# Establishing New Levels of Interoperability
# for
# Web-Based Scholarship



Cartoon by:
Patrick Hochstenbach

Herbert Van de Sompel
Los Alamos National Laboratory
@hvdsomp

# Reminiscing About 15 Years of Interoperability Efforts

Herbert Van de Sompel
Los Alamos National Laboratory
herbertv@lanl.gov

Michael L. Nelson
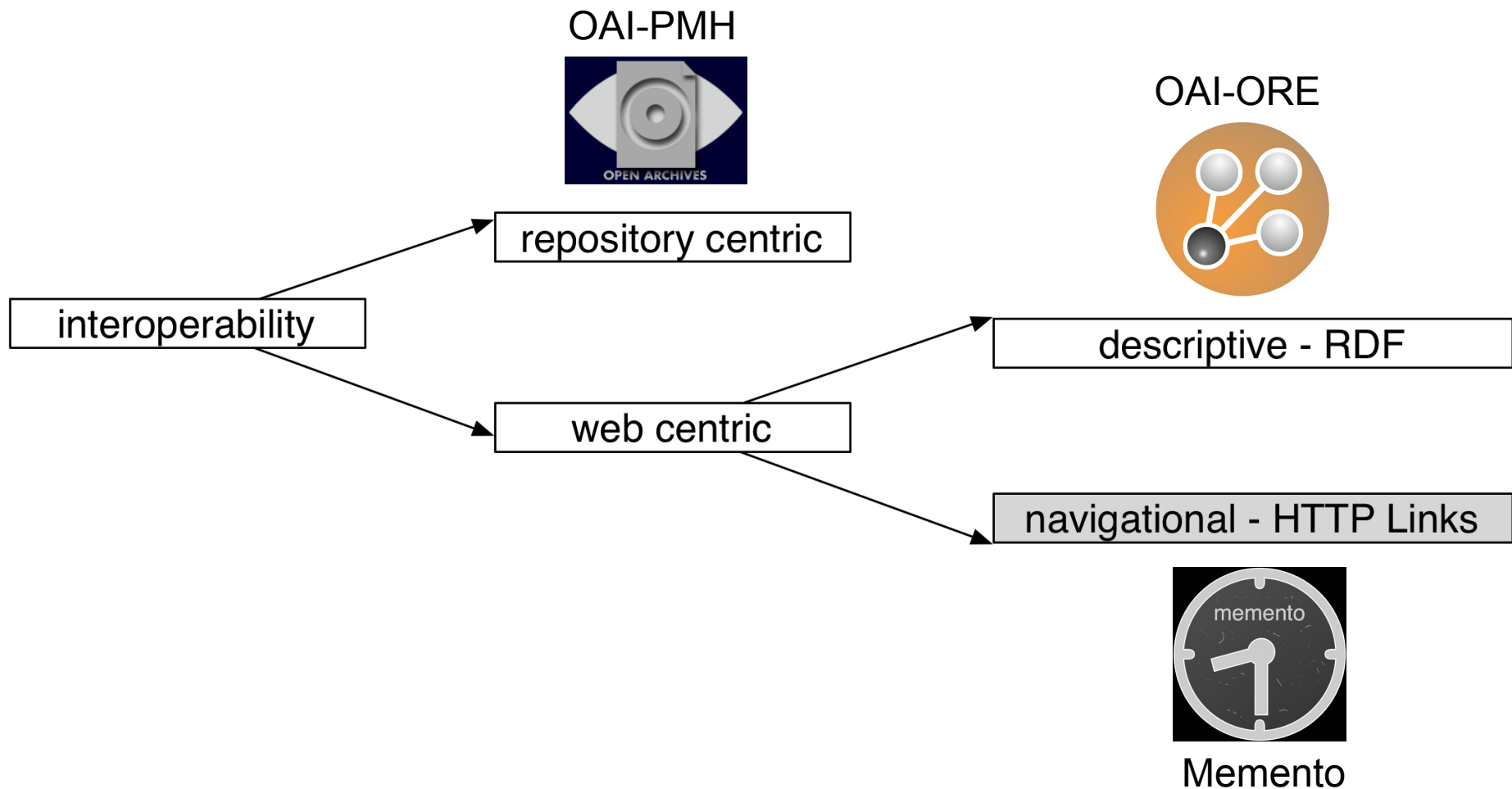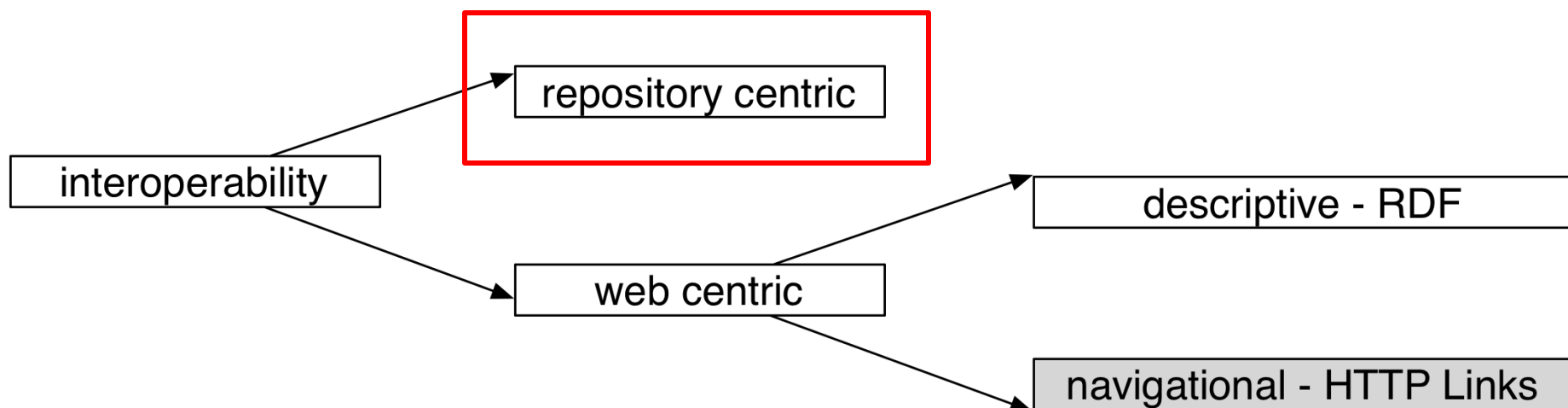Old Dominion University
mln@cs.odu.edu

