

Identifiers

CS431 - Web Information Systems

Carl Lagoze - Cornell University - Feb. 6 2008

BEWARE

Most discussions and work on web
information involves (degenerates to)
discussions about what is the information
unit and how is it identified!

Acknowledgments

- Stuart Weibel - OCLC
- Herbert Van de Sompel - LANL
- Andy Powell - EduServ
- Norman Paskin - International DOI Foundation

Identifiers

- Provide a key or **handle** linking abstract concepts to physical or perceptible entities
- Provide us with a necessary figment of persistence
- They are perhaps the one **essential** and common form of **metadata**
- Why bother?
 - Finding things
 - Comparing things
 - Referring to things (Citations)
 - Asserting ownership over things

Identity <-> Change <-> Persistence

- Paradox: reality contains things that persist and change over time
 - Heraclitus and Plato: can you step into the same river twice?
 - Ship of Theseus: over the years, the Athenians replaced each plank in the original ship of Theseus as it decayed, thereby keeping it in good repair. Eventually, there was not a single plank left of the original ship. So, did the Athenians still have one and the same ship that used to belong to Theseus

Identity <-> Change <-> Persistence



I have lots of identifiers for different roles and applications

- Carl Jay Lagoze, Dad, Hey you
- 123-456-7890 (SSN)
- 1234-5678-1234-1234 (Visa Card)
- FZBMLH (US Airways locator on January 18 flight to San Diego)

Lots of (non-digital) Identifier Standards

- ISBN (International Standard Book Number)
 - Origin 1966 U.K.
 - ISO 2108 1970
 - Uniquely identifies each edition and variation of a book
 - Number is semantically meaningful (components)
 - prefix/country code/pub code/item #/checksum
 - International administration (>150 countries)
- ISSN (International Standard Serial Number)
 - Uniquely identifies every serial (not issue or volume)
 - Semantically meaningless (anonymous)
 - International administration
- Lots of others
 - Recording Code, Tech Report, Audiovisual

<http://www.collectionscanada.ca/iso/tc46sc9/index.htm>

Some overarching comments

- Identification is complex “even” in the physical world
 - Librarians have dealt with it via Name/Authority records
- Identification has many non-technical dimensions
- The Web Architecture through URIs provides a simple uniform “technical” solution
- There are many more complex solutions that interleave architecture with policy
- Experience has shown that:
 - Simplicity often wins
 - Separation of concerns makes sense

What do we want from identifiers?

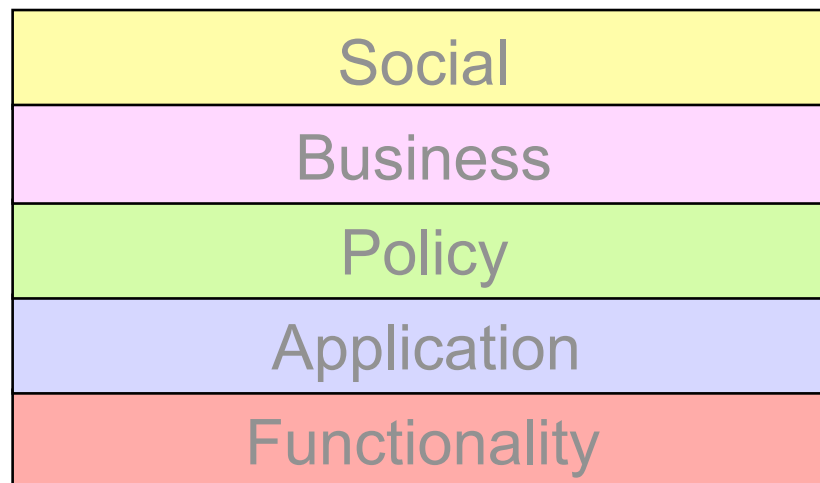
- Global uniqueness
- Authority/Reliability
- **Appropriate** functionality
 - Resolution
 - Other services
- Persistence

Identifier Issues

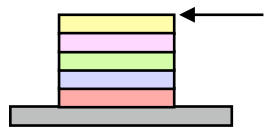
- Object granularity
- Identifier Context
 - Object atomicity
 - Part/whole relationships
- Location independence
 - Multiple location resolution
- Administration (centralized vs. decentralized)
- Human vs. machine processing

The Identifier Layer Cake

- In the digital world identification has lots of dimensions, only some of which are technical

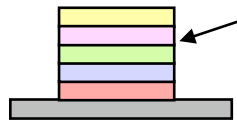


The Web: http...TCP/IP...future infrastructure?



