# Introduction to Web Archiving

**Marc Spaniol**

Saarbrücken, May 28, 2009

**Introduction to Web Archiving**

Marc Spaniol

Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-1/50

# Agenda

**Introduction to Web Archiving**

Marc Spaniol

Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-2/50

# Indexing vs. Archiving

- Indexing
  - Completeness
  - Access to content
  - Scalability (speed)
  - Efficiency
  - Freshness

$\Rightarrow$   "Taking a Photo"

- Archiving
  - Completeness
  - Access to content
  - Scalability (coverage)
  - Authenticity
  - Coherence
  - Durability

$\Rightarrow$   "Shooting a Movie"

# The Challenge of Web Archiving

- Digital library
    - Organized
    - Groomed content
    - Lots of metadata
    - Structured changes
    - Active preservation policies

- World Wide Web
    - A disorganized free-for-all
    - Very little metadata
    - Unpredictable additions, deletions, modifications
    - No (coordinated) preservation strategy

Databases and
Information Systems
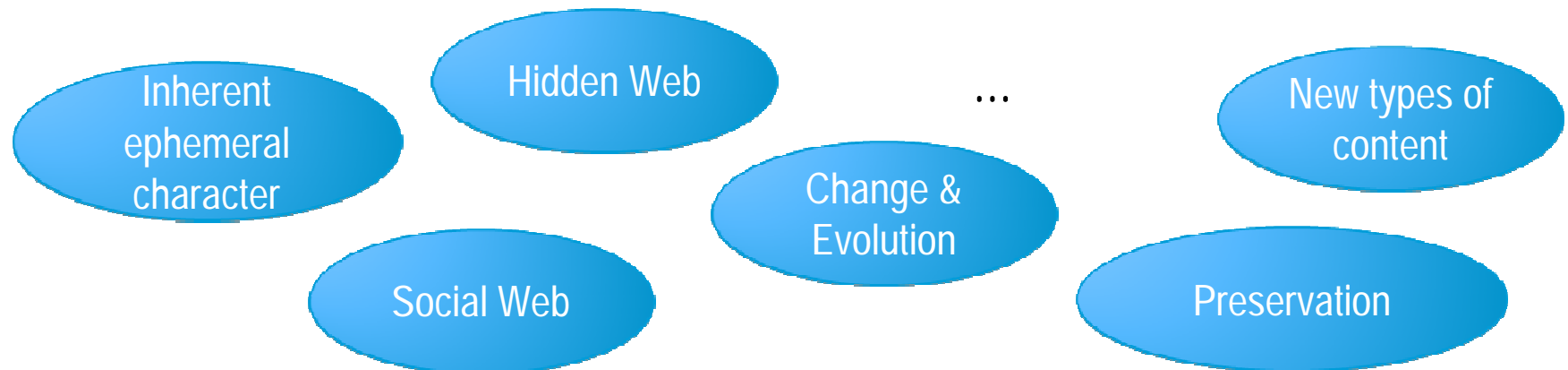Prof. Dr. G. Weikum

MPII-Sp-0509-4/50

# Goals of Web Archiving

- Role of Web
  - Providing information and services for seemingly all domains
  - Reflecting all types of events, opinions, and developments within society, science, politics, environment, business, etc.
  - Giving room for the articulation for a multitude of stakeholders
- $\Rightarrow$ Archiving this quickly changing multifaceted information space has becomes a relevant issue for cultural heritage

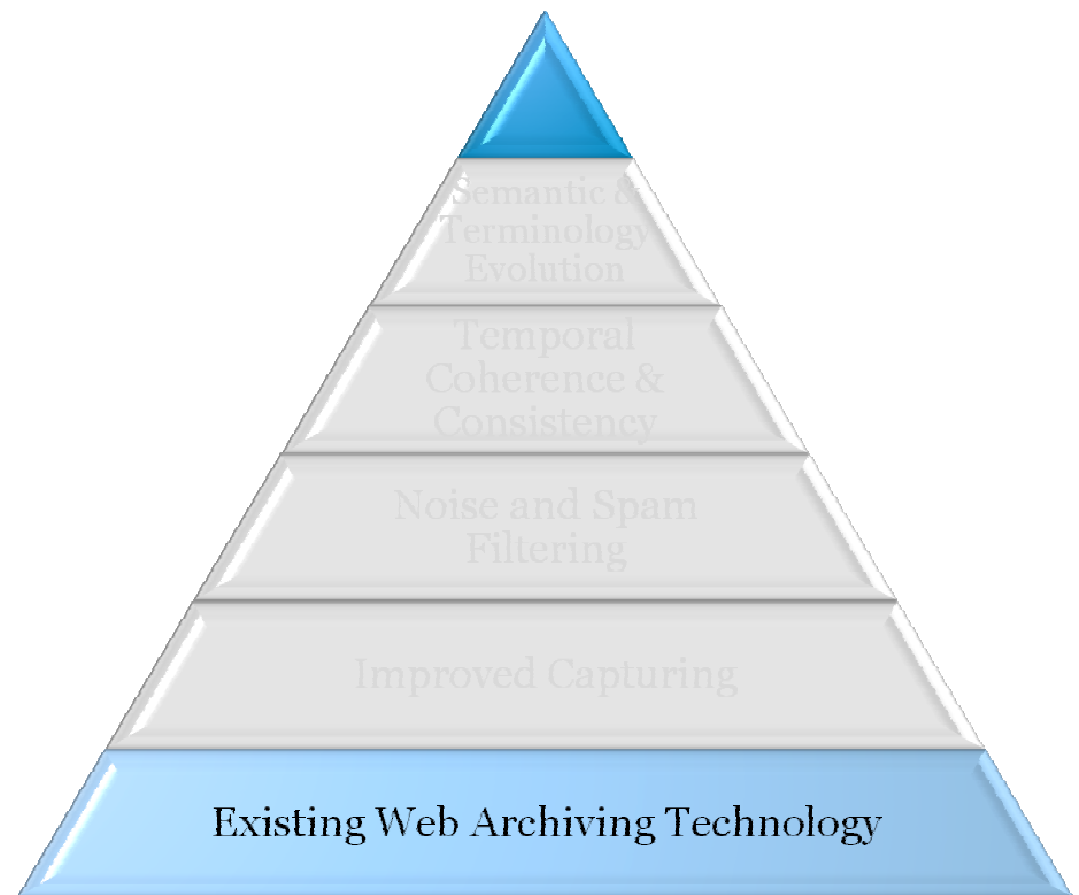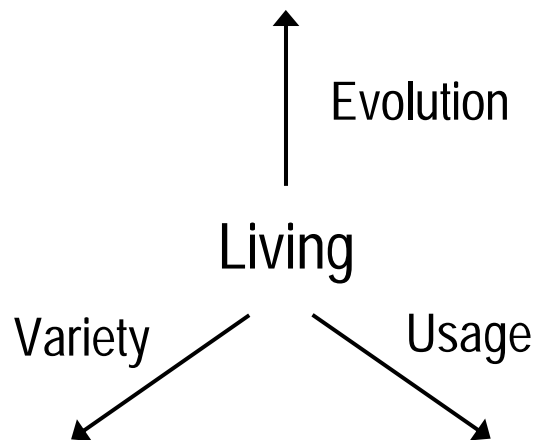- Web archiving imposes various challenges:

Inherent ephemeral character

Hidden Web

...

New types of content

Change & Evolution

Social Web

Preservation

# Next Generation Web Archiving

Development of Web archiving technology for

- High quality Web archives
- Long-term archive usability

$\Rightarrow$ From Web page storage
to "Living Web Archives"

Evolution

Living

Variety        Usage



Semantic &
Terminology
Evolution

Temporal
Coherence &
Consistency

Noise and Spam
Filtering
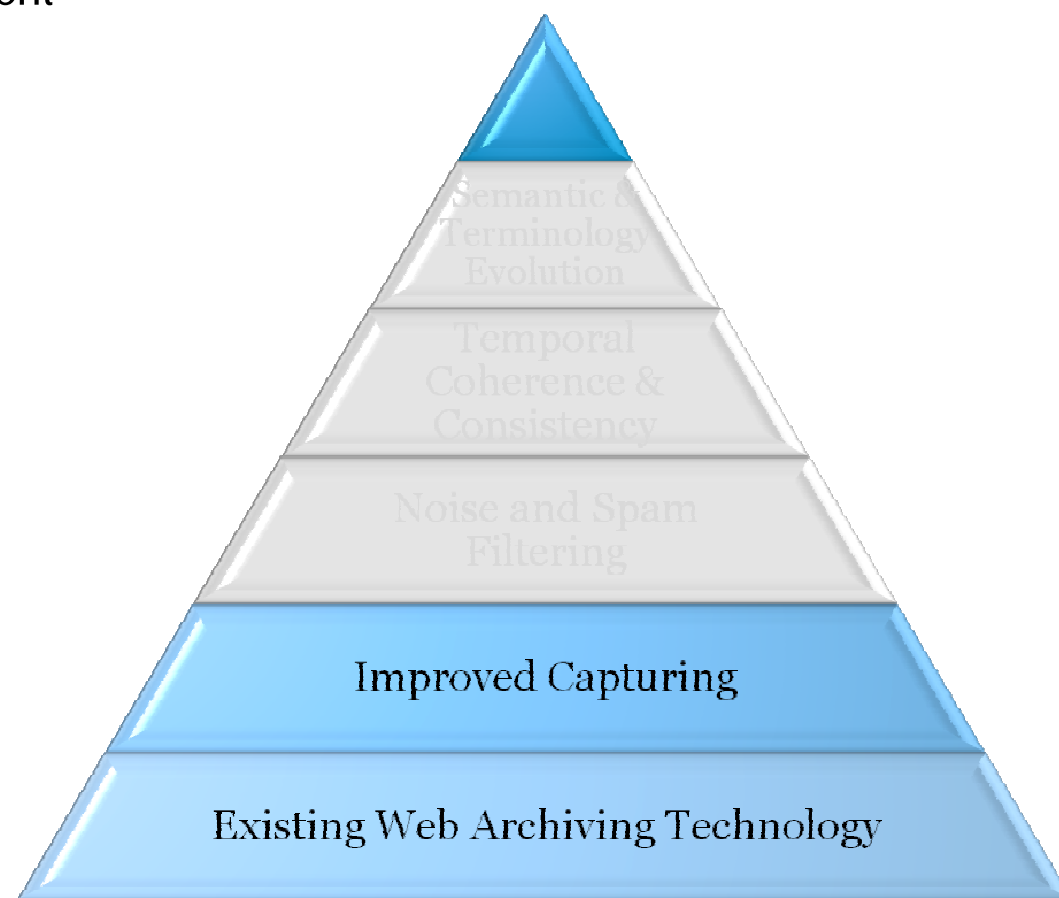
Improved Capturing