# Research Communication & Research Process on the Web

- A highly distributed activity

- Turning this distributed activity from a gathering of silo-ed nodes into an ecology of collaborating nodes, requires establishing interoperability
  - In the web context, this seems like a rather unique challenge: Most web enterprises do not want interoperability they want dominance, monopoly

- To a large extent, interoperability <u>across</u> this distributed activity remains restricted to persistent identification of communicated objects and contributors
  - Which results in added-value – services can be created

# Evolution of Our Thinking about Interoperability

OAI-PMH

OAI-ORE

repository centric

interoperability

descriptive - RDF

web centric

navigational - HTTP Links

Memento

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

# Evolution of Our Thinking about Interoperability

interoperability

repository centric

web centric

descriptive - RDF

navigational - HTTP Links

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

1999

- OAI was a heroic effort to fundamentally transform scholarly communication
  - By promoting communication via preprints, non-peer-reviewed papers

- The OAI took a technical approach to achieve the goal
  - Make preprints easier to discover, access

The OAMH protocol is a low-barrier interoperability specification for the recurrent exchange of metadata between systems

# the Metadata Harvesting protocol

service provider

data provider

harvester

repository

**6** Requests

Replies

# Those Were the Days

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

## 3.1.1.1 Encoding an OAI-PMH request in a URL for an HTTP GET

Don't trust HTTP

## 3.6 Error and Exception Conditions

In event of an error or exception condition, repositories **must** indicate OAI-PMH errors, distinguished from HTTP Status-Codes, by including one or more error elements in the response. While one

```
http://an.oa.org/OAI-script?
verb=GetRecord&identifier=oai:arXiv.org:hep-
th/9901001&metadataPrefix=oai_dc
```

HTTP GET with GetRecord verb

A repository replies to a request with an *incomplete list* and a `resumptionToken;`

An HTTP link

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

# Repository-Centric Interoperability Paradigm

Address interoperability challenges from the perspective of a <u>node</u>, e.g. an IR, a publisher, a web-based authoring portal, a software repository, …

- **The node at the center of the universe**

- Define a machine interface for your node, expect others to use it

- Piggybacking on the web without truly embracing its core technologies

- The node resembles a brick & mortar library that can be visited subject to well-intended yet idiosyncratic policies – the interface

# Launching the OAI -  Luce, Van de Sompel, Ginsparg (1999)



Repositories still use OAI-PMH, created in the olden days when I looked like this

# Evolution of Our Thinking about Interoperability

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

# Web-Centric, Resource-Centric Interoperability Paradigm

Address interoperability challenges from the perspective of the web

- **The resource at the center of the universe**
  - The notion of a node, a repository, not even of a web server exists in the architecture of the web

- The tools of the interoperability trade are the primitives of the web

# Tools of the Web-Centric Interoperability Trade

- Resource
- URI
- HTTP as the API: HEAD/GET, POST, PUT, DELETE
- Representation
- Media Type
- Link
- Content Negotiation

<span style="color:red">**W3C Architecture of the World Wide Web**</span>

# Evolution of Our Thinking about Interoperability

repository centric

interoperability

web centric

descriptive - RDF

navigational - HTTP Links

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

- OAI-ORE observation: Scholarly assets are rapidly becoming *compound,* consisting of multiple resources with various:
  - Relationships
  - Interdependencies

- How to convey this compound-ness in an interoperable manner so that applications can access, consume such assets?



2006

**Aiming for New Levels of Cross-Repository Functionality**
TICER, Digital Libraries a la Carte, Tilburg, The Netherlands, August 22 2006
Herbert Van de Sompel

