**SHAMAN and Storage Virtualization**
Adil Hasan
Univ. of Liverpool
Toulouse Workshop on Storage Virtualization
2nd June 2009

**SHAMAN Collaborators:**

# What is SHAMAN?

- Sustaining HeritAge through Multivalent ArchiviNg.

- FP7 EU Integrated Project start Dec/07 finish Dec/11.

- Aim to investigate long-term preservation of data-sets.

# What is SHAMAN?

- Issues:
  - Accommodate unknown changes to hardware (infrastructure).
  - Accommodate unknown changes to preservation tools (processes).
  - Accommodate unknown changes to format and description of data (content).

- Hardware will change:

  - Provide infrastructure layer on-top of hardware.

  - Encapsulate changes, only hardware layer worries about changes.

  - Uniform interface to hardware; driver maps from hardware interface to infrastructure interface.

  - Logical-to-physical object mapping; insulates from changes in data location.

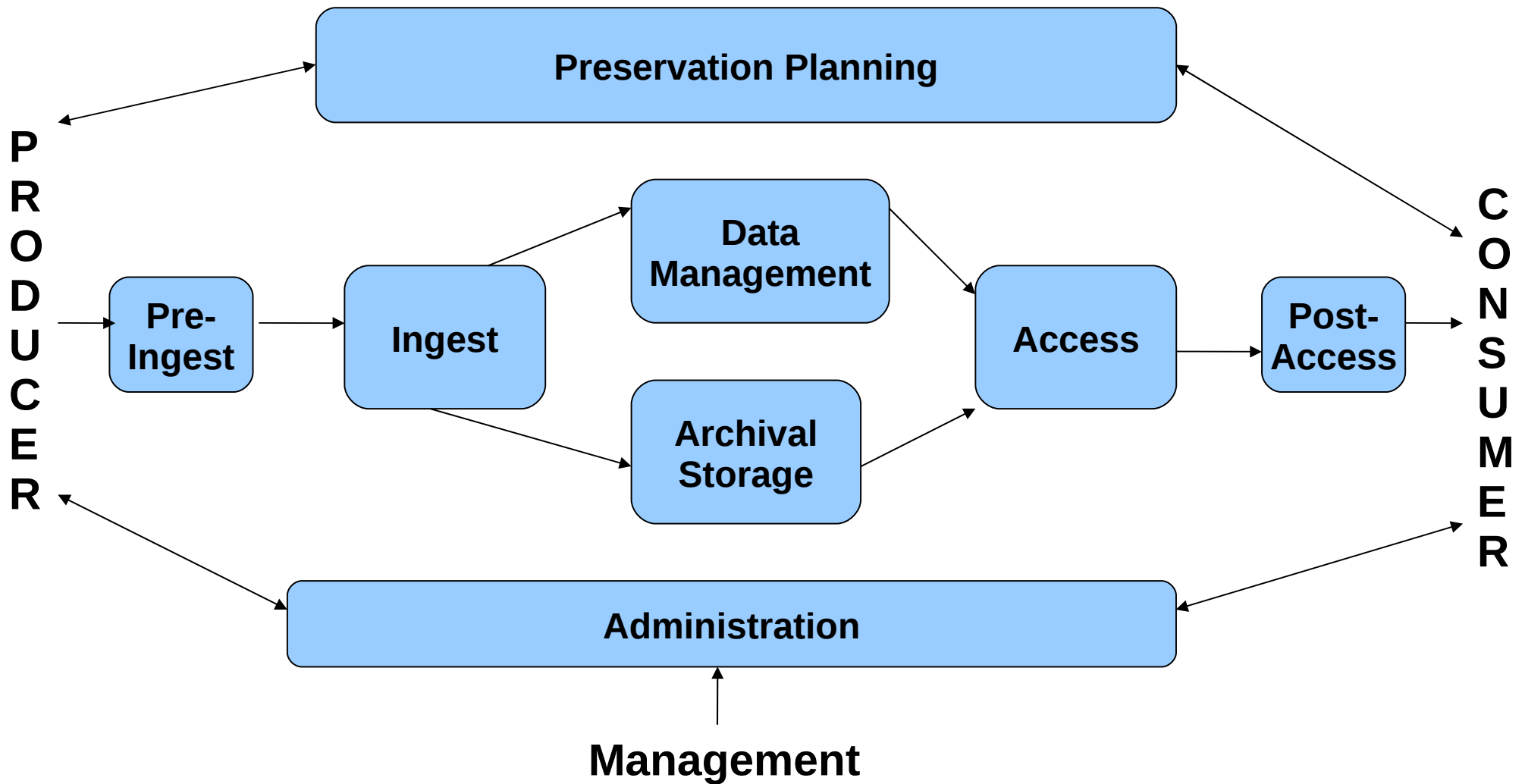  - Record characteristics of hardware (for debugging).

