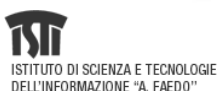


Approaches for Large Scale Digital Library



Introduction & Motivation

Tutorial at RCDL 2007
October, 15th 2007

Thomas Risse
L3S Research Center, Hannover, Germany
Email: risse@L3S.de

In cooperation with Claudia Niederee (L3S), Carlo Meghini (CNR-ISTI), Heiko Schuldt (Uni Basel)

Agenda





1. Introduction & Motivation
2. Challenges of bringing DL to distributed Infrastructures
3. Underlying Technologies and their promises (SOA, P2P)
4. Solutions for decentralized DL infrastructures (with BRICKS Demos)
5. Conclusions and future directions

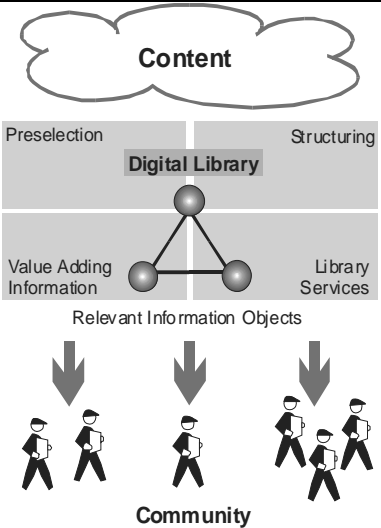



Page 2

Introduction

Digital Libraries:
Content-to-Community
Mediation





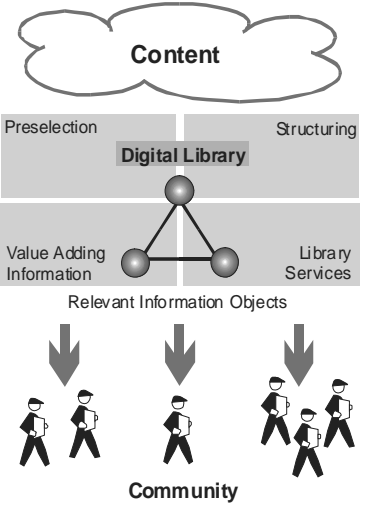



Page 3

Digital Libraries: That's where we are

- Core functionality of a DL is well-understood (Content-to-Community Mediation)
- Prototypical as well as operational Digital Libraries for various kinds of content have been build (Maps, Music, Historical Artifacts, ...)
- Standards for metadata, search, metadata harvesting, etc. have been established and are (more or less consequently) used fostering interoperation and integration
- Underlying technologies like information retrieval, have been advanced





Page 4

Digital Libraries: That's where we are (2)

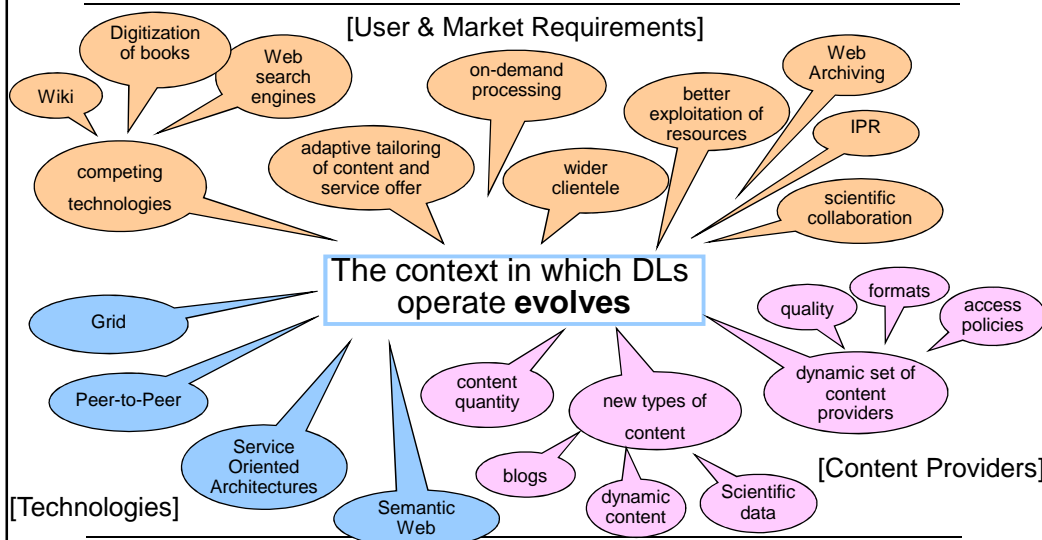
- Digital Library Management Systems in operation (e.g. Greenstone, DSpace, OpenDLib, FEDORA)

Example: DSpace

- Open source system developed by MIT Libraries and Hewlett-Packard Labs
- Repository for the digital research and educational material produced by members of a research university or organization
- Addresses the need of organizations to collect, preserve, index and distribute their material
- Is designed to make participation by depositors easy
- Information model is built around the idea of organizational "Communities"



However ...



Scenario 1 – Finds¹ identification today



- The public is an important source for understanding the past and archeological sites
- Thousands of objects are discovered by people during walking, traveling, gardening, etc.
- One day-to-day job of Archeological Museum Curators is to identify these objects
- Traditional way of working
 - Examining objects
 - Preliminary identification
 - Comparison to specialist reference collections, e.g. roman coins
 - Sometime object descriptions are enough
 - Sometime it is the starting point for a scientific analysis process
 - Comparison to objects in the museum collection
 - Record object in detail



¹ With *Finds* we mean *archeological finds*



Page 7

Scenario – Finds identification today



Disadvantage of the traditional way

- The reference collections have to be known by the curator
- Curator needs information about new collections, e.g. when a museum opens its archive
- Curators need to access them one by one
- Collections are heterogeneous from point of view of data, search facilities, user interface, languages, etc.

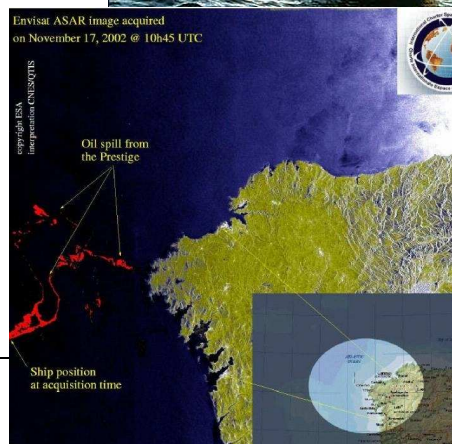
As a consequence curators invest a lot of time in the identification process



Page 8

Scenario 2 – Environmental Incidents (ImpECt¹)

- ESA ENVISAT satellite monitored an oil spill caused by a tanker off the Galician coast
- An expert group from different domains has to be coordinated and work on the problem (e.g. oil-fighting bodies, ministries of the involved countries, ...)
- different resources tailored to situation are required (relevant ecological and environmental data, historical dossiers on tankers accidents;
- On demand simulations on oil distributions have to be computed based on weather forecast, wind, temperature, currents, ...
- Postprocessing: reports have to be generated, environmental effects have to be analysed

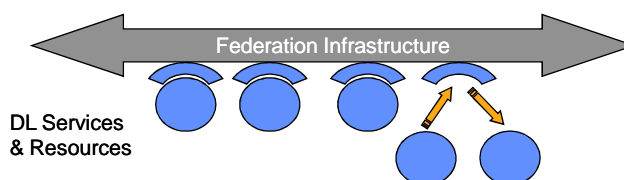
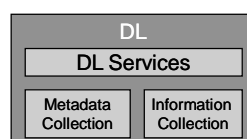
¹ImpECt= Implementation of Environmental Conventions

Next Generation Digital Libraries (NGDL)



Current plans for next generation DL architectures are aiming for a **transition**

- from the DL as an integrated, centrally controlled system
- to a **dynamic configurable federation** of DL services and information collections for enabling large-scale DLs.



Page 10

In the Tutorial ...



Closer look

- on these challenges and the resulting requirements
- Selected solutions for these challenges developed in the European Project BRICKS
- technologies for implementing next generation digital libraries (SOA, Peer-to-Peer)
- lessons learned from implementing large-scale digital libraries on top of distributed infrastructures
- small demos (if time is available)



Page 11

Approaches for Large Scale Digital Library



ISTITUTO DI SCIENZA E TECNOLOGIE
DELL'INFORMAZIONE "A. FAEDO"



Challenges of bringing DL to distributed Infrastructures

Tutorial at RCDL 2007
October, 15th 2007

Thomas Risse
L3S Research Center, Hannover, Germany
Email: risse@L3S.de



In cooperation with Claudia Niederee (L3S), Carlo Meghini (CNR-ISTI), Heiko Schuldt (Uni Basel)

Agenda

1. Introduction & Motivation
2. Challenges of bringing DL to distributed Infrastructures
3. Underlying Technologies and their promises (SOA, P2P)
4. Solutions for decentralized DL infrastructures (with BRICKS Demos)
5. Conclusions and future directions

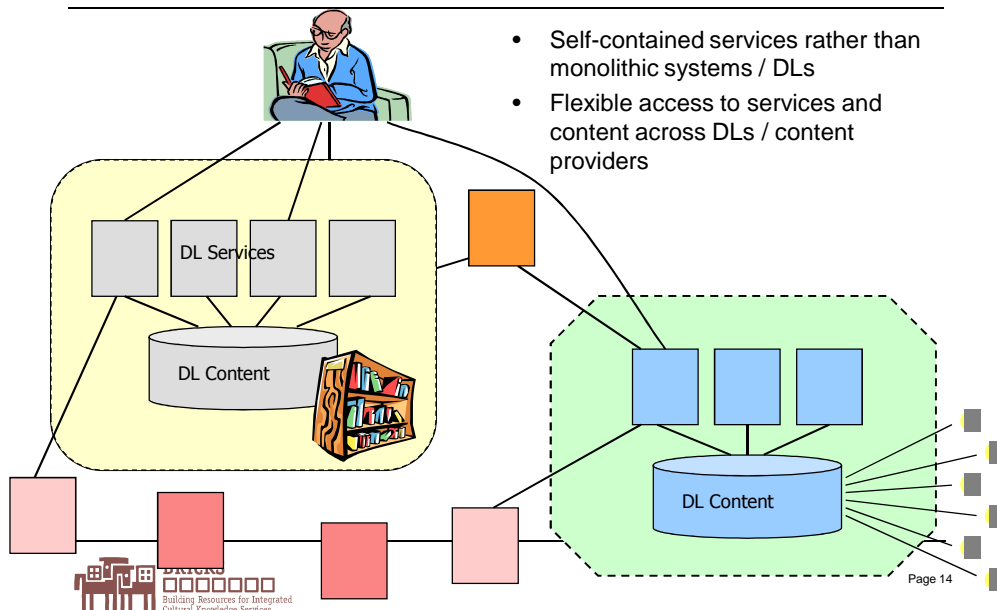


Page 13

Digital Libraries – The Future





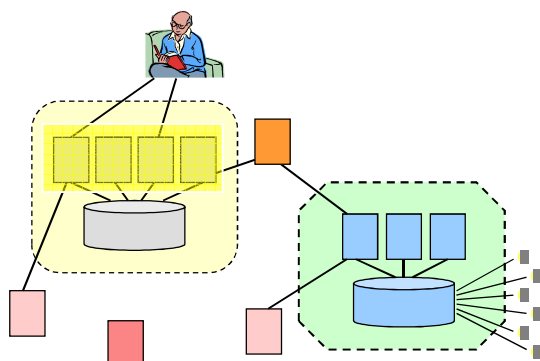
- Self-contained services rather than monolithic systems / DLs
- Flexible access to services and content across DLs / content providers




Page 14

Requirements ...