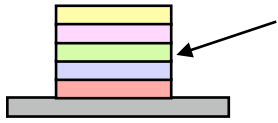
Social Layer

- The only guarantee of the usefulness and persistence of identifier systems is the commitment of the organizations which assign, manage, and resolve identifiers
- Whom do you trust?
 - Governments?
 - NGOs?
 - Cultural heritage institutions?
 - Commercial entities?
 - Non-profit consortia?
- We trust different agencies for different purposes at different times



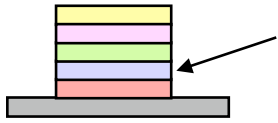
Business layer

- Who pays the cost?
- How, and how much?
- Who decides (see governance model)?



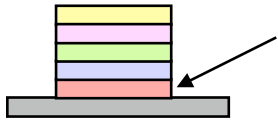
Policy Layer

- Who has the 'right' to assign or distribute Identifiers?
- Who has the 'right' to resolve them or offer services against them?
- What are appropriate assets for which identifiers can be assigned, and at what granularity?
- Can identifiers be recycled?
- Can ID-Asset bindings be changed?



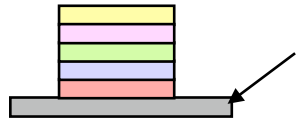
Application Layer

- What underlying dependencies are assumed?
 - http... tcp/ip...(bar code|RFID) scanners...
- What is the nature of the systems that support assignment, maintenance, resolution of identifiers?
- Are servers centralized? federated? peer to peer?
- How is uniqueness assured?



Functional Layer: Operational characteristics of Identifiers

- Is it globally unique? (easy)
- How does it 'behave'? What applications recognize it and act on it appropriately?
- Do identifiers need to be matched to the characteristics of the assets they identify?
- Do humans need to read and transcribe them?



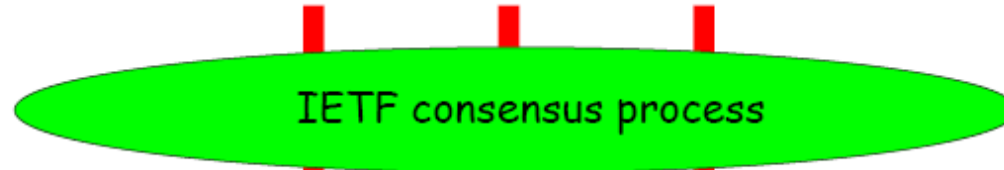
Technology layer: The Web

Some fundamental questions:

- Must our identifiers be URIs?
- Must they be universally actionable?
- If so, what is the desired action?
- Is there ever a reason to use a URI other than an http-URI as an identifier?

Identifiers in the web architecture

1992: Berners-Lee: "universal document identifier"



1994: RFC 1738 : Uniform Resource Locator



1995: RFC 1808 : Relative Uniform Resource Locators



1998: RFC 2396 URI Generic Syntax ("replaces 1738 and 1808")



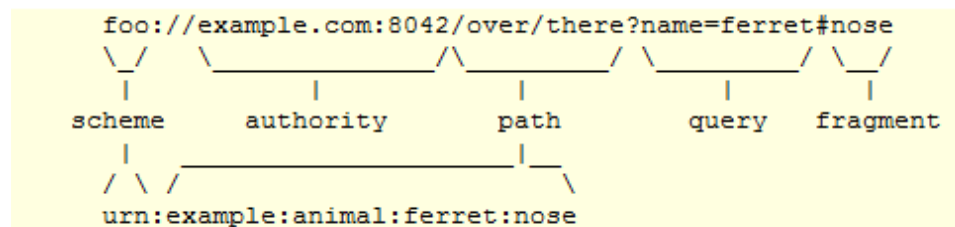
2004: RFC 2396 bis (revision) ?