- Components of process will change:

  - Use abstract language to describe preservation processes (these are the policies).

  - Translate from abstract language to actual work-flows.

  - Insulates from changes to work-flow engines or services.

  - Abstract description of processes must be preserved.

- Format and description of data will change:

  - Try to keep data in original format (save space), render to user with migratable tool.

  - Migrate data-sets that cannot be rendered reliably (closed source format, complex, etc).

  - Use standards for description of data.

  - Ensure dictionary of terminology (domain-specific terms) archived.

- Preservation of a digital object requires:

  - data-format, data-description, processes, hardware

- Can regard this collection as the context of the data.

- Want to capture this context in a single unit; the Preservation Description Information.

- PDI must be related to data.

- PDI itself must be preserved.

- OAIS describes long-term archive system.

- Used in many projects. Is the basis for SHAMAN.

- Needs improvements:

  - Preserving process information necessary.

  - Pre-Ingest phase required (recognise importance of assembling data for ingest).

  - Post-Access phase required (recognise importance of further processing after extraction of data).

- Further observations:

  - Preservation Planning must encompass Ingest and Access.

  - Preservation Planning requires input from Producer, Consumer and Management.

  - Roles must be further refined to understand better mapping to existing business roles.

- Essentially 3-types of storage:

  - Ingest Storage – holds ingested data where AIP is created. Temporary Storage

  - Dissemination Storage – holds data for access where DIP is assembled. Temporary Storage.

  - Archival Storage – holds archived data and description (i.e. AIP). Permanent Storage.

- Ingest and Dissemination are most-likely frequent-access systems.

  - Viewed as cache-systems for the Archival Storage.

  - System must have good performance.

- Archival storage less-frequently access, must be reliable.

- Amount of digital data increasing rapidly.

- Many institutions collaborate.

- Makes sense to make use of collaborators storage.

- Allows storage to scale with data volume.

- Geographically distributed data guards against storage failure.

- Data Grids provide a means of combining distributed resources into logical resource.

- Insulate storage from access to storage.

  - Provide uniform access to resources through drivers.

- Removes dependence on physical location through logical-physical file mapping.

- Virtualization of storage.

- Observations:

  - Data Grid resource may not be part of the collaboration. May be third-party.

  - Need SLA on provision of resource such that it is possible to replace resource with different type (e.g. cloud).

  - Need to ensure SLA is implemented (evidence).

  - Need to make sure preservation processes are sufficiently encapsulated to run on collaborator/external resources.

- Data Grid provides distributed storage.

- Need to manage that storage. To apply policies to that storage.

- IRODS – policy-driven data management system.

- Allows virtualization of storage AND virtualization of policies.