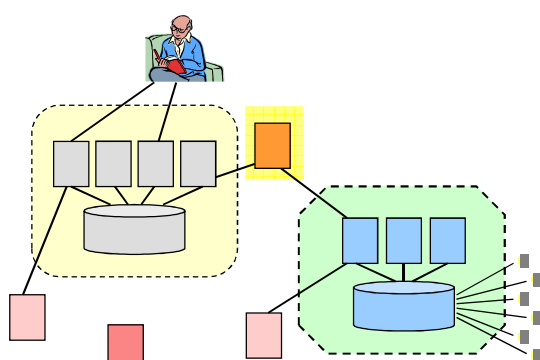
- Specialized services, **local** to a content provider
 - Search
 - Different media types
 - Content-based
 - Multi-object, multi-feature
 - Multilingual Access
 - Relevance feedback
 - ...
 - Indexing
 - Annotation
 - Metadata management
 - Content management
 - Resource Management
 - ...




Page 15

... Requirements ...

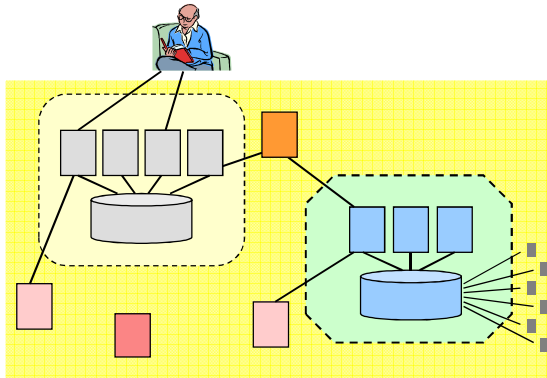


- Specialized services, **across** different providers
 - Search
 - Different media types
 - Content-based
 - Multi-object, multi-feature
 - Multilingual Access
 - Relevance feedback
 - ...
 - Metadata management
 - Content management
 - Resource Management
 - Indexing
- without central control (no censorship)



Page 16

... Requirements ...



Virtual DL (virtual collection)

- Easy to extend by integrating new content providers, services, etc.
- Example: collaboration in eScience applications

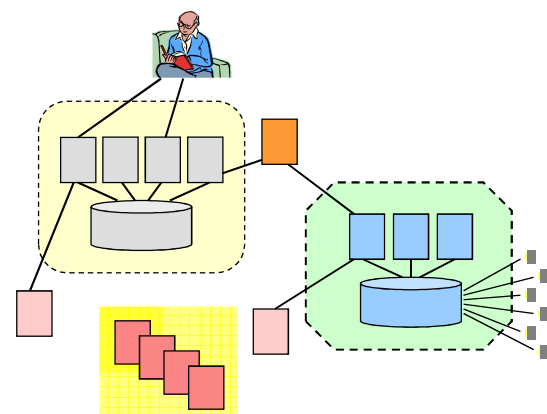
Management of services which are

- Distributed
- Heterogeneous
- Autonomous



Page 17

... Requirements ...

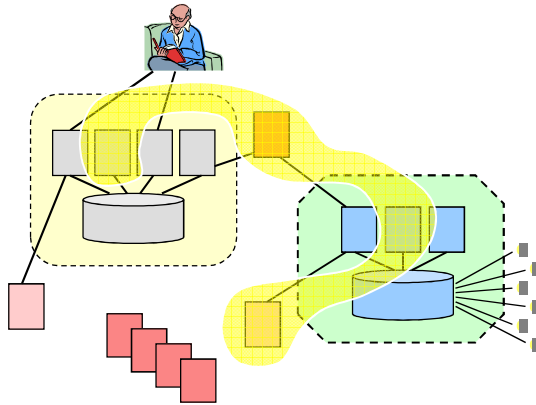


- Management of computationally intensive services
 - Scale-out
 - Load balancing
- Example: services for extracting characteristic features from image/audio/video documents



Page 18

... Requirements ...



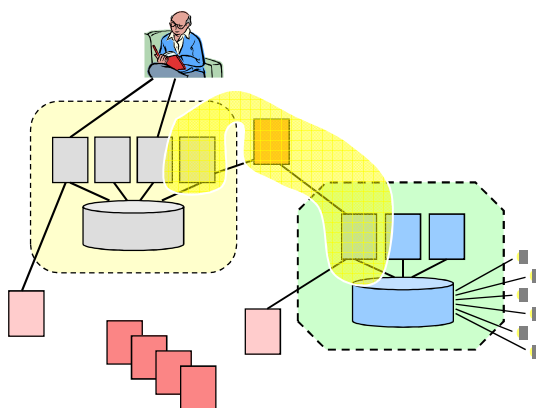
Composition of services

- Defining complex services / processes / workflows on the basis of existing services
- Flexibility: automatic adaptation
- Example: complex processes for automated storage and replication of data, generation of meta data (content features), ...



Page 19

... Requirements ...

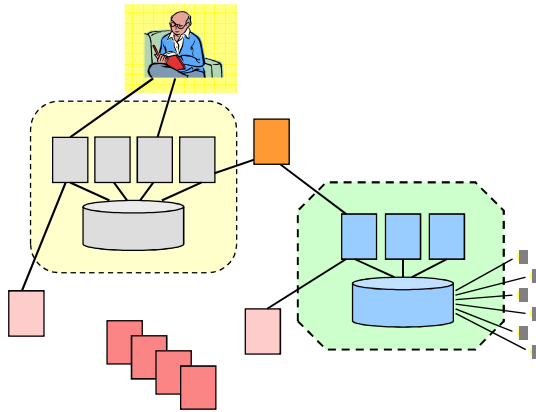


- Notification of changes
- Guaranteed consistency of derived data
- Example: automated updates of complex reports



Page 20

... Requirements ...

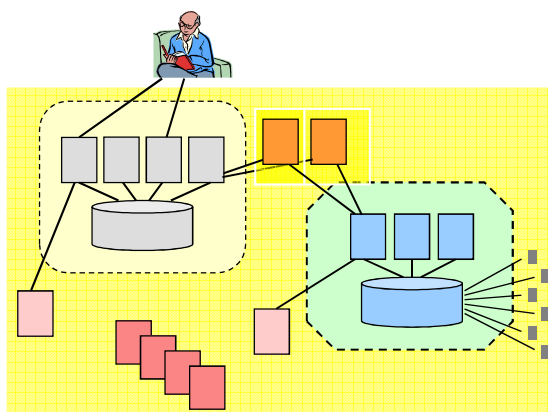


- Personalization
- Visualization
- Access anywhere
 - Mobile devices
- Context- and location-aware services
- Authentication & authorization



Page 21

... Requirements ...

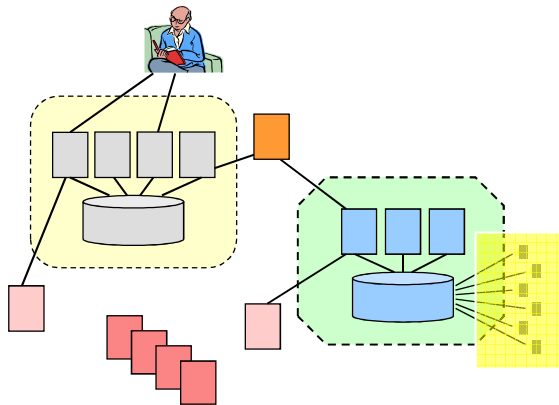


- High degree of availability:
 - Access anytime
 - Replication
- High degree of dependability / reliability
 - Systems their users can count on
- High degree of scalability



Page 22

... Requirements



Continuously generated data

- Sensor networks
- Sensor data streams (hardware / software sensors)
- Example: data created by certain instruments in eScience

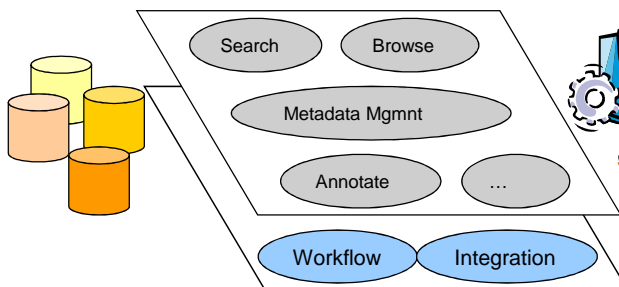


Page 23

DL Architecture – an Abstract View



DL users ...



... have access to DL functionality that ...

- operates on distributed content
- is implemented by distributed services
- is supported by system services



Page 24

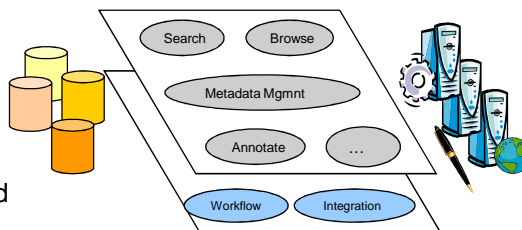
Summary



- The challenges and requirements of new DLs and their applications cannot be met by current DL Systems



- New technologies are needed that allow to bring DLs to **Distributed Infrastructures**



Page 25

Questions & Discussions



Page 26

Approaches for Large Scale Digital Library



ISTITUTO DI SCIENZA E TECNOLOGIE
DELL'INFORMAZIONE "A. FAEDO"



Underlying Technologies

Tutorial at RCDL 2007
October, 15th 2007

Thomas Risse
L3S Research Center, Hannover, Germany
Email: risse@L3S.de

In cooperation with Claudia Niederee (L3S), Carlo
Meghini (CNR-ISTI), Heiko Schuldt (Uni Basel)

Agenda



1. Introduction & Motivation
2. Challenges of bringing DL to distributed Infrastructures
3. Underlying Technologies and their promises (SOA, P2P)
4. Solutions for decentralized DL infrastructures (with BRICKS Demos)
5. Conclusions and future directions



Page 28

Overview



- Service-oriented Architectures
 - Definition
 - Service Model
 - Web Service Stack
 - Example: Supply Chain
- Peer to Peer Infrastructures
 - The nature of peer-to-peer systems
 - Some history
 - Application domains
 - The Data access challenge
- Summary



Page 29

Service Oriented Architectures (SOA)



Page 30

Definitions



Gartner Group (technical definition)

- Web Services are loosely coupled software components that interact with one another dynamically via standard Internet technologies.

