# Pdist-RIA Crawler: A P2P architecture to crawl RIAs

SOFTWARE SECURITY
RESEARCH GROUP

In Collaboration With IBM

**Seyed M. Mirtaheri, Gregor v. Bochmann, Guy-Vincent Jourdan, Iosif Viorel Onut**
**School of Information Technology and Engineering  - University of Ottawa**

## Introduction

Rich Internet Applications (RIAs) allow better user interaction and responsiveness than traditional web applications.

Thanks to new technologies like AJAX (Asynchronous JavaScript and XML), Rich Internet Applications can communicate with the server asynchronously. This allows continuous user interactions.
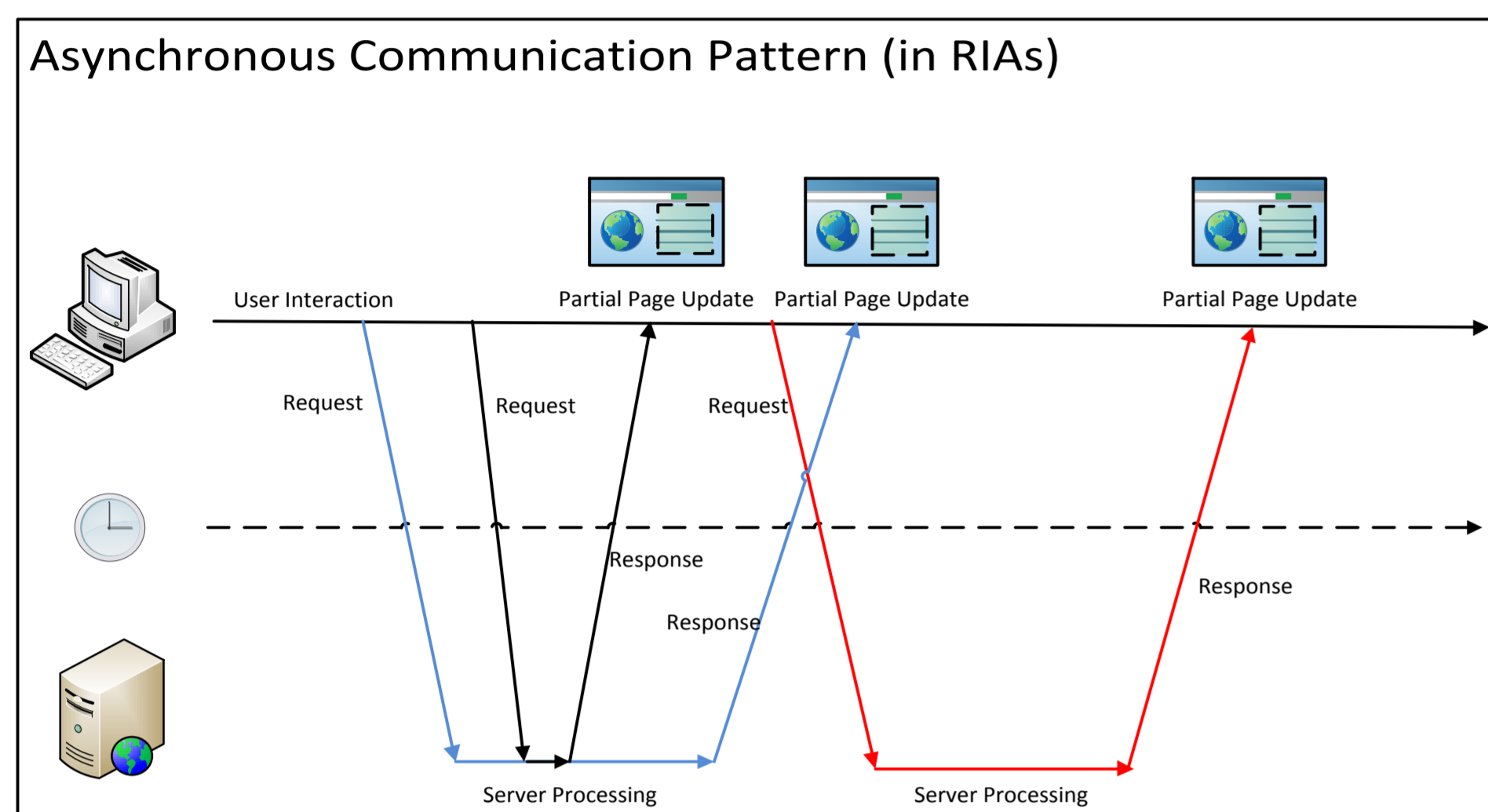


Figure 1: AJAX enabled RIA communication pattern.

