

# Home Archiving: Moving Digital Preservation Capabilities from Large Institutions to SMEs and Home Users

Andreas Rauber

Department of Software Technology and  
Interactive Systems

Vienna University of Technology

[www.ifs.tuwien.ac.at/~andi](http://www.ifs.tuwien.ac.at/~andi)

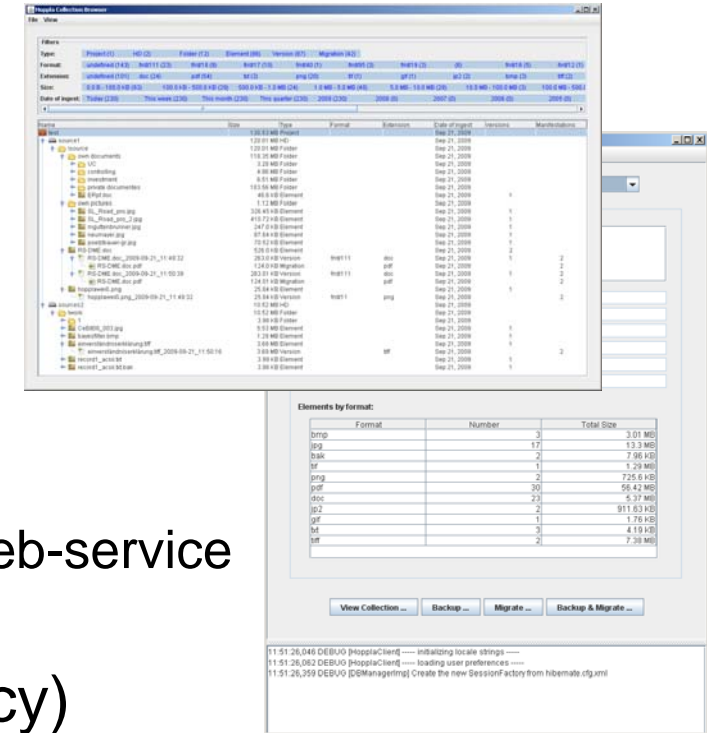
- Digital Preservation R&D predominantly for large, professional institutions (libraries, archives, museums)
- Need for Digital Preservation (DP) solutions in
  - smaller institutions
  - SMEs
  - SOHO: Small Office / Home Office
  - individuals
- Currently hardly any solutions available
- CMS, no real preservation support apart from back-up
- Goal:  
A solution automating preservation activities so that they can be deployed in non-expert settings

# Requirements

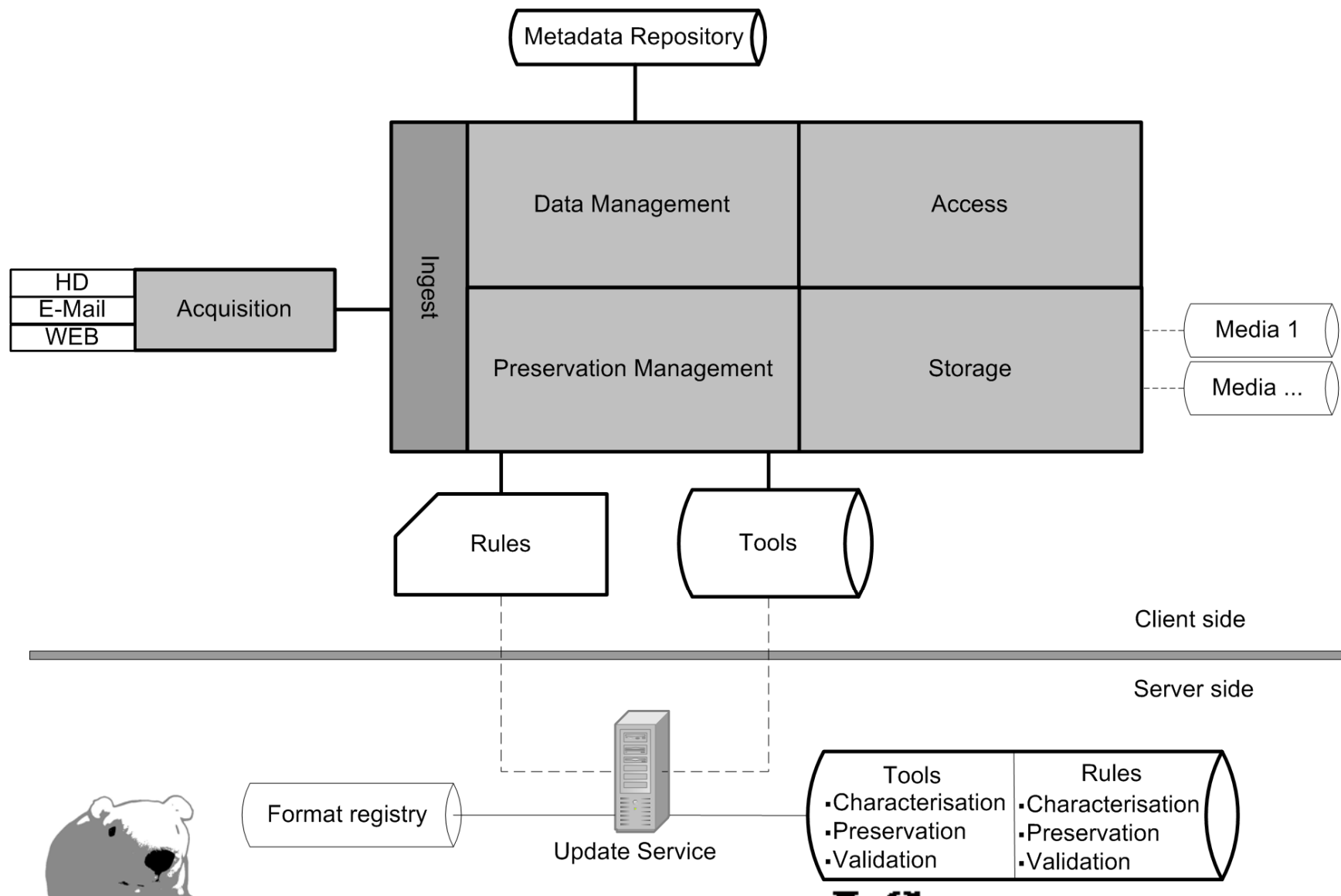
- Operate with little or no DP expertise
- Operate with little or no IT expertise
- Effortless operations: not core business, infrastructure support
- Highly automated
- High level of robustness
  - non-professional interaction with system
  - recovery on complete loss of core management system
- Support for common (also non-archival) media for storage
- “Easier” challenges, more focused collections
- Probably weaker requirements on quality levels