Forrester Research (business definition)

- Web Services are automated connections between people, systems and applications that expose elements of business functionality as a software service and create new business value.



Page 31

Service Oriented Computing



Goals of Service Oriented Computing

- Distributed interoperable Systems
- Cost reduction due to reuse of services
- Flexibility
- Easy adaptation to future developments

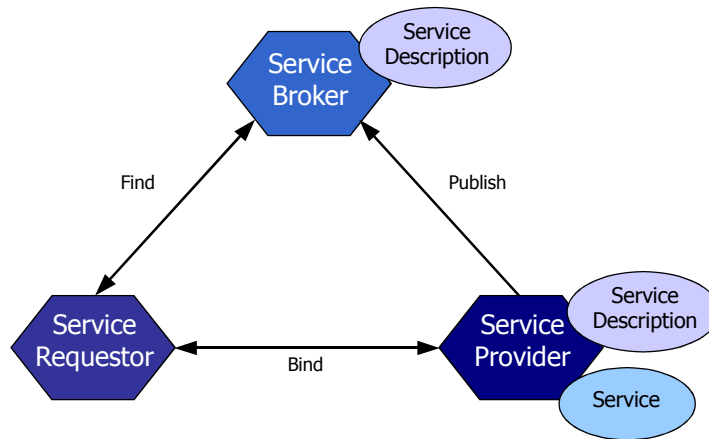
Properties of services

- Self describing
- Accessed using a well-defined interface
- Autonomous
- Loosely coupled
- Independent from the Platform / Operating System
- Independent from the Programming language



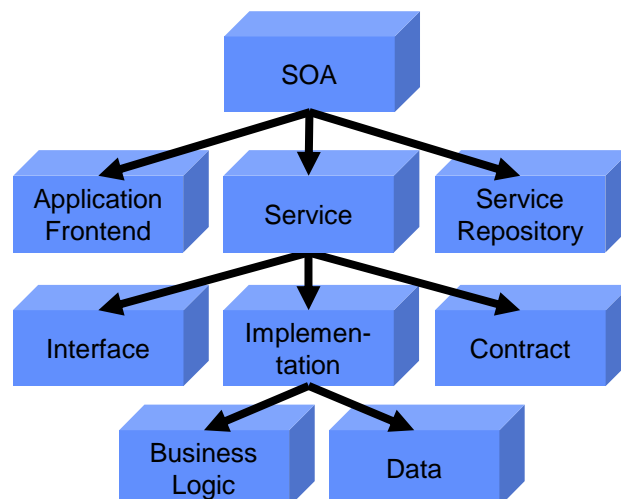
Page 32

Service Model

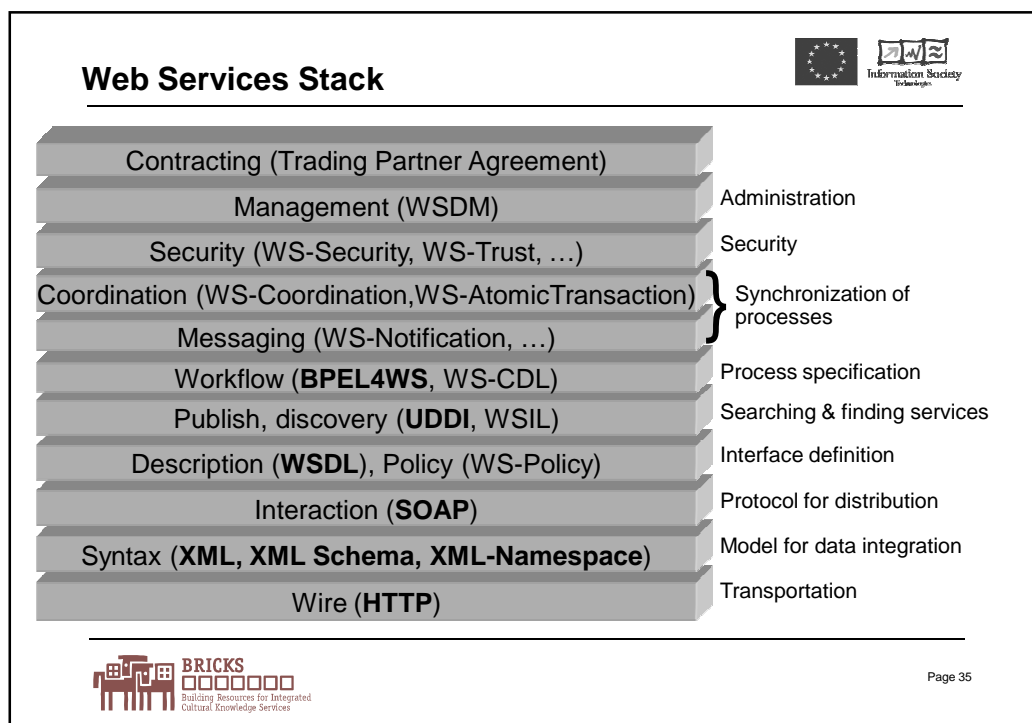


Page 33



Elements of SOA




Page 34

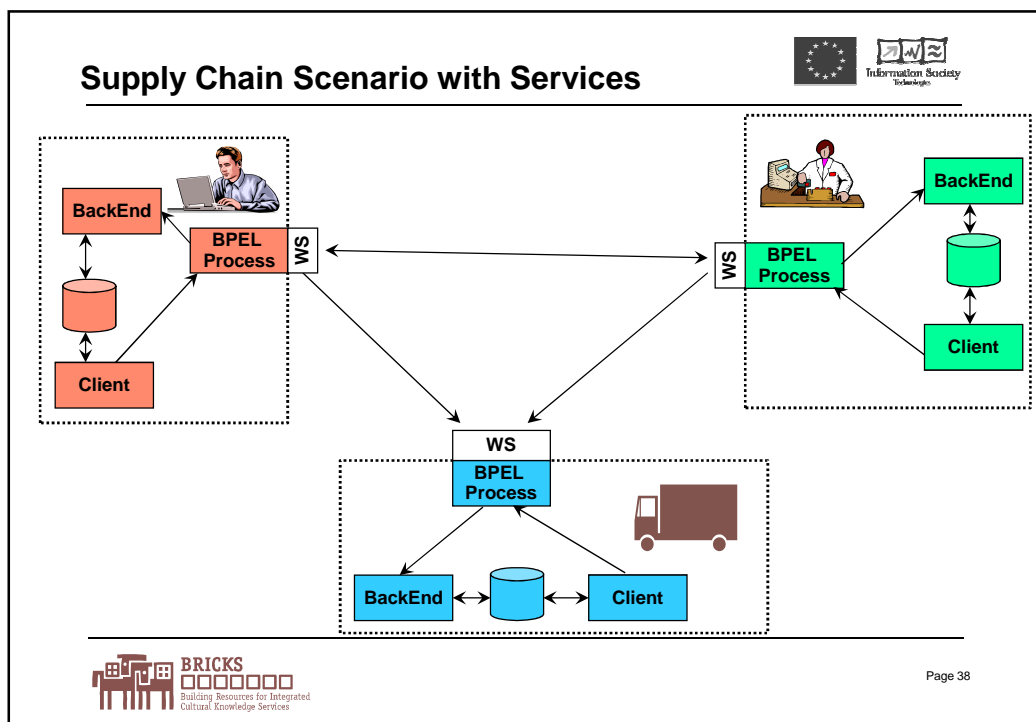
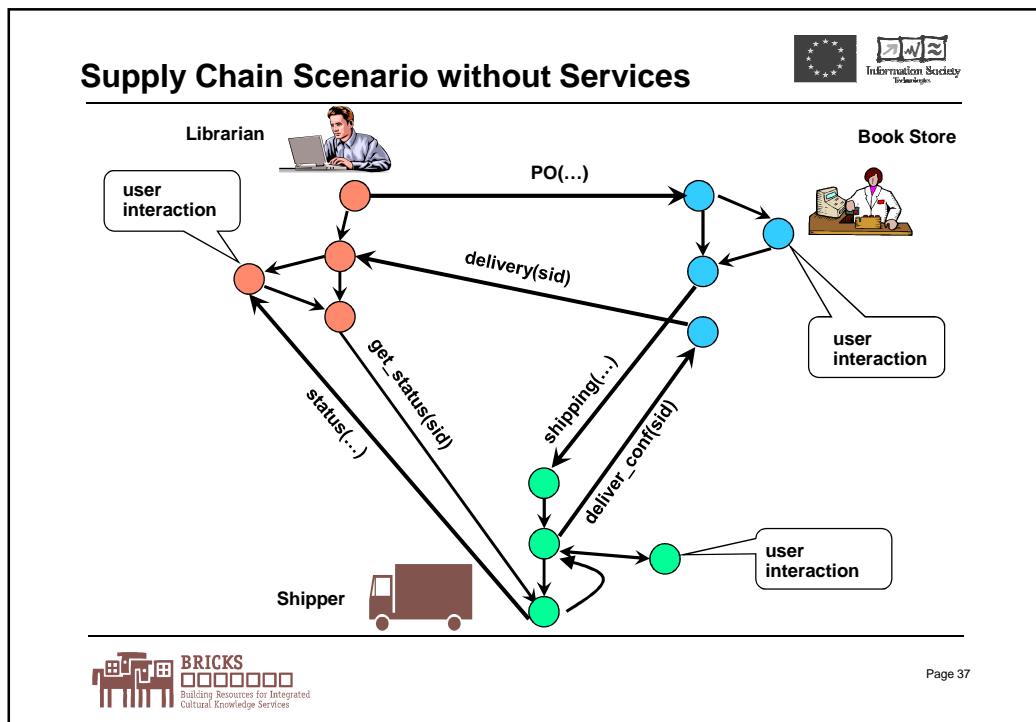


Relevant Standards Overview

XML	A universal model for data exchange and data integration
XML Schema	Defines the schema of a XML document, makes syntactical restrictions, defines structural patterns, defines data types
XML Namespace	Provide means to avoid naming conflicts in XML documents
SOAP	Simple Object Access Protocol - A universal communication protocol
WSDL	WS Description Language - Description of the WS interfaces, parameters, etc.
WS-Policy	General-purpose model to describe the policies of a Web service
UDDI	Universal Description, Discovery and Integration A standard for publication and discovery of information
BPEL4WS	Business Process Execution Language Specification of workflow for the composition of services
WS-Notification	Defines the publish/subscribe pattern for message oriented systems
WS-Coordination	Coordination of distributed actions. Includes transaction management.
WS-Security	Secure SOAP messages
WS-Trust	Management of trust relationships
WSDM	Distributed Management of Web Services


Page 36



Advantages and Disadvantages



Advantages

- Clear Description of Services and Interfaces
- Transparent access to Legacy Systems
- Higher Flexibility and Dynamics
- Widely accepted Web Service standards
- Various software is available

Disadvantages

- Semantic standards are still under development
- Several complementing standards are in development
- No Resource Sharing
- Depends partly on central management services