Security of RIA and automating security testing are important, ongoing, and growing concerns. One important aspect of this automation is the crawling of RIAs i.e. reaching all possible states of the application from the initial state. Being able to do so automatically is also valuable for search engines and accessibility assessment.

Crawling of RIA applications is an expensive and time consuming process due to their large number of states. To accelerate this operation we distribute the operation over many nodes in an elastic cloud environment.

## Background

* **Crawling Strategy**:

  ➤ Breath-First Search
  ➤ Bounded Depth-First Search
  ➤ Based on page weight
  ➤ Model Based Crawling
    ➤Hypercube Model
    ➤Menu Model
  ➤ Greedy algorithm
  ➤ Probabilistic model

* **Partitioning strategies**: Mostly use server related matrix as primary tool to partition search space:
  ➤ Page URL
  ➤ Server IP address
  ➤ Server geographical location
  ➤ [Loo 2004] describes distributed web crawling by hashing the URL
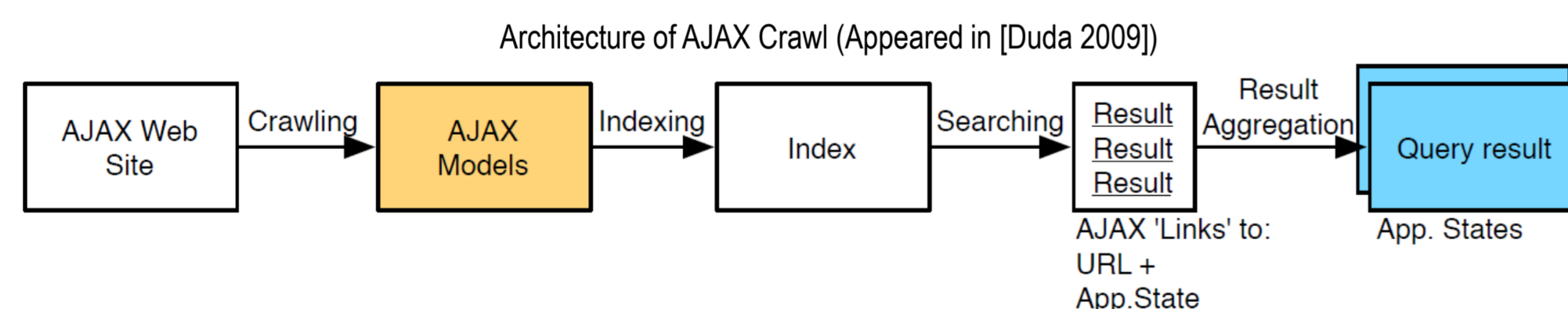


Figure 2: [Loo 2004] System Architecture

References:
* [Amalfitano 2010] D. Amalfitano, A.R. Fasoline, P. Tramontana, Techniques and tools for Rich Internet Application testing, in Proc. WSE, 2010, pp.63-72.
* [Mesbah 2008]  A. Mesbah and A. van Deursen, A Component- and Push-based Architectural Style for Ajax Applications. 2008, Journal of Systems and Software (JSS) 81(12):2194-2209
* [Duda 2009] C. Duda, G. Frey, D. Kossmann, R. Matter, C. Zhou, AJAX Crawl: Making AJAX Applications Searchable, Proc. IEEE Intern. Conf. on Data Engineering, Shanghai, 2009, pp. 78-
* [Loo 2004] B.T. Loo, O. Cooper, S.Krishnamurthy, Distributed Web Crawling over DHTs, Technical report, University of California, Berkeley, 2004, http://www.eecs.berkeley.edu/Pubs/TechRpts/2004/5370.html
* [Exposto 2008] J. Exposto, J. Macedo, A. Pina, A. Alves and J. Rufino, Efficient partitioning strategies for distributed web crawling, in Information Networking: Towards Ubiquitous Networking and Services, Springer LNCS 5200 (2008), pp. 544-553.
* [Boldi 2003] P. Boldi, B. Codenotti, M. Santini, S.Vigna, UbiCrawler: A Scalable Fully Distributed Web Crawler, Software: Practice & Experience, Vol. 34, 2003, p. 2004.
* [Dincturk 2014] Dincturk, Mustafa Emre and Jourdan, Guy-Vincent and Bochmann, Gregor von and Onut, Iosif Viorel: Model Based Crawling, ACM Transactions on the WEB
* [Dincturk 2012] Dincturk, Mustafa Emre  and Choudhary, Suryakant and Bochmann, Gregor von and Jourdan, Guy-Vincent and Onut, Iosif Viorel: A Statistical Approach for Efficient Crawling of Rich Internet Applications, Proceedings of the 12th international conference on Web engineering