- **HOPPLA:**  
**Home Office Painless Persistent Long-term Archiving**
- Client-server system
- Inspired by: antivirus-SW and software firewall solutions
  - data resides with client
  - server provides DP expertise and solutions
- Support for
  - ingest of data from different sources (home, e-mail, on-line)
  - multiple back-ups on (also low-end) storage (DVD, ext. HDD, RAID systems, on-line storage)
  - recovery on loss of system data
- Focus on robustness and automation
- Meet requirements of audit and certification initiatives

- Developed in Java (plattform-independent)
- Combination of back-up and DP
- Outsourcing of expertise
- Flexible client – server architecture
  - rules
  - tools (mostly plattform-dependent, plus web-service based solutions on request)
- Data remain only on client side (privacy)
- Metadata provided by external experts and automated tools/services



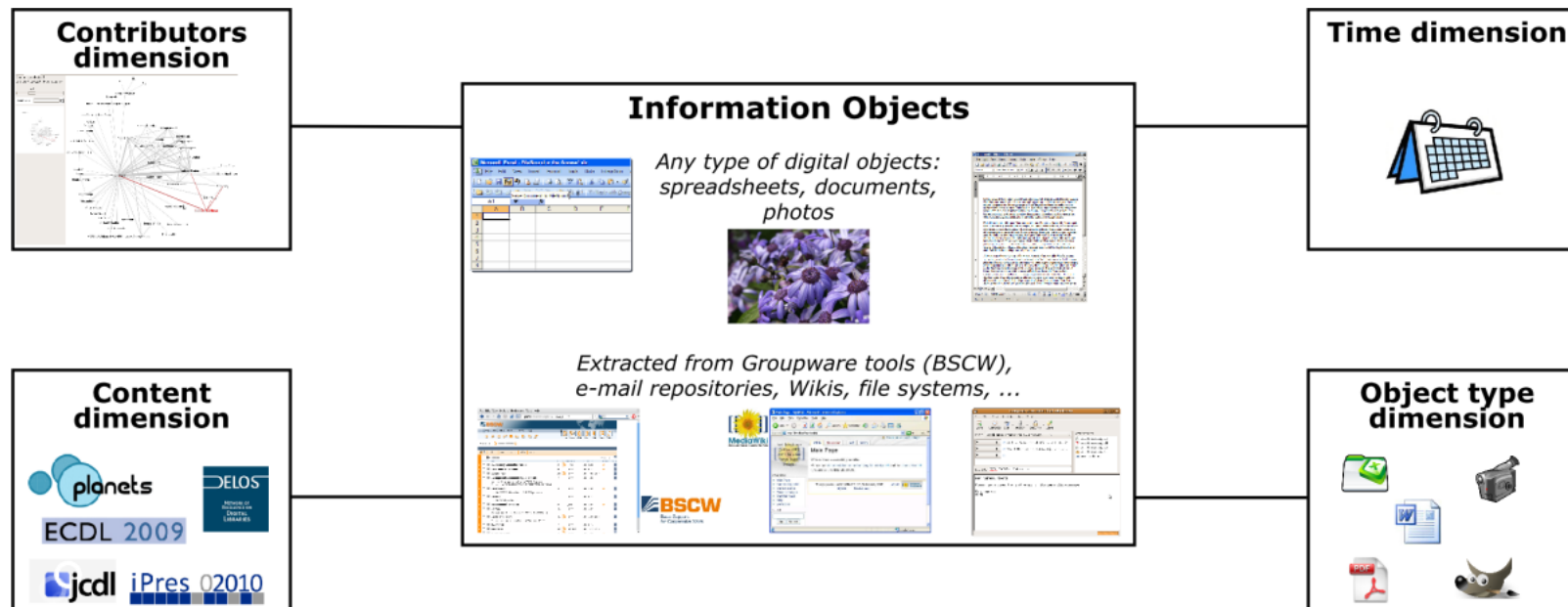
# HOPPLA



## Core Workflow

- Ingest from different source media
  - Creation of collection profile (technical metadata)
  - Extract / estimate contextual information from source info
- Collection profile is sent to server  