RESEARCH
LIBRARY
Los Alamos
NATIONAL LABORATORY

# ORE Will Allow Web Crawlers to Unambiguously Recover CDO Structure from the Web Graph

# Express the `ore:describes` relationship

# Tools of the Web-Centric Interoperability Trade – RDF Stack

- Resource
- URI
- HTTP as the API
- Representation
- Media Type
- Link
- Content Negotiation, e.g. for preferred Media Type
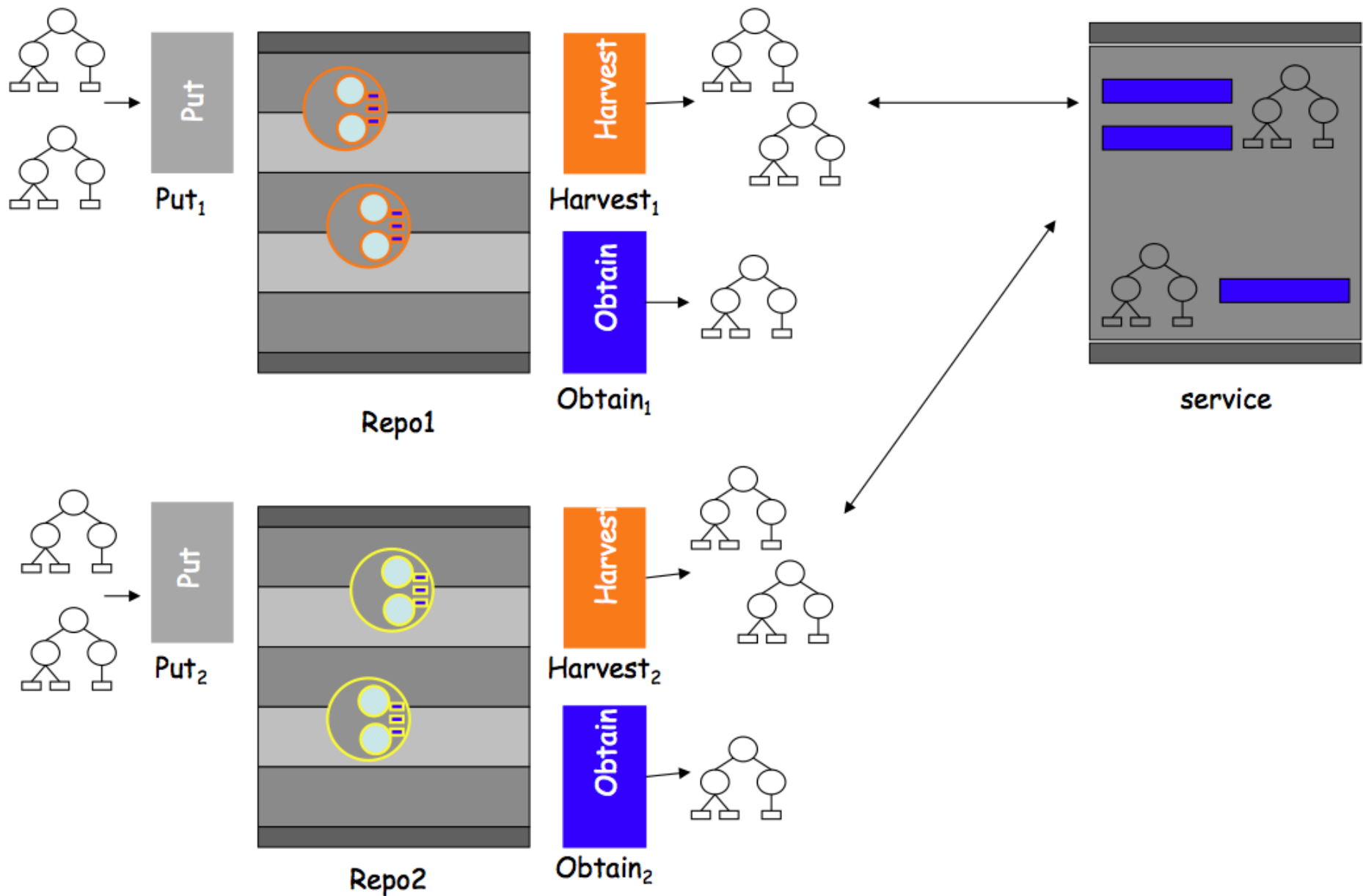
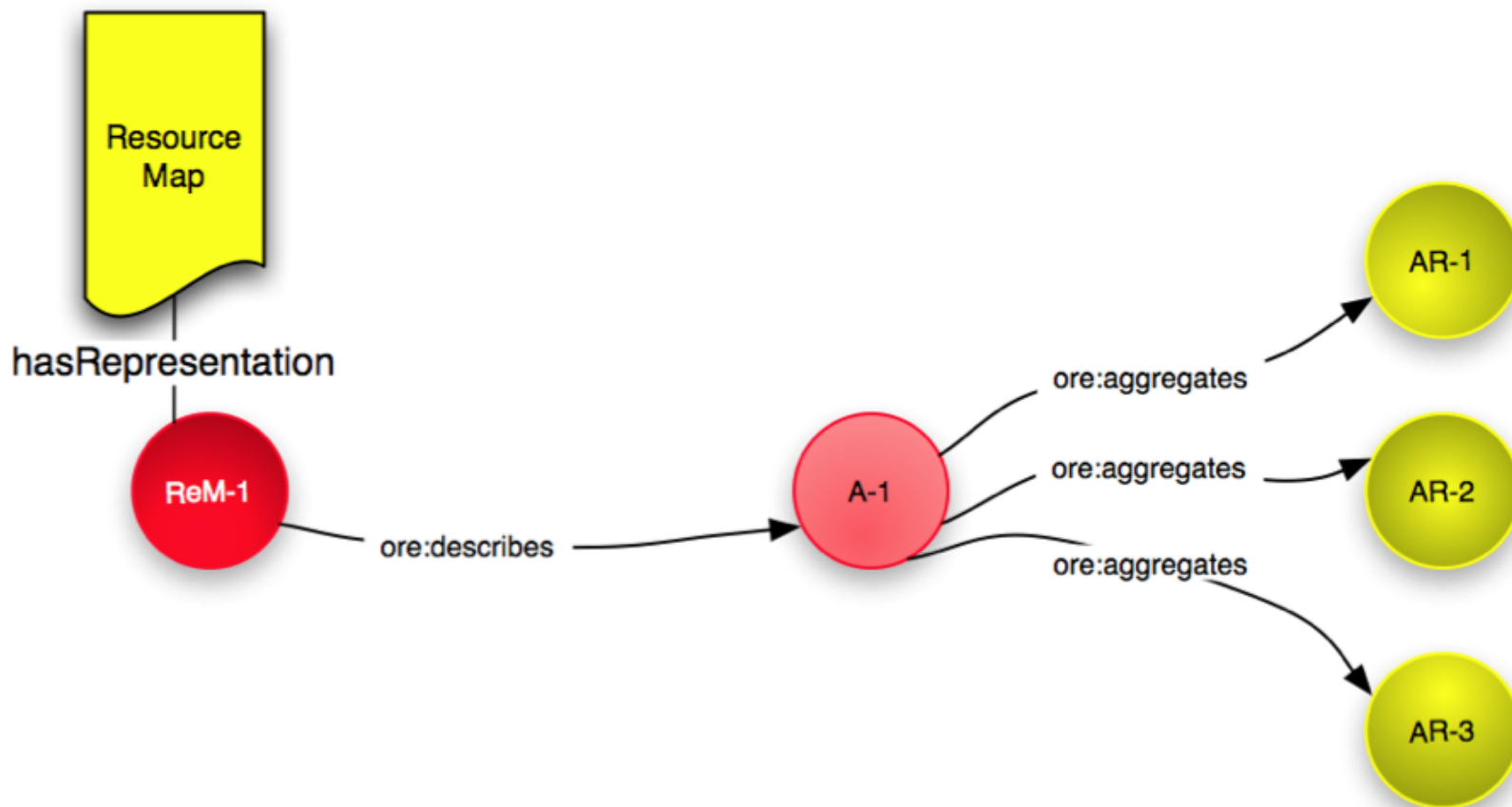W3C Architecture of the World Wide Web