Existing Web Archiving Technology

# Archive Fidelity

Next generation Web archiving methods and tools

• Enhance archive fidelity and authenticity by

- Capturing all types of content
- Capturing of hidden Web
- Detecting traps



Pyramid levels from bottom to top:
- Existing Web Archiving Technology
- Improved Capturing
- Noise and Spam Filtering
- Temporal Coherence & Consistency
- Semantic & Terminology Evolution

Databases and
Information Systems
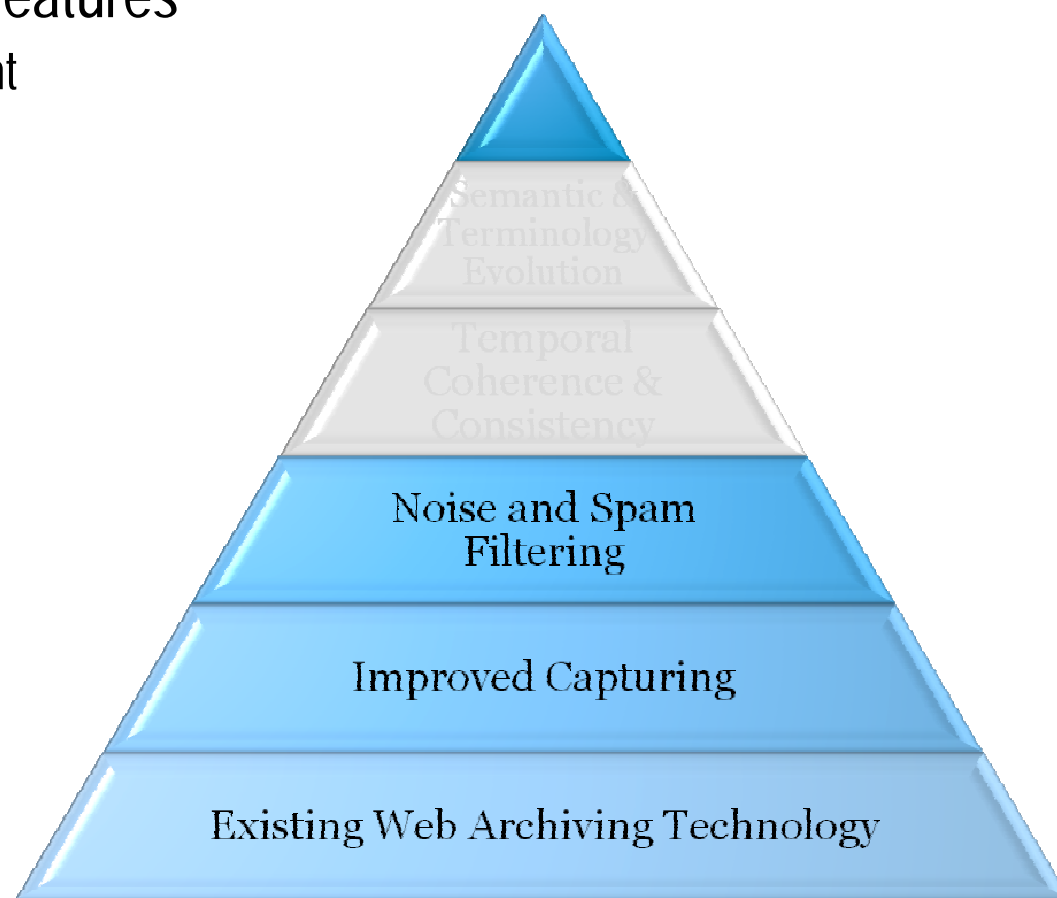Prof. Dr. G. Weikum

MPII-Sp-0509-7/50

# Advanced Filtering

Next generation Web archiving methods and tools:

- Enhance archive fidelity and authenticity

- Provide advanced filtering features

    - Capture all types of content

    - Detect traps

    - Filtering Web spam

    - Filtering noise



Pyramid levels (top to bottom):
Semantic & Terminology Evolution
Temporal Coherence & Consistency
Noise and Spam Filtering
Improved Capturing
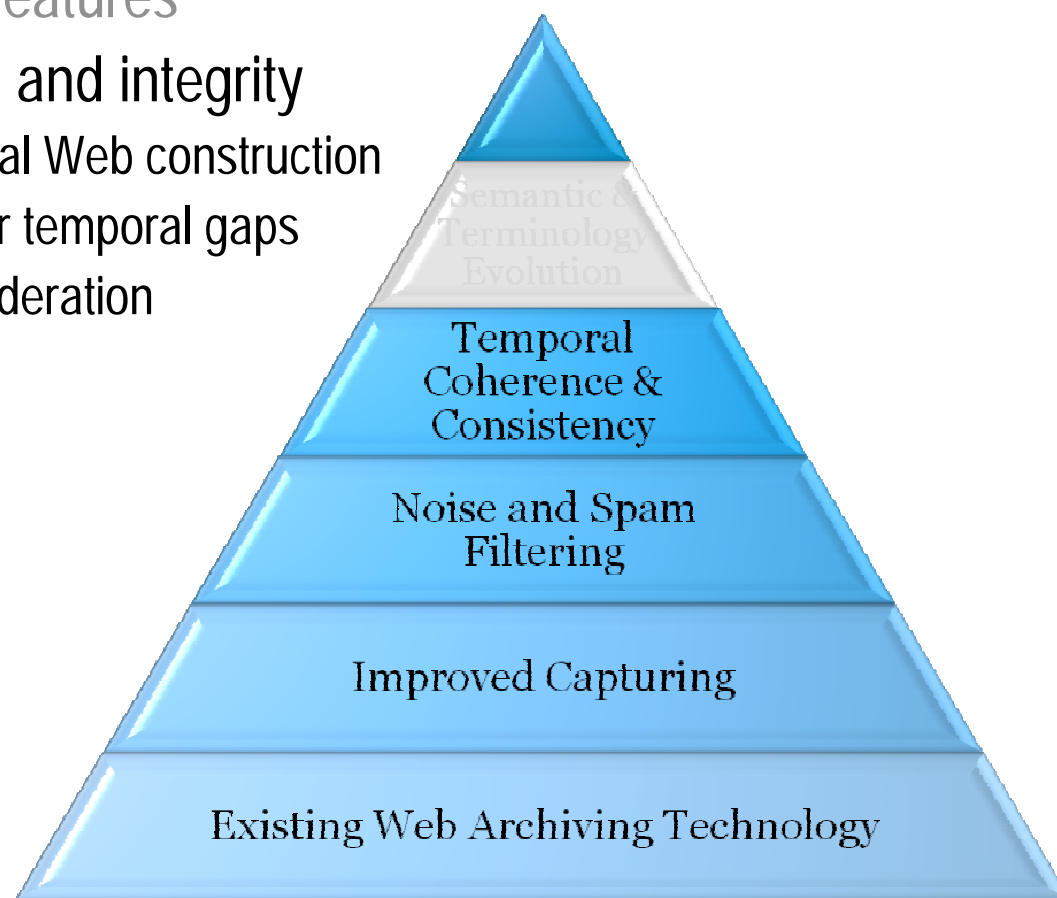Existing Web Archiving Technology

# Archive Coherence

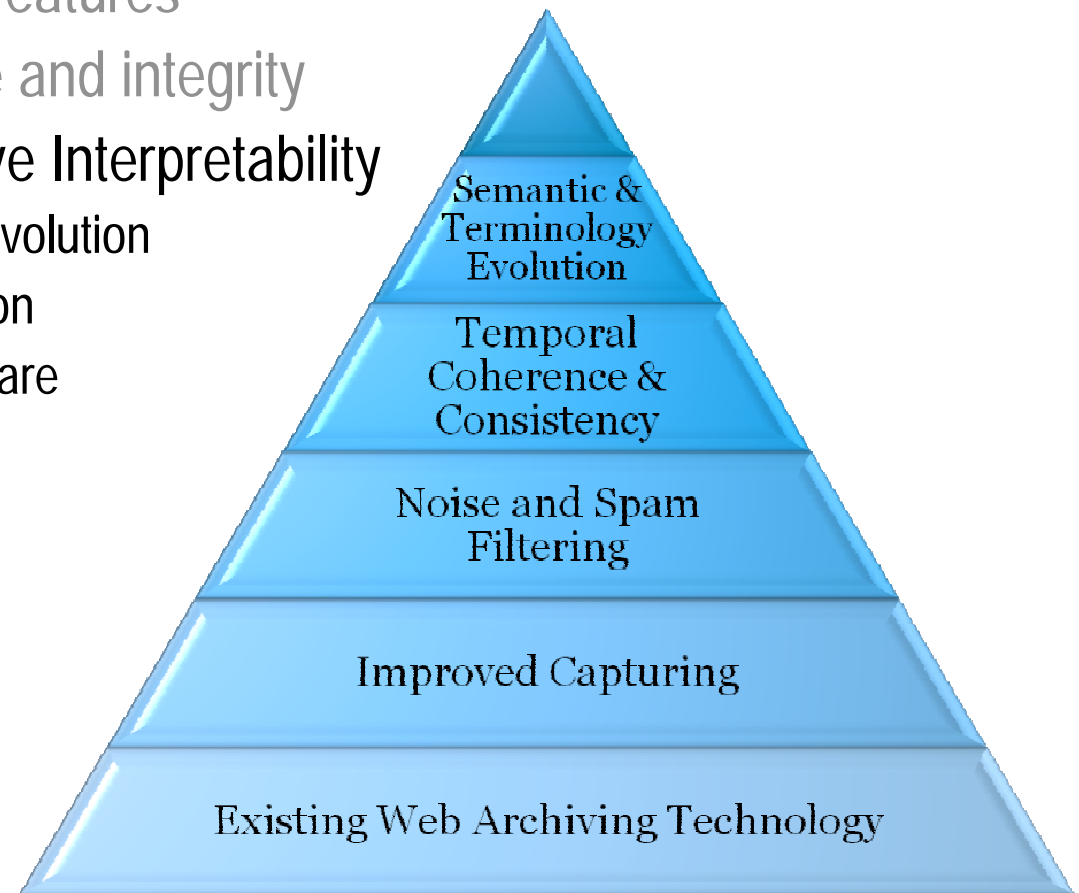## Next generation Web archiving methods and tools