Page 39

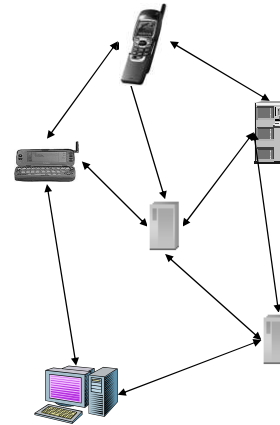
Peer-to-Peer Computing



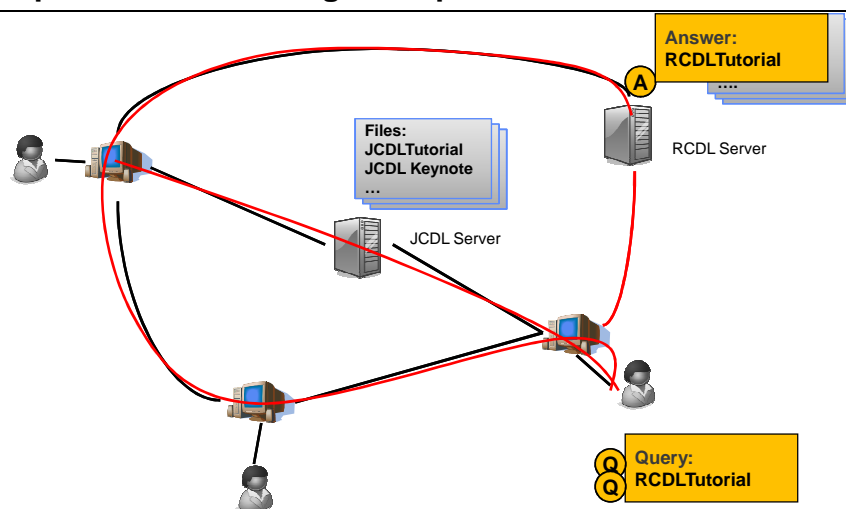
Page 40

The Nature of P2P

- Peers act as client and server
 - Peers communicate directly
 - Both provide services, e.g. storage, calculations
- Peers are autonomous
 - Owner decide e.g. when a service is active
- No central administration
- No global system view
 - Peers have only some knowledge about their neighborhood
- Highly dynamic
 - Peers appear and disappear whenever they like
- Unreliable connections
- Mostly self organizing
 - Services and clients are actively looking for appropriate partners



Simple P2P File Sharing Example



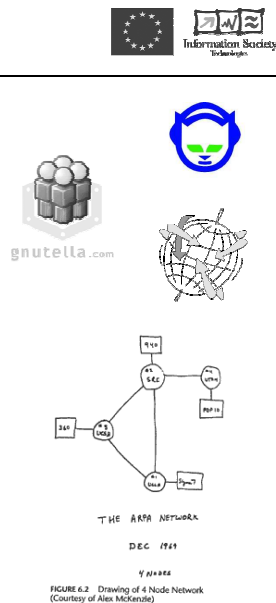
Some History

P2P receives a lot of attention ...

- File sharing systems
 - Gnutella, Freenet, Napster, ...

... but P2P is nothing new

- Telephone systems
- Federated databases
- ARPAnet
 - First networks were P2P systems
- Usenet, Domain Name Service
 - Servers communicate in a P2P manner
- Service Oriented Architectures



Page 43

Advantages of Peer-to-Peer

- Scalability
 - P2P solutions should be scalable by definition
- Availability of information, services, ...
 - Redundancy increases the availability
- Flexibility
 - Flexible reaction on changes, e.g. leaving of peers, changing from wire to wireless communication
- Low administration costs
 - Self organization requires no central administration



Page 44

Disadvantages



- Advantages grow with the number of peers
 - To accomplish the advantages many peers are necessary
- Software has to deal with P2P properties
 - Some P2P properties (e.g. dynamicity) are challenging for the development
 - Software components can provide functions to other peers as a service (e.g. Microsoft .Net)



Page 45

Application Domains



- Resource sharing
 - Joined usage of unutilized resources (e.g. CPU, storage)
- Content sharing
 - Sharing of files (music, video), calendar, spam sender data, ...
 - News distribution
 - Searching for information
 - Examples: Gnutella, Freenet, Napster, Groove, Hive, Cloudmark SpamNet
- Collaboration
 - Peers work together
 - Collaboration within the project team, e.g. document creation
 - Instant Messaging
 - Benefits
 - Autonomous
 - Offline working
 - No administration
 - Examples: Groove (Collaboration), Skype (Messaging), ...



Page 46

Challenge: Data Access



Data Access Structures

- Methods to query and access „nomadic data”
- Partial knowledge

Simple approaches

- Breadth-first-search (Gnutella)
- Depth-first-search (Freenet)

Disadvantage

- High accumulated bandwidth

More sophisticated approaches

- Use global unique identifier (GUID) of object for routing, (DHT, Plaxton Routing, Chord, Tapestry, Pastry, CAN, P-Grid)



Page 47

Summary



Web Services

- Well established standard for interoperable services
- Lightweight communication protocol
- Easy to use
- Depends partly on central management services
→ Fully decentralized Web Services are developed in BRICKS

Peer-to-Peer

- Describes the way of communication between entities
- Self-organizing system of autonomous entities
- Highly flexible and scalable



Page 48

Questions & Discussions



Page 49

Approaches for Large Scale Digital Library



ISTITUTO DI SCIENZA E TECNOLOGIE
DELL'INFORMAZIONE "A. FAEDO"



Solutions for decentralized DL
infrastructures: issues, solutions and
advantages in the framework of BRICKS

Tutorial at RCDL 2007
October, 15th 2007



Thomas Risse
L3S Research Center, Hannover, Germany
Email: risse@L3S.de

In cooperation with Claudia Niederee (L3S), Carlo
Meghini (CNR-ISTI), Heiko Schuldt (Uni Basel)

Agenda

1. Introduction & Motivation
2. Challenges of bringing DL to distributed Infrastructures
3. Underlying Technologies and their promises (SOA, P2P)
4. Solutions for decentralized DL infrastructures (with BRICKS Demos)
5. Conclusions and future directions



Page 51



The BRICKS Project
<http://www.brickcommunity.org/>



Page 52

BRICKS - Project Identity Card



- Project Acronym: BRICKS - Building Resources for Integrated Cultural Knowledge Services
- Consortium: 24 organizations from 9 countries
- Thematic area: Digital Libraries Services
- Duration: 42 months
- Started in January 2004
- Budget: 12,2 Mill. Euro
- Homepage: <http://www.brickscommunity.org/>
- Technical Partners: Engineering, Italy; Fraunhofer IPSI; CNR-ISTI, Italy; Metaware, Italy; ARC Seibersdorf Research GmbH, Austria; EPFL, Switzerland; Canoo, Switzerland; Uni. of Athens, Greece; Uni. of Sheffield, UK; Scuola Normale Superiore di Pisa, Italy; Uni. of Florence, Italy; Oxford ArchDigital, UK; PolyDisplay, Norway
- User Partners: Italian Ministry of Culture, Italy; Vatican Secret Archives; Uffizi Gallery, Italy; Schönbrunn Palace, Austria; Austrian National Library, Austria; European Museum Forum, UK; Museum of Cycladic Art, Greece; Russian Cultural Heritage Network, Russia; Museums, Libraries and Archives Council, UK; Studio Azzurro, Italy



Page 53

Access to professional cultural resources today



Situation

- Large number of independent distributed digital information sources
- Often professional sources are not public
- Still a large number of interesting sources are missing, e.g. small museums

Find and access information sources (examples)

- Search engines, e.g. CiteSeer, DBLP, Google
- Z39.50 for some libraries
- Personal Knowledge
- ...

Disadvantages

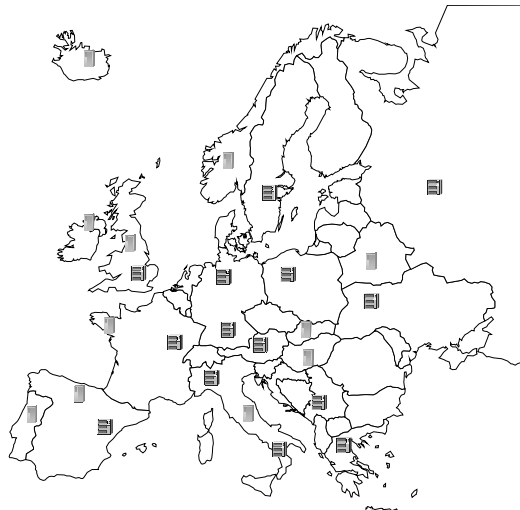
- Time consuming
- Limited openness
- Seldom collaboration support
- Seldom language support



Page 54

Overall Goals of BRICKS

- Transparent access to distributed available information sources
- Retrieval of information with knowledge support
- Multi-lingual
- Protection of intellectual properties
- Low Cost
- Easy installation and maintenance
- Platform-independent



Page 55

BRICKS Approach

Service oriented

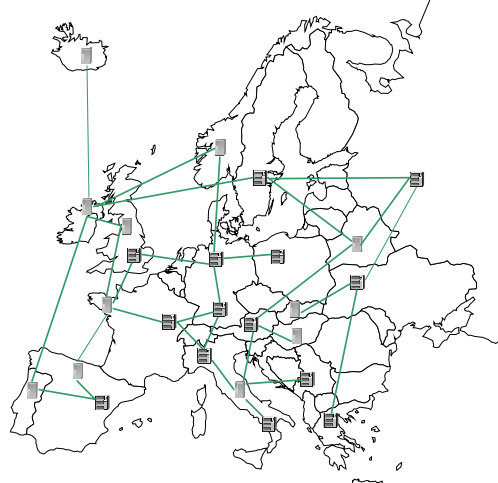
- Standardized interface descriptions based on Web Services
- Platform independent
- Flexible composition of services