Identifiers vs. Locators in the Web Architecture

- But **locators** and **identifiers** are not the same
- Not every identifier is a locator, but every locator is an identifier
- There is no deterministic way to distinguish if an identifier is a locator
 - Remember an HTTP GET returns the “state” of the respective resource at the time of request.
 - In this manner we can think of the web graph that is presented as the result of a GET as a state machine
 - REST: Later in the course

URI: Universal Resource Identifier

- Generic **syntax** for identifiers of resources
- Defined by [RFC 2396](#)
- Syntax: <scheme>:<scheme-specific-part>
 - ftp://ftp.is.co.za/rfc/rfc1808.txt
 - http://www.ietf.org/rfc/rfc2396.txt
 - mailto:John.Doe@example.com
 - urn:oasis:names:specification:docbook:dtd:xml:4.1.2
- Hierarchically-organized, components in order of decreasing significance



Mixing identifier syntax and semantics: Opaque versus Identifiers with Meaning

- DOI:10.1045/3451/13x.4
- <http://store.apple.com/1-800-MY-APPLE/WebObjects/AppleStore>
- Should identifiers carry semantics?
 - People like semantic identifiers
 - Semantic Drift can be a problem
 - Words and names change meaning over time
 - Semantics can compromise persistence
 - Organizations/People/Concepts change over time
 - Semantics is culturally laden

Varieties of semantics

- Opaque
 - Nothing can be inferred, including sequence
 - Cannot be reverse-engineered (feature or bug?)
- Low-resolution date semantics
 - LCCN 99-087253
- Encoded semantics
 - ISBN 1-58080-046-7
 - Country codes... agency codes... checksums...
- Sequential Semantics
 - OCLC numbers
- Name/Word Semantics
 - Work Name/Chapter Name/Section Name

URI Schemes (as of 2005 06 03)

<http://www.iana.org/assignments/uri-schemes>

ftp	File Transfer Protocol	modem	modem
http	Hypertext Transfer Protocol	ldap	Lightweight Directory Access Protocol
gopher	The Gopher Protocol	https	Hypertext Transfer Protocol Secure
mailto	Electronic mail address	soap.beep	soap.beep
news	USENET news	soap.beeps	soap.beeps
nntp	USENET news using NNTP access	xmlrpc.beep	xmlrpc.beeps
telnet	Reference to interactive sessions	xmlrpc.beeps	xmlrpc.beeps
waits	Wide Area Information	urn	Uniform Resource Names
prospero	Prospero Directory	go	go
z39.50s	Z39.50	h323	H.323
z39.50r	Z39.50 Retrieval	ipp	Internet Printing Protocol
cid	content identifier	tftp	Trivial File Transfer Protocol
mid	message identifier	mupdate	Mailbox Update (MUPDATE) Protocol
vemmi	versatile multimedia	pres	Presence
Interfaceservice	service location	im	Instant Messaging
imap	internet message access protocol	mtqp	Message Tracking Query Protocol
nfs	network file system protocol	iris.beep	iris.beep
acap	application configuration access	dict	dictionary service protocol
protocoltsp	real time streaming protocol	snmp	Simple Network Management Protocol
tip	Transaction Internet Protocol	crld	TV-Anytime Content Reference Identifier
pop	Post Office Protocol v3	tag	tag
data	data		
dav	dav		
opaquelocktoken opaquelocktoken		Reserved URI Scheme Names:	
sip	session initiation protocol	afs	Andrew File System global file names
sips	secure session intitiacion protocol	tn3270	Interactive 3270 emulation sessions
tel	telephone	mailserver	Access to data available from mail servers
fax	fax		

Why is RFC 2396 so big?