- Rule-engine manages the policies.
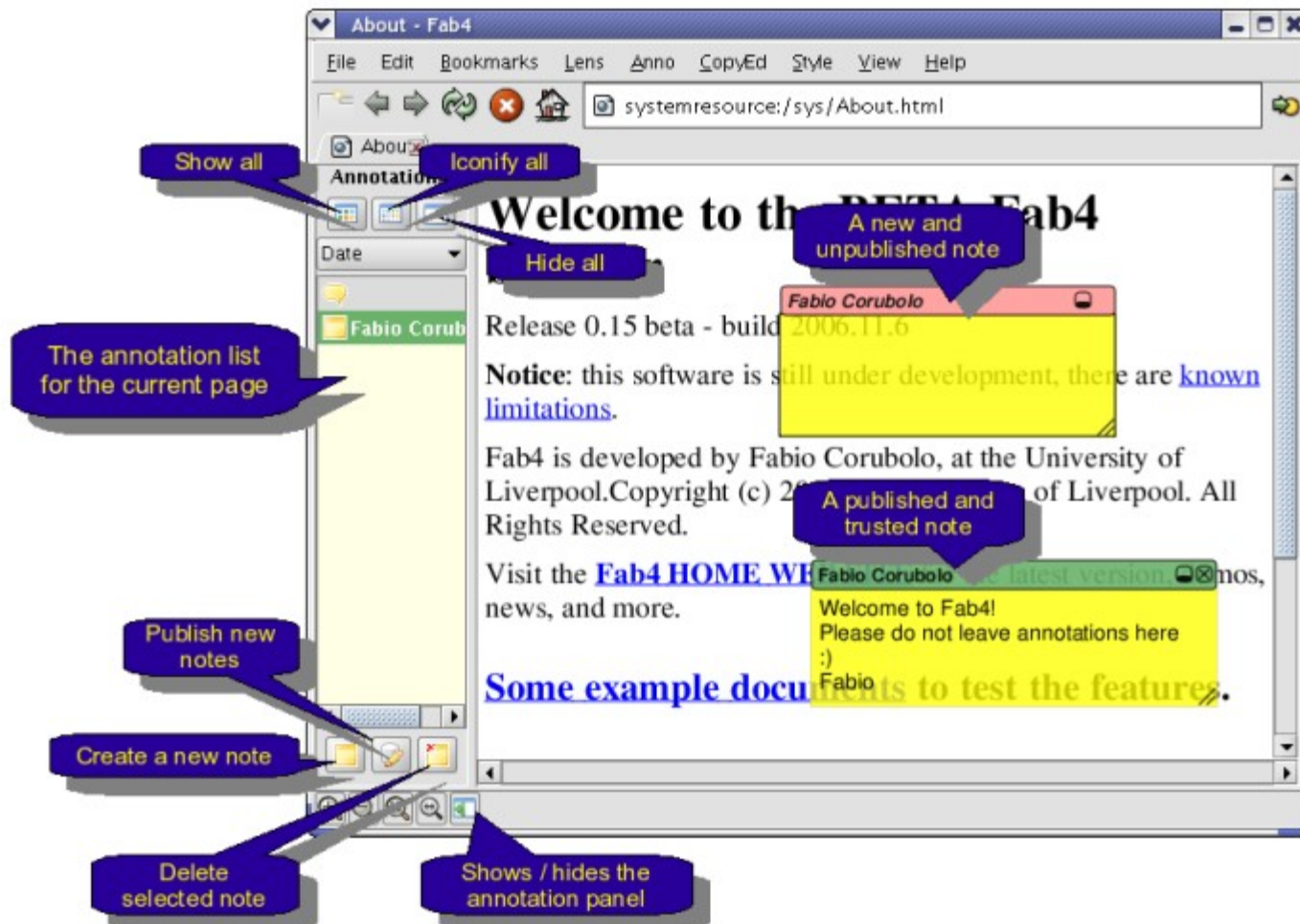
- Policies implemented as rules.

- Can view iRODS rule-engine as managing the preservation processes:

  - Replicate data, checksum data, transform data, extract metadata etc.

- Require tool to transform from abstract policies to rules (in preparation).

- Rule execution must be logged (which rule, microservice, result, when run) for auditing.