Decentralized Organization

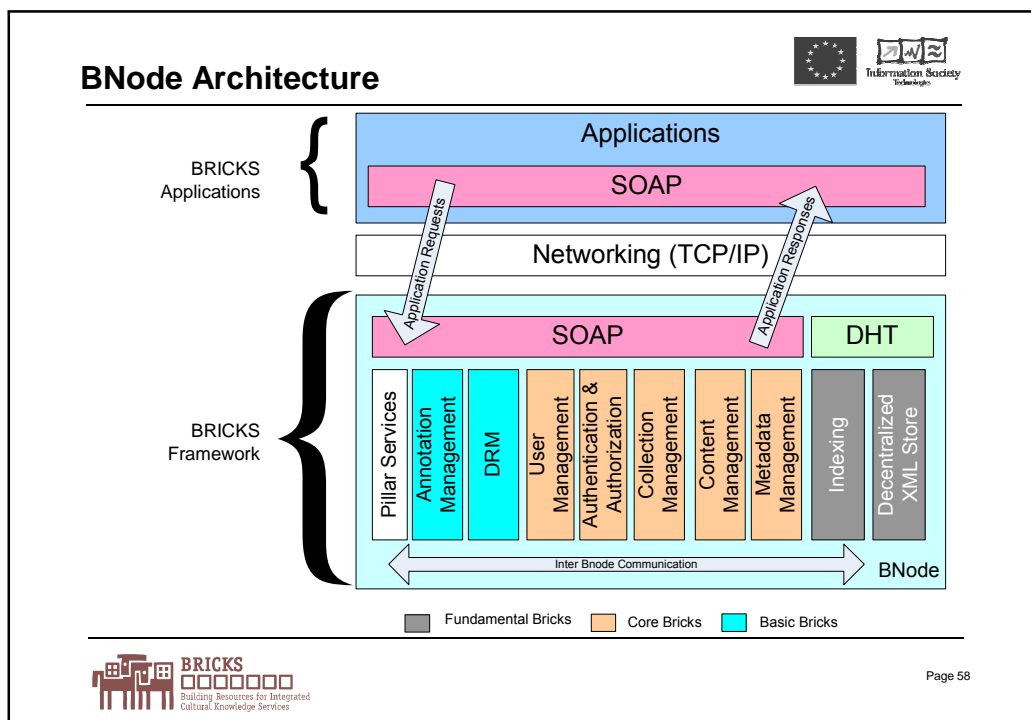
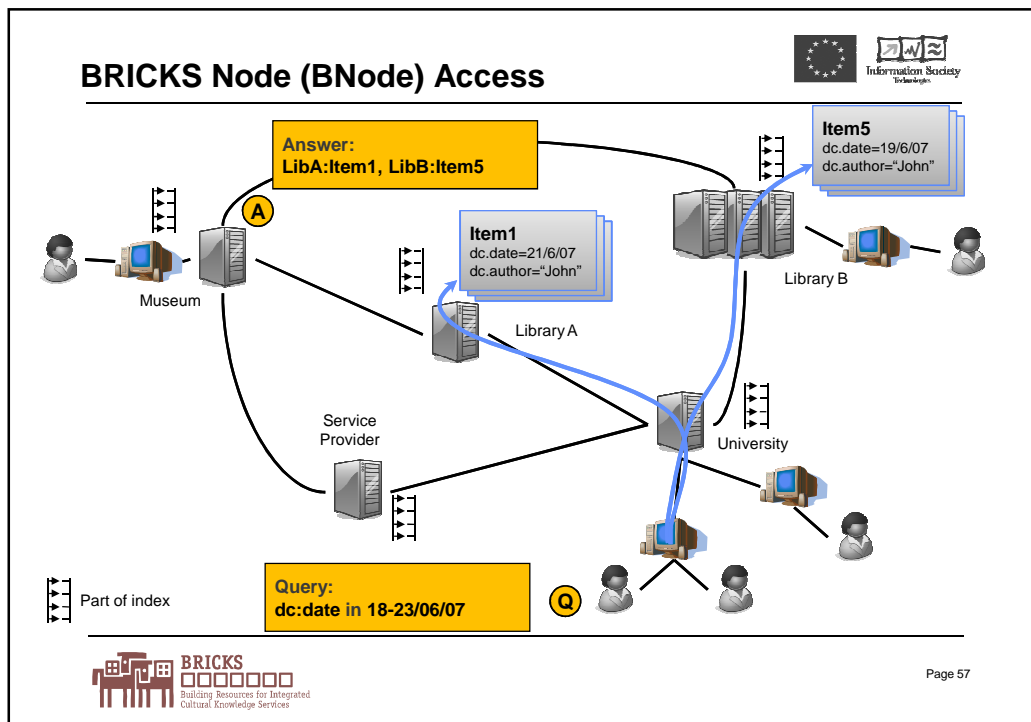
- Avoid central coordination
→ No Single Point of Failure
→ No central maintenance organization necessary
- Highly scalable
- Increased reliability
- Minimized maintenance cost

Usage & Costs

- Easy installation
- Open Source lowers the barrier to use BRICKS
- Allows tailoring to user needs



Page 56








BRICKS Workspace Demo

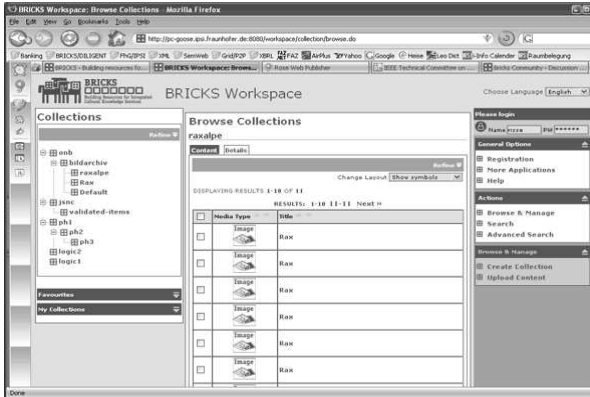



Page 59

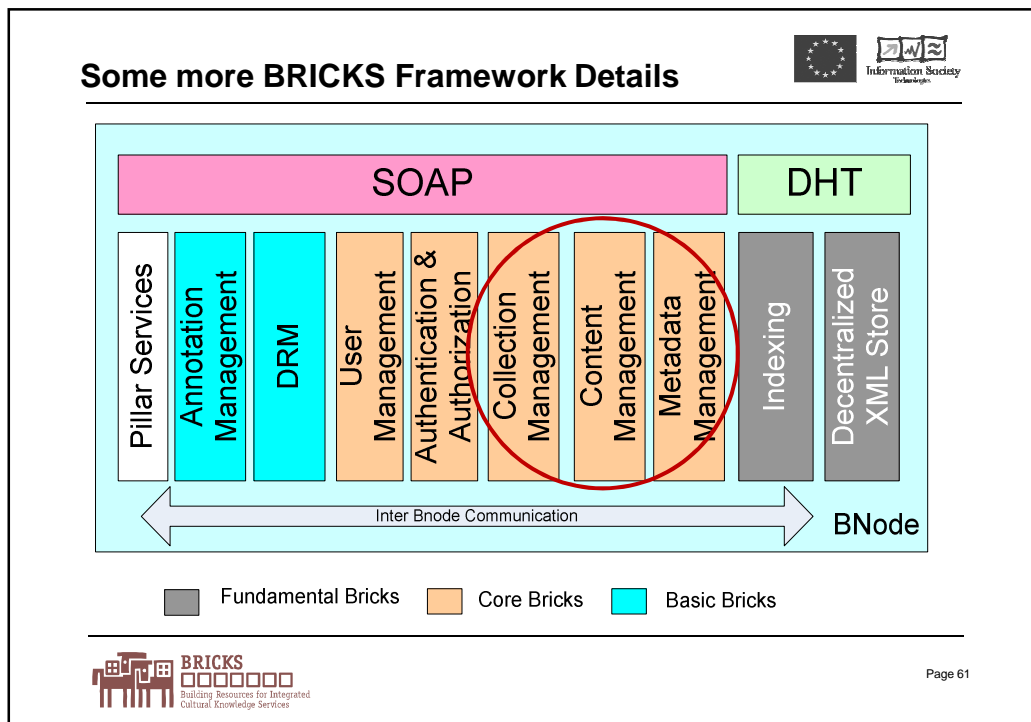
The BRICKS Workspace



- What does it demonstrate?
 - A web application (thin client) accessing BRICKS Foundation services
 - Navigation through the BRICKS information space
- Target audience
 - General end-users
 - Application developers (as an example)
- Alternatives
 - BRICKS Desktop
An ECLIPSE based native user interface
 - Domain & Application specific interfaces





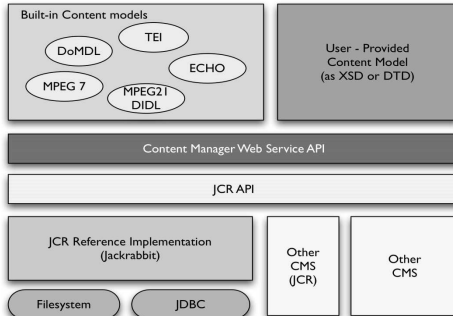
Page 60






BRICKS Content Management

- BRICKS is based on the Java Content Repository
- Content Repository for Java™ Technology API (JSR-170)
 - Born from the need to standardize the proprietary content repositories
 - Supports a wide range of applications
 - Provides a unified API
 - Creation and access
 - Versioning
 - Access control
 - Full text searching
 - API is independent from the Back-End storage systems, e.g. file system, WebDAV repository, XML database, SQL database
 - JCR based systems: Apache-Jackrabbit, Magnolia, Alfresco, ...
- BRICKS extends JCR with a Web Service interface
- Provide some meta-content models, e.g. DoMDL Model, MPEG-21, MPEG-7, TEI-Lite, ...






Page 63

BRICKS

Metadata Management



Page 64

Metadata



Minerva classification

A lot of definitions → Always a good reason for long discussions

Professional Users (e.g. Librarians) have a concrete view

Metadata are value-added information that professional users create to arrange, describe, track and access information objects

Different Types of Metadata

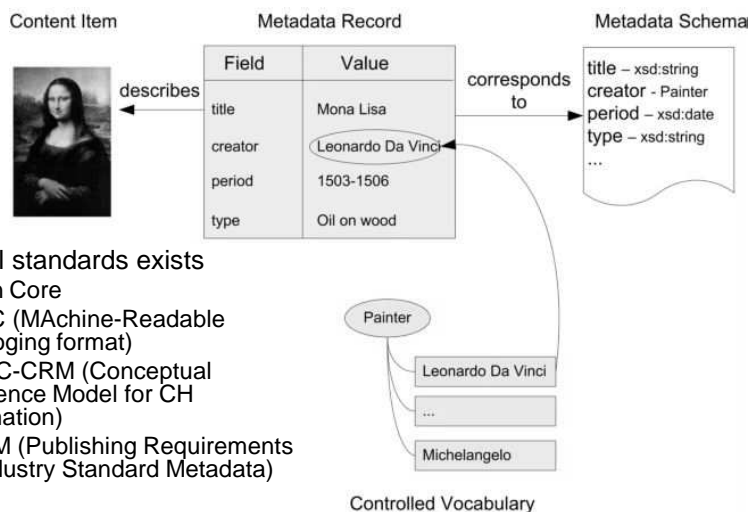


Typ	Definition
Administrative	Metadata used in managing and administering information resources
<i>Descriptive</i>	<i>Metadata used to describe or identify information resources</i>
Preservation	Metadata related to the preservation management of information resources
Technical	Metadata related to how a system functions or metadata behave
Use	Metadata related to the level and type of use of information resources



Page 65

Descriptive Metadata



- Several standards exist
 - Dublin Core
 - MARC (Machine-Readable Cataloging format)
 - CIDOC-CRM (Conceptual Reference Model for CH information)
 - PRISM (Publishing Requirements for Industry Standard Metadata)