- Typed Link
- Controlled Vocabularies for Typed Links

RDF, RDFS, OWL

# Interoperability via RDF, RDFS, OWL Stack

Used by various interoperability efforts, e.g. OAI-ORE, Open Annotation, W3C PROV, Research Objects, …

- Provides <u>extensive expressiveness for description</u>
- Typically based on publishing documents that adhere to a certain "profile" and reveal relations, properties, …
- Non-Trivial barrier to entry as illustrated by slow adoption, likely related to unfamiliar technology stack

# Evolution of Our Thinking about Interoperability

Herbert Van de Sompel

- Memento is about the Web and time:
    - Resources evolve over time
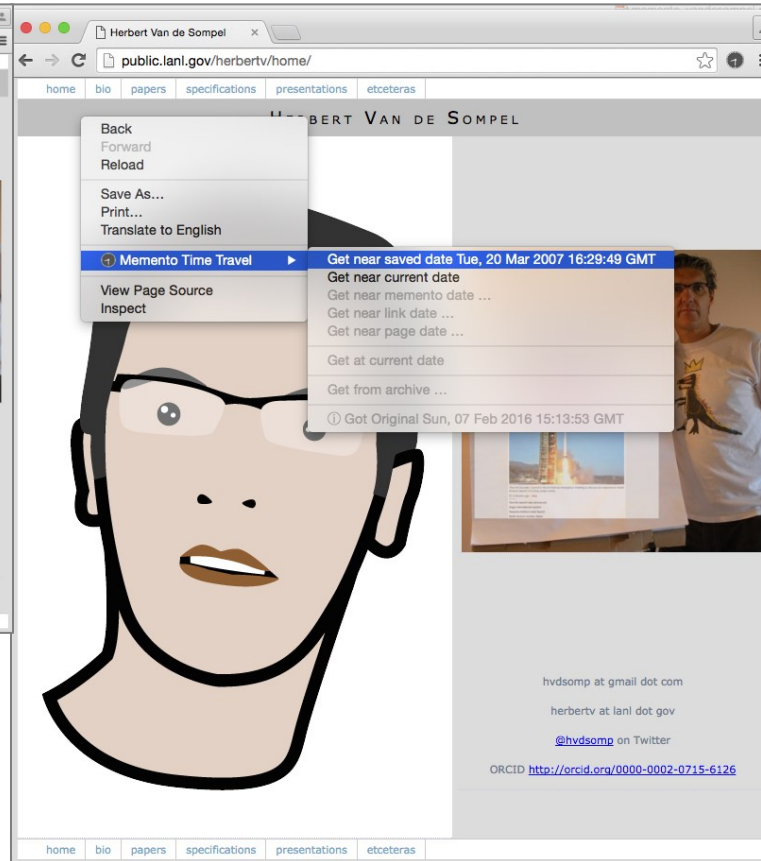    - Only the current resource version is available from a resource's URI
    - How to seamlessly access prior versions, if they exist, using the resource's URI and a version datetime

- Memento looks at this problem for the Web, in general:
    - Time-Based access to resource versions across web archives, resource versioning systems