- Character encodings
- Escaping Characters
- Partial and relative URIs
 - e.g. chap2/start.html, /top/next/part.html, #head1
 - Algorithms for establishing base URL and attaching relative reference to it
- URI Equivalence

Mixing Identifiers with Resolution

URL: Universal Resource Locator

- Deprecated term but still in common use
- String representation of the location for a resource that is available via the Internet
- Use URI syntax
- Scheme has function of defining the access (protocol) method. Used by client to determine the protocol to "speak".
 - `http://an.org/index.html` - open socket to an.org on port 80 and issue a GET for index.html
 - `ftp://an.org/index.html` - open socket to an.org on port 21, open ftp session, issue ftp get for index.html....

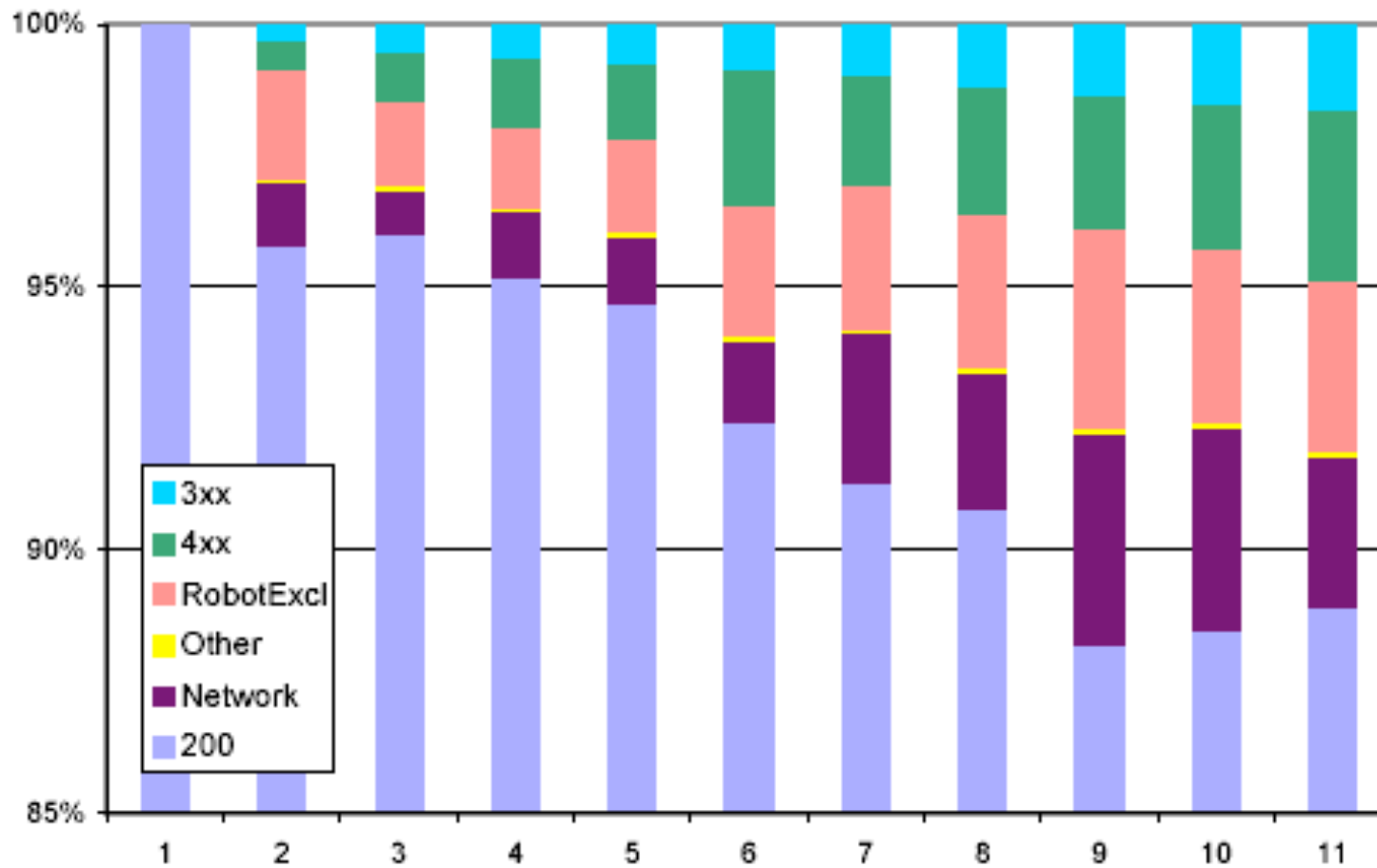
Identification vs. Location Again (URI fragments)