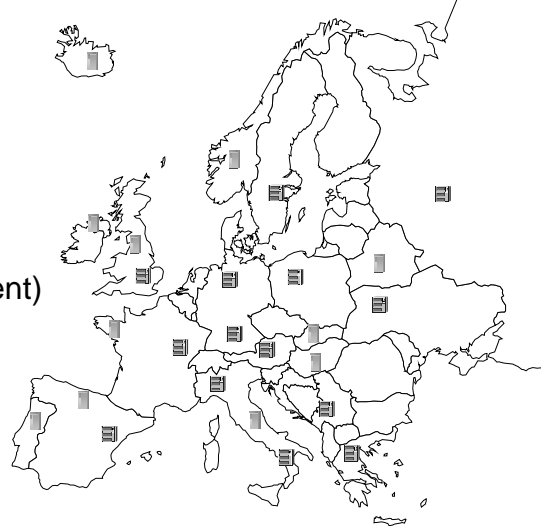


Page 66

Metadata Storage Requirements



- Durability
- Availability / Accessibility
- Scalability
- Consistency
(des. metadata \leftrightarrow content)



Page 67

Metadata Storage Modes (1/2)



- Local Storage
 - Advantage: simple, well-known, ...
 - Disadvantage: no guarantee of 100% availability

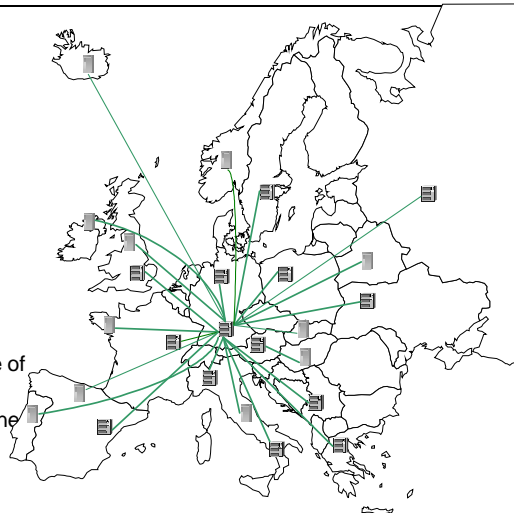


Page 68

Metadata Storage Modes (1/2)



- Local Storage
 - Advantage: simple, well-known, ...
 - Disadvantage: no guarantee of 100% availability
- Central Storage
 - Advantage
 - Well-known
 - Easy to implement transaction guarantees
 - Disadvantage
 - Scalability problems for large volume of data and requests
 - High concentration of resources at one place (bandwidth, space, CPU)
 - High costs
 - Central point of failure

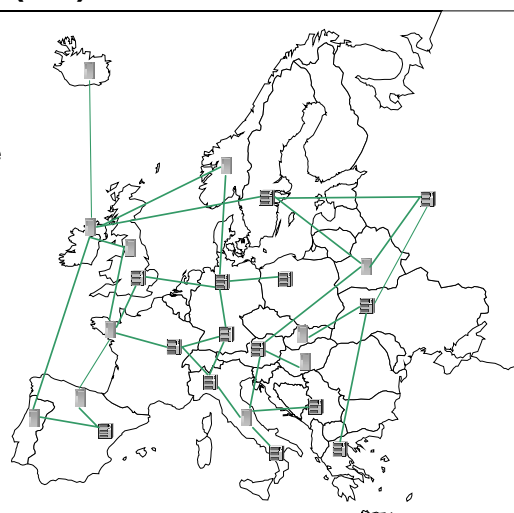


Page 69

Metadata Storage Modes (2/2)



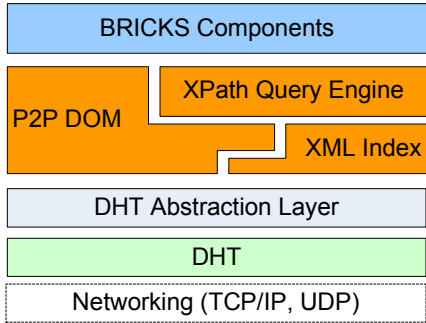
- Decentralized Storage
- Advantage
 - High scalability
 - Fault tolerance, no central point of failure
 - Better usage of available resources
 - Decentralized architecture → decentral storage
 - Transparent access
 - Disadvantage
 - Slower access → an index speed it up in many cases
 - Complicated logic behind
 - No transaction guarantees




Page 70

Decentralized Storage in BRICKS

- Used for administrative metadata
- Data spread within BRICKS community, but transparent access to users
- Uses P2P layer for BNode communication
 - Discovering neighboring peers
 - Routing and processing messages within BRICKS community, but without global topology knowledge
- Implements
 - Well-known W3C DOM API for creating and accessing XML documents
 - XPath language for querying XML documents
- Built-in protocols for
 - Maintaining high data availability through replication
 - Concurrency and consistency control



The diagram shows a stack of components for BRICKS. At the top is a blue box labeled 'BRICKS Components'. Below it are two orange boxes: 'P2P DOM' on the left and 'XPath Query Engine' on the right, with 'XML Index' below the XPath Query Engine. These are followed by a light blue box 'DHT Abstraction Layer', a green box 'DHT', and a dashed box 'Networking (TCP/IP, UDP)' at the bottom.





BRICKS
Building Resources for Integrated
Cultural Knowledge Services

Page 71

Storage Locations of Metadata Types in BRICKS

- Local storage
 - Descriptive metadata (with decentralized index)
 - Technical (EXIF, etc.)
 - Security metadata (ACL, etc.)
 - Annotations
 - Ontologies (with decentralized discovery)
- Decentralized
 - Service descriptions
 - Collection descriptions





BRICKS
Building Resources for Integrated
Cultural Knowledge Services

Page 72

Descriptive Metadata in BRICKS



- The BRICKS **Metadata Manager**
 - is responsible for managing cultural assets from various institutions
 - must support arbitrary metadata formats from various institutions
- Schema/Format definition
 - the semantics of records is modeled and exposed in OWL-DL
- Bibliographic records
 - are internally stored and represented in RDF
- Controlled vocabularies, Thesauri, etc.
 - supported if they are represented in RDF, RDFS or OWL

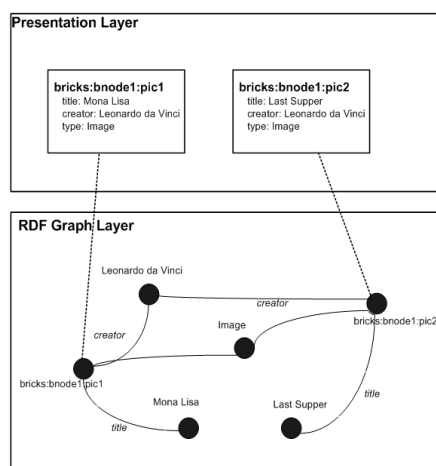


Page 73

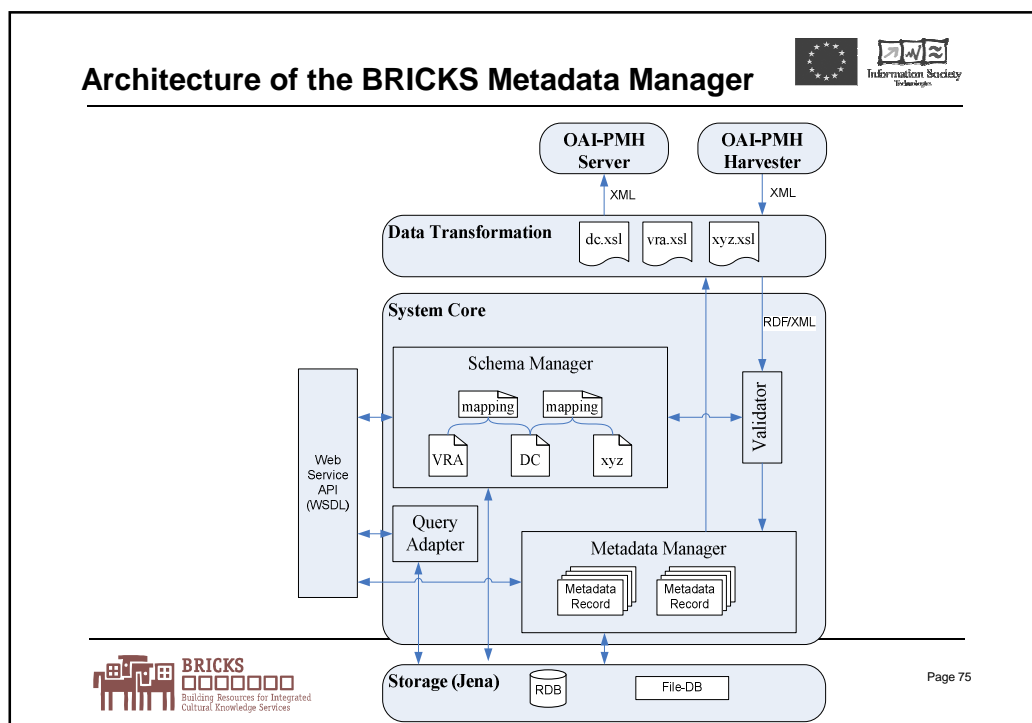
Metadata Manager - Design Decisions



- RDF and OWL
 - are well-suited for handling heterogeneous metadata internally (!)
- But:
 - do NOT show triples or anything else related to RDF or OWL to the user or front end developers!
 - nobody wants to obtain triples in a query result
- In BRICKS we've decided to
 - completely hide RDF and OWL from users and high-level applications
 - expose it only for machines



Page 74



BRICKS and Existing Systems

- The BRICKS idea is not to replace but to integrate with existing systems
- Rely on already accepted protocols in the Digital Libraries domain to tap existing metadata and content databases
- Currently available:
 - Import of “existing” metadata via the OAI protocol for metadata harvesting (OAI-PMH)
 - easy to understand
 - easy to use and minimal implementation effort for institutions
 - gives us a minimum level of interoperability (Simple Dublin Core)
 - Wrapping of SRU sources (Search/Retrieval via URL)
 - Integration of SRU sources as collections
 - Providing BRICKS metadata via SRU interface



BRICKS Collection Management

or

How to navigate through the BRICKS information space?



Page 77



Outline

- What is a collection manager
- Requirements on collections
- Implementing the collection management
- Searching a distributed DL
- Implementing the search mechanism



Page 78

Outline



- What is a collection manager
- Requirements on collections
- Implementing the collection management
- Searching a distributed DL
- Implementing the search mechanism



Page 79

What is a Collection Manager



- A Collection Manager is a mediator between:
 - the applications
 - and the DL contents.
- Applications may range from GUIs to arbitrarily complex programs for performing domain specific tasks.
- DL contents are the digital objects stored in the DL and their associated information.
- The job of the CM is:
 - to represent the DL contents via a conceptual model that is understandable and intuitive from an application viewpoint
 - to offer primitives for the manipulation of this model in support of the application tasks.