2009

Today  Select Date Mar 20 2007  Apr 03 2007

From Internet Archive

Memento for Chrome at http://bit.ly/memento-for-chrome

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

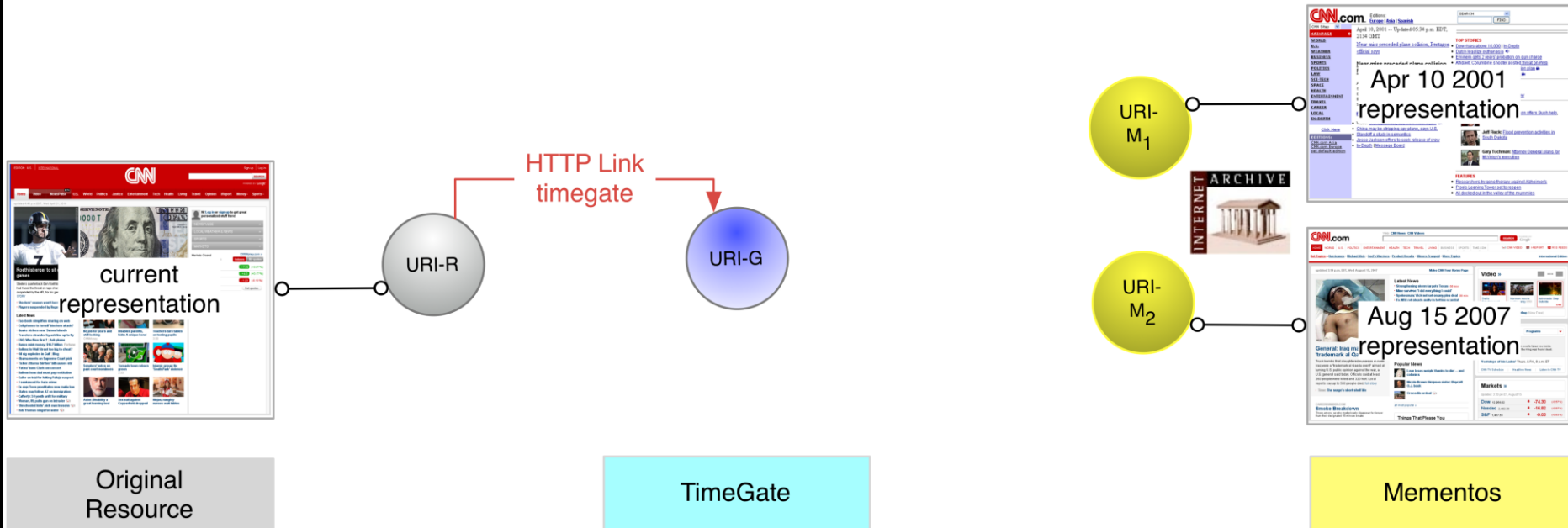# Original Resource and Mementos



current
representation

URI-R

Original
Resource

URI-M₁

Apr 10 2001
representation

ARCHIVE
INTERNET

URI-M₂

Aug 15 2007
representation

Mementos

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

# Bridge from Present to Past



**HTTP Link timegate**

URI-R

URI-G

URI-M₁ — Apr 10 2001 representation

URI-M₂ — Aug 15 2007 representation

current representation

Original Resource

TimeGate

Mementos

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

Los Alamos
NATIONAL LABORATORY

# Bridge from Present to Past



HTTP Link
timegate

Datetime Negotiation
Accept-Datetime = $T_i$

Datetime Negotiation
Accept-Datetime = $T_j$

URI-R

URI-G

URI-M$_1$

URI-M$_2$

current
representation

Apr 10 2001
representation

Aug 15 2007
representation

Original
Resource

TimeGate

Mementos

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

Los Alamos
NATIONAL LABORATORY

# Bridge from Past to Present



HTTP Link
original

URI-M₁

Apr 10 2001
representation

URI-R

current
representation

URI-M₂

Aug 15 2007
representation

HTTP Link
original

Original
Resource

Mementos

Los Alamos
NATIONAL LABORATORY

# Tools of the Web-Centric Interoperability Trade – HTTP Stack

- Resource
- URI
- HTTP as the API
- Representation
- Media Types
- Link
- Content Negotiation, e.g. for Media Type, Time

- Typed Link
- Controlled Vocabularies for Typed Links

W3C Architecture of the World Wide Web

HTTP Links, IANA link relation registry, community link relation types

HATEOAS – Hypermedia As The Engine Of Application State

http://en.wikipedia.org/wiki/HATEOAS

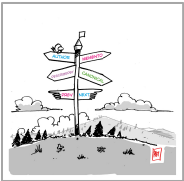# Interoperability via HTTP Links, IANA Link Relation Types

Used by Memento, ResourceSync, Signposting the Scholarly Web:

- Provides <u>coarse expressiveness for navigation</u> via IANA registered relation types (expressed as reserved terms)
  - Finer grained expressiveness via community-defined relation types (expressed as HTTP URIs)
- Typically based on publishing typed links that support a client to navigate among resources in an informed manner
- Low implementation barrier because of familiar technology stack

# Establishing New Levels of Interoperability: Examples

 ResourceSync

 Signposting the Scholarly Web

 Robust Links

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

# ResourceSync

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

# Anurag Acharya Told Us Why We Need ResourceSync

## What does indexing need?

- List of all article urls

  *Web search Scholar*

- Ability to fetch article urls

- What we index is what the user sees

- Identify scholarly articles

  *Scholar*

- Determine article metadata

Anurag Acharya. Indexing Repositories: Pitfalls & Best Practices. Open Repositories 2015 Keynote.
http://www.or2015.net/wp-content/uploads/2015/06/or-2015-anurag-google-scholar.pdf

# Anurag Acharya Told Us Why We Need ResourceSync

## List of articles - IV

- Best practice: Year-month browse
  - Linked from homepage - EPrints
  - Helps crawlers as well as users

- Best practice: Article sitemap
  - Include urls for ALL articles
  - Linked from robots.txt or homepage
  - DSpace if sitemaps are enabled

Anurag Acharya. Indexing Repositories: Pitfalls & Best Practices. Open Repositories 2015 Keynote. http://www.or2015.net/wp-content/uploads/2015/06/or-2015-anurag-google-scholar.pdf

# ResourceSync is Based on Sitemaps

- Sitemap is the document format used throughout the framework
  - Used widely by web servers to advertise their resources to search engines

```
<urlset xmlns="http://www.sitemaps.org/schemas/sitemap/0.9">

  <url>
    <loc>http://example.com/res1</loc>
    <lastmod>2013-01-02T13:00:00Z</lastmod>
  </url>

  <url>
    <loc>http://example.com/res2</loc>
    <lastmod>2013-01-02T14:00:00Z</lastmod>
  </url>
  …
</urlset>
```
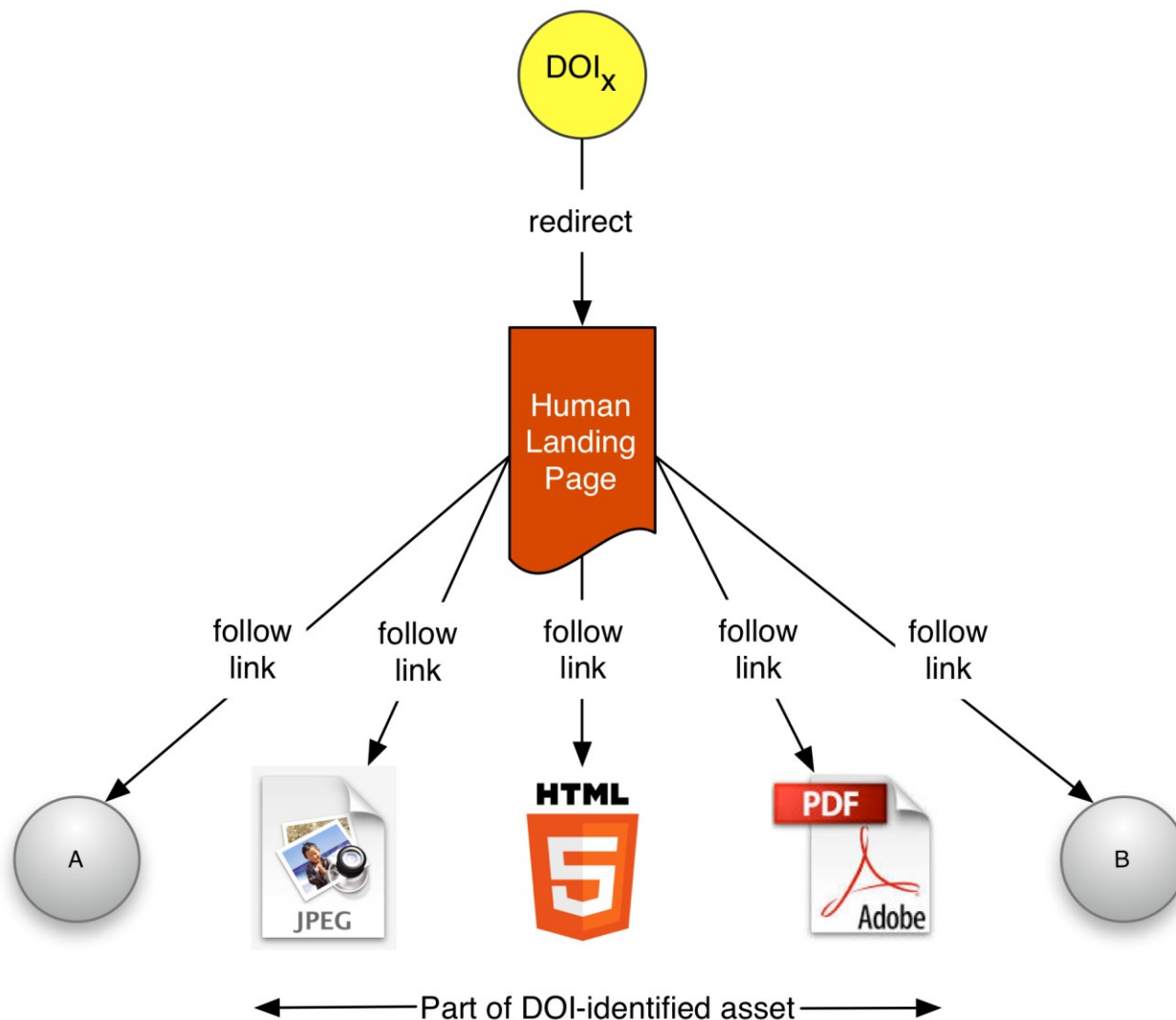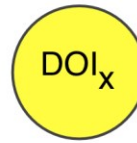
# ResourceSync, ANSI/NISO Z39.99-2014