(adjustable level of detail – not implemented yet)
  - Experts on server side provide registry with preservation plans (Plato)
  - Appropriate preservation plan is chosen according to user profile (data volumes, risk level, cost/benefit settings)
  - Preservation action plan sent to client
- Client performs migration activities
  - Data stored redundantly
  - Media refresh

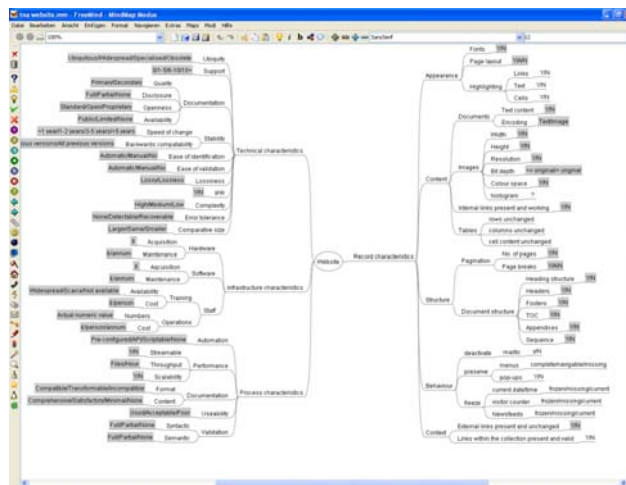
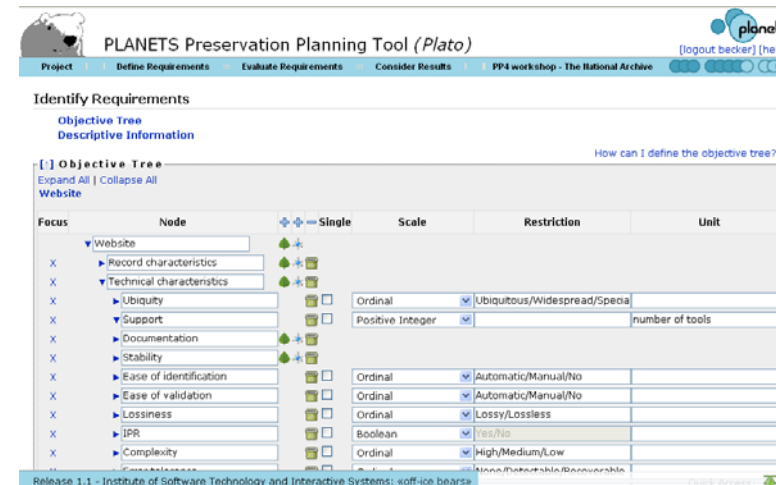
- Technical Metadata: JHOVE, Pronom, XCDL, ...
- Semantic metadata: utilising context of objects
  - Extract context using IR, IE and NLP techniques
  - Organise objects along multiple dimensions (DWH-inspired)
  - Finding groups of related objects (semi-automatic)





## Server-side: Preservation Planning

- Plato Preservation Planning Tool
- Implements Planets Preservation Planning Workflow
- Allows creation of objective tree
  - within application or via import of mindmaps
- Allows the selection of Preservation action tools to be evaluated

PLANETS Preservation Planning Tool (Plato)

Project: Define Requirements Evaluate Requirements Consider Results PP4 workshop - The National Archive

