# Client-Side Preservation Techniques for ORE Aggregations

Michael L. Nelson & Sudhir Koneru

Old Dominion University, Norfolk VA

OAI-ORE Specification Roll-Out

Open Repositories 2008

Southampton UK, April 4, 2008

# Outline

- Background: Let the "Web Infrastructure" preserve your information
- Premise: ReMs are critical for preservation purposes
- Client-side vs. Server-side approaches to preservation
- Sketch of a possible framework for client-side preservation techniques

# Web Infrastructure

**preservation = refreshing + migration**

# Web Repository Contributions

Frank McCown, Joan A. Smith, Michael L. Nelson, Johan Bollen, Lazy Preservation: Reconstructing Websites by Crawling the Crawlers, Proceedings of WIDM 2006,pp. 67-74.  http://www.cs.odu.edu/~mln/pubs/widm-2006/lazyp-widm06.pdf

# Overlap with Internet Archive

| SE | In IA | | Not in IA | |
| --- | --- | --- | --- | --- |
| | Cached (II) | No cache (I) | Cached (III) | No cache (IV) |
| Ask | 9.2% | 36.0% | 0.3% | 54.5% |
| Google | 40.7% | 3.7% | 50.3% | 5.3% |
| MSN | 51.1% | 1.1% | 43.7% | 4.1% |
| Yahoo | 39.3% | 1.8% | 47.7% | 11.2% |



Frank McCown, Michael L. Nelson, Characterization of Search Engine Caches, Proceedings of IS&T Archiving 2007, pp. 48-52. http://arxiv.org/abs/cs.DL/0703083

# Warrick -- A Service to Recover Lost Websites
## warrick.cs.odu.edu

# How Much Did We Reconstruct?

"Lost" web site

Reconstructed web site

Four categories of
recovered resources:

1) Identical: A, E
2) Changed: B, C
3) Missing: D, F
4) Added: G

Missing link to D;
points to old
resource G

F can't
be found

# Resource Maps Unambiguously Define an Aggregation

- The "manifest" nature of ReMs allow us to know "if we got it all"
  - "known knowns"
  - "known unknowns"
  - "unknown unknowns"

- Assuming the ReM is recovered, the implications for preservation are clear:
  - defines members of the aggregations
  - defines relationships between them

# Server-Side Techniques

- Repository A uses ReMs for their aggregations.
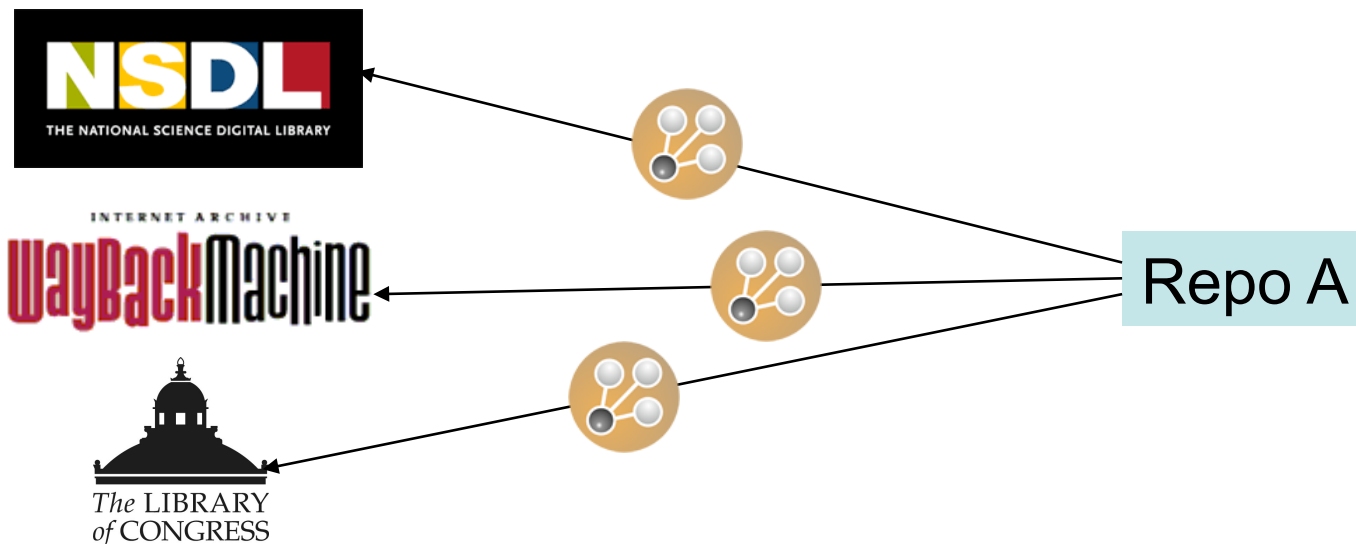- Repository B harvests ReMs to ensure total coverage of Repository A.
- Repository A can use its ReMs to validate transfer to Repository B.
- Third parties use ReMs to audit B's preservation of A.
- New ReMs created to reflect migration, refreshing of aggregations.