Open Archives Initiative ResourceSync Framework Specification

**ResourceSync Framework Specification (ANSI/NISO Z39.99-2014)**

21 April 2014

**This version:**
    http://www.openarchives.org/rs/1.0/resourcesync
**Latest version:**
    http://www.openarchives.org/rs/resourcesync
**Previous version:**
    http://www.openarchives.org/rs/0.9.1/resourcesync

## Abstract

This ResourceSync specification describes a synchronization framework for the web consisting of various capabilities that allow third-party systems to remain synchronized with a server's evolving resources. The capabilities may be combined in a modular manner to meet local or community requirements. This specification also describes how a server should advertise the synchronization capabilities it supports and how third-party systems may discover this information. The specification repurposes the document formats defined by the Sitemap protocol and introduces extensions for them.

- Synchronization of resources from a Source to Destinations
  - Includes exposing repository content to aggregators, search engines

- Applies to any resource with an HTTP URI

- Leverages key ingredients of web interoperability, follow your nose, existing Search Engine Optimization practice

http://www.openarchives.org/rs/toc

# Publish Inventory, Changes, Notifications

- Repository communicates about the state of its resources:

  - <u>Publish inventory</u>: snapshot of the state of resources at a moment in time

  - <u>Publish changes</u>: enumeration of resource changes that occurred during a temporal interval

  - <u>Notify about changes</u>: send notifications as changes occur

# Payload for Inventory, Changes, Notifications

- A repository may communicate additional information pertaining to each resource:

  - <u>Technical metadata about a resource</u>: content encoding, content length, mime type, content-based hash

  - <u>Links to related resources</u>: mirror copies, alternate representations, resource versions, diff between current and previous version, metadata-to-content link, content-to-metadata link, collection membership, persistent identifier, etc.

# ResourceSync is Based on Sitemaps

- Extensions to Sitemaps:
  - ○ &lt;rs:ln&gt; for links
  - ○ &lt;rs:md&gt; for metadata

```
<urlset xmlns="http://www.sitemaps.org/schemas/sitemap/0.9">
        xmlns:rs="http://www.openarchives.org/rs/terms/">
  <rs:ln …/>
  <rs:md …/>

  <url>
    <loc>http://example.com/res1</loc>
    <lastmod>2013-01-02T13:00:00Z</lastmod>
    <rs:ln …/>
    <rs:md …/>
  </url>
  …
</urlset>
```

# Signposting the Scholarly Web



## Example pattern: The PID, the Landing Page, the Stuff

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

Los Alamos
NATIONAL LABORATORY

# Response to HTTP HEAD on
## http://dx.doi.org/10.2218/ijdc.v9i1.320

HTTP/1.1 303 See Other
Server: Apache-Coyote/1.1
Date: Fri, 9 Jan 2015 16:31:46 GMT
Vary: Accept
Location: http://www.ijdc.net/index.php/ijdc/article/view/320
Link: <http://www.ijdc.net/index.php/ijdc/article/view/320>
 ; rel=" describedby"
 ; type="text/html"
Content-Length: 188

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016

# This Allows a Machine Agent …

- To understand that the splash page describes the DOI-identified asset

- To determine that resource A is not part of the DOI-identified asset

- To navigate towards the profile of the authors of the asset when landing on any of the constituent resources of the DOI-identified asset

- To understand that a DOI is associated with the PDF, HTML, and JPEG resources and that this DOI should preferably be used to refer to those resources

- To associate annotations made to the HTML page with the DOI

# Signposting: Work in Progress

# Signposting: Work in Progress

Herbert Van de Sompel
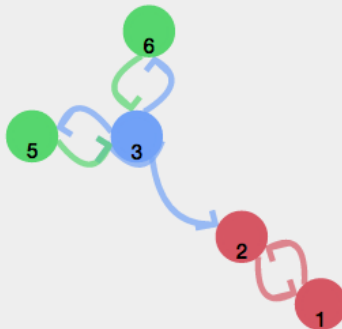COAR Annual Meeting, Vienna, Austria, 12/04/2016