## Motivation and Aim

* **Non-URL-Based Crawling strategy**:
  ➤ In a RIA one URL corresponds to many states of DOM. Unlike traditional websites in which every call to server would change the whole DOM and the page URL, RIA relies on small AJAX updates that does not necessarily modify the page URL:

  ➤ Traditional distributed crawlers rely heavily on URL in order to partition the search space. Underlying assumption for this strategy is a one to one correspondence between the URL and the state of DOM which does not hold in RIA.

  ➤ Therefore we propose to partition the search space based on events.

* **Crawling Strategy:**  Reduce the workload by choosing the events to execute using Greedy algorithm.
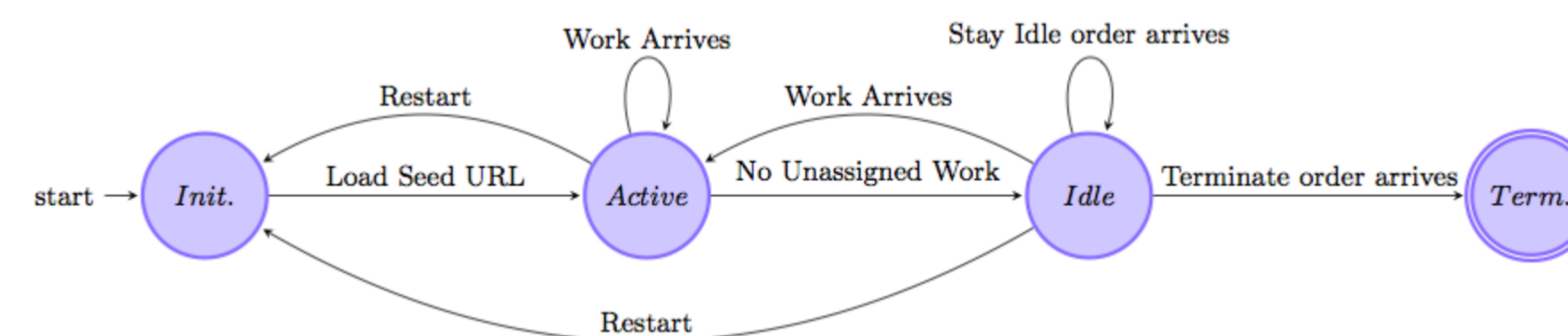
* **Crawling Efficiency**: Discover states as soon as possible, using Probabilistic model.

Reference:
* [Benjamin 2010] K. Benjamin, G. v. Bochmann, G.-V. Jourdan and V. Onut, Some modeling challenges when testing Rich Internet Applications for security, First Intern. Workshop on Modeling and Detection of Vulnerabilities (MDV 2010), Paris, France, April 2010. 8 pages.

## Proposed Architecture

* Nodes act autonomously and independently. Each node starts at Init state, when get tasks go to Active state, when has nothing to do goes to idle state, and finally terminates when termination order arrives.