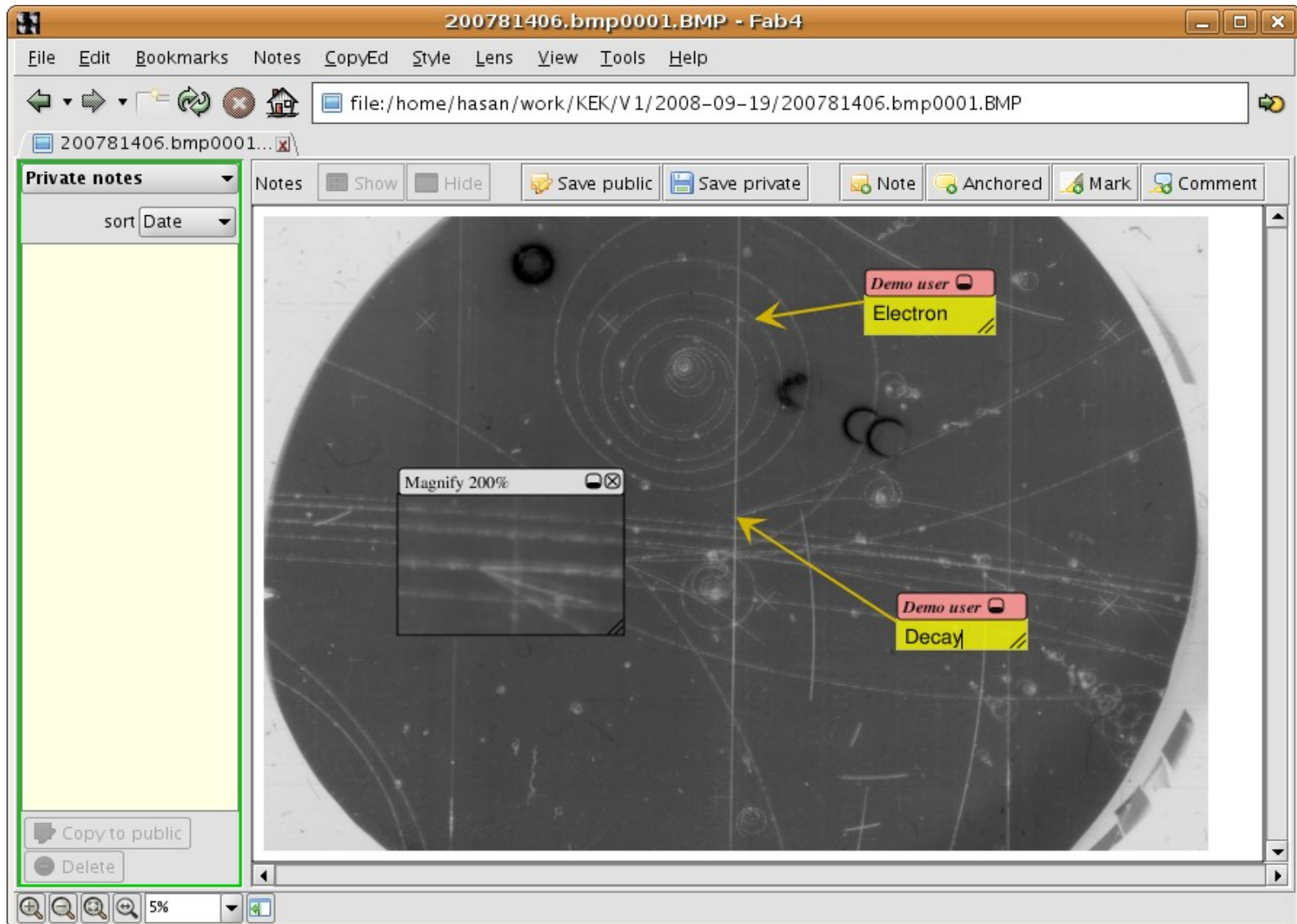- Validation of each microservice for auditing.

- Cloud, such as Amazon S3 can be viewed as another storage resource within the data grid.

- SLA would ensure Cloud resources provide necessary functionality.

  - Would need to ensure cloud providers are capable of providing mechanism to ensure SLA enforced.

  - Checks can be implemented as rules in iRODS.

- Long-term access can be achieved by:

- Migration to new formats

  - Requires more storage to hold original and migrated copy. Intensive process.

  - May be only option for some closed formats.

- Rendering of old format

  - Tool capable of reading old format means 'migration' on demand.

  - Need to ensure all the required properties of data are accessible.

- SHAMAN make use of Multivalent tool:

  - Renders data in old format to user.

  - Written in Java.

  - Capable of reading different formats through different drivers ('media engines').

- Plan to archive tool and package-up tool and data in original format for download to allow access to data.

- Fab4 browser interface to Multivalent.

- Allows annotations to be made on document.

- Annotations saved in separate file from document.

- Data not altered.

