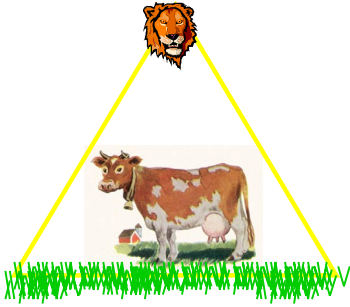


Information Carnivores



05-Feb-01 15:36

1

Outline

- MetaCrawler (briefly)
- Wrappers
- Meta-bots
 - Jango
- The Story of Netbot

05-Feb-01 15:36

2



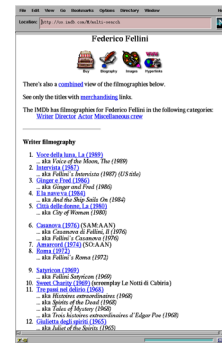
- Built by Selberg & Etzioni
- Release in Jun '95
- Currently aggregates 12 search engines:
 - LookSmart, About, Infoseek,
 - GoTo, Google, DirectHit,
 - RealNames, Webcrawler, AltaVista,
 - Excite, Lycos, Thunderstone

05-Feb-01 15:36

3

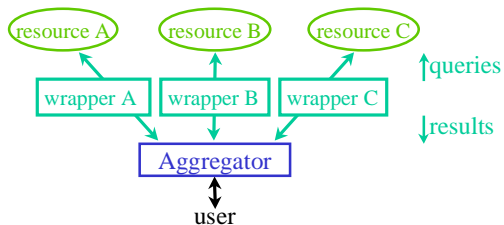
The Need for Wrappers

lots of
information
but
computers don't
understand
much of it



05-Feb-01 15:36

Strategy: Wrappers



05-Feb-01 15:36

5

Scaling issues

Need custom wrapper for *each* resource.

```
<HTML><BODY BGCOLOR="#FFFFFF" LINK=
"00009C" ALINK="00009C" VLINK="00009C"
TEXT="000000"><center><table><tr><td><NO
BR><NOBR><A HREF="/bin/cgi.qd.d?MEM
=1" TARGET="" top"></A><A HREF="/bin/cgi.dir.d?MEM=1"
TARGET="" top"></A><A HREF="" TARGET="" top"><img src
="/ypimages/b_r_hd_3.gif" border=0 ALT="Home"
width=47 height=20 align=top></A></NOBR></tr>
</td></tr></table></center><center><table border=0
width=576><tr><td colspan=2 align=center><center>
```

→ useful
information

But hand-coding is *tedious*.

Especially since sites frequently change format

05-Feb-01 15:36

6

Wrapper Approaches

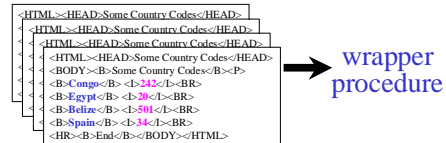
- **Perl-like languages**
 - Simple and effective (if tedious)
- **Proprietary languages & tools**
 - Click and generalize
- **Conversion to tree form**
 - Use XML as intermediate representation
 - Extract children of specified node
- **Machine Learning**
 - Promising, but not yet fielded

05-Feb-01 15:36

7

Kushmerick Contribution

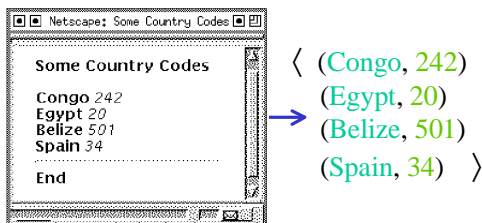
machine learning techniques
to automatically construct
wrappers from examples



05-Feb-01 15:36

8

Example

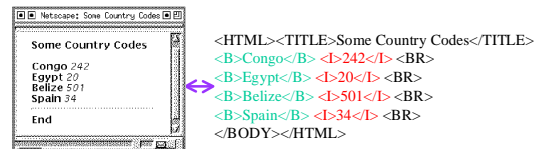


05-Feb-01 15:36

9

LR wrappers: The basic idea

exploit *fortuitous non-linguistic regularity*



Use ****, ****, **<I>**, **</I>** for parsing

05-Feb-01 15:36

10

Country/Code LR wrapper

procedure **ExtractCountryCodes**
while there are more occurrences of ****
1. extract **Country** between **** and ****
2. extract **Code** between **<I>** and **</I>**