* A virtual ring is created based on breath-first-search traversal of the nodes. A termination token goes around this wring that keeps the list of states IDs and the number of states visited by each node. When all states are visited on all nodes, a termination order is broadcasted.
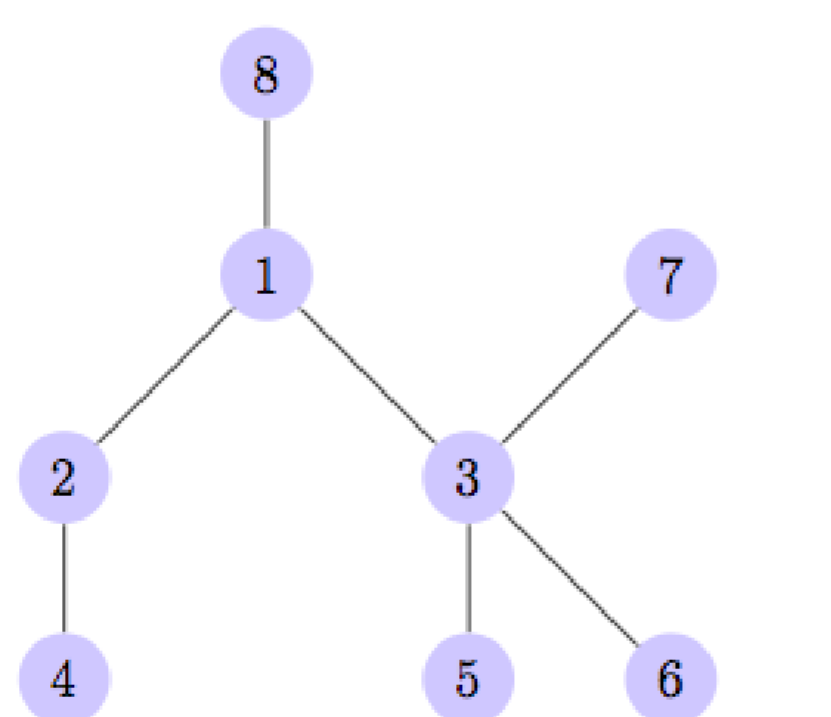


