically.



lecture of Internet-based IE technolgies

# Application 4



## Track price changes

We rarely remember prices, but Diff-IE does. Here, the prices of these HP workstations dropped.

### » Personal workstations

#### » Affordable power

Enhanced performance in an affordable package.

» **HP Z400 Workstation** NEW!
**Starting at: $ 929.00*
As low as $27/mo.**

- Up to **16 GB** of system memory
- Up to **4.5 TB** of internal storage
- Up to NVIDIA Quadro FX4800 or **dual** NVIDIA Quadro FX1800 graphics

» **HP xw4600 Workstation**
**Starting at: $ 679.00*
As low as $20/mo.**

#### » Compact power

Eight core performance in a compact footprint.

» **HP Z600 Workstation** NEW!
**Starting at: $ 1,589.00*
As low as $46/mo.**

- Up to **24 GB** of system memory
- Up to **4.5 TB** of internal storage
- Up to **eight 2D** displays
- Up to **dual** NVIDIA Quadro FX1800 graphics

» **HP xw6600 Workstation**
**Starting at: $ 1,219.00*
As low as $35/mo.**

#### » Extreme power

Ultimate performance with extreme expandability.

» **HP Z800 Workstation** NEW!
**Starting at: $ 1,839.00*
As low as $53/mo.**

- Up to **192 GB** of system memory
- Up to **7.5 TB** of internal storage
- Up to **dual** Quadro FX5800 graphics

» **HP xw8600 Workstation**
**Starting at: $ 1,339.00*
As low as $39/mo.**

» **HP xw9400 Workstation**
**Starting at: $ 2,599.00*
As low as $74/mo.**

lecture of Internet-based IE technolgies

# References

- Kushmerick, N. (2000) [Wrapper induction: Efficiency and expressiveness](). Artificial Intelligence J. 118(1-2):15-68 (special issue on Intelligent Internet Systems).

- Chun-Nan Hsu and Ming-Tzung Dung. [Generating finite-state transducers for semistructured data extraction from the web.]() *Information Systems*, 23(8):521-538, Special Issue on Semistructured Data, 1998.

- **Ion Muslea**, Steve Minton, Craig Knoblock. [Hierarchical Wrapper Induction for Semistructured Information Sources,]() *Journal of Autonomous Agents and Multi-Agent Systems*, *4:93-114, 2001* .

# References sites

- Repository of online information sources used in information extraction task: http://www.isi.edu/info-agents/RISE/index.html

- Chia-Hui Chang, et al, " A survey of Web Information Extraction Systems" in IEEE Transactions on Knowledge and Data Engineering.

- Papers, tutorials, lectures, code
  - http://www.cs.cmu.edu/~wcohen/10-707

# Summarization

- What is wrapper?

- How to do wrapper induction?

- How to maintenance wrapper?