Architecture of the World Wide Web Web Information Systems

CS/INFO 431

January 31, 2007 Carl Lagoze - Spring 2007



Acknowledgments

• Erik Wilde – UC Berkeley

– http://dret.net/lectures/infosys-ws06/http



Information Management: A Proposal

Tim Berners-Lee, CERN March 1989, May 1990

http://www.w3.org/History/1989/proposal.html



Web as a graph

We can call the circles nodes, and the arrows links. Suppose each node is like a small note, summary article, or comment. I'm not over concerned here with whether it has text or graphics or both. Ideally, it represents or describes one particular person or object. Examples of nodes can be

- People
- Software modules
- Groups of people
- Projects
- Concepts
- Documents
- Types of hardware
- Specific hardware objects

Web as a graph

The arrows which links circle A to circle B can mean, for example, that A...

- depends on B
- is part of B
- made B
- refers to B
- uses B
- is an example of B

W3C*

Architecture of the World Wide Web, Volume One

W3C Recommendation 15 December 2004

This version: <u>http://www.w3.org/TR/2004/REC-webarch-20041215/</u> Latest version: <u>http://www.w3.org/TR/webarch/</u> Previous version: <u>http://www.w3.org/TR/2004/PR-webarch-20041105/</u>	Make Note: Three URLs for the "same" Intellectual object
Editors:	

lan Jacobs, W3C Norman Walsh, Sun Microsystems, Inc.

http://www.w3.org/TR/webarch/

Naïve view of web graph



Think for a second...

- When I access google.com on my cell phone it looks different than on my desktop
- When I access google.com from Paris it looks different than when I access it from Ithaca

Architectural Components of the Web

- Participants
 - Servers
 - Web Agents
 - General agents such as robots e.g., google crawler
 - specialized as User Agents
- Identification
 - URIs to identify Resources
 - Some URIs resolve URLs
 - Some do not resolve INFO URIs
- Interaction
 - Standardized protocols with exchange of messages
 - One example HTTP
 - Requests result in return of Representations
- Formats
 - Representation is in form of sequence of bytes with a media type that provides hints to the agent about processing
 - One example of a format is a MIME type

A resource is...

- An entity that has an identity (a URI)
 - Some resources are digital URI <-> URL
 - Some are non-digital people, institutions, etc.
- Abstract you can't examine/touch/see a resource
 - Information hiding
- A service point for initiating protocol (HTTP) actions
- Target of links (you make hyperlinks to a URI) —

A representation is...

- The result of applying a service request upon a resource
- What the server determines to be the state of the resource
 - Parameterized
 - Time, space
 - Request parameters
 - Many to one mapping from resource to representation
- A package:
 - Metadata about request, server actions, agent
 - Data the "content"
- The entity that is processed by a web agent
 - In the case of a browser (user agent) rendered and displayed
 - Note that many agents such as crawlers make extensive use of metadata (lastmodified)
- The entity that is the source of links
 -

Use cases of representations in increasing complexity

- Static transcription of a file
 - foo.html is disseminated from http://my.org/foo.html
- Dynamic dissemination of data based on static file
 - Translation from base jpeg to thumbnail
 - Result of PHP program
- Multiple representations from resource based on user agent or other parameters
 - google for cell phones and desktops
 - Content negotiation
- Totally time dependent representations
 - http://cnn.com

Real nature of web graph



Think about it...

The web graph is completely ephemeral: Based on representations that are time and agent context

Closer look at Resource/Representation Relationship



Content Negotiation



HTTP (Hypertext Transfer Protocol)

- HTTP 1.1 RFC 2516 <u>ftp://ftp.isi.edu/in-notes/rfc2616.txt</u>
- Basis of interaction between web agents and servers
- Layered on top of TCP
- Uses DNS
- Text-based
 - All messages and requests are human readable
- Stateless
 - No persistent client/server connection
 - All state carried in protocol request/reply (cookies)



HTTP Request Types

- GET initiate a retrieval action based on the URI
- POST transmit information package to URI at server
- HEAD same as GET but return only header (metadata) information
- PUT store the request content at the URI
- **DELETE** Remove resource at URI

HTTP Example



HTTP Request

- Start line
 - Consists of method, path, version
 - GET index.html HTTP/1.1
 - Valid methods include:
 - GET, POST, HEAD, PUT, DELETE
- Headers
 - HTTP/1.1 requires a Host: header Host: cs.cornell.edu
 - Many other headers
- Optional body content

HTTP Response

- Start line
 - consists of HTTP version, status code, and description
 - HTTP/1.1 200 OK
 - HTTP/1.1 404 Not Found
- Headers (Many header definitions) Content-type: text/html
- Content

HTTP Response Codes

- Response coded by first digit
 - 1xx: informational, request received
 - 2xx: success, request accepted
 - 3xx: redirection
 - 4xx: client error
 - 5xx: server error

Simple HTTP GET - Response

GET /path/file.html HTTP/1.1 Host: cs.cornell.edu User-Agent: Mozilla/3.0 [Blank line]

HTTP/1.1 200 OK Content-Type: text/html Date: Wed, 31 Jan 2007 14:58:57 GMT Content-Length: 1354

<html> <head>

•••

Simple HTTP Post

POST /path/script.php HTTP/1.1 Host: cs.cornell.edu User-Agent: Mozilla/3.0 Content-Type: application/x-www-form-urlencoded Content-Length: 32

home=Cosby&favorite+flavor=flies [Blank line]

HTTP Content Negotiation (serverside)

- Resources may have multiple dimensions
- General idea
 - Web agent makes HTTP requests stating constraints
 - Using constraints server decides on "best" representation
- HTTP defined constraints are language, encoding, format, character encoding
 - Accept, Accept-Charset, Accept-Encoding, Accept-Language
- Server may also use:
 - User-agent: device specificity
 - Host: localization
 - Cookies

Header request examples for content negotiation

Other types of content negotiation

- Client-side
 - Server response with list of different representations
 - Client (or user) makes a choice
- Transparent
 - Cache plays a role in client-side negotiation