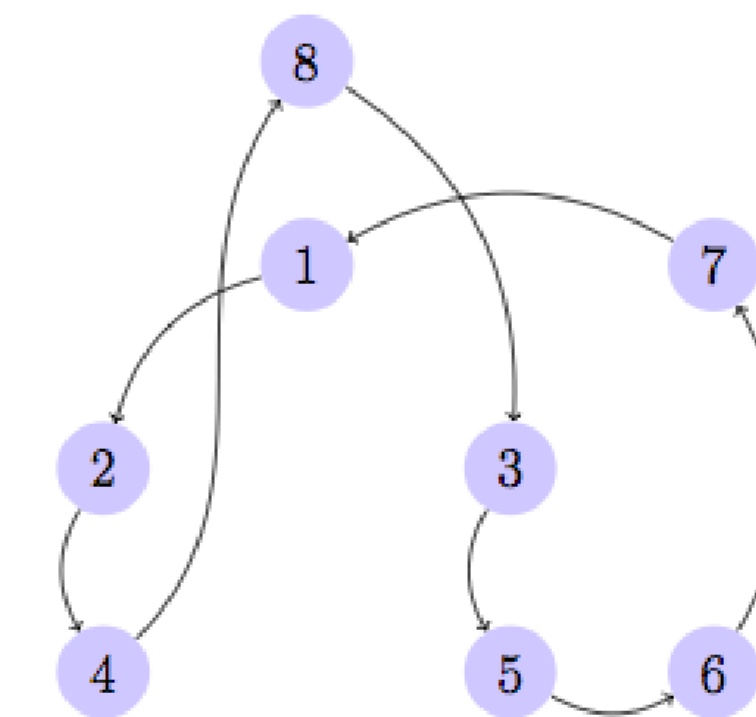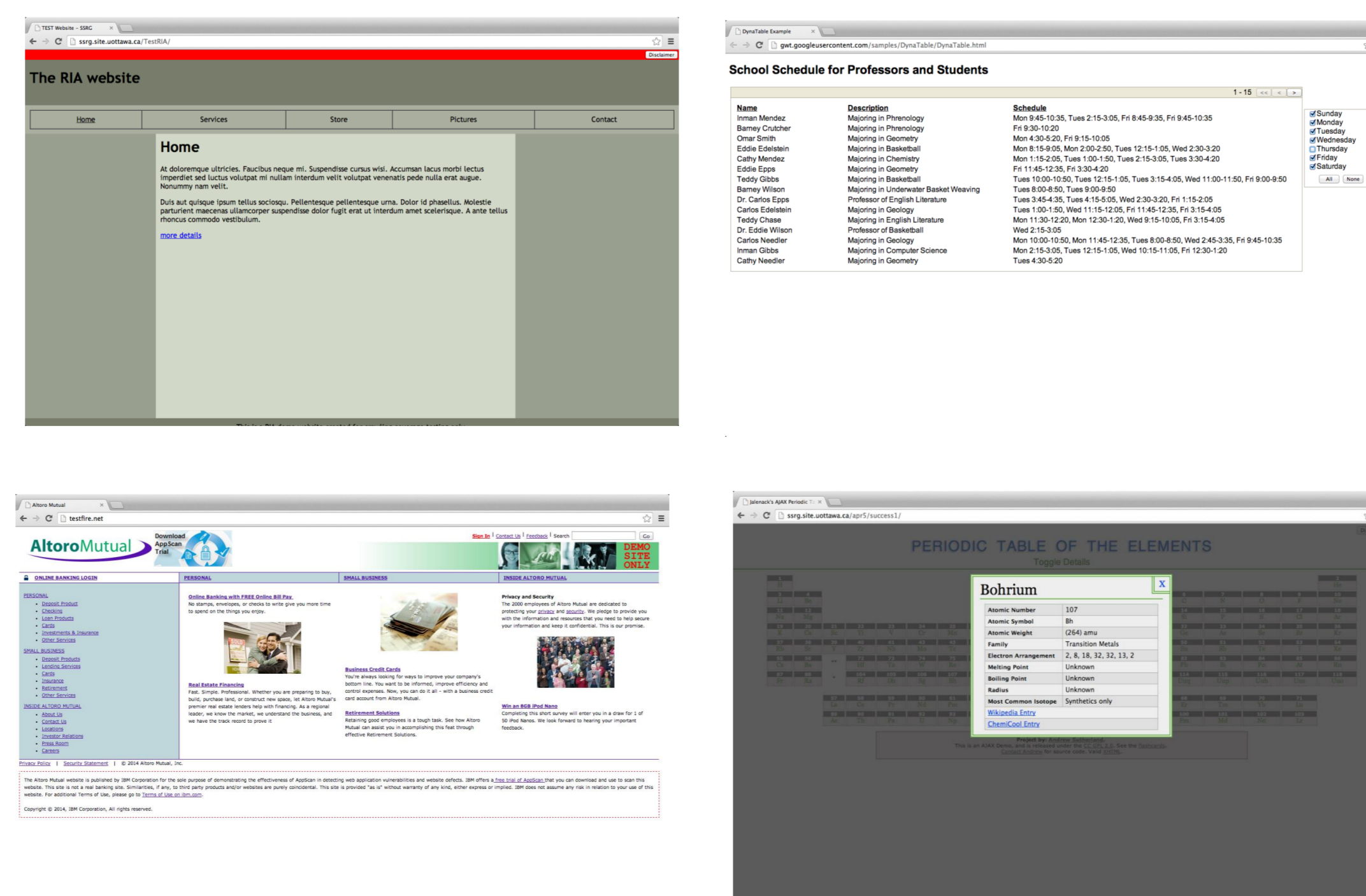Figure 5.1: Sample virtual spanning tree    Figure 5.2: The termination virtual ring