Page 80

The DL Conceptual Model



- We can view the contents of a DL as a set of documents, that is **multimedia complex objects**, with associated information, which supports the basic services offered on these objects.
- Distinctive features:
 - the size is typically $O(10^3)$ - $O(10^6)$
 - the members are very heterogeneous, e.g. in
 - Structure
 - Types
 - Language
 - Format
 - Contents
- *How do we go about taming such a set?*



Page 81

The DL Conceptual Model



- The problem is not new: organization!

"Hence we view the organization of a digital library network as basically an abstraction mechanism in terms of which details from a lower level of representation are suppressed. This is a crucial issue when dealing with large digital libraries — just as abstracting techniques are important in the development of programs"

(Mylopoulos & Levesque, 1979)

- Digital Library Collections are an abstraction mechanism!



Page 82

Outline

- What is a collection manager
- Requirements on collections
- Implementing the collection management
- Searching a distributed DL
- Implementing the search mechanism



Page 83

Requirements: Content Providers

Content providers structure its object space in sets of items termed "collections"
→ Natural to mirror real-world collections and to hold the primary content of the DL.

BRICKS

- These containers are called Physical Collections
- A physical collection is a set of content items which belong to the same content provider and are homogeneous *from the provider point of view*:
 - items are of the same kind
 - items are described by the same metadata format(s)
 - Items have same digital rights
 - ...
- In BRICKS the contribution of a content provider to the DL is always defined in terms of one or more physical collections, thus when content items are added to (removed from) physical collections, they are added to (removed from) the digital library.



Page 84

Requirements: Content Providers



- Conversely: **a content item exists in the digital library only if there is a physical collection holding it.** Physical collections partition the DL information space: they are pair wise disjoint and their union makes up the content of the DL.
- **Physical collections are structured in (physical) sub-collections,** a notion that has been found a useful organizational mechanism by content providers.
- **A physical collection can have an arbitrary number of sub-collections** but only one parent collection. The graph of the sub-collection relation is a forest, each tree of which has a physical collection at the root. Within this tree, a content item of the physical collection belongs to exactly one sub-collection.
- **Physical collections are a central notion in the BRICKS content model:** not only content is organized by physical collection, but so is the discovery of resources and the definition of logical collections.



Page 85

Requirements: Content Consumers



- People go to libraries to acquire knowledge for carrying out their own tasks.
- Typically, they search the whole library and end up with a subset of the library contents, consisting of the items that are relevant to them.
- This subset may be regarded as the *consumer view of the library*.
- The view is never static:
 - Consumers may evolve it manually, by adding newly discovered items or removing no longer useful ones
 - The view may evolve automatically: consumers describe their needs in some language and the items that satisfy these needs are added to the view:
 - *pull mode*: the consumer initiates the process (BRICKS)
 - *push mode*: someone else initiates the process (e.g. publish-subscribe)



Page 86

Requirements: Content Consumers



- This notion of *view* is captured in BRICKS by Logical Collections.
- A Logical Collection in BRICKS is a set of *references* to DL items, with identity.
- Any registered BRICKS user can create and operate upon Logical Collections.
- By creating logical collections, BRICKS users organize the information space and at the same time enrich the Digital Library content: once a logical collection is created, it becomes part of the digital library information space and can be e.g. searched by other users as any other digital library resource.
- Operations on Logical Collections affect *references* not content items!

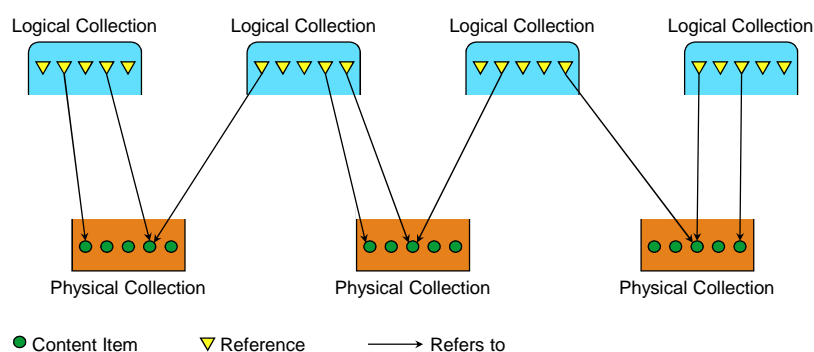


Page 87

Requirements: Content Consumers



- A BRICKS content item may be referenced in many Logical Collections.
- Conversely, a Logical Collection may contain references to many items, coming from many different Physical Collections.



Page 88

Static Collections



- BRICKS logical collections come in two flavors:
 - Static Collections
 - Dynamic Collections (a.k.a. Stored Queries)
- A static collection can be evolved manually by its owner, by inserting or removing references.
 - If the collection has been created with the intent of holding, for instance, all works of Maurolico on conics (intention), there is no guarantee that at any moment the collection indeed contains (extension) references to all Maurolico works on conics available in the digital library.
 - A static collection is a static container which communicates with the rest of the digital library only via the intervention of the users authorized to modify its extension.



Page 89

Dynamic Collections



- A dynamic collection is defined by a **query** over the DL content, including other collections.
- Any time the dynamic collection is accessed, either in browsing or via a query, the defining query is evaluated.
- Dynamic collections evolve automatically as the DL content evolves.
- Dynamic collections evolve **only** automatically: no operation is offered to add or remove references from a dynamic collection.
- Static collections are **extension-driven**:
 - for consumers who cannot describe their needs other than extensionally, i.e. by pointing at the relevant items
- Dynamic collections are **intension-driven**:
 - for consumers who can describe their needs as a BRICKS query
 - BRICKS supports only the pull-mode (for now)



Page 90

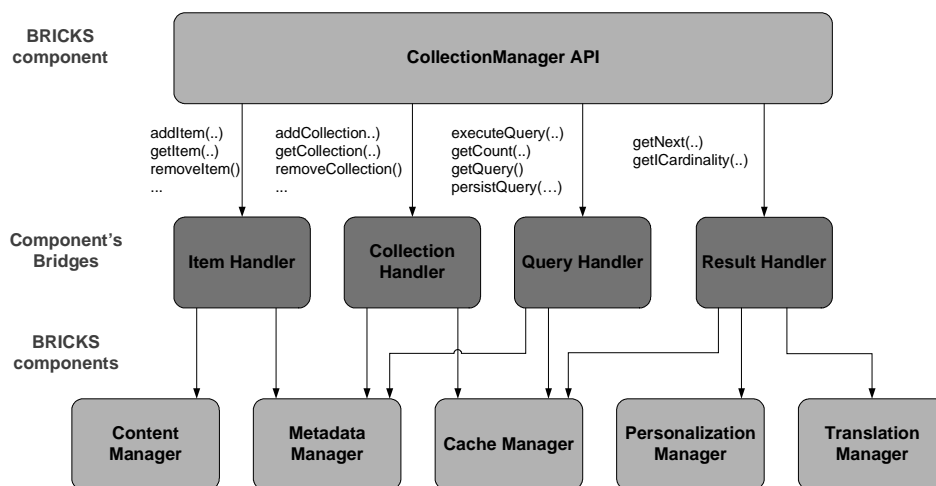
Outline

- What is a collection manager
- Requirements on collections
- Implementing the collection management
- Searching a distributed DL
- Implementing the search mechanism



Page 91

Architecture

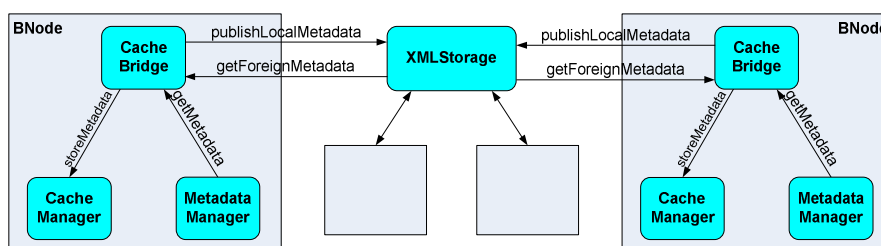


Page 92

Cache Bridge: Start-up



- In every BNode, the metadata of the locally defined collections are stored in the local Metadata Manager
- At BNode start-up, the Cache Bridge:
 1. Publishes the local collection metadata from the local Metadata Manager
 - into the local Cache and
 - into the decentralized XML Storage (to be retrieved by foreign BNodes)
 2. gets foreign collection metadata from the XML Storage into the local Cache



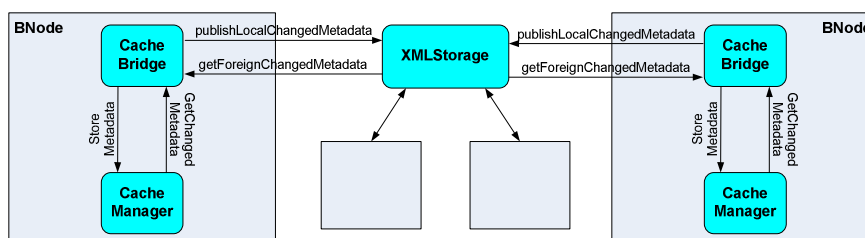
Page 93

Cache Bridge : Runtime



At runtime, whenever the local collections metadata are modified (by some operation)

1. the Cache Bridge gets the the modified local metadata from the Cache Manager and publishes them to the XML Storage for the foreign BNodes to retrieve them upon synchronization
2. Periodically synchronizes the foreign collections metadata with the XML Storage



Page 94

Outline

- What is a collection manager
- Requirements on collections
- Implementing the collection management
- **Searching a distributed DL**
- Implementing the search mechanism



Page 95

Searching a distributed DL

- The users of a DL form a very heterogeneous class, ranging from casual, computer-illiterate people to highly specialized scholars.
- All these users expect the DL system help them discover the DL objects of interest, independently from:
 - Location
 - Language
 - Type
 - Format
 - Structure
- in a way that reflects their preferences, if any, amongst which language plays a fundamental role.



Page 96

Requirements



- An important distinction: application-dependent vs application-independent discovery
- Application-dependent: based on criteria highly specialized with respect to a class of applications
 - x-ray images similar to a given one
 - brilliant executions of a piece of music
 - video sequences giving a feeling of anguish
 - spectacular takeovers in Formula 1 races
- Application independent: objective criteria
 - movies directed by Woody Allen
 - recent books on bio-informatics
 - Chopin's preludes played by Pollini



Page 97

Requirements



The query facility of a DL must be:

- extensible: it must be possible to plug-in it specialized search engines, devoted to capture the semantics of application-dependent search criteria
- flexible: it must be possible to express different kinds of queries, each addressing a different level of skill or knowledge of the DL contents
- user-sensitive: it must adapt to the preferences of the user
- efficient: it must respond as fast as possible



Page 98

Outline

- What is a collection manager
- Requirements on collections
- Implementing the collection management
- Searching a distributed DL
- Implementing the search mechanism



Page 99

The BRICKS query language

Types of queries:

- *Collection queries*: allow to discover collections by stating a Boolean condition on the collection metadata
- *DL objects queries*:
 - allow to discover content objects
 - may be contextualized to certain BNodes or collections
 - may be personalized
 - come in 3 flavors:
 - Simple
 - Advanced
 - Ontology



Page 100

Simple queries



- A simple query is the simplest form of query and most popular
- It is meant to serve casual users or users having a rather vague idea of the desired resources.
- These users typically do not want to address their search to specific metadata attributes, or operators.
- The language for simple queries allows to express sequence of unconstrained terms, very much in the style of e.g. Web search engines. In their search, users will be able to use:
 - wild cards: `ca?lo databas*`
 - phrases: `"jakarta apache"`
 - proximity operators: `"jakarta apache" ~10`
 - Boolean operators: `"jakarta apache" AND "jakarta lucene"`
- In a simple search, metadata records are seen as texts whose words are the metadata attribute values.