- Different resources:
 - <http://blatz.org/grotz>
 - <http://blatz.org/grotz#remblat>
- HTTP treats them the same
 - Strips off "#remblat"
 - User agent processes fragment

UR(I)L Issues

- **Persistence**
 - “link rot”
- Location dependence
- Valid only at the item level
 - What about works, expressions, manifestations
- Multiple resolution
 - “get the one that is cheapest, most reliable, most recent, most appropriate for my hardware, etc.”
- Non-digital resources?
- How about identifying representations?

Link-rot



crawls ran consecutively, starting on 5 Dec. 2002 and ending on 12 Feb. 2003

<http://www2003.org/cdrom/papers/refereed/p097/P97%20sources/p97-fetterly.html>

The identifier persistence myth

“No scheme or syntax guarantees persistence of any kind”

John Kunze, California Digital Library

URI's - The Web Gurus View

Henry Thompson W3C

- The web works because you can
 - View source
 - Follow your nose
 - Write URIs on the side of a bus
 - Use generic tools
 - Redirect, cache and proxy
- The Web is hands-down the most successful distributed name-based system the world has yet seen
 - Hmm... Postal addresses, phone #'s?
- Ergo anyone designing a persistent identifier system should start from the assumption that http URIs are sufficient for their **technology** needs.
 - Remember there are non-technology issues that need to be dealt with otherwise

Cool URIs don't change

Tim Berners-Lee 1998

<http://www.w3.org/Provider/Style/URI>

What makes a cool URI?

A cool URI is one which does not change.

What sorts of URI change?

URIs don't change: people change them

Other community/application specific "persistent" identifier mechanisms

- Digital Object Identifier(DOI)
- Technology and social infrastructure for naming
- Established by publishers for persistent naming of entities (articles, journals, conference proceedings)
- Cognizant of FRBR elements
- Underlying technology is handle system
 - Resolution server
 - Governance mechanism to establish "persistent"
 - Multiple resolution
 - Registration/mechanism has metadata associated with it
- Used in Crossref - citation linking
 - <http://www.nature.com/nature/journal/v451/n7178/full/nature06496.html>

Other community identifiers

- Astrophysics Data Service (ADS) bibcode
 - http://adsdoc.harvard.edu/abs_doc/bibcodes_help.html
 - <http://adswww.harvard.edu/>
 - Useful for linking among multiple sources of information in a reliable manner
- PubMed Identifier (PMID)
 - unique number assigned to each PubMed citation of life sciences and biomedical scientific journal articles.

Why haven't URNs caught on beyond certain communities?

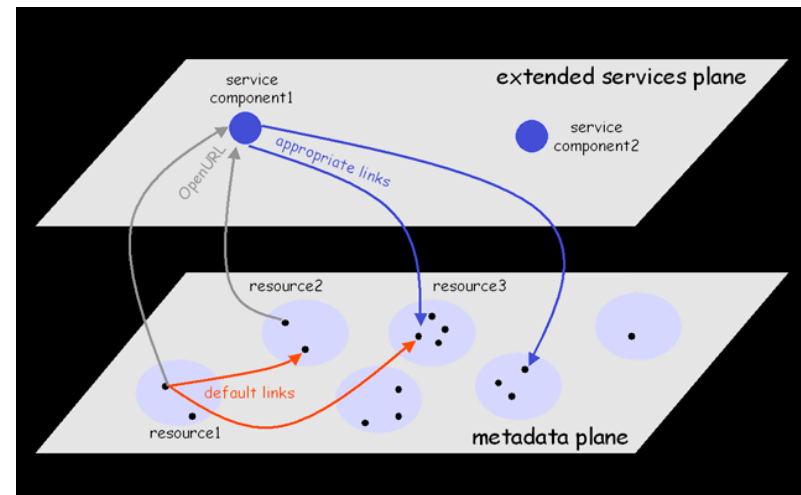
- Complexity of systems
- One size does not fit all - special purpose URN schemes have been successful, e.g., PubMed ID, Astrophysics BibCode
- No guarantee of persistence - longevity is an organizational not technical issue
- Requires well-regulated administrative systems
- Absence of "killer" applications - although reference linking is emerging