## Test Beds

* Four target web applications used to measure the performance of the distributed web crawler:
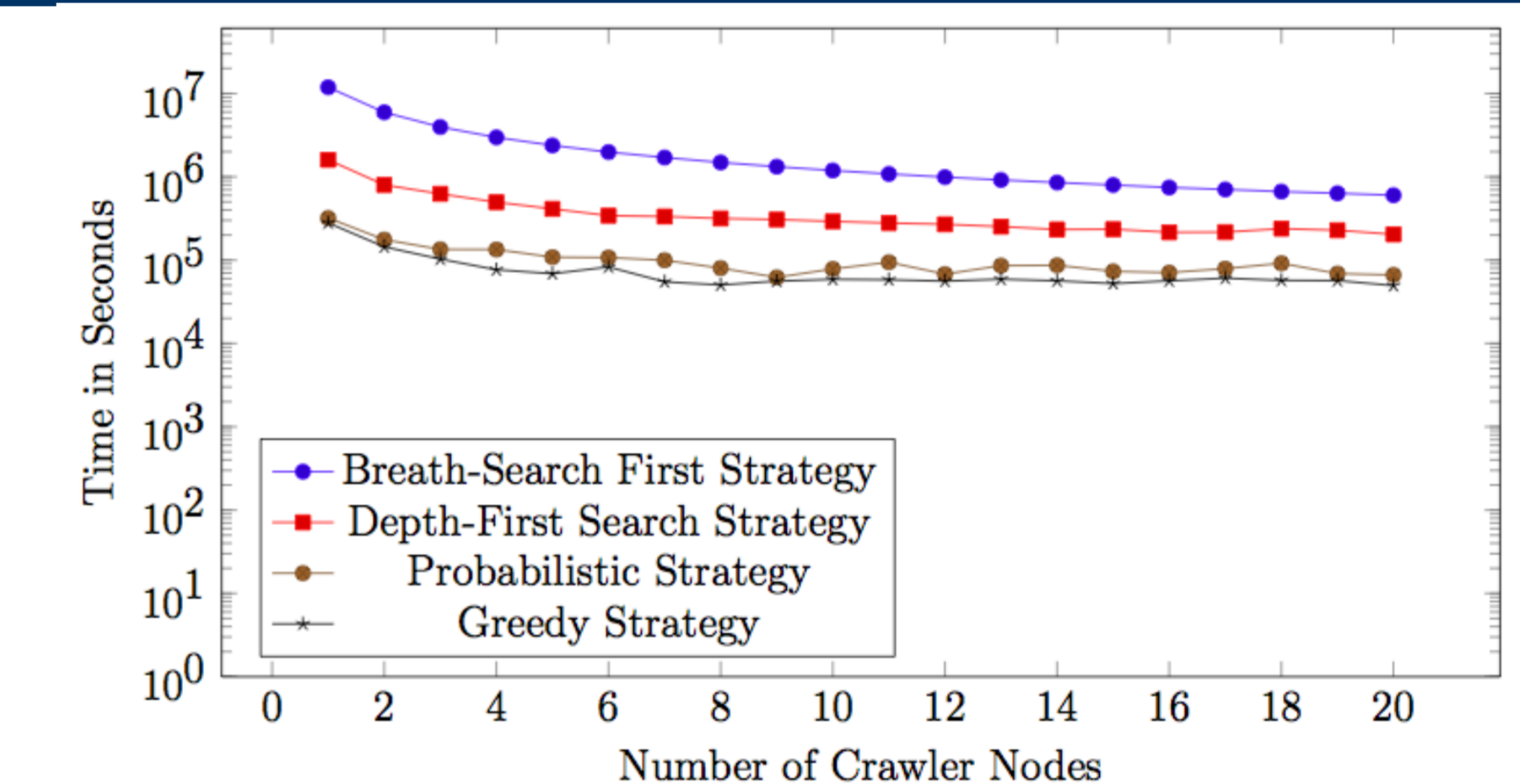


## Algorithm



## Experimental Results



## Crawling Efficiency

* Crawling efficiency measures how early in the crawl new states are discovered:
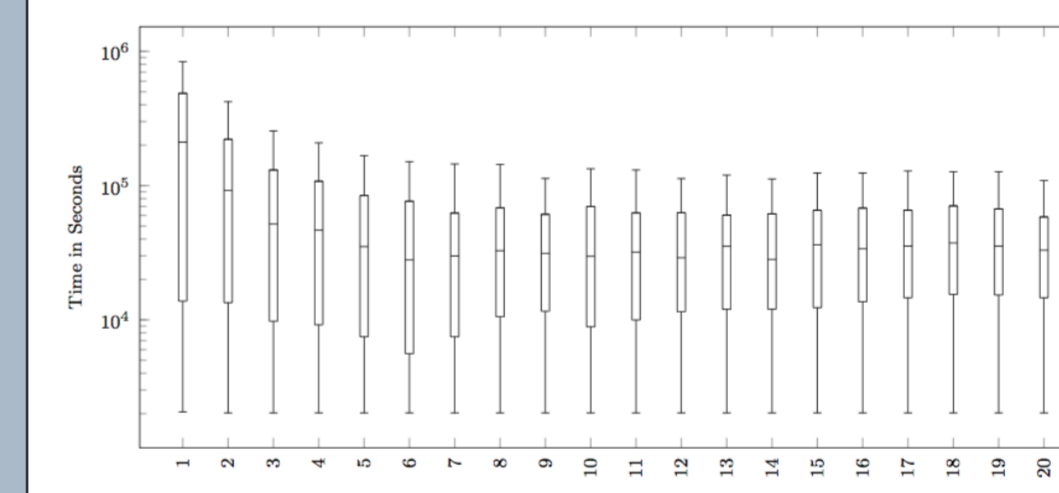