- Enhance archive fidelity and authenticity
- Provide advanced filtering features
- Improve archive coherence and integrity
    - Deal with issues of temporal Web construction
    - Identify, analyze and repair temporal gaps
    - Consistent Web archive federation

Semantic &
Terminology
Evolution

Temporal
Coherence &
Consistency

Noise and Spam
Filtering

Improved Capturing

Existing Web Archiving Technology

# Archive Interpretability

Next generation Web archiving methods and tools

- Enhance archive fidelity and authenticity
- Provide advanced filtering features
- Improve archive coherence and integrity
- Facilitate (long-term) archive Interpretability
  - Dealing with terminology evolution
  - Handling semantic evolution
  - Preparing for evolution aware access support

Semantic & Terminology Evolution

Temporal Coherence & Consistency

Noise and Spam Filtering

Improved Capturing

Existing Web Archiving Technology

# Goals of Web Archiving Summarized

- Archiving function $\alpha$ applied to website $W$ produces a capture $C_W$ of the web site's resources and related metadata:

$$\alpha(W) \rightarrow C_W$$

- Restoration function $\rho$ "unpacks" the capture $C_W$ and reproduces the original site:

$$\rho(C_W) \rightarrow W$$

- Transformation function $\tau$ "unpacks" the capture $C_{W'}$ converts the components to the modern-day equivalent, and reproduces the original site within a new environment:

$$\tau(C_W) \rightarrow W_\Delta$$
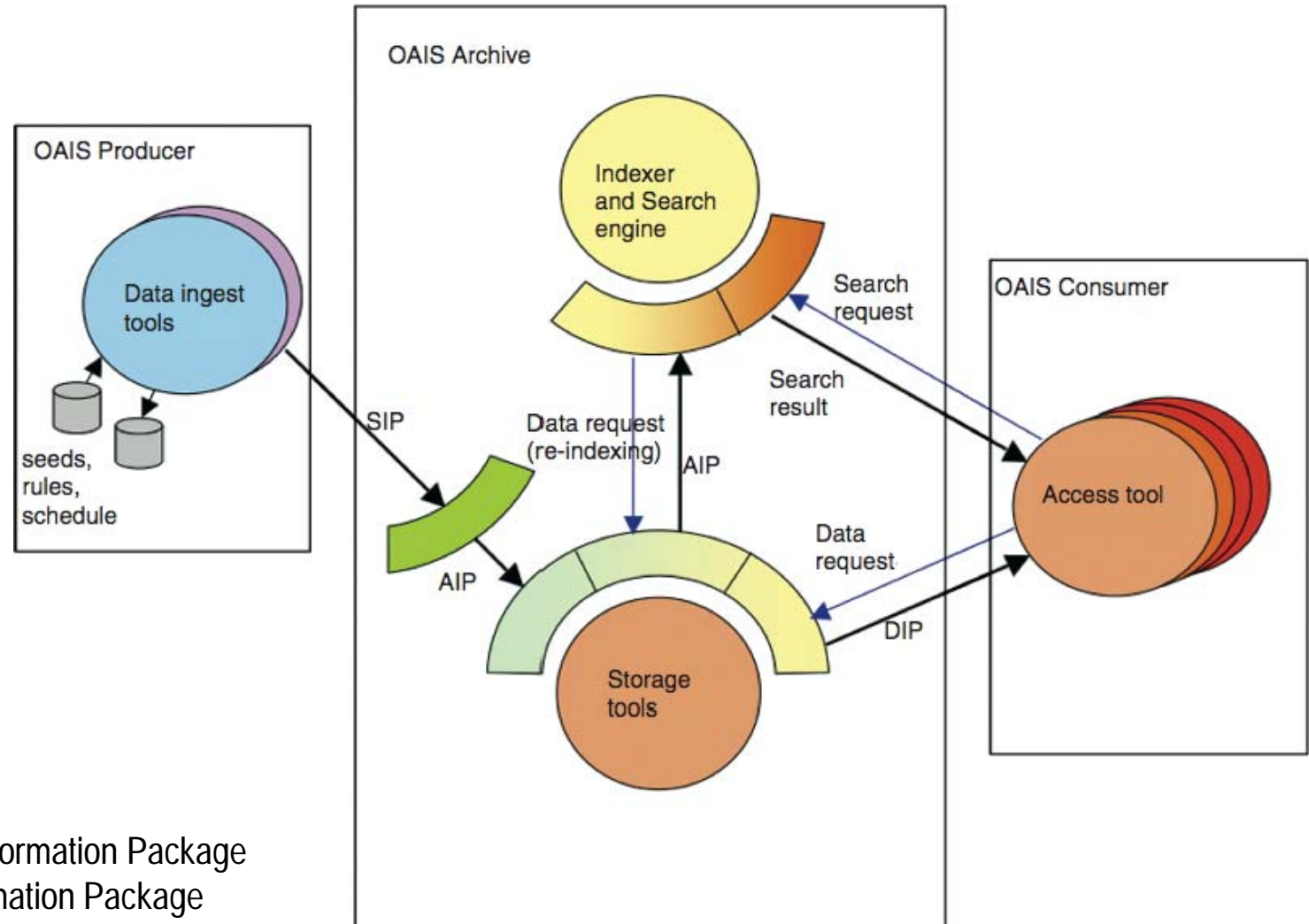
# Aspects of Web Archiving

**Introduction to Web Archiving**

Marc Spaniol

Databases and
Information Systems
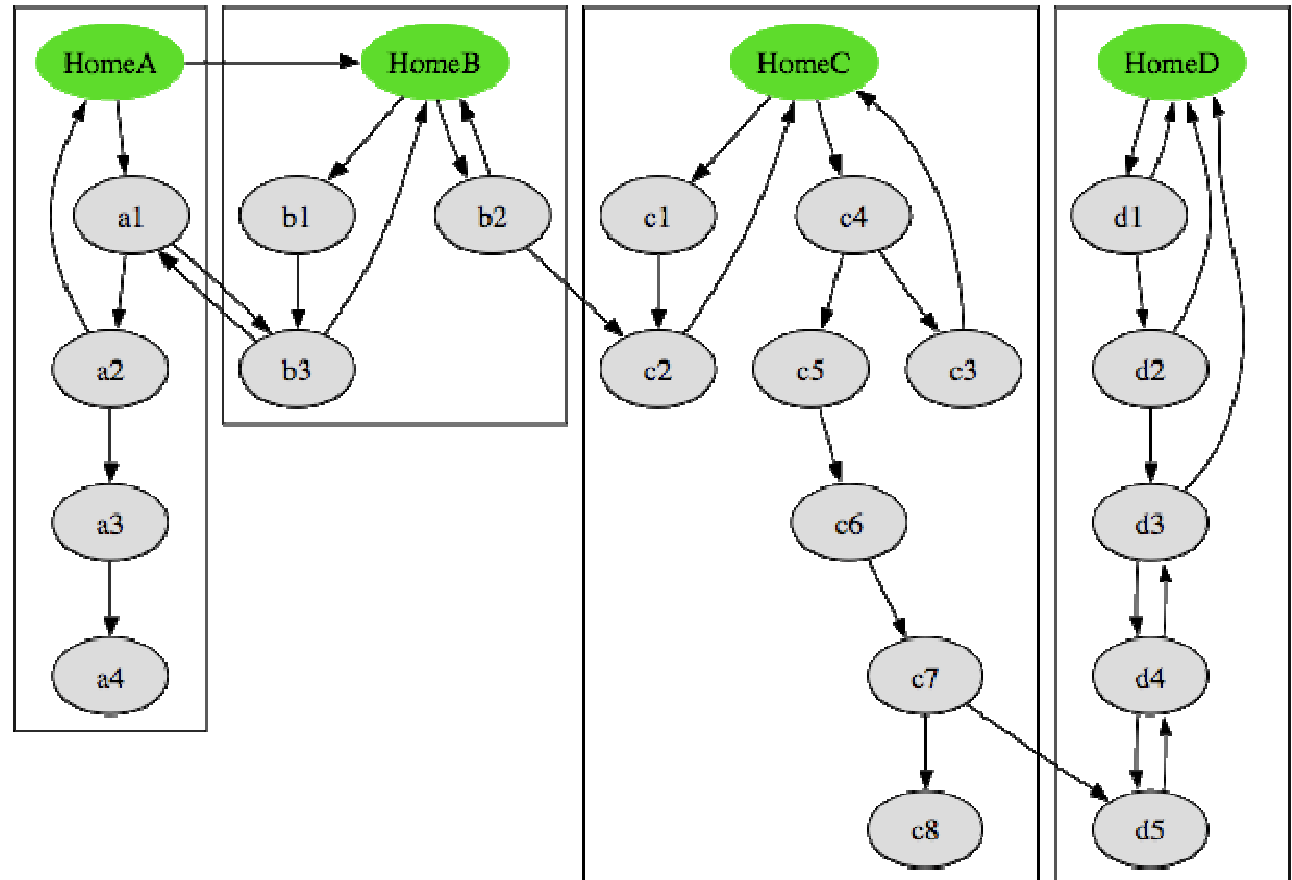Prof. Dr. G. Weikum