Left-Right wrappers

05-Feb-01 15:36

11

Value of Meta-Information

- **TV guide makes more than NBC**
- **American Airlines makes more on SABRE than it does on tickets**
- **Product Reviews, Advice = \$\$\$**
 - (Likewise the DB of info adapters)
- **Optimal strategy for info access**
 - (Etzioni FOCS 96)

05-Feb-01 15:36

12

Meta-is-'beta'

- Web Search
- Shopping
- Product Reviews
- Chat Finder
- Columnists (e.g. jokes, sports,)
- Email Lookup
- Event Finder
- People Finder
- Restaurant Reviews
- Job Listings
- Classifieds
- Apartment + Real Estate

05-Feb-01 15:36

13

Personal Shopping Assistant

- You Name the Product
- Agent Visits Stores, Review Sites...
- Makes Summary Report...
- If Requested, Buys Items for You

Inevitable New Category

05-Feb-01 15:36

14

The Netbot Story: Part I

- **Bargain Finder (1995)**
 - Proof of concept; CD stores only; hardcoded
- **UW Shopbot Prototype (Jan '96)**
 - Focus on scalability to many stores, categories
- **Netbot Founded (May '96)**
 - Five Developers, 4 licenses, Seed funding

05-Feb-01 15:36

15

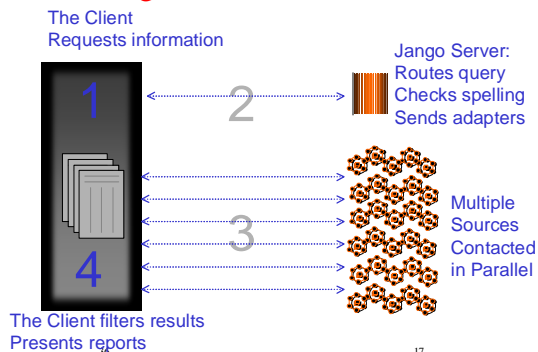
The Netbot Story: Continued

- **Hired Eric Zocher (Sept '96)**
 - Weld/Etzioni back to school
- **Second Round Financing (Jan '97)**
 - Netbot now at 25 people
- **Jango Client Launch**
 - Beta (April '97)
 - 1.0 (July '97)