# Signposting: Work in Progress

**Demo**

Input any HTTP URI of a scholarly article, and hit Get Headers to see its corresponding signposting headers.

http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0115253.pdf  ← **URI of PDF file**

Get Headers ⤨

## Signposting Headers for the Landing Page Pattern

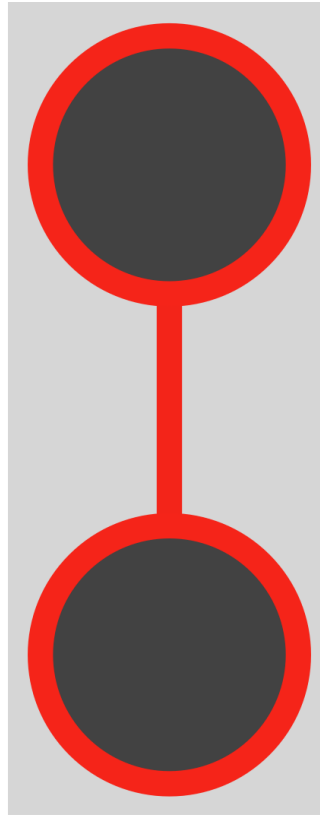| Resource | Link Header |
|---|---|
| **Node: 5**<br>http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0115253.PDF | <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0115253>; rel="collection" |
| **Node: 4**<br>http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0115253 | <http://dx.plos.org/10.1371/journal.pone.0115253>; rel="canonical",<br><http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0115253.PDF>; rel="item"; type="application/pdf",<br><http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0115253>; rel="item"; type="text/html",<br><http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0115253.XML>; rel="item"; type="text/xml" |
| **Node: 6**<br>http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0115253.XML | <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0115253>; rel="collection" |
| **Node: 2**<br>http://dx.plos.org/10.1371/journal.pone.0115253 | <http://dx.doi.org/10.1371/journal.pone.0115253>; rel="pid" |

Los Alamos
NATIONAL LABORATORY

# Robust Links

Herbert Van de Sompel
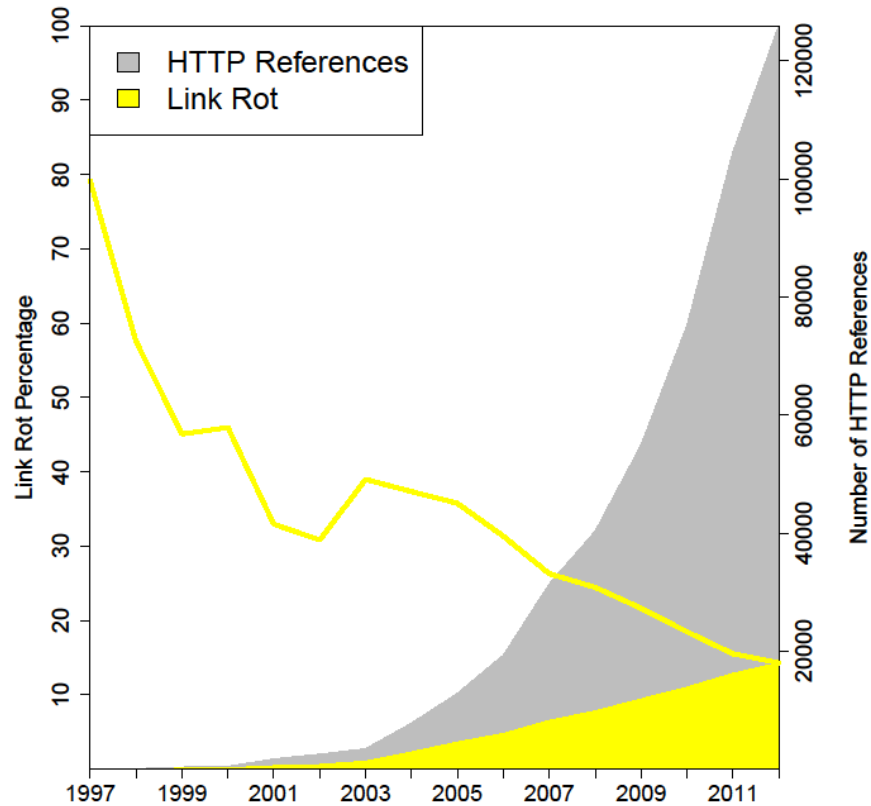COAR Annual Meeting, Vienna, Austria, 12/04/2016
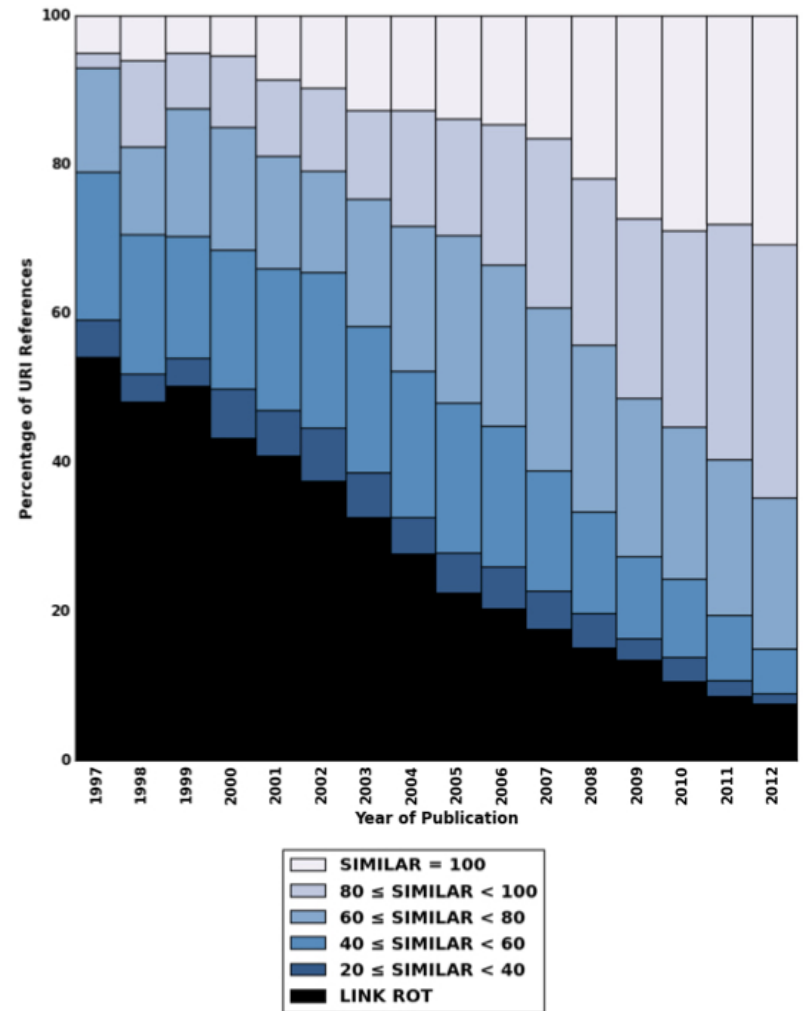
# Reference Rot