Conclusions

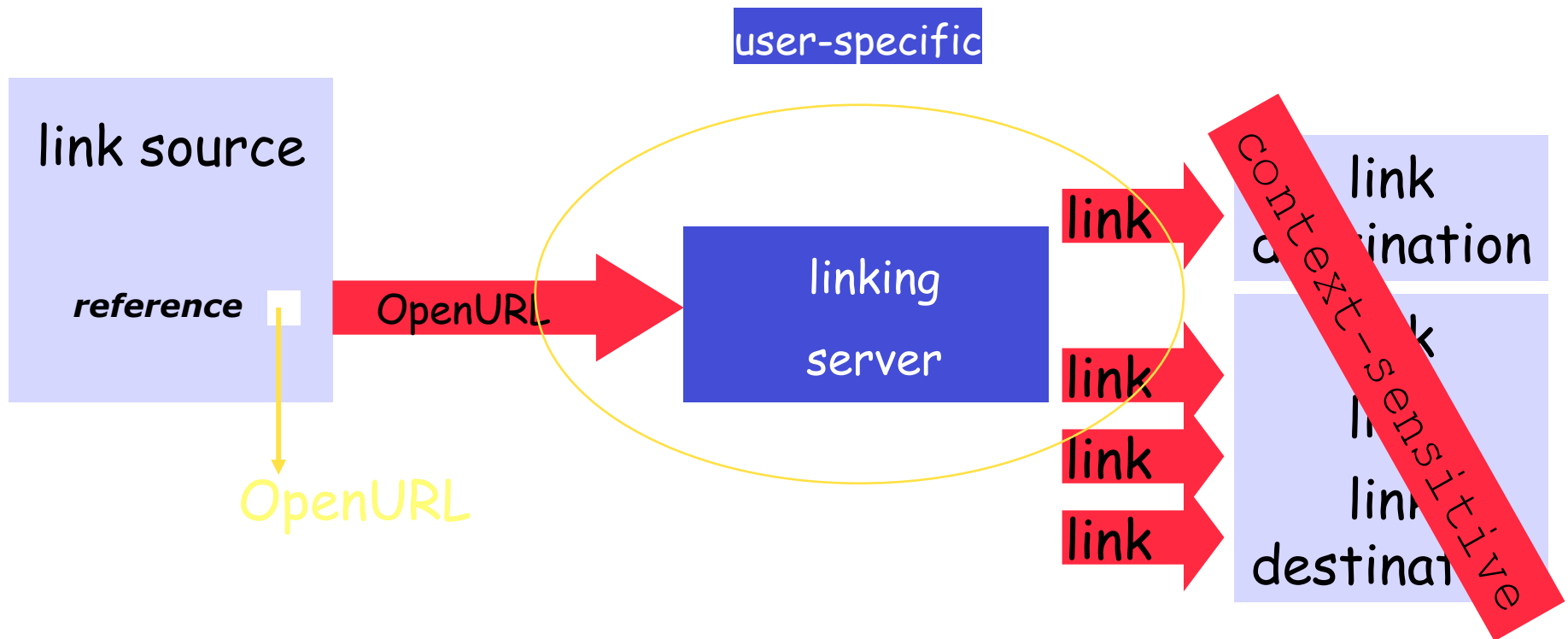
- There is no established "answer" the identification problem
 - Lots of identify wars
 - Turf protecting
- In reality there are different needs with different appropriate solutions
- URIs do work as an appropriate technological solution and must always be considered.

openURL: Making links context sensitive

- Why?
 - "Appropriate item" differs for each user
 - Licensing locality
 - Some users may want a choice (abstract, full text, etc.)
- Conceptualize link as service rather than object targeted.
- OpenURL
 - Transports metadata about the work to...
 - A localized service that interprets the metadata and provides contextualized choices to the user.