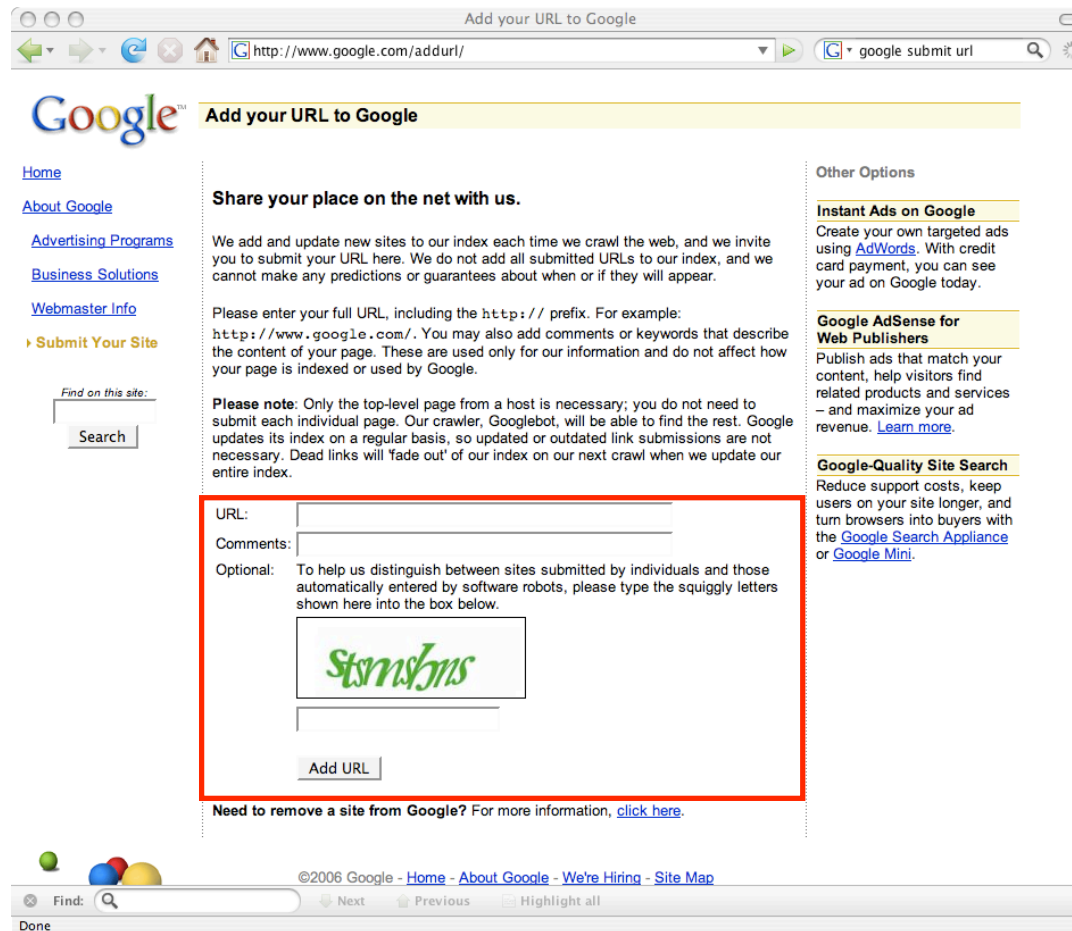# Can We Involve End-Users in the Preservation Process?

- Leverage the actions of end users?
  - "people helping robots…"

- Make preservation more accessible?
  - light-weight and easy like Google Analytics and reCAPTCHA?

```html
<html>
…
<h1>hello world</h1>
…
<script type="text/javascript">
resourcemap="http://www.foo.edu/repo/helloworld/rem.atom";
webReposToCheck="google,yahoo,internetArchive";
checkMirrors="yes";
writeBack="http://www.bar.org/wiki/"; </script>
<script type="text/javascript"
src="http://ore.cs.odu.edu/preserve.js"> </script>
…
</html>
```

# Client-Side Techniques

- Operations on the ReM and Aggregated Resources (ARs)
  - validation, http status, ReM visualization, etc.

- Interacting with the Web Infrastructure
  - checking for ReM, ARs in Internet Archive, search engine caches, etc.
  - reconstructing aggregation for a given time interval
  - submitting ReM, ARs to WI

- Inter-client communication
  - my client updates/repairs ReM -- how to communicate that to other clients and servers?

# One Reason Why We Need Humans in the Loop

# A Possible Scenario…



wiki.somewhere.org

```
<script type="text/javascript">
resourcemap="http://www.foo.edu/repo/helloworld/rem.atom";
webReposToCheck="google,yahoo,internetArchive";
checkMirrors="yes";
writeBack="http://www.bar.org/wiki/"; </script>
<script type="text/javascript"
src="http://ore.cs.odu.edu/preserve.js"> </script>
```

ore.cs.odu.edu

# Wikis Would Make a Nice Inter-Client Message Store



Function as a publicly (computers + humans) readable revision control system for ReMs

# "Help Preserve This Object"

# Current Status

- Hierarchical view of ReM
- Finds copies of Aggregated Resources in Internet Archive, Google, Yahoo
- Next up:
  - use Simile time line software (http://simile.mit.edu/timeline/) to display ARs in time
  - chose a time interval for reconstruction
  - send edited ReMs to a wiki or public email service
  - write a program to read & vet edited ReMs from public store