



-
-
-
-
-
-
-

Crawling the Hidden Web


by S. Raghavan & H.G. Molina



Presenter: Jack Ng
Oct 30, 2002



Background Info

- 
- **Hidden Web** - databases whose content is accessible only through search forms
 - Why is it important to tap into the hidden Web?



Background Info

- According to "*The Deep Web: Surfacing Hidden Value*", 2001:

- ƒ 500 billion documents; 500 times > PIW
- ƒ 7500 TB of data; 19 TB for PIW
- ƒ grows much faster than the PIW
- ƒ High quality, topic-specific information
- ƒ 95 % is publicly accessible - no fees or subscriptions

3



Background Info

- Challenges faced by crawlers to extract content from the hidden Web:
 - ƒ size of hidden Web is enormous!
 - ƒ content not reachable by following hypertext links
 - ƒ "form-filling" is a human activity

“Training” a crawler is very difficult!!

4



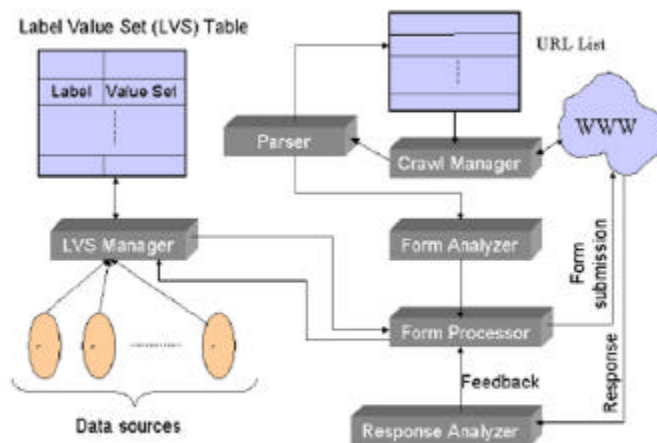
Background Info

- Authors' approach to address the challenges:
 - f* task-specificity
 - f* human-assistance
- Propose:
 - f* model of hidden Web crawler
 - f* model of form page
 - f* LITE (Layout-based Information Extraction Technique) for content extraction
- Implementation - HiWE (Hidden Web Exposer)

5



HiWE Architecture



6



HiWE Data Structures - LVS Table

- Task-specific DB
- Organized by concepts
- Vocabularies for filling out forms

Fuzzy set: membership function assigns 'confidence' to each value

Task: search for game reviews

Platform	{Xbox, PS2, GameCube, PC}
Genre	{action, RPG, strategy, sports}
Developer	{EA, Sega, Squaresoft, Bioware}
Release date	{1999, 2000, 2001, 2002}

7



HiWE - Form Processing Strategy

- Given internal form representation:

$$F = (\{E_1, E_2, \dots, E_n\}, S, M)$$

- For each infinite domain element, label matching algorithm finds closest match in LVS table and assigns value set to it
- Rank value assignments to ensure quality submission
 - f Fuzzy conjunction --> conservative
 - f Average
 - f Probabilistic --> aggressive
- Submit only if rank is greater than threshold

8



What is LITE?

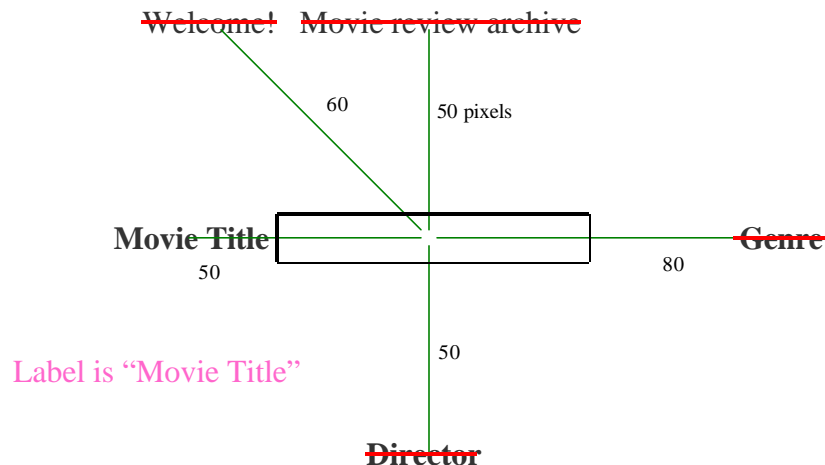
- Label extraction *heuristics* based on how page is laid out for *human* viewing
- *Idea*: label is often visually adjacent to widget (e.g., textbox) and obvious to viewer
- Partial layout is sufficient to determine adjacency --> prune unnecessary elements (see Figure 4 in paper)
- Applications in HiWE:
 - f* form page analysis
 - f* response page analysis

9



LITE Application - Form Page Analysis

How LITE heuristics identify label of form element



10



LITE Application - Response Page Analysis

- Based on idea that results must be obvious to viewers
- Prune page to find visually *center-most* portion & interpret it as results location
- To identify error pages:
 - f*search center portion for common error text (e.g., "No results")
 - f*compute hash value for center portion
 - common hash values = error pages

11



Experimental Results

- Value assignment ranking
 - f**fuzzy conj.* --> best submission efficiency
 - f**average* ✓ --> most successful submissions
 - f**probabilistic*--> poor performance
- LITE outperforms other label extraction techniques; overall **93 %** accuracy

12



Thoughts...

- Strength & novelty of solution
 - +flexible framework
 - +works with non-cooperative DBs
 - +crawler has learning capability
 - +crawls both PIW and hidden Web
 - +‘mines’ visual layout info for semantics

13



Thoughts...

- *Implementation* limitations
 - LVS table - how to handle semantically ambiguous labels?
 - what about image labels?
 - doesn't consider relationships among elements when assigning values
 - 'all-or-none' form submission policy

14



Thoughts...

- Presentation of paper
 - easy to follow and understand
 - right level of details
 - goals & pre-conditions clearly defined

- overall, a good paper!