OpenURL linking

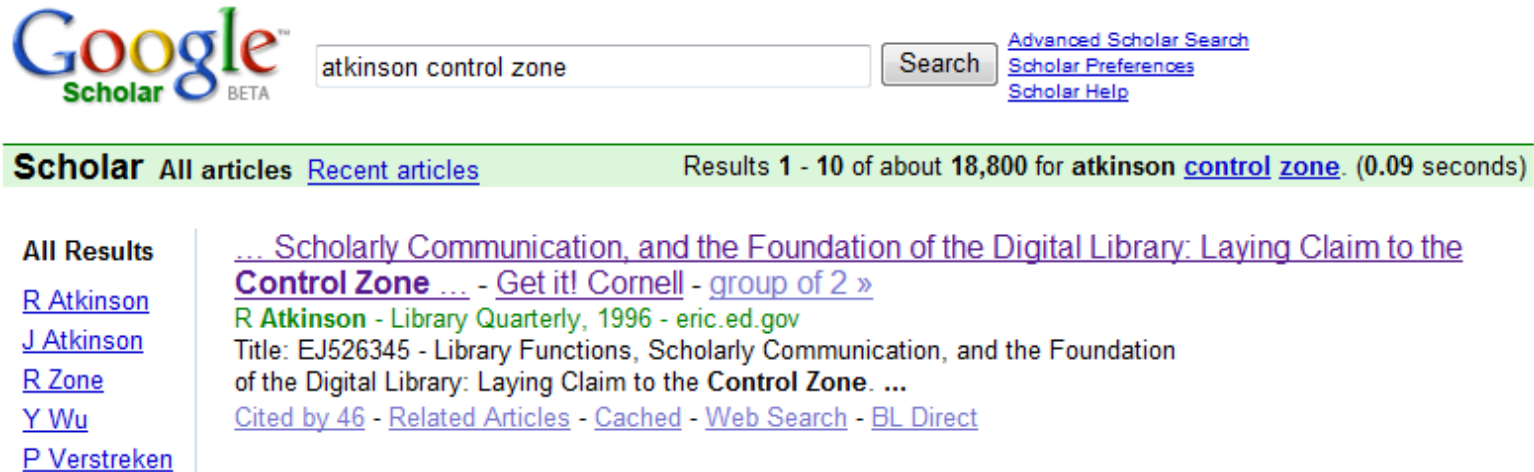


Components of an OpenURL

- Base-URL - Service component that accepts the openURL
- Object Description - Identifying information about an object (e.g., the identifier of a resource, metadata about the resource)
- Origin Description - Identifying information about origin of request.

<http://www.ukoln.ac.uk/distributed-systems/openurl/>

Google Scholar and OpenURL



The screenshot shows the Google Scholar search interface. The search bar contains the text "atkinson control zone" and a "Search" button. To the right of the search bar are links for "Advanced Scholar Search", "Scholar Preferences", and "Scholar Help". Below the search bar, a green banner displays "Scholar All articles Recent articles Results 1 - 10 of about 18,800 for atkinson control zone. (0.09 seconds)". On the left side, there is a list of authors: "All Results", "R Atkinson", "J Atkinson", "R Zone", "Y Wu", and "P Verstreken". The main content area shows a search result for "Control Zone" by R Atkinson, published in "Library Quarterly" in 1996. The title is "EJ526345 - Library Functions, Scholarly Communication, and the Foundation of the Digital Library: Laying Claim to the Control Zone. ...". Below the title are links for "Cited by 46", "Related Articles", "Cached", "Web Search", and "BL Direct".

<http://scholar.google.com/scholar?hl=en&lr=&q=atkinson+control+zone>