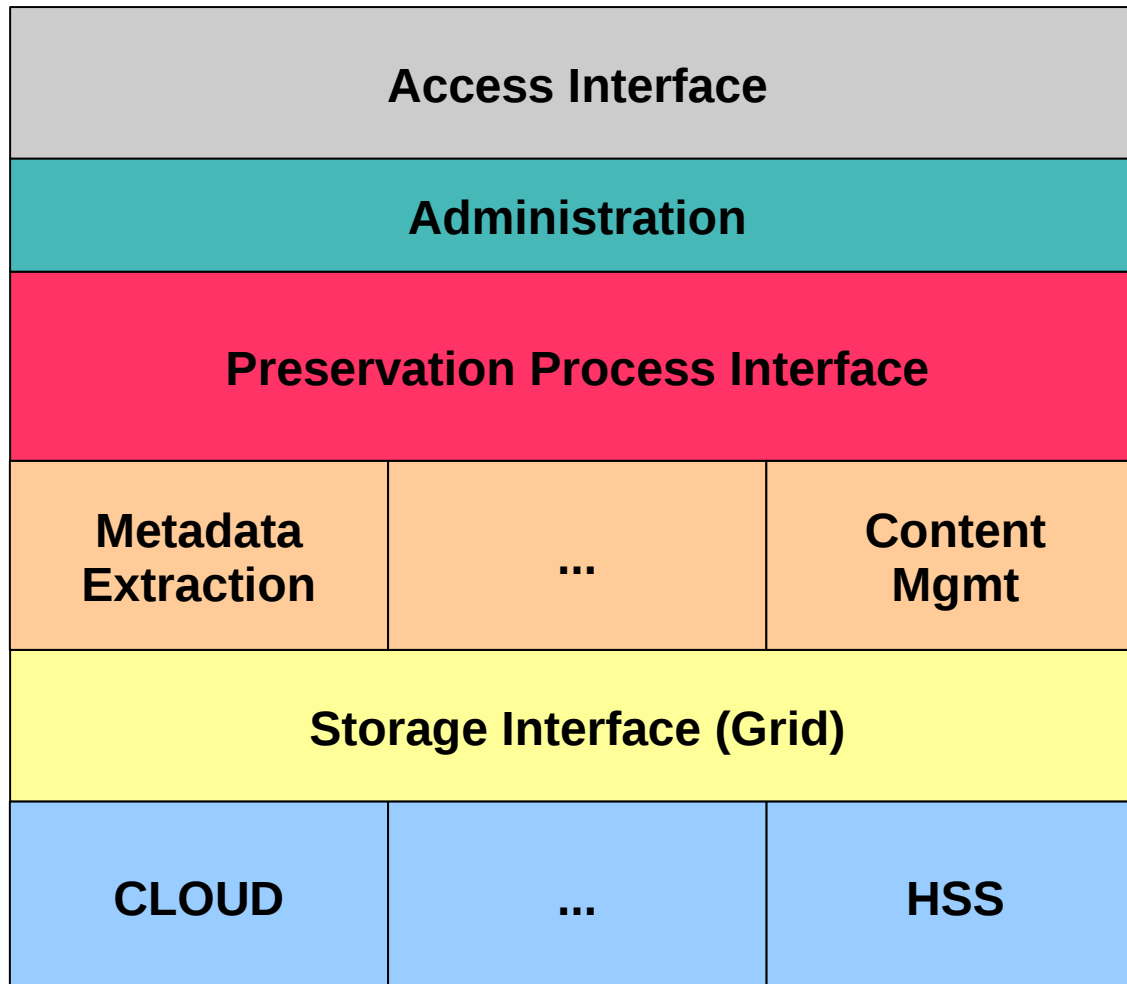- Also apply different behaviours such as magnifying lens

- Annotations semantically attached to document.

  - For images will need to be re-worked.

- Allows annotations and behavours to appear in correct position in document in different formats.

- Successful discovery requires as much information about the digital object as possible is supplied during ingest.

- Also requires information is extracted and indexed.

- SHAMAN make use of powerful Cheshire digital library system to extract information such that it is discoverable.

- Processes involved in extraction must be preserved.

# Preservation Layers

| Access Interface |
|---|

Provides uniform access to data.

| Administration |
|---|

Manages access & system, Transforms pres processes from abstract

| Preservation Process Interface |
|---|

Interface provides uniform access to different preservation processes.

| Metadata Extraction | ... | Content Mgmt |
|---|---|---|

| Storage Interface (Grid) |
|---|

Grid interfaces to different Types of storage. Provides Uniform Interface

| CLOUD | ... | HSS |
|---|---|---|

SHAMAN
Sustaining Heritage Access through Multivalent ArchiviNg

- To ensure data usable in the long-term:

  - Insulate from hardware changes.

  - Insulate from changes to processes.

  - Insulate from changes to data format.

  - Insulate from changes to description.

  - Ensure as much information as possible about data is captured.

    - Ideally test data is understandable without ANY external dependencies.

- SHAMAN aims to provide a framework that accounts for these issues.