

Web-scale discovery indexes

Eric Lease Morgan
University of Notre Dame

August 18, 2009

Outline

- Introduction
- The problems to be solved
- What versus how; two case studies
- Indexes, not databases
- Next steps; a do-it-yourself recipe
- Plan B
- Do-it-yourself and Plan B compared
- Web-scale indexes and “next generation” library catalogs

Problems to be solved

What is the purpose of our libraries? Why do they exist? What are the problems they are trying to solve?

The answers are not definitive, especially when applied to individual libraries...







The whats of librarianship

Libraries collect, preserve, organize, and disseminate data, information, and knowledge for the purposes of making the work of their respective communities easier.

To one degree or another, just about everything us librarians do can be associated with one of these processes.

These things are the whats of librarianship, and they *change very slowly*.









The hows of librarianship

The hows of librarianship are the things of our everyday work, our day-to-day operations, the specific workflows within each of our libraries.

The hows of librarianship change at a much faster pace, and these *changes are usually driven by technology*.

Catalogs are a good example

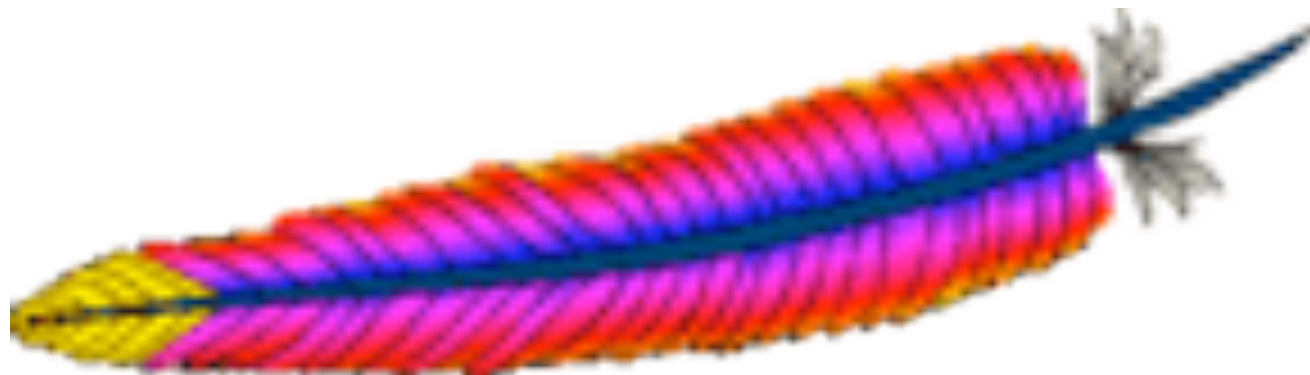


Journal indexes are another



Indexes, not databases

With the advent of freely available, industrial strength indexers – not databases – we are beginning to see another evolutionary step in the development of the library catalog and journal article index.