Identify Requirements

Objective Tree  
Descriptive Information

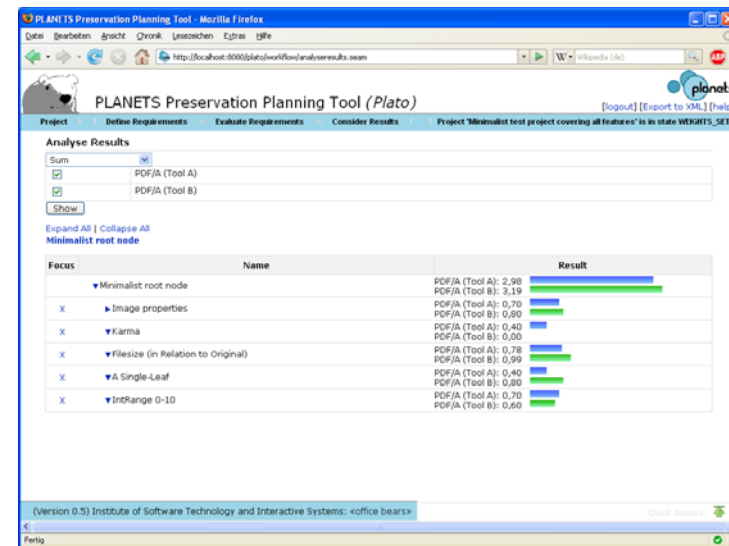
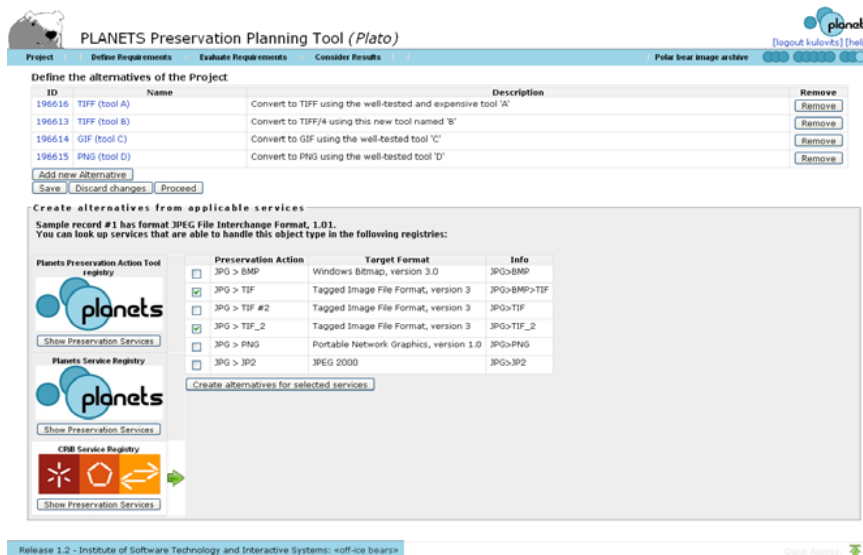
How can I define the objective tree?

[+] Objective Tree  
Expand All | Collapse All  
Website

Focus	Node	Single	Scale	Restriction	Unit
	Website				
X	Record characteristics				
X	Technical characteristics				
X	Ubiquity	<input type="checkbox"/>	Ordinal	Ubiquitous/Widespread/Special	
X	Support	<input type="checkbox"/>	Positive Integer		number of tools
X	Documentation				
X	Stability				
X	Ease of identification	<input type="checkbox"/>	Ordinal	Automatic/Manual/No	
X	Ease of validation	<input type="checkbox"/>	Ordinal	Automatic/Manual/No	
X	Lossiness	<input type="checkbox"/>	Ordinal	Lossy/Lossless	
X	IPR	<input type="checkbox"/>	Boolean	Yes/No	
X	Complexity	<input type="checkbox"/>	Ordinal	High/Medium/Low	

Release 1.1 - Institute of Software Technology and Interactive Systems: koff-ice bears

- Runs preservation action experiments, documents results
- Allows definition of transformation rules, weightings
- Performs evaluation, sensitivity analysis,
- Provides recommendation (ranks solutions)
- Stores/exports preservation plan with evidence, triggers,...



- Preservation Planning to ensure “optimal” preservation
- Operated by experts on server side
- A simple, methodologically sound model to specify and document requirements
- Repeatable and documented evaluation
- Basis for well-informed, accountable decisions
- Concretization of OAIS model
- Follows recommendations of TRAC and nestor
- Plato:
  - tool support to perform solid, well-documented analyses
  - creates core preservation plan
- Rules to match preservation requirements with plans

- Server side has risk profile for object types
- Client side has preferences with
  - user / institution implicit preferences (e.g. obj. importance, ...)
  - degree of risk avoidance based on object type, file size, ...
  - preferences in terms of storage space availability / cost
- Objects identified at some risk level are matched with preservation plans
- List of potential preservation actions:
  - tools installed on client side
  - tools wrapped as plug-ins
  - external tools potentially to be installed at client side (license / cost / willingness to install)
  - external webservice (if client is willing to send data)
- Identify most suitable solution

- Server receives form client
  - collection profile (technical metadata)
  - preferences settings
  - list of locally available tools
  - level of detail will be configurable to meet privacy requirements
  
- Server provides to client
  - Preservation action plan: recommended preservation actions
  - Preservation plan for documentation / evidence
  - wrapped migration tools or install packages (potentially also emulators)
  
- Client
  - performs migration actions and redundant storage

## Preservation Action Plan

[...]

<migrationToolId>1</migrationToolId>

<constraint>

<metadataConstraints><metadataConstraint>

<characterisationRuleId>1</characterisationRuleId>

<comparisonMode>Smaller</comparisonMode>

<metadataId>height</metadataId>

<constraintValue>10000</constraintValue>

</metadataConstraint><metadataConstraint>

[...]

<minSize>3</minSize>

<maxSize>10</maxSize>

[...]