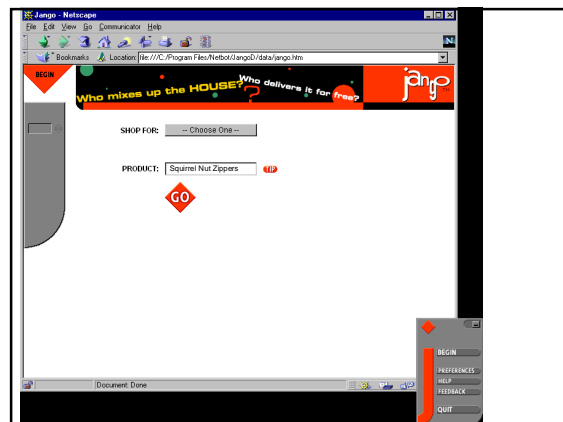
05-Feb-01 15:36

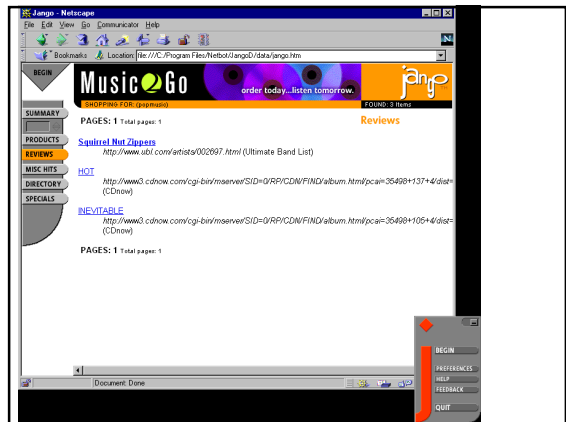
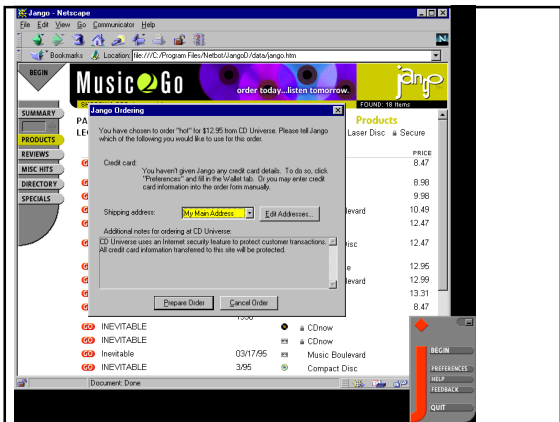
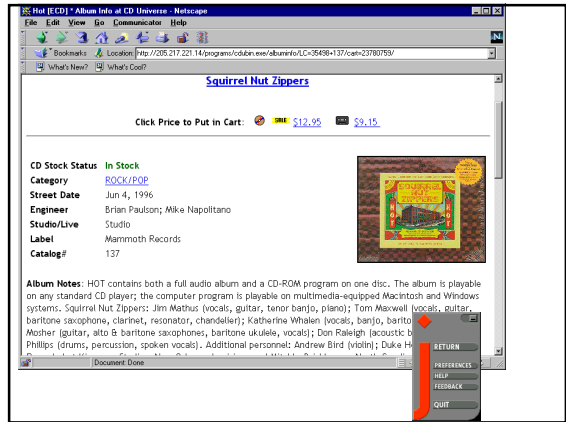
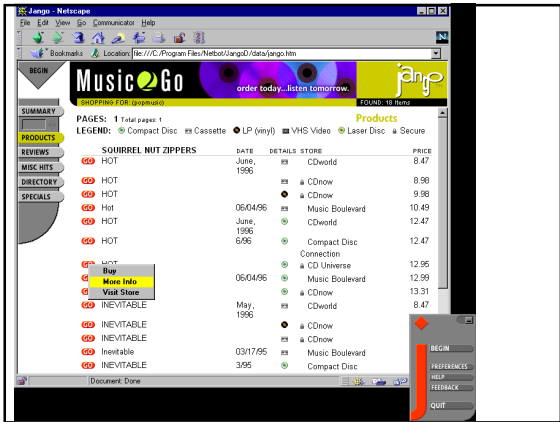
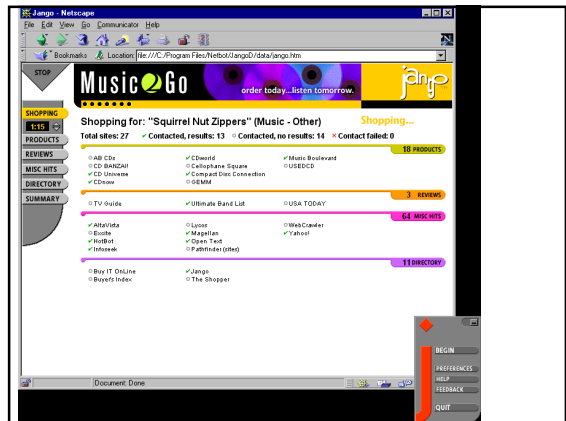
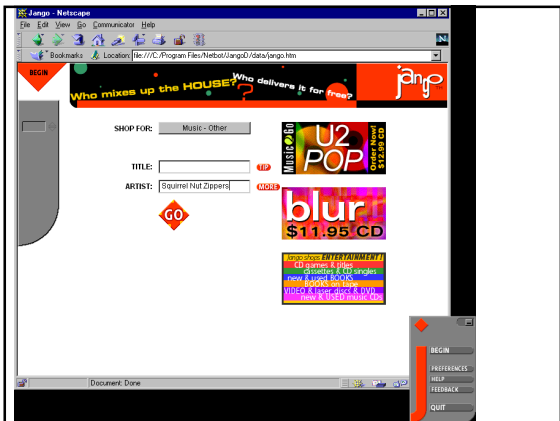
16