MPII-Sp-0509-13/50

# Web Archiving Tools



AIP: Archival Information Package
DIP: Data Information Package
SIP: Submission Information Package
OAIS: Open Archival Information System

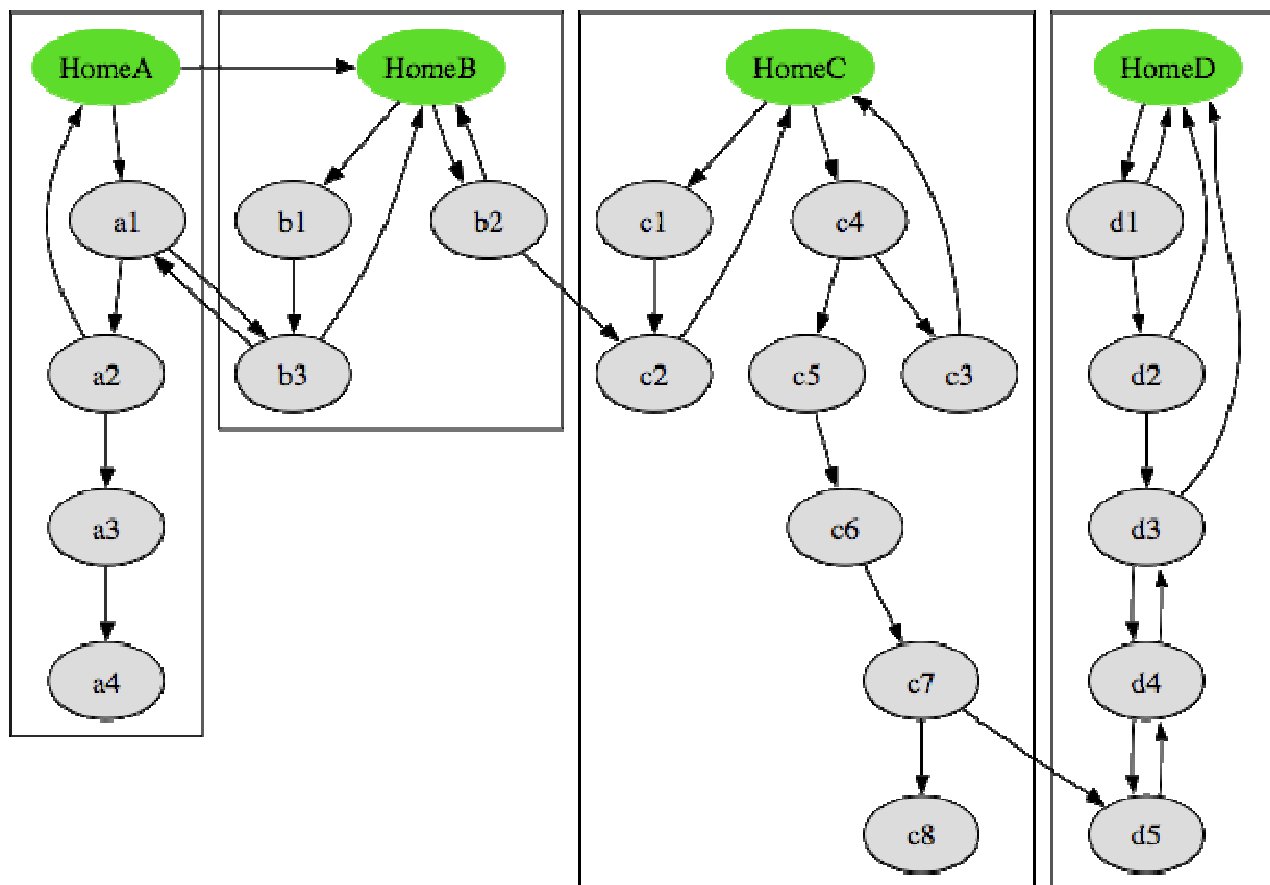# Selection of Seed(s) and Scope

- Entry point / seed: Where the capturing process (crawl) starts. Top of the hypertext path that will be followed.

- Scope: The extent of the area that will be included in the gathering, as defined by criteria applicable to each node.
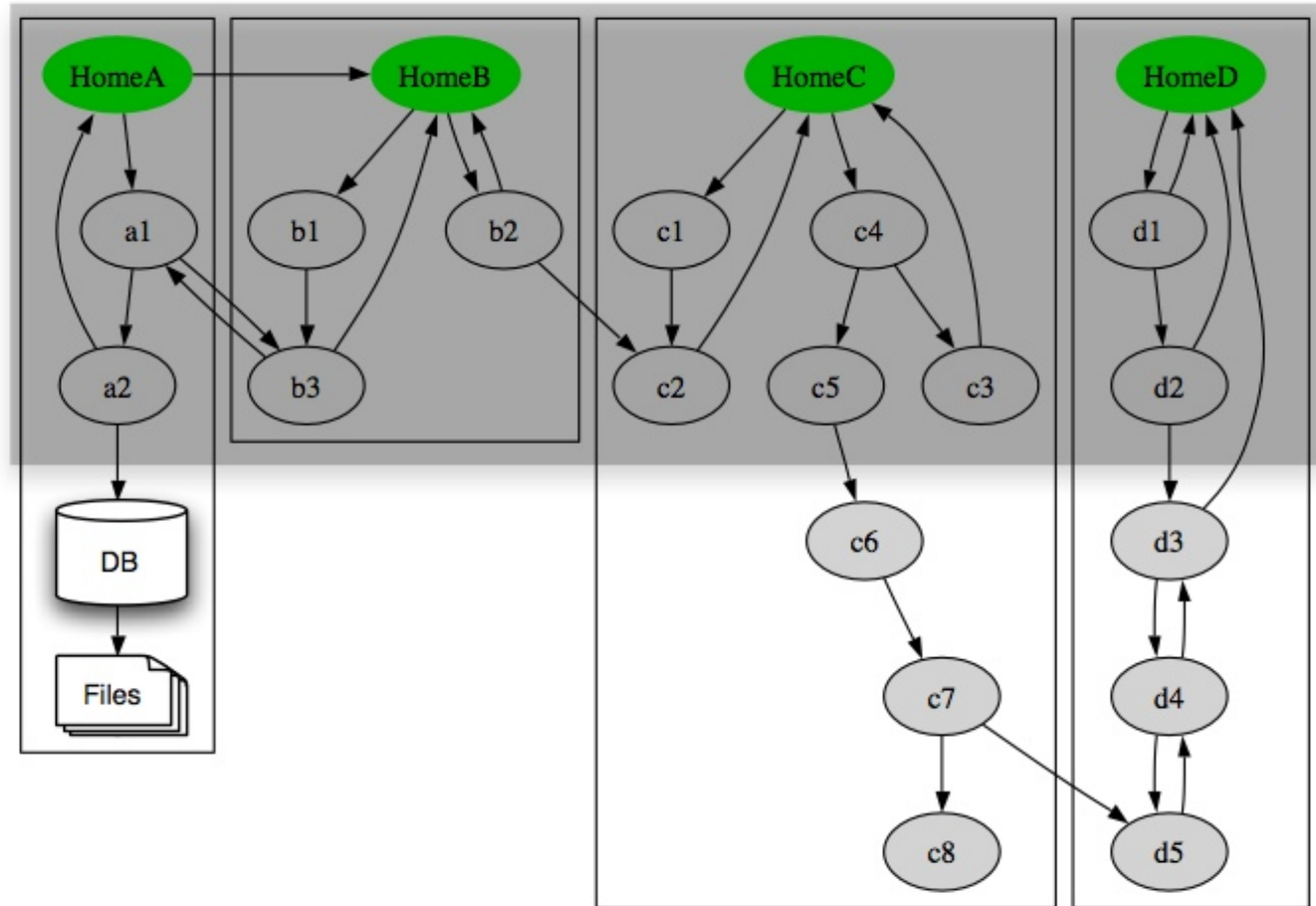
# Completeness

- Vertically: Number of relevant nodes found from entry point.

- Horizontally: Number of relevant entry points found within the designated perimeter.