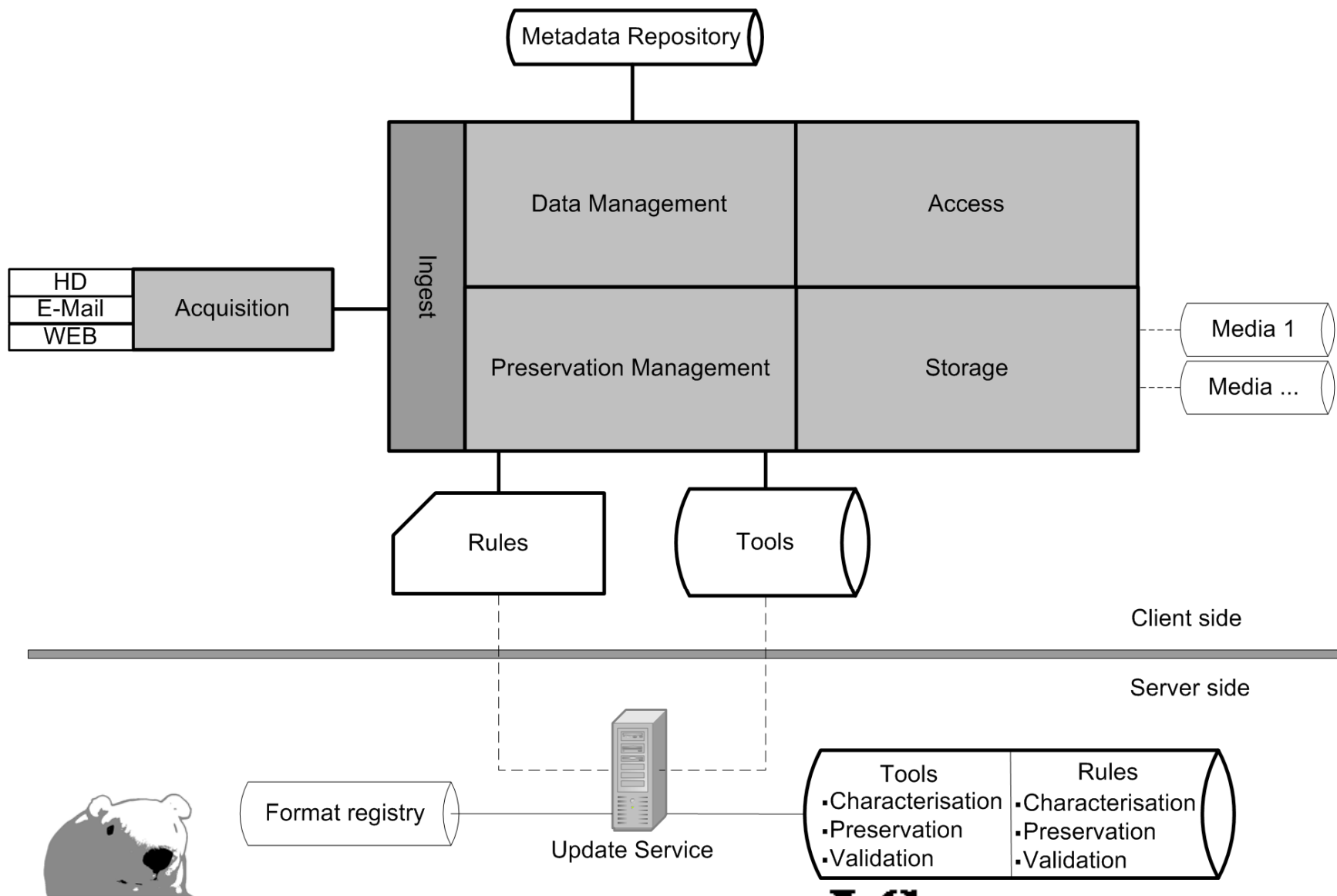
<riskScore>bestpractice</ riskScore>

<estChangeInSize>50%</estChangeInSize>

<estDurationPerMb>5</estDurationPerMb>

[...]