Figure 8.12: Dyna-Table: Cost of Discovering States using Depth-First Search Strategy
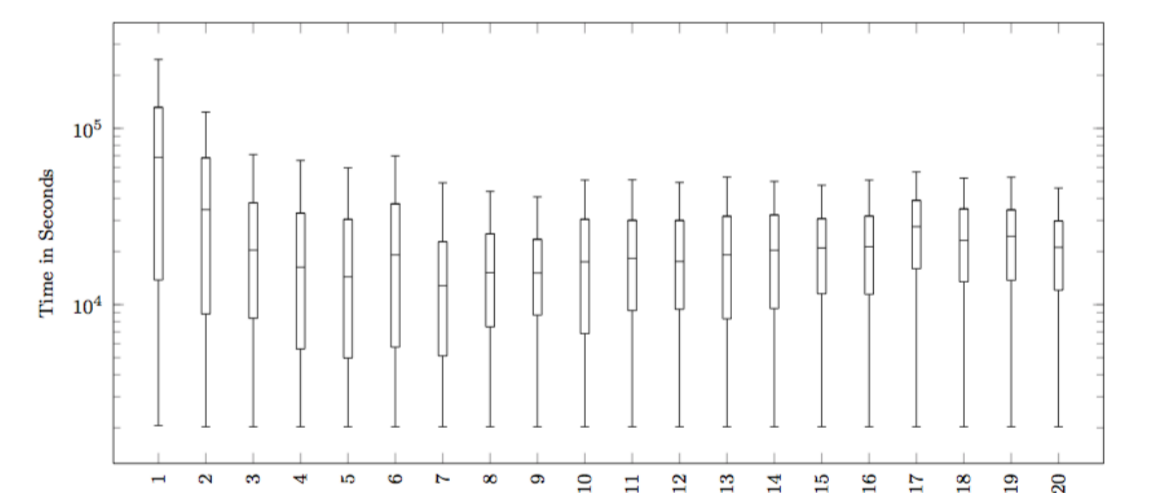
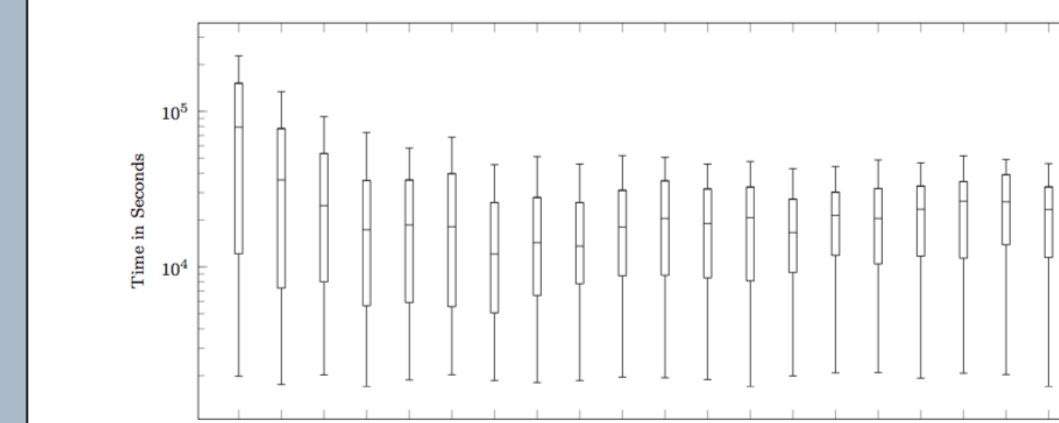Figure 8.13: Dyna-Table: Cost of Discovering States using Greedy Strategy

Figure 8.14: Dyna-Table: Cost of Discovering States using Probabilistic Strategy

## Conclusion

* Distributed Greedy algorithm has the best performance in terms to total time it takes to crawl a website.
* Distributed Probabilistic model is the most efficient algorithm and discovers states early in the crawl.

## Future Work

* We are currently working on distributed crawling of RIA in a cloud environment.
* We plan to add fault tolerance to our strategy so that if some of the nodes crash rest of the nodes continue without interruption.
* Once we have a working implementation of the system we plan to optimize it based on different infrastructure parameters such as cost of communication or the processing power available to different nodes.

## Acknowledgments