# Extensive Collection

- Horizontal completeness is preferred to vertical completeness

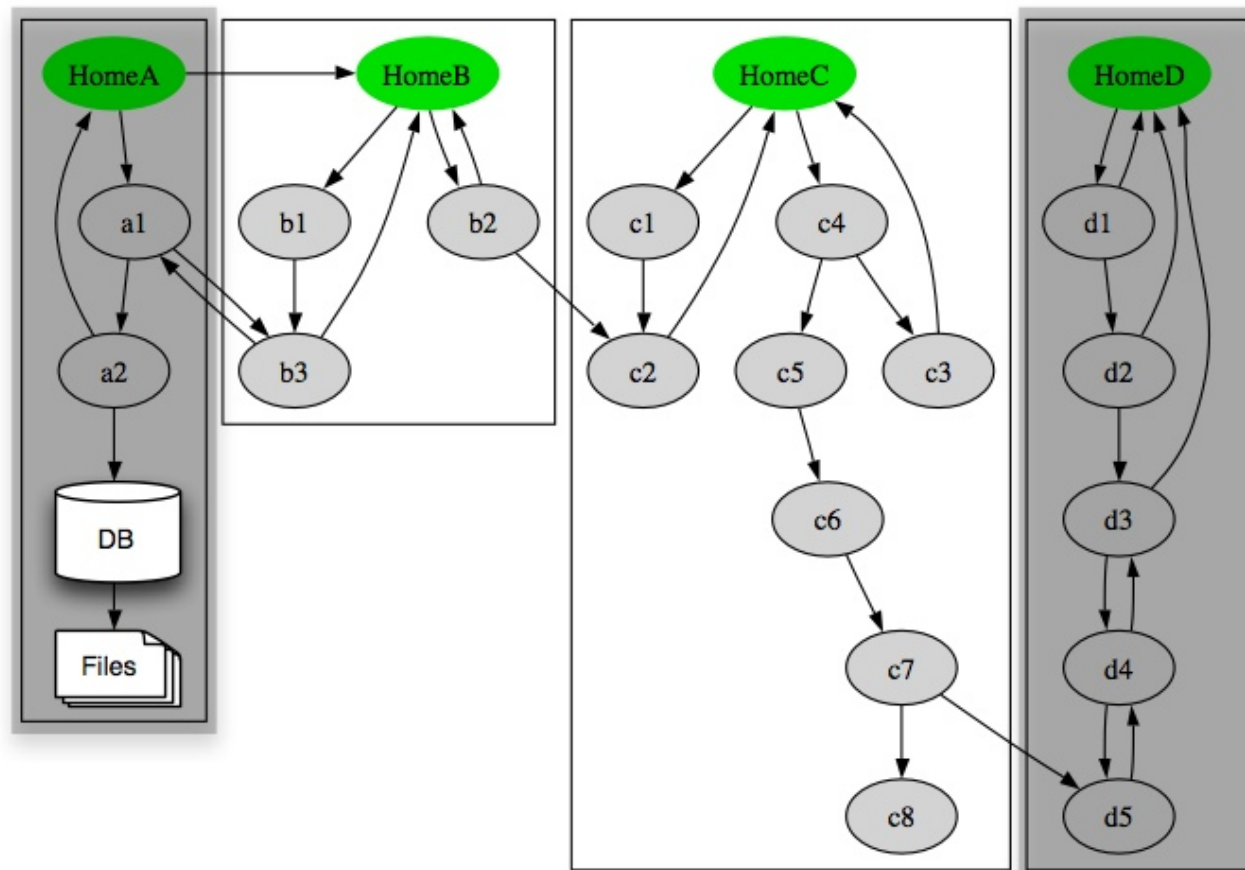- Holistic, domain based, or topic-centric archiving

# Intensive Collection

- Vertical completeness is preferred to horizontal completeness

- Site-based archiving

- Defines the high level target of a collection

- Explicit exclusion to avoid duplicate content with other collections

# The Challenge of Web Archiving

- HTTP cannot ask for only new or modified contents
  - *Timestamps* have limited benefit
  - No list of pages that have been deleted, changed, and added
  - *Each* content must be requested, one at a time, *by name*

- There is no "SELECT *" in HTTP
  - Crawlers can only GET one resource at a time, by name
  - HTTP cannot give a crawler a list of all URLs for the site

$\Rightarrow$ Undiscovered or hidden resources will not be captured or refreshed

$\Rightarrow$ "Strategy" required

# Server Side Archiving
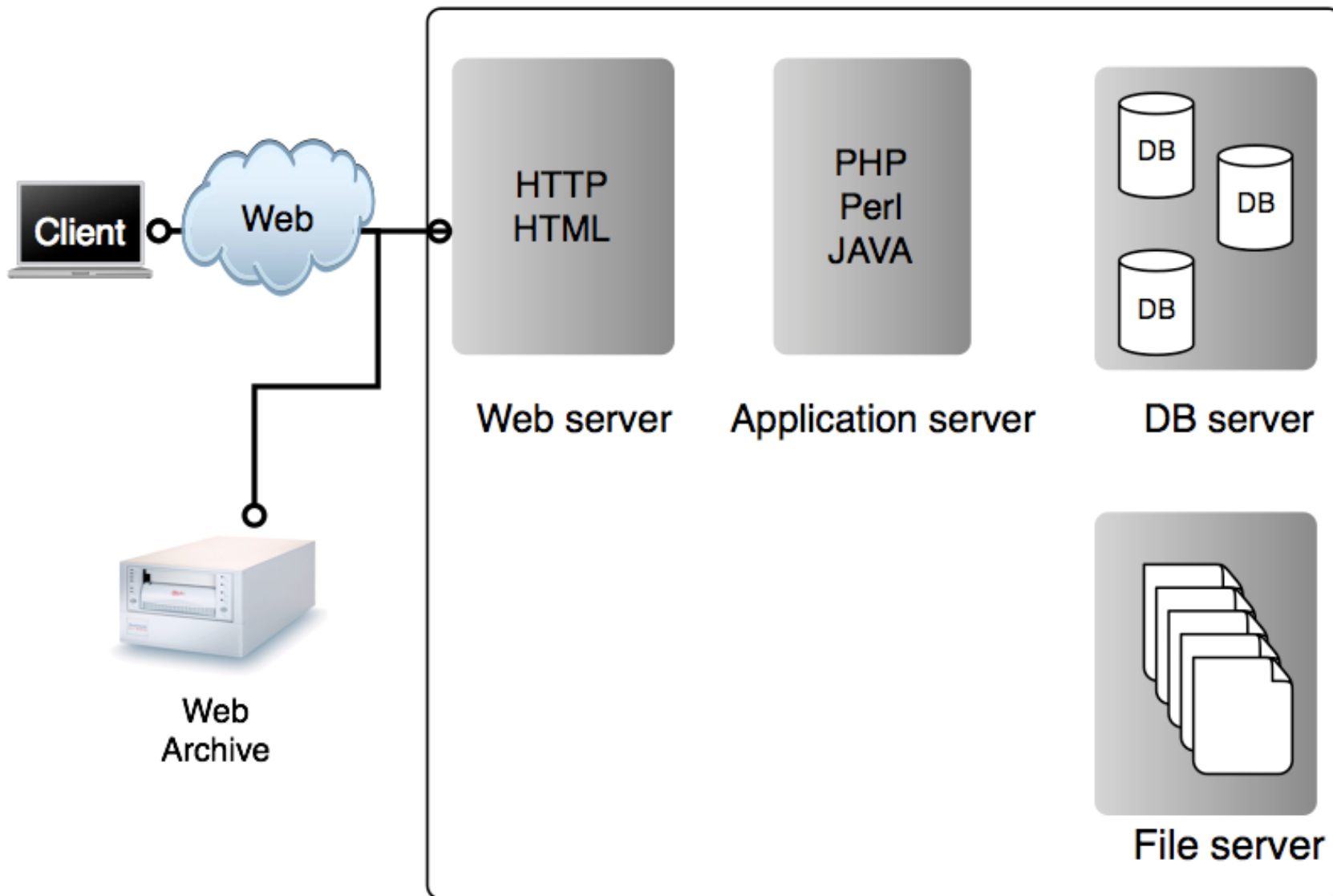
# Server Side Archiving Revisited

- Benefits

  + Extremely comprehensive

  + Changes are fully traceable (if budget permits)

  + Instantaneous snapshots possible

  + No network latency or limitations

  + Deep Web compliant

- Drawbacks

  - Change monitoring may decrease server performance

  - Needs sophisticated set-up

  - Requires server access

# Transaction based Archiving

# Transaction based Archiving Revisited

- Benefits
  - + Comes for "free"
  - + "Smart" coverage achieved by human interaction
  - + Simple maintenance
  - + No server collaboration/manipulation required


- Drawbacks
  - - Unsystematic
  - - Data quality is potentially poor
  - - Needs traffic monitoring
  - - Privacy issues
  - - Potential network latency or limitations
  - - Requires constant traffic
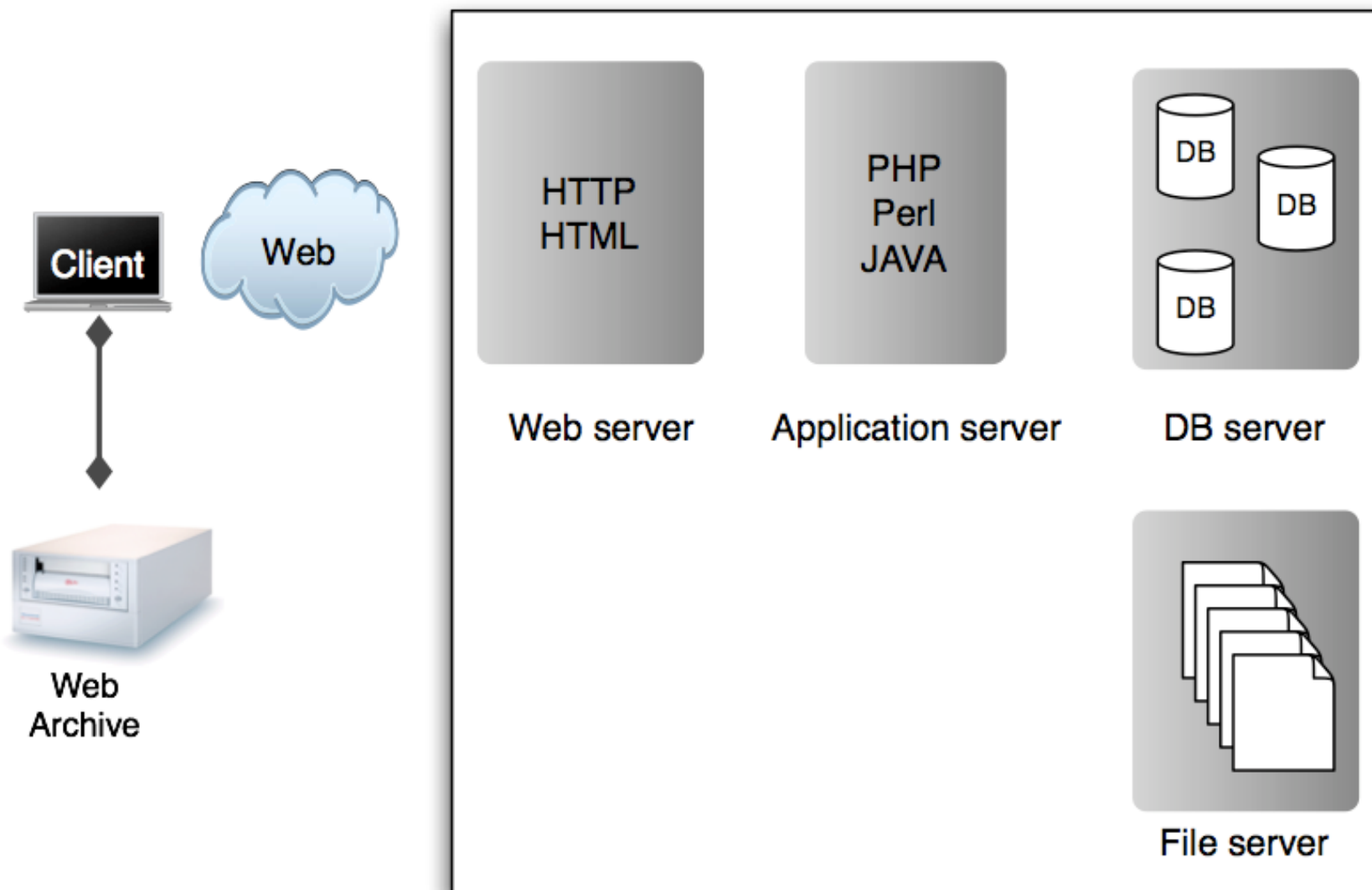
# Client Side Archiving