# HOPPLA



## Storage

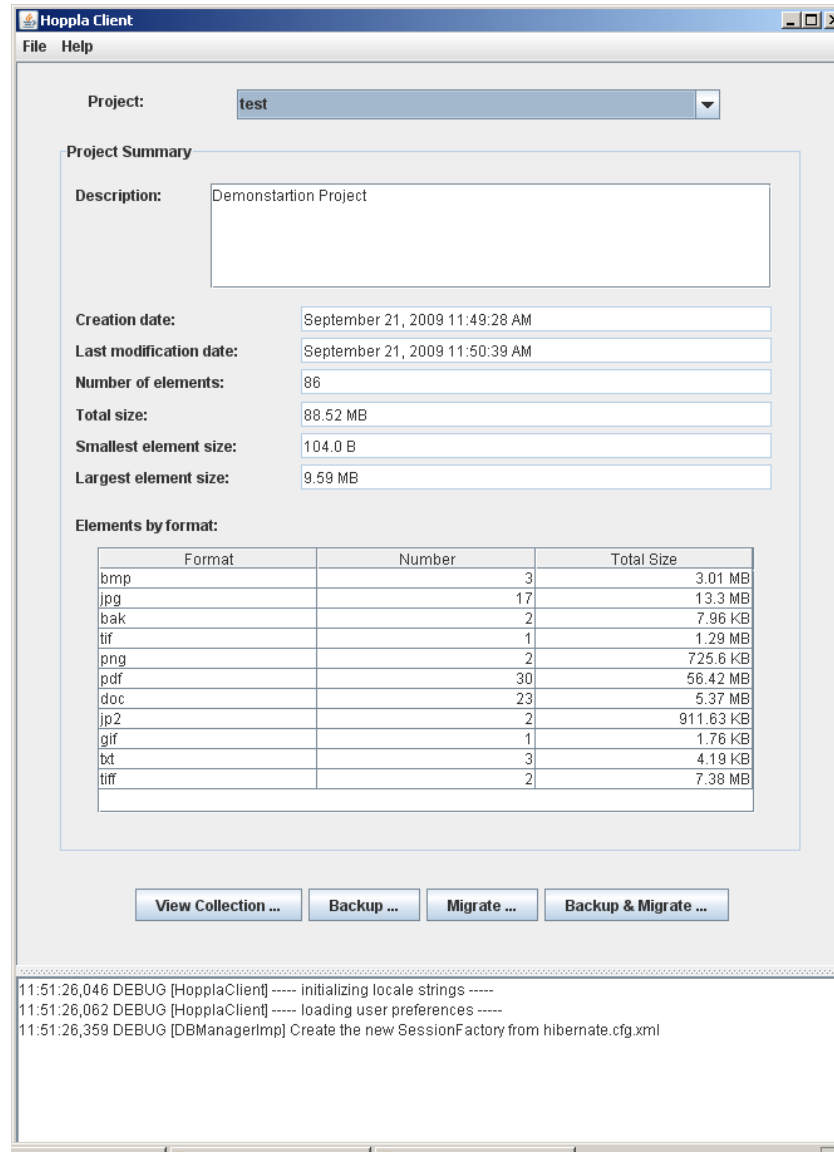
- Re-create source hierarchy
  - directory structure
  - mail folder structure
  - complex objects with sub-directories
- Locate each object in respective location
- Embed some metadata in filename (version, timestamp)
- Add XML file per directory with all object metadata in data management
- Bit-level preservation:  
configure degree of redundancy and media types
- media refreshment via reminders / automatic
- Can be used directly (even if somewhat limited) without an system



## Access

- Via collection browser
  - representation of directory structure
  - archived versions and time-stamps
  - migrated versions and time stamps
  - operates via data management, faceted browser
- Can be used directly (even if somewhat limited) without any HOPPLA system
  - needs to be able to mount file system (FAT, ISO9660, on-line)
  - use physical directory structure
  - limited set of metadata in filenames
  - detailed information (basis for recovery) in XML file
- Preservation Plans as regular objects in dedicated directory

# Hoppla - Software Prototype



The screenshot shows the Hoppla Client application window. At the top, there is a menu bar with 'File' and 'Help'. Below the menu bar, a 'Project:' dropdown menu is set to 'test'. The main area is titled 'Project Summary' and contains several fields: 'Description' (text area with 'Demonstartion Project'), 'Creation date' (September 21, 2009 11:49:28 AM), 'Last modification date' (September 21, 2009 11:50:39 AM), 'Number of elements' (86), 'Total size' (88.52 MB), 'Smallest element size' (104.0 B), and 'Largest element size' (9.59 MB). Below these fields is a table titled 'Elements by format' with columns for 'Format', 'Number', and 'Total Size'. At the bottom of the main area, there are four buttons: 'View Collection ...', 'Backup ...', 'Migrate ...', and 'Backup & Migrate ...'. The bottom of the window shows a log window with the following text:

```

11:51:26,046 DEBUG [HopplaClient] ----- initializing locale strings -----
11:51:26,062 DEBUG [HopplaClient] ----- loading user preferences -----
11:51:26,359 DEBUG [DBManagerImp] Create the new SessionFactory from hibernate.cfg.xml

```



# Hoppla - Software Prototype



**Hoppla Collection Browser**

File View

**Filters**

**Type:** Project (1) HD (2) Folder (12) Element (86) Version (87) Migration (42)

**Format:** undefined (143) fmt/111 (23) fmt/18 (9) fmt/17 (10) fmt/40 (1) fmt/95 (3) fmt/19 (3) (6) fmt/16 (5) fmt/12 (1)

**Extension:** undefined (101) doc (24) pdf (54) txt (3) png (20) tif (1) gif (1) jpeg (2) bmp (3) tiff (2)

**Size:** 0.0 B - 100.0 KB (93) 100.0 KB - 500.0 KB (29) 500.0 KB - 1.0 MB (24) 1.0 MB - 5.0 MB (48) 5.0 MB - 10.0 MB (28) 10.0 MB - 100.0 MB (3) 100.0 MB - 500.0 MB (0)

**Date of ingest:** Today (230) This week (230) This month (230) This quarter (230) 2009 (230) 2008 (0) 2007 (0) 2006 (0) 2005 (0)