Page 101

Advanced queries



- For users who can characterize their information needs very precisely in terms of a metadata schema, such as librarians or expert library users.
- An advanced query is a
 - Boolean combination of simple conditions
 - on metadata fields, possibly coming from different schemas:

`DC.creator = "Bob" AND XY.date > 01.01.2000`

- In an advanced search, metadata records are seen as sets of (attribute value) pairs, exactly like database records, which may or may not satisfy a query stating simple conditions on such attributes.



Page 102

Ontology queries



- For highly specialized users, who are familiar with an ontology formally conceptualizing their domain of interest, and want to retrieve documents annotated by descriptions derived from that ontology.
- An ontology query is a SPARQL expression:

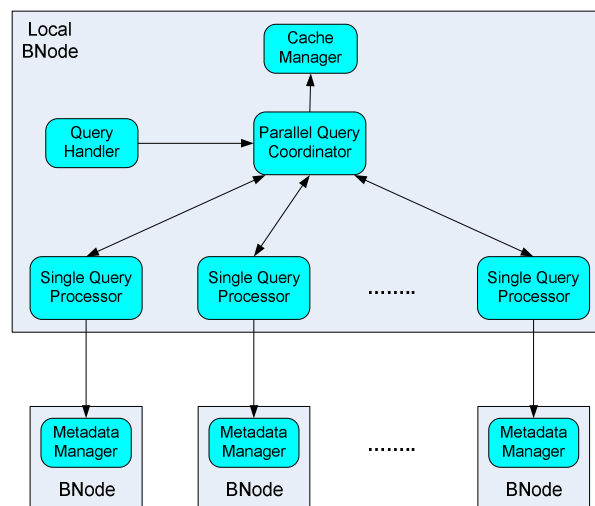
```

PREFIX table: <http://example.org/book/book1#>
SELECT ?title, ?author
FROM <http://example.org/book/book1.owl>
WHERE {
    ?element table:title ?title
    ?element table:author ?author
}
  
```



Page 103

Query Handler: Overall Architecture



Page 104

Archaeological Sites Application Prototype

Finds¹ identification today

- The public is an important source for understanding the past and archeological sites
- Thousands of objects are discovered by people during walking, traveling, gardening, etc.
- Museum curators, archaeologists, students, amateurs all need to identify these finds
- Traditional way of working
 - Examining objects
 - Preliminary identification
 - Comparison to specialist reference collections, e.g. roman coins
 - Sometime object descriptions are enough
 - Sometime it is the starting point for a scientific analysis process
 - Comparison to objects in the museum collection
 - Record object in detail



¹ With *Finds* we mean *archeological finds*

Finds identification today



- Disadvantage of the traditional way
 - The reference collections have to be known by the curator
 - Curator needs information about new collections, e.g. when a museum opens its archive
 - Curators need to access them one by one
 - Collections are heterogeneous from point of view of data, search facilities, user interface, languages, etc.
- As a consequence curators invest a lot of time in the identification process



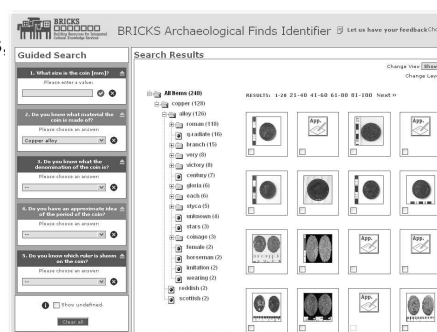
107

Page 107

Finds identification with BRICKS



- Examining objects
- Preliminary identification
- Comparison to specialist reference collections
 - Easy access to all relevant archeological collections (Even to unknown collection)
 - Advanced search facilities
 - Support of ontologies, translation services, different types of search
 - One application to access all collections
 - Collaborative identification processes
- Comparison to objects in the museum collection
- Record object in detail



108

Page 108



Archaeological Sites Application Prototype

DEMO



Page 109



Lessons Learned

- Developers and system designers have a longer learning curve about service oriented technologies
- A service oriented design is different from traditional design, e.g. the number of function should be limited and well considered
- Communication costs between services are often underestimated
- The communication of the concept of distributed architectures to the end user is a hard process
- An early prototype is helpful for the communication between users and developers



Page 110

Questions & Discussions



Page 111

Approaches for Large Scale Digital Library



ISTITUTO DI SCIENZA E TECNOLOGIE
DELL'INFORMAZIONE "A. FAEDO"



Conclusions & Future Directions

Tutorial at RCDL 2007
October, 15th 2007



Thomas Risse
L3S Research Center, Hannover, Germany
Email: risse@L3S.de

In cooperation with Claudia Niederee (L3S), Carlo Meghini (CNR-ISTI), Heiko Schuldt (Uni Basel)

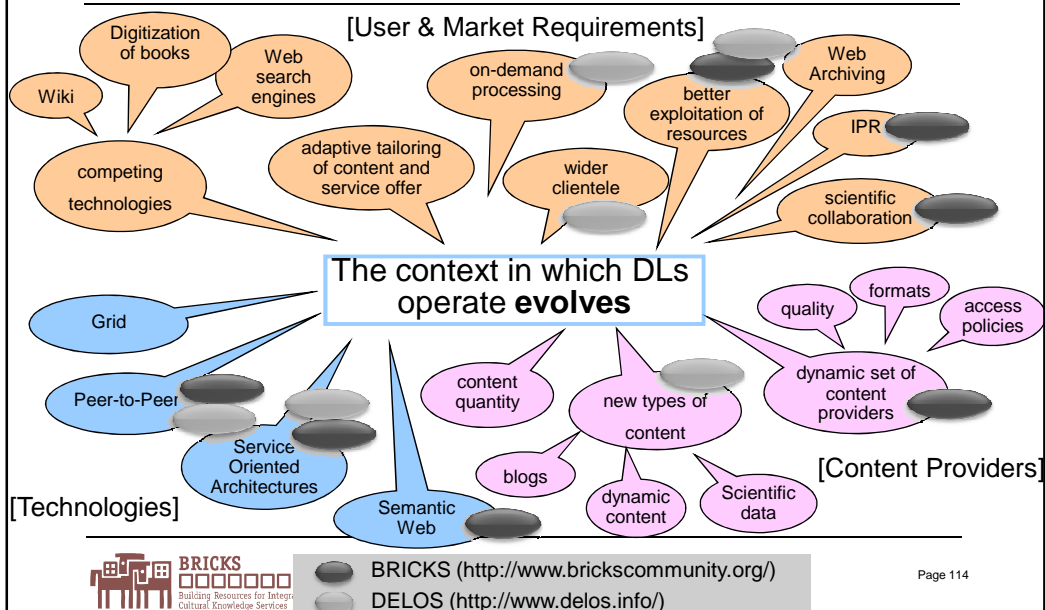
Agenda

1. Introduction & Motivation
2. Challenges of bringing DL to distributed Infrastructures
3. Underlying Technologies and their promises (SOA, P2P)
4. Solutions for decentralized DL infrastructures (with BRICKS Demos)
5. Conclusions and future directions



Page 113

Summary



Page 114

Next Steps for Establishment & Take Up

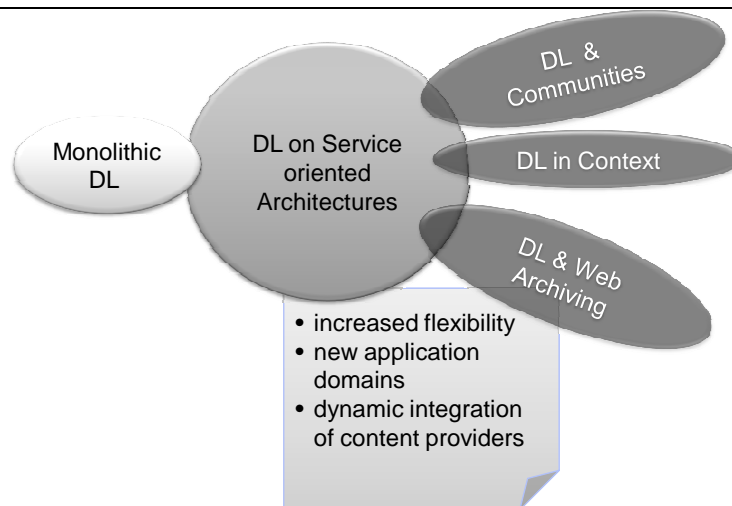


- Consolidation & Standardization
 - compatibility
 - exchangeable modules
- Attraction of a wider community
 - cultural heritage related
 - new application areas (e.g. e-Science)
- Solving of open challenges
 - IPR
 - quality assurance
 - long-term preservation
 - information integration
 - interoperability



Page 115

Future Directions: DL in Context



Page 116

DL & Communities



- Strong trend: Involvement of communities
 - community based content creation (e.g. Wikipedia)
 - community-based content collection (e.g. YouTube, Flickr)
 - community-based recommendation + rating (e.g. online Bookstores)
 - community based tagging (folksonomies, ...)
- DL Opportunities
 - use the intelligence of communities to improve DL content and service offer (e.g. content enrichment and interlinking)
 - exploit DL technologies in community-based systems
 - create even richer content by combining community created content with DL content in seamless ways.



Page 117

DL in Context



- Support for the individual so far:
- standard query based access
 - Personalization support (e.g. considering interests and preferences)
 - Support for personal digital libraries (e.g. Daffodil) and virtual libraries (DILIGENT)

DL Opportunities

- Improved personalization support
- integration of DL functionality in working processes
- targeted pro-active information provision
- requires process- and context modelling
- application example: e-Science



Page 118

DL and Web Archiving



- Increased role of Web in information access, content creation, reflecting opinions, performing transactions etc. in daily life activities, society, culture, politics, education, etc.
- → increasing importance of sound Web archiving
- first generation of Web archiving technology in place
- but: poor access methods for Web archives e.g. archive.org
- little integration with other information sources
- special challenges due to inherent structure w.r.t access
- DL opportunities:
 - apply experience of DL community in metadata, search, etc.
 - integrate Web archive collection with DL collection for more comprehensive information provision
 - develop new types of access methods (time, evolution) and application for targeted communities



Page 119

Thank you for your attention



Page 120