# Client Side Archiving Revisited

- Benefits

    + No server collaboration/manipulation needed

    + Only crawler set-up required

    + Mostly automated process (daily/weekly/monthly)


- Drawbacks

    - Changes might get lost

    - Good data quality requires sophisticated crawling strategies

    - Potential network latency or limitations
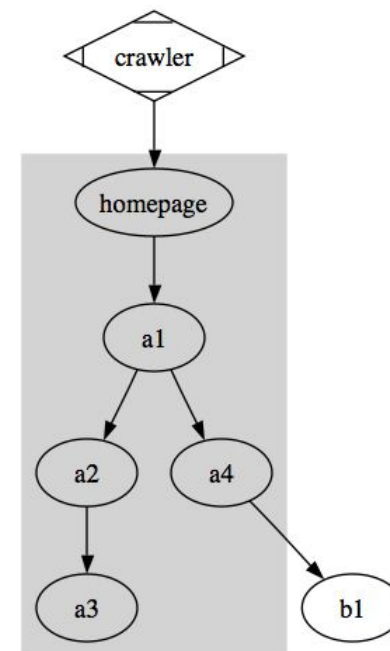
    - Computational "expensive"


Next week's lecture: "Data Quality in Web Archiving"

# Web Capturing with Heritrix

- Internet Archive's crawler

    - Open source java implementation

    - Web-scale archiving crawler

    - Extensible

- Key components

    - Scope

    - Frontier

    - Processor chains

- Configuration options include

    - Crawl scope, e.g. via

        - SURT expression: +http://(de,mpi-inf.mpg,www,)/

        - Regular expression: ^(http|https|dns):(//)?[a-zA-Z0-9\.]*mpi-inf.mpg.de/.*

    - Lot of fine-tuning features options

        - delay-factor

        - max-delay-ms

        - min-delay-ms

        - max-retries

        - retry-delay-seconds

        - etc.

# SURT
# Sort-friendly URI Reordering Transform

- Transformation applied to URIs
  - Left-to-right representation matching the natural hierarchy of domain names
  - Useful when comparing or sorting URIs

- Converting URIs according to SURT
  - Make all characters lowercase
  - Change the 'https' scheme to 'http'
  - '/' after a URI authority component only appear in the SURT form if it appeared in the plain URI form

- SURT form URIs are typically not used to specify exact URIs for fetching

- Less expressive than regular expressions → Exercises

# SURT Prefix

- Used for crawl scope specification in Heritrix

- Conversion to SURT prefix:

  1. Convert the URI to its SURT form.

  2. If there are $\geq$ 3 slashes ('/') in the SURT form, remove everything after the last slash

     <http://(org,example,www,)/main/subsection/> ✔

     <http://(org,example,www,)/main/subsection> $\rightarrow$ <http://(org,example,www,)/main/>

     <http://(org,example,www,)/> ✔

     <http://(org,example,www,)> ✔

  3. If the resulting form ends in an off-parenthesis ')', remove the off-parenthesis

     <http://(org,example,www,)> $\rightarrow$ <http://(org,example,www,>

Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-27/50

# Heritrix Output

- ARC/WARC files (Web ARChive) ~ 500 MB – 1 GB each
- "ZIP files" of content(s) and some metadata

Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-28/50

# A Webmaster's Omniscient View

**MySQL**

1. Data1
2. User.abc
3. Fred.foo

**Dynamic**

Authenticated

Entry point / seed

Tagged:
No robots

**httpd**

1. file1
2. /dir/wwx
3. Foo.html

Orphaned

Deep

Unknown/not visible

# Web Server's View of a Web Site

Require authentication

Entry point / seed

Generated on-the-fly (e.g. by CGI)

Tagged: No robots

Unknown/not visible

**Introduction to Web Archiving**

Marc Spaniol

# A Crawler's View of a Web Site

Not crawled
(protected)

Entry point / seed

Not crawled
(generated on-the-fly,
e.g. by CGI)

Not crawled
robots.txt or
robots META tag

Crawled pages

Not crawled
(unadvertised & unlinked)

Not crawled
(remote link only)

Not crawled
(too deep)

Remote web site

Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-31/50

# Streaming Media Capturing

**Introduction to Web Archiving**

Marc Spaniol

# Web Information Systems

Dynamic Web sites



Hidden Web



- Each interaction with a Web information system can potentially generate a unique customized response

$\Rightarrow$ Document the context of this interaction, or pseudo-transaction

# Hidden Web Archiving

- Procedure:

  1. Detect it

  2. Try to crawl it by automatic query generation

  3. <u>Or</u>  encourage site producer to be more friendly to crawlers

- Special crawl for documentary gateways

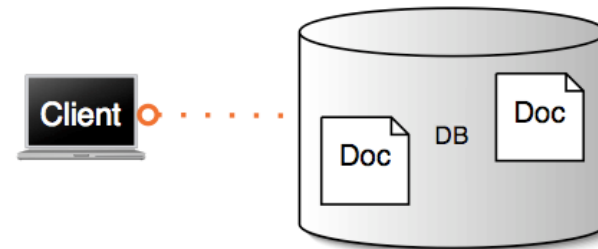  - Find patterns

  - Feed the fields

    - Finding forms is easy, filling them is not

    - Proximity of text near fields (beyond and on the left)

    - Tokenization and analyze

  - Reconstruct navigation or access logic

# Hidden Web Archiving: HTML Form extraction



[Fontes & Soares Silva, WIDM 2004]

# Crawler-Server Collaboration

- Open Archives Initiative (OAI) Protocol for Metadata Harvesting

- Provided flat list (maybe hidden for public)

- RSS feeds

- OAI server

  - Pushed by search-engines

  - Yahoo content acquisition program, google

$\Rightarrow$ The *sitemap* standard is intended to list the resources at a site

Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-36/50

# HTTP GET vs. OAI-PMH GetRecord

**Introduction to Web Archiving**

Marc Spaniol

Databases and Information Systems Prof. Dr. G. Weikum

MPII-Sp-0509-37/50

HTTP GetRecord
**Machine-readable**

HTTP GET
**Human-readable**

JHOVE METADATA

MD-5 LS

**Complex Object**

mod_oai

**Apache Web Server**

"GET /headlines.html HTTP1.1"

"GET /modoai/?verb=GetRecord&identifier= headlines.html&metadaprefix=oai_didl"

WEB SITE

# OAI-PMH data model

Adobe PDF ← resource

| OAI-PMH identifier<br>= entry point to all records pertaining to the resource | ← item |

Metadata

| Dublin Core<br>metadata | MARCXML<br>metadata | MPEG-21<br>DIDL | METS | ← records |

**simple**   **more expressive**   **highly expressive**   **highly expressive**
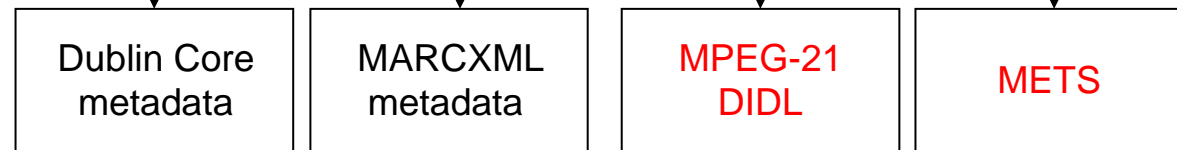
# OAI-PMH Syntax

| Verb | Function |
|------|----------|
| Identify | Repository description |
| ListMetadataFormats | Supported metadata formats |
| ListSets | Sets defined by repository |
| ListIdentifiers | Unique IDs contained in repository |
| ListRecords | Listing of $n$ records |
| GetRecord | Listing of a single record |

Repository metadata { Identify, ListMetadataFormats, ListSets }

Harvesting calls { ListIdentifiers, ListRecords, GetRecord }

**http://www.sample.edu/modoai?verb=ListIdentifiers&metdataPrefix=oai_dc&from=2004-09-15&set=mime:video:mpeg**

Human-readable Web site

OAI-PMH invocation

Give me a list of all resources, include Dublin Core metadata, dating from 9/15/2004 through today, and that are MIME type video-MPEG.