Name	Size	Type	Format	Extension	Date of ingest	Versions	Manifestations
test	130.53 MB	Project			Sep 21, 2009		
source1	120.01 MB	HD			Sep 21, 2009		
source1\source	120.01 MB	Folder			Sep 21, 2009		
source1\source\own documents	118.35 MB	Folder			Sep 21, 2009		
source1\source\own documents\UC	3.28 MB	Folder			Sep 21, 2009		
source1\source\own documents\controlling	4.96 MB	Folder			Sep 21, 2009		
source1\source\own documents\investment	6.51 MB	Folder			Sep 21, 2009		
source1\source\own documents\private documentes	103.56 MB	Folder			Sep 21, 2009		
source1\source\own documents\ERpf.doc	46.6 KB	Element			Sep 21, 2009	1	
source1\source\own pictures	1.12 MB	Folder			Sep 21, 2009		
source1\source\own pictures\SL_Road_pro.jpg	326.45 KB	Element			Sep 21, 2009	1	
source1\source\own pictures\SL_Road_pro_2.jpg	410.72 KB	Element			Sep 21, 2009	1	
source1\source\own pictures\mguttenbrunner.jpg	247.0 KB	Element			Sep 21, 2009	1	
source1\source\own pictures\neumayer.jpg	87.64 KB	Element			Sep 21, 2009	1	
source1\source\own pictures\poelzlbauer-gr.jpg	70.52 KB	Element			Sep 21, 2009	1	
source1\source\RS-DME.doc	526.0 KB	Element			Sep 21, 2009	2	
source1\source\RS-DME.doc_2009-09-21_11:49:32	263.0 KB	Version	fmt/111	doc	Sep 21, 2009	1	2
source1\source\RS-DME.doc.pdf	124.0 KB	Migration		pdf	Sep 21, 2009		2
source1\source\RS-DME.doc_2009-09-21_11:50:39	263.01 KB	Version	fmt/111	doc	Sep 21, 2009	1	2
source1\source\RS-DME.doc.pdf	124.01 KB	Migration		pdf	Sep 21, 2009		2
source1\source\hopplaweiß.png	25.84 KB	Element			Sep 21, 2009	1	
source1\source\hopplaweiß.png_2009-09-21_11:49:32	25.84 KB	Version	fmt/111	png	Sep 21, 2009		2
sources2	10.52 MB	HD			Sep 21, 2009		
sources2\work	10.52 MB	Folder			Sep 21, 2009		
sources2\work\1	3.98 KB	Folder			Sep 21, 2009		
sources2\work\1\CeBit08_003.jpg	5.53 MB	Element			Sep 21, 2009	1	
sources2\work\1\bayesfilter.bmp	1.29 MB	Element			Sep 21, 2009	1	
sources2\work\1\einverständniserklärung.tiff	3.69 MB	Element			Sep 21, 2009	1	
sources2\work\1\einverständniserklärung.tiff_2009-09-21_11:50:16	3.69 MB	Version		tiff	Sep 21, 2009		2
sources2\work\1\record1_acsii.txt	3.99 KB	Element			Sep 21, 2009	1	
sources2\work\1\record1_acsii.txt.bak	3.98 KB	Element			Sep 21, 2009	1	

**Hoppla Collection Browser**

File View

**Filters**

Type: Project (1) HD (2) Folder (12) Element (86) Version (87) Migration (42)

Format: undefined (143) fmt/111 (23) fmt/18 (9) fmt/17 (10) fmt/40 (1) fmt/95 (3) fmt/19 (3) (6) fmt/16 (5) fmt/12 (1)

Extension: undefined (101) doc (24) pdf (54) txt (3) png (20) tif (1) gif (1) jp2 (2) bmp (3) tiff (2)

Size: 0.0 B - 100.0 KB (93) 100.0 KB - 500.0 KB (29) 500.0 KB - 1.0 MB (24) 1.0 MB - 5.0 MB (48) 5.0 MB - 10.0 MB (28) 10.0 MB - 100.0 MB (3) 100.0 MB - 500.0 MB (0)

Date of ingest: Today (230) This week (230) This month (230) This quarter (230) 2009 (230) 2008 (0) 2007 (0) 2006 (0) 2005 (0)

