# PURDUE
UNIVERSITY®

# CS47300:  Web Information Search and Management
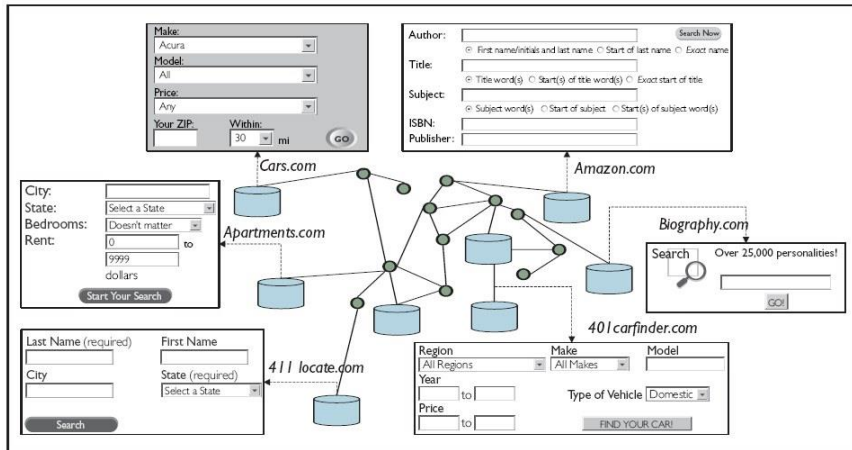
*Deep Web*
Prof. Chris Clifton
17 November 2017

Indiana
Center for
Database
Systems
TM

---

# PURDUE Deep Web vs. Dark Web
UNIVERSITY®

- Dark Web:  Hidden intentionally
  - Largely to support illegal or socially unacceptable activity
  - *But legality and acceptability vary, web is trans-national and trans-cultural*
  - We won't go here…
- Deep Web:  Data hidden behind interfaces
  - Federated search is one answer
  - But can we crawl this data?

2

---

# Conceptual View
## *(He, Patel, Zhang, Chang '07)*
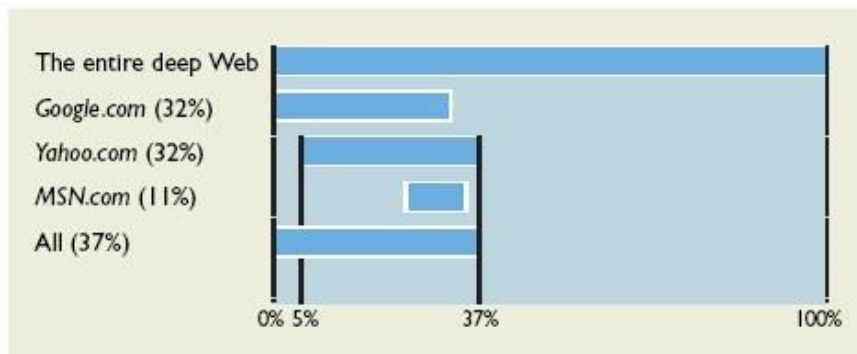


3

---

# Does the Deep Web Matter?

- Where are the entry points?
- What is the scale?
- How "structured" is the data?
- What topics are covered?
- How well do search engines already cover this?
- Wat about existing specialized portals?
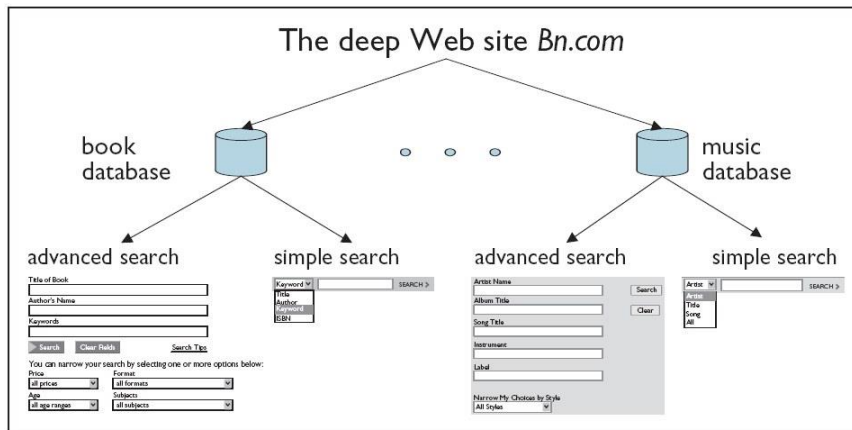
4

## Size Estimate of the Deep Web

| | Sampling Results | Total Estimate | 99% Confidence Interval |
|---|---|---|---|
| Deep Web sites | 126 | 307,000 | 236,000 - 377,000 |
| Web databases | 190 | 450,000 | 366,000 - 535,000 |
| —unstructured | 43 | 102,000 | 62,000 - 142,000 |
| —structured | 147 | 348,000 | 275,000 - 423,000 |
| Query interfaces | 406 | 1,258,000 | 1,097,000 - 1,419,000 |

5

## Search Engine Coverage



The entire deep Web
Google.com (32%)
Yahoo.com (32%)
MSN.com (11%)
All (37%)

0% 5%    37%    100%

6

3

# Deep Web Components
## *(He, Patel, Zhang, Chang '07)*

The deep Web site *Bn.com*

book database        •  •  •        music database

advanced search    simple search    advanced search    simple search

7

---

# Challenges

- How do we know what is in a database?
  - Sample queries?
  - Search page
    - Descriptive information
    - Form fields
- How do we query it?
- How do we process results?

8

4

# Can this be real?

- "General" search
  - See Google, etc.
- "Specialized" search
  - Metaquerier
  - Cazoodle

9