# Exemplary Application of OAI-PMH

Three from Tivoli



**FANtastisch**

ZU JEDEM HEIMSPIEL:
Gewinnen auf
www.takeda.de

Ebenso einzigartig wie die Alemannia sind ihre Fans. Zum Beispiel Boris, Marc und Markus: Spätestens seit sie im Pokalfinale 2004 die Originalhose von George Mbwando ergattern konnten, kennt man die Clique im X-Block. Was die drei und alle anderen Fans verbindet, sind die Liebe und die Leidenschaft für Schwarz-Gelb. Auch wir von Takeda Pharma teilen diese Leidenschaft. Und stellen deshalb hier echte Alemannia-Fans vor. Weitere Fan-Porträts sowie Infos über unser Gewinnspiel finden Sie auf www.takeda.de.

EUREGIO PARTNER

Takeda Pharma

• Official Alemannia Aachen fan leaflet
• No. 8, Season 2005/2006
• …

**http://www.takeda.de/unternehmen/pdf/fantastisch/pdf8_17.pdf
encoded as an MPEG-21 DIDL**

- DC metadata
- Jhove metadata
- Checksum
- • • •
- Provenance

Adobe PDF

```
<didl> <metadata source="jhove">...</metadata>
    <metadata source="file">...</metadata>
    <metadata source="essence">...</metadata>
    <metadata source="grep">...</metadata> ...
    <resource mimeType="application/pdf"
    identifier="http://www.takeda.de/unternehmen/
    pdf/fantastisch/pdf8_17.pdf"
    encoding="base64">
    SADLFJSALDJF...SLDKFJASLDJ </resource>
</didl>
```
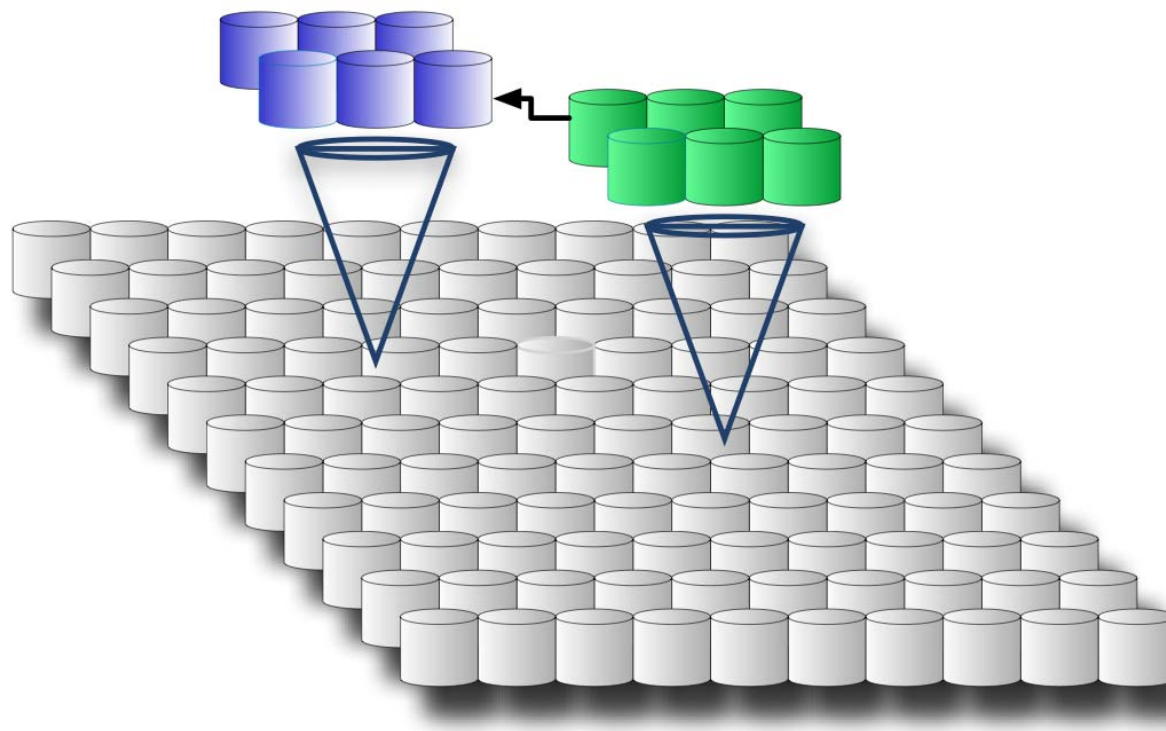
- Resource and metadata packaged together as a complex digital object represented via XML wrapper
- Uniform solution for simple & compound objects
- Unambiguous expression of locator of datastream
- Disambiguation between locators & identifiers
- OAI-PMH datestamp changes whenever the resource (datastreams & secondary information) changes
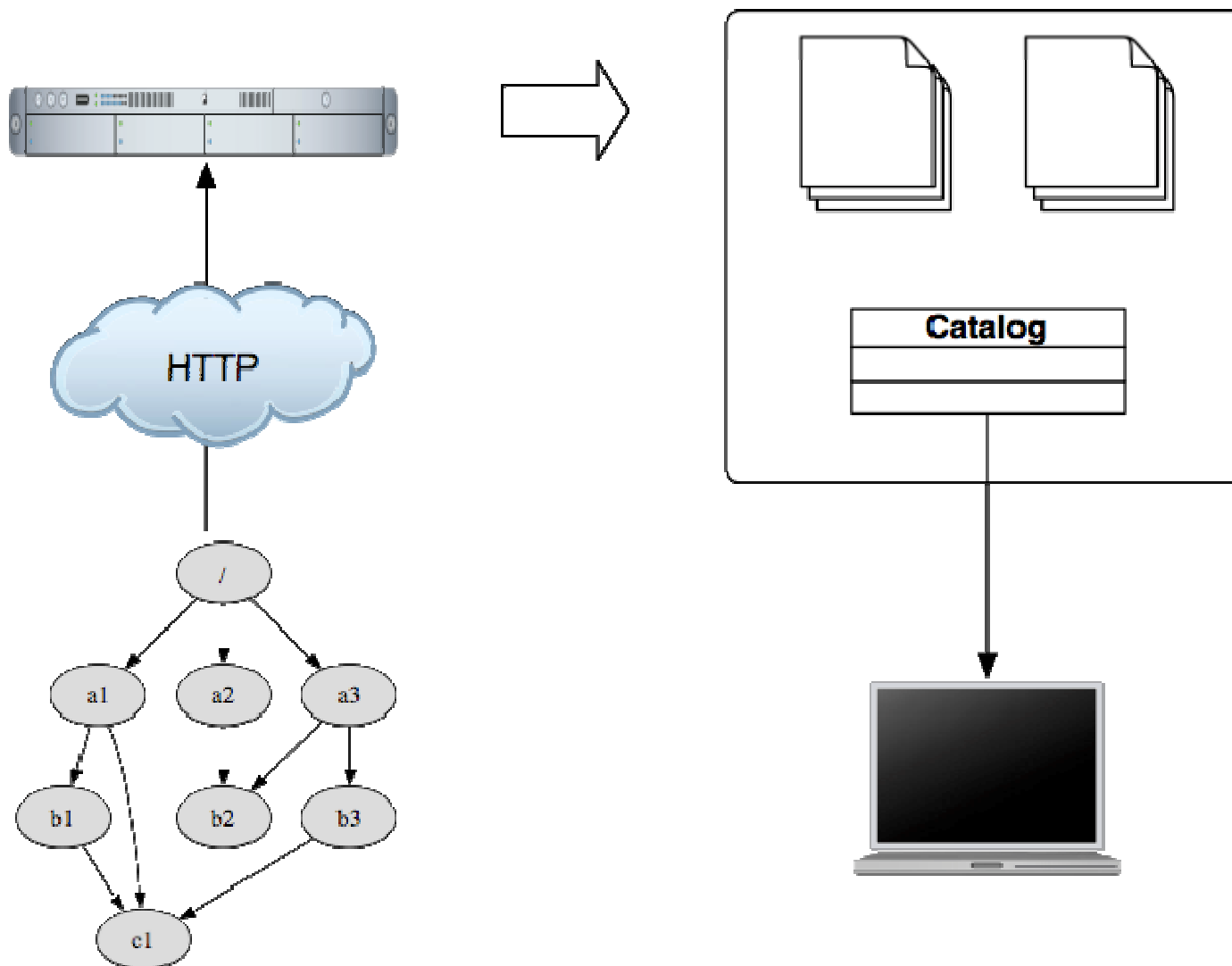
Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-40/50

# Archiving: Web Archives Grid

- Many "connected" servers
- WARC files spread among several servers
- Indexing of WARC files for access by URL and date