- Links to Web at Large resources are subject to <u>Reference Rot</u>:

  - Link Rot: Link stops working, e.g. HTTP 404 "Not Found"

  - Content Drift: Linked content changes over time
    - Possibly to the extent that it becomes no longer representative of the content that was initially referenced



hiberlink

---

# Link Rot

# Content Drift



Martin Klein et al. (2014) Scholarly context not found. In: PLOS ONE
http://dx.doi.org/10.1371/journal.pone.0115253

# Combating Reference Rot

① Create a snapshot of the referenced resource in one of many web archives that support on-demand archiving (manual, API):
- o archive.today
- o Internet Archive
- o perma.cc
- o webcitation.org

② Reference snapshots actionably by using:
- o Original URI
- o Snapshot URI
- o Date/Time of snapshot

in order to maximize link robustness

# Reference Resources Actionably

- When referencing resources, use Link Decorations to convey Original URI, Snapshot URI, Date/Time

```
<a href="http://hiberlink.org"
    data-versionurl="https://archive.is/drFFu"
    data-versiondate="2015-11-16" >
```

```
<a href="https://archive.is/drFFu"
    data-originalurl="http://hiberlink.org"
    data-versiondate="2015-11-16" >
```

- Legitimate in HTML5
- Can be made actionable with JavaScript, e.g. robustlinks.js

Herbert Van de Sompel et al. (2015) Robust Links - Link Decorations
http://robustlinks.mementoweb.org/spec/

# See Robust Links at Work

# See Robust Links at Work

ed" with actionable attributes as per the Robust Links🔗 specification.
snapshots of the referenced resources. These robust links can be
linked content was around the time the original link was put in place.
esults from the Mellon-funded Hiberlink project🔗. It leverages the

### Robust Links

Get near page creation date 2015-11-16

Get near link date 2015-10-06

Get from archive.is

— OAI-PMH (1999)

rm scholarly communicat
d to be working towards
y available peer-reviewe
e and commercial servic

bility as a way to break ground for a universal adoption of e-print
he discoverability of e-prints, actually making them easier to discover
vas OAI-PMH🔗, a protocol for the recurrent exchange of metadata
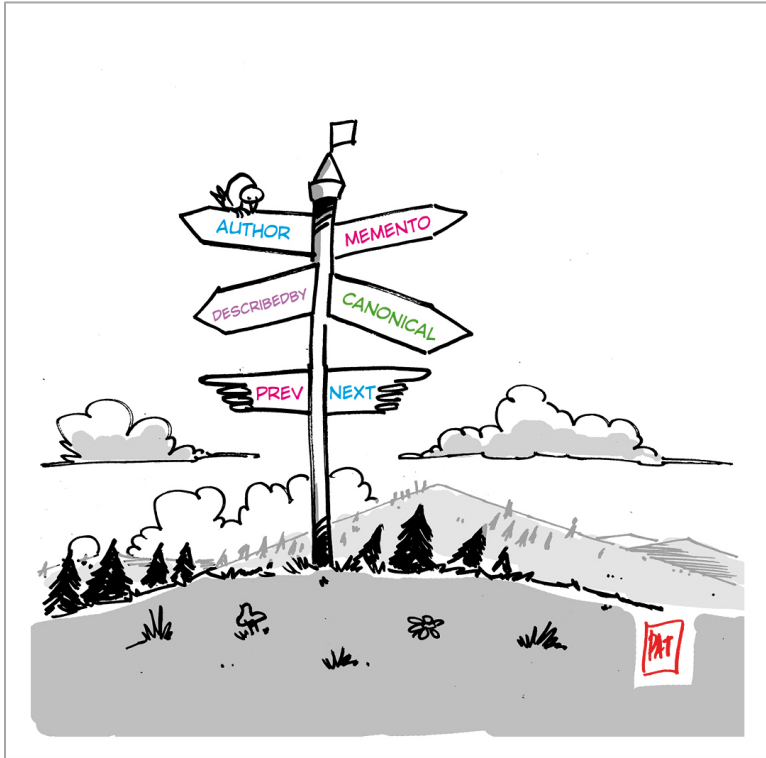which was to an extent inspired by the Dienst protocol🔗.

# Conclusion

There is no real conclusion. There are insights:

- One doesn't do interoperability because of interoperability but to enable cross-node applications that add value

- Establishing interoperability across a vast amount of nodes is a huge challenge. But meaningful levels of interoperability can be achieved via really basic approaches.

- Unfortunately, not even discovery is a solved problem (although the solution is available)
  - Anurag's keynote is a real embarrassment for our community

Leading organizations and projects should promote web-centric cross-repository interoperability

# Establishing New Levels of Interoperability
# for
# Web-Based Scholarship



Cartoon by:
Patrick Hochstenbach

Herbert Van de Sompel
Los Alamos National Laboratory
@hvdsomp

Herbert Van de Sompel
COAR Annual Meeting, Vienna, Austria, 12/04/2016