Name	Full Path Name	Size	Type	Format	Extension	Date of ingest	Versions	Manifestati...
bayesfilter.bmp_2009-09-21_11:49:32	testsource1\source\down ...	1.29 MB	Version	fmt/116	bmp	Sep 21, 2009		2
bayesfilter.bmp_2009-09-21_11:50:16	testsources2\work\bayes...	1.29 MB	Version	fmt/116	bmp	Sep 21, 2009		2
BroschuereJobcard_02.pdf_2009-09-21_11:49:32	testsource1\source\down ...	1.41 MB	Version	fmt/17	pdf	Sep 21, 2009		2
CeBit08_001.jpg.png	testsource1\source\down ...	4.53 MB	Migration		png	Sep 21, 2009		2
CeBit08_002.jpg.png	testsource1\source\down ...	4.24 MB	Migration		png	Sep 21, 2009		2
CeBit08_003.jpg.png	testsource1\source\down ...	4.19 MB	Migration		png	Sep 21, 2009		2
CeBit08_003.jpg.png	testsources2\work\CeBit...	4.19 MB	Migration		png	Sep 21, 2009		2
CeBit08_004.jpg.png	testsource1\source\down ...	4.16 MB	Migration		png	Sep 21, 2009		2
CeBit08_006.jpg.png	testsource1\source\down ...	4.41 MB	Migration		png	Sep 21, 2009		2
CeBit08_007.jpg.png	testsource1\source\down ...	2.61 MB	Migration		png	Sep 21, 2009		2
CeBit08_007.jpg_2009-09-21_11:49:32	testsource1\source\down ...	3.7 MB	Version	fmt/43	jpg	Sep 21, 2009	1	2
CeBit08_008.jpg.png	testsource1\source\down ...	3.77 MB	Migration		png	Sep 21, 2009		2
CeBit08_009.jpg.png	testsource1\source\down ...	3.03 MB	Migration		png	Sep 21, 2009		2
CeBit08_009.jpg_2009-09-21_11:49:32	testsource1\source\down ...	4.14 MB	Version	fmt/43	jpg	Sep 21, 2009	1	2
ControllingU0610.pdf_2009-09-21_11:49:32	testsource1\source\down ...	1.26 MB	Version	fmt/17	pdf	Sep 21, 2009		2
EinfuehrungSoftwareArchitekturen.pdf_2009-09-21_11:49:32	testsource1\source\down ...	2.58 MB	Version	fmt/18	pdf	Sep 21, 2009		2
einverstaendniserklaerung.tiff_2009-09-21_11:49:32	testsource1\source\down ...	3.69 MB	Version		tiff	Sep 21, 2009		2
einverstaendniserklaerung.tiff_2009-09-21_11:50:16	testsources2\work\keinver...	3.69 MB	Version		tiff	Sep 21, 2009		2
Finanzwirtschaft0511.pdf_2009-09-21_11:49:32	testsource1\source\down ...	2.2 MB	Version	fmt/17	pdf	Sep 21, 2009		2
KerzenEWF0606.pdf_2009-09-21_11:49:32	testsource1\source\down ...	1.03 MB	Version	fmt/17	pdf	Sep 21, 2009		2
Kopie von CeBit08_008.jpg.png	testsource1\source\down ...	3.77 MB	Migration		png	Sep 21, 2009		2
MDA.pdf_2009-09-21_11:49:32	testsource1\source\down ...	1.49 MB	Version	fmt/18	pdf	Sep 21, 2009		2
MySQL&PHP0601.pdf_2009-09-21_11:49:32	testsource1\source\down ...	2.25 MB	Version	fmt/17	pdf	Sep 21, 2009		2
Planets_Implementation_Plan_Fina_I.doc_2009-09-21_11:49:...	testsource1\source\down ...	3.26 MB	Version	fmt/111	doc	Sep 21, 2009	1	2
ReWe0510.pdf_2009-09-21_11:49:32	testsource1\source\down ...	2.22 MB	Version	fmt/17	pdf	Sep 21, 2009		2
rufzeichen_tu_4c_uk.tif_2009-09-21_11:49:32	testsource1\source\down ...	1.29 MB	Version		tif	Sep 21, 2009		2
Testbed Use Cases v0.9.doc_2009-09-21_11:49:32	testsource1\source\down ...	1.87 MB	Version	fmt/111	doc	Sep 21, 2009	1	2

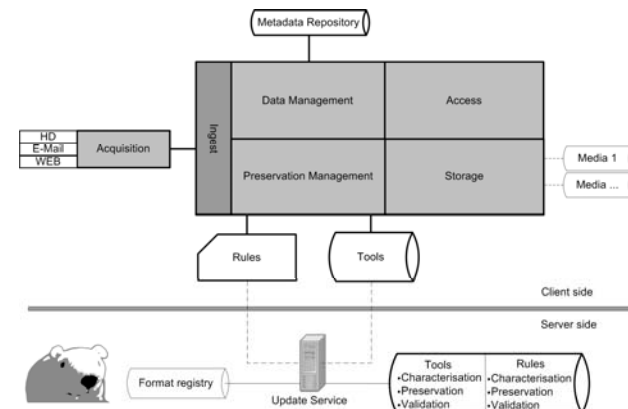
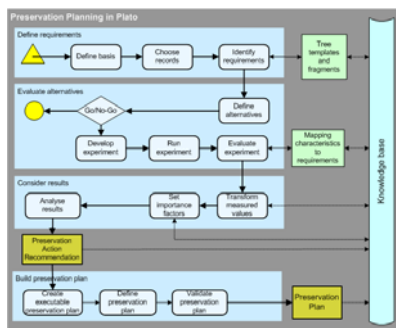
## Next steps

- Prototype development via DME-funded project
- Co-funding for adapting and integrating prototype in partner systems
- First functional prototype winter 2009 (internal, partner)
- Tighter coupling of preservation actions on client side and preservation planning on server side
- Design adaptations to allow more flexible integration into different systems
  - repository solutions
  - back-up systems
- Eventually better support for audit trails

- Digital Preservation is a challenge for everybody
- Lack of solutions for small institutions / individuals
- Digital Preservation as a service
- Automation: metadata creation, preservation actions
- Server side:
  - Preservation planning
  - Tool provisioning
- Covering bit preservation and logical preservation
- Flexible adaption to needs via rule-based mappings
- Outsourcing of expertise

<http://www.ifs.tuwien.ac.at/dp/hoppla>

# Thank you!



<http://www.ifs.tuwien.ac.at/hoppla>  
<http://www.ifs.tuwien.ac.at/dp/plato>