# Hosting: Non-Web Archive
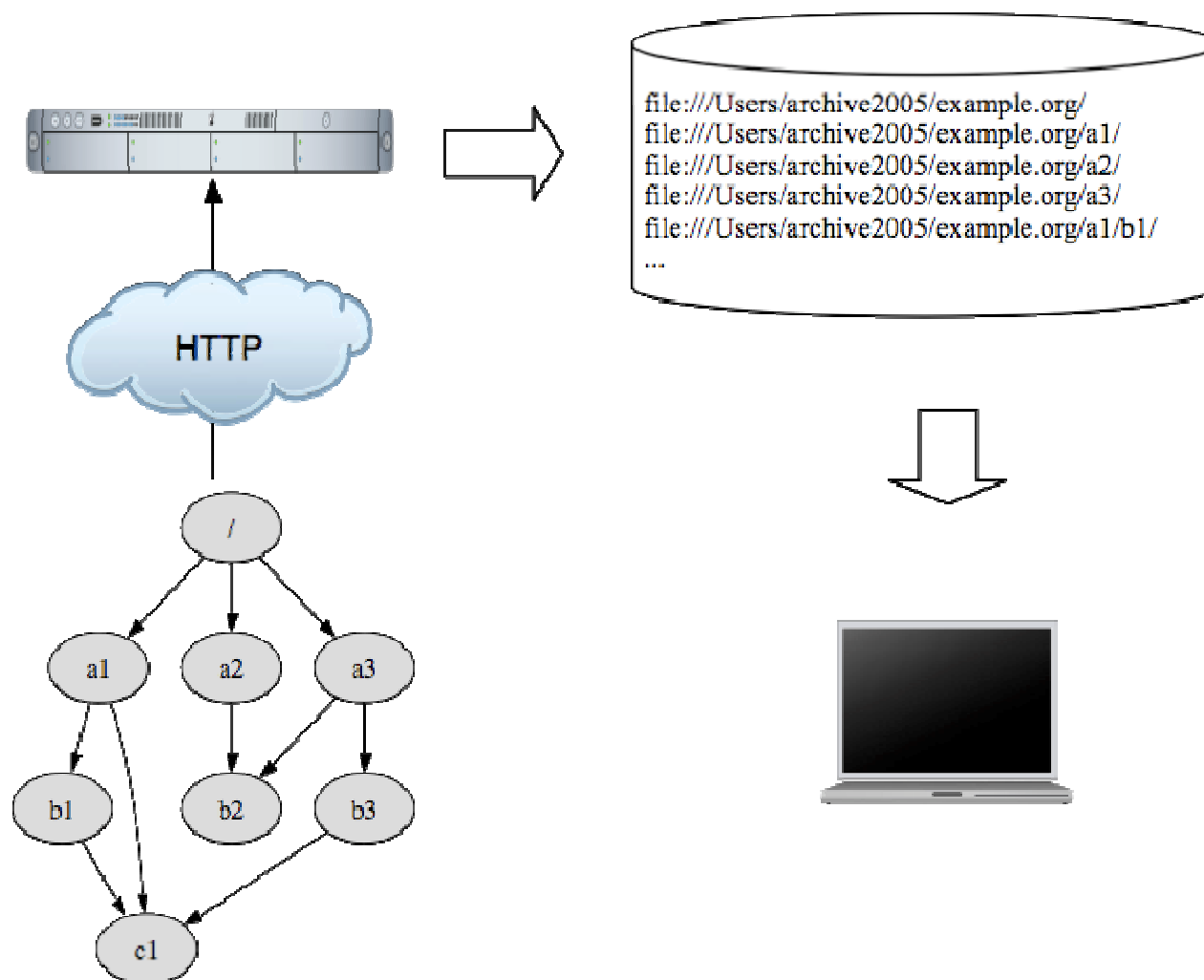
**Introduction to Web Archiving**

Marc Spaniol

Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-42/50

# Non-Web Archive Summary

- Benefits

  + Designed for archiving of specific (non-Web) collections

  + Potentially fast data access


- Drawbacks

  - Cataloging (usually) does not resemble hyperlink structure

  - Implementation cost for cataloging logic

  - Special search interface required

Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-43/50

# Hosting: Local File Navigation



file:///Users/archive2005/example.org/
file:///Users/archive2005/example.org/a1/
file:///Users/archive2005/example.org/a2/
file:///Users/archive2005/example.org/a3/
file:///Users/archive2005/example.org/a1/b1/
...

# Local File Navigation Summary

- Benefits
    - + Cheap
    - + Simple
    - + No additional infrastructure needed
    - + Fast

- Drawbacks
    - - Limited accessibility
    - - Small scale only
    - - Links are converted in relative ones
    - - Copying only

# Hosting: Web-served Archive

# Web-served Archive Summary

**Introduction to Web Archiving**

Marc Spaniol

- Benefits

    + Realistic "look&feel"

    + Convenient navigation

    + Time-travel also for non-technical experienced users possible


- Drawbacks

    - Web server needed

    - WARC/ARC file access required

    - Indexing tool for WARC/ARC files necessary

    - Time consuming sequential reads of WARC/ARC files
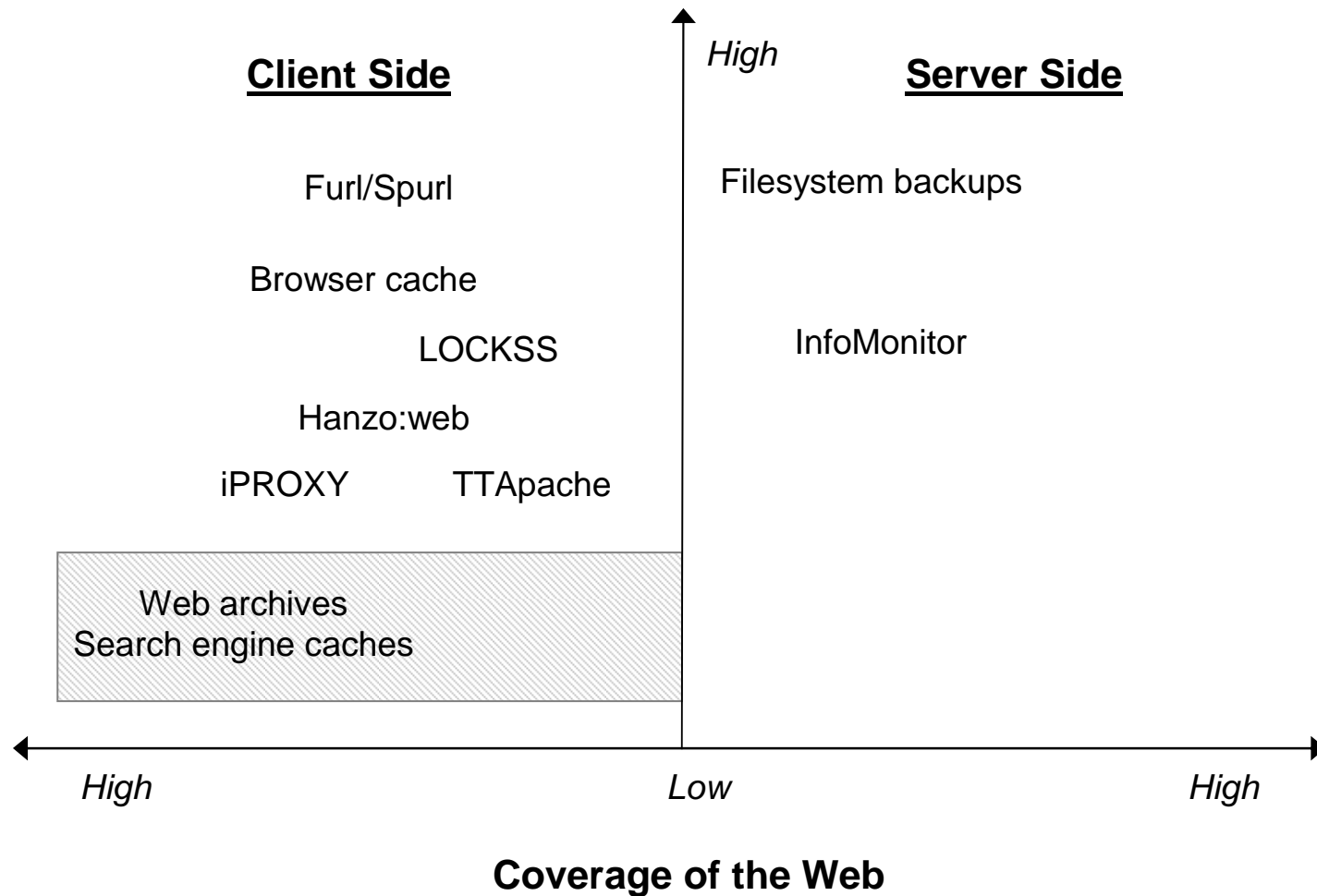
Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-47/50

# Cost of Web Archiving

**Publisher's cost**
(time, equipment, knowledge)

*High*

**Client Side**                    **Server Side**

Furl/Spurl                         Filesystem backups

Browser cache

LOCKSS                             InfoMonitor

Hanzo:web

iPROXY          TTApache

Web archives
Search engine caches

*High*                *Low*                *High*

**Coverage of the Web**

# Summary

- Web archiving is different from Web indexing

- Archiving crawlers
    - Do not aim at efficiency or freshness
    - Target at authenticity, coherence and durability

- Important aspects of Web archiving
    - Scope of archiving requires a clear definition
    - Seeds need to be carefully selected
    - Capturing of all URIs on a site and streaming media is hard
    - Preservation of hidden or dynamically generated contents is almost impossible
    - Pages may be orphaned intentionally or accidentally
    - Sitemaps rarely exist
    - WARC file processing is the bottleneck in retrieval
    - Capturing takes a long time (!!!) and contents may not fit to each other

Databases and
Information Systems
Prof. Dr. G. Weikum

MPII-Sp-0509-49/50

# References

[Heri09]     Heritrix: "Glossary".
             http://crawler.archive.org/articles/user_manual/glossary.html
             [last access: May 27, 2009]

[Masa06]     J. Masanès: "Web Archiving". Springer, New York, Inc., Secaucus, NJ,
             2006.

[MKSR04]     G. Mohr, M. Kimpton, M. Stack and I. Ranitovic: "Introduction to Heritrix, an
             archival quality Web crawler". 4th International Web Archiving Workshop
             (IWAW'04), 2004.
             http://crawler.archive.org/Mohr-et-al-2004.pdf
             [last access: May 27, 2009]

[NCSm06]     M. Nelson, F. McCown, J. Smith Thinking: "Differently About Web Page
             Preservation".
             http://www.loc.gov/today/cyberlc/feature_wdesc.php?rec=3896
             [last access: May 27, 2009]