Lucerne



Powered by
Swish-e

“Smart” computer indexes

```
# calculate term frequency/inverse document frequency
sub tfidf {

    my $n = shift; # number of times found in document
    my $t = shift; # total number of words in document
    my $d = shift; # total number of documents
    my $h = shift; # number of hits in the corpus

    my $tfidf = 0;

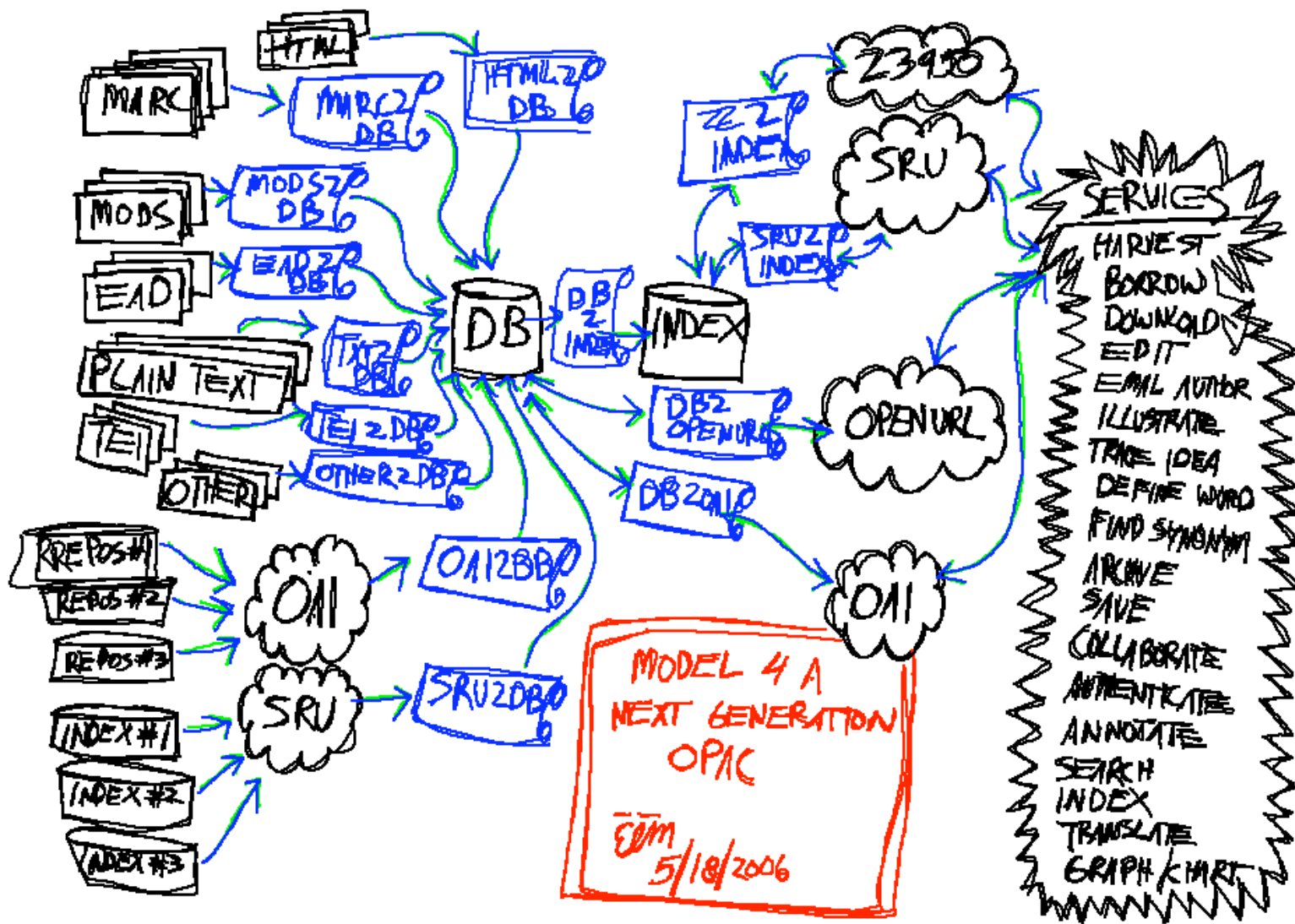
    if ( $d == $h ) { $tfidf = ( $n / $t ) }
    else { $tfidf = ( $n / $t ) * log( $d / $h ) }

    return $tfidf;

}
```

Next steps

What process can be used to take advantage of this environment? What are the next steps?
The creation and maintenance of a combined book/journal article index is an excellent example...



A recipe

1. Allocate resources
2. Charge the group
3. Wait three months
 - ☐ Dump MARC records
 - ☐ Export metadata from repositories
 - ☐ Harvest and/or mirror article and etext content and metadata
 - ☐ Feed all metadata to indexer
 - ☐ Design simple user interface to search the index
4. Ask for an update
5. Go to Step #3 four times
6. Evaluate
7. Share your experience
8. If the process was successful, then go to Step #1
9. Otherwise, consider Plan B





3 Eggs And Toast 5.99
3 Eggs And Toast 5.99
Add Bacon or Sausage 1.00
4 strips of bacon, 3 patties or 3 saus

Big Breakfast
3 Eggs, Toast, Coffee, Home Fries or Hash
One Ham Steak or Sausage And
2 buttermilk pancakes or 2 French

OMELETS (3 LG EGGS)
SERVED WITH BUTTER AND TOAST
1 Item 5.99
2 Item 6.99
3 Item 7.99
Choice of:
Sausage, Bacon, Turkey, Ham, or Sausage 50 per Extra Item

Waffle Pancakes Chocolate
3 Stack 3.75

Waffle Pancakes available

THE LITE SIDE
1/2 Pancake or 1 French To
Bacon, 2 Links or 1 Patty
INCLUDES COFFEE

KIDS, TOAST AND COFFEE
FRENCH TOAST AND COFFEE

KIDS CORNER
1 EGG, 1 SLICE TOAST
2pcz Bacon, 1 patty or 1 lb

OR
1 PANCAKE OF CHOC
1 FRENCH TOAST
Includes SMALL dr

FOR CHILDREN UN
NO SUBSTIT

Waffle Cones
1.50 CHU
15% C



© infomotions.com



Plan B

For any number of reasons the do-it-yourself approach is not feasible for you and your library...







© infomotions.com



© infomotions.com

Plan B

...If this is the case, then you might have someone else do the work for you, and they will:

1. Decide what content to include (collections)
2. Collect it (acquisitions)
3. Normalize it (cataloging)
4. Index it (systems)
5. Provide access to the index (public service)

Do-it-yourself and Plan B compared

Neither the do-it-yourself nor the Plan B approach are perfect. There are strengths, opportunities, weaknesses, and threats in both.









Web-scale indexes – the right direction



But search is not enough

Search and access are not really the problems to be solved. Instead we need to be devising ways to make content more useful and placing it into context: annotate, apply concordance, compare & contrast, create flip book, create tag cloud, graph, highlight, plot on map, rate, review, summarize, trace author and citation, translate, etc.

Thank you

<http://www.library.nd.edu/daiad/morgan/musings/web-scale/>