Jango Client Architecture



17





Core Technology I: Information Adapters

- **Low cost, flexible, “glue”**
 - Connects Jango to 100s of sites
- **Development & QA**
 - Proprietary language & compiler
 - Semiautomatic construction (learning)
 - Continuous monitoring

05-Feb-01 15:36

25

Core Technology II: Query Routing

- **Selects WWW sites given query**
- **Limited natural language abilities**
- **Spell checking**
- **Taxonomic knowledge base**

05-Feb-01 15:36

26

Core Technology III: Parallel Pull

- **Conventional Pull**
 - Typing a URL yields a single web page
- **Parallel Pull**
 - Typing high-level request yields info compiled from dozens of web pages
- **Parallel Aggregation Engine**
 - Efficient parallel access: aggregation, relevance, duplicate elimination

05-Feb-01 15:36

27

Consumer Benefits of Jango

- **Comprehensive.** Knows about 100s of stores
- **Quick.** Contacts stores simultaneously
- **Current.** Gathers up-to-the-minute details
- **Convenient.** Creates custom shopping guide
 - Product prices & availability...
 - Reviews...
 - Links to each online store...

05-Feb-01 15:36

28

Jango Benefits for Retailers

- **More Traffic.** Jango pulls from store with every query to a category.
- **Targeted Shoppers.** Jango users are looking for specific products.
- **One click buy.** Jango speeds sales process.
- **Reduced cost of customer acquisition.** Jango will catalyze more shopping
- **Shopping Reports.** What's selling

05-Feb-01 15:36

29

Jango Business Model

- **Give client away for free**
- **Get 1,000,000 users by Xmas '97**
- **Sell to merchants:**
 - Advertising (targeted users in act of buying!)
 - Real time promotions
 - Eventually aim for transaction cut

05-Feb-01 15:36

30

Reaction to Jango

- Early adopters loved it
- Awards
- Substantial “buzz” in the press
- Two imitators (Sept '97)
- Contract with ATT Worldcom
- But only **30,000** downloads
- All efforts shifted to server implementation 8/97

05-Feb-01 15:36

31

New Business Model

- **Shopping Infrastructure Company**
- **License Server to High-traffic Sites**
 - Search Engines: Yahoo, Excite, ...
 - ISPs: ATT Worldcom, AOL, MSN, ...
 - Specialty sites: ESPN, ...
- **Split revenue from merchants**

05-Feb-01 15:36

32

Enter Excite

- **Very Aggressive Company**
 - #7 -> #2 search engine in 18 months
 - 26 Million pages views / day
- **Believed in Internet Shopping**
 - Already had leading shopping channel
 - More users/day than we got in 6 months
 - Sales force of 40
- **Acquires Netbot for \$35 M**
 - Net Present Value of \$300M in @Home stock

05-Feb-01 15:36

33

Copycats

- **Jungle**
 - Snagged Yahoo contract for Xmas '97
 - Bought by Amazon for \$185M in Aug '98
 - Net Present Value of \$1.11B
 - ZShops?
 - Eliminate technology?
- **C2B**
 - Bought by Inktomi for \$90M in Sept '98
 - Net Present Value of \$630M
- **Infospace**
 - Activeshopper June '99
- **MySimon**
 - Bought by Cnet for \$700M in Jan '00
- **Cadabra**
 - Bought by Goto for \$250M in Feb '00

05-Feb-01 15:36

34

The Future

- **Price Tracking**
- **Auctions**
- **Category Wizards**
 - E.g. cars, long distance carriers
- **Cross Selling**
 - E.g. batteries, concert tickets
- **Gift Advisor**
- **Virtual Registry**
- **Cross Vendor Loyalty Program**

05-Feb-01 15:36

35