

Database preservation toolkit:

a flexible tool to normalize
and give access to databases

DLM Forum: “Making the Information Governance Landscape in Europe”

**José Carlos Ramalho jcr@keep.pt
KEEP SOLUTIONS // University of Minho**

Nov. 12-14, 2014, Lisbon, Portugal

KEEP SOLUTIONS: Projects



- DigitArq, CRAV (2003..[2008-2012])
- RODA (2006..[2008-...])
- RCAAP (2008-...)
- PPA (2009)
- Open source: RODA, KOHA, DSpace, Moodle, etc.
- Scientific research (european projects)
 - **SCAPE**: Large scale preservation
 - **4C**: the cost for curation
 - **e-ark**: presented by Kuldar



<http://www.keep.pt>

Database Preservation Toolkit

Developed within RODA project

Now stand-alone open-source project

<http://keeps.github.io/db-preservation-toolkit/>

Imports and exports between DBMS and DB formats

Supports preservation formats: DBML, SIARD

The Past: 2006 - 2009

RODA: Repositório de Objectos Digitais Autênticos

José Carlos Ramalho
jcr@di.uminho.pt

Francisco Barbedo
frbarbedo@iantt.pt

Miguel Ferreira
mferreira@dsi.uminho.pt

Cecília Henriques
chenriques@iantt.pt

Rui Castro
Rcastro@iantt.pt

Luis Corujo
lcorujo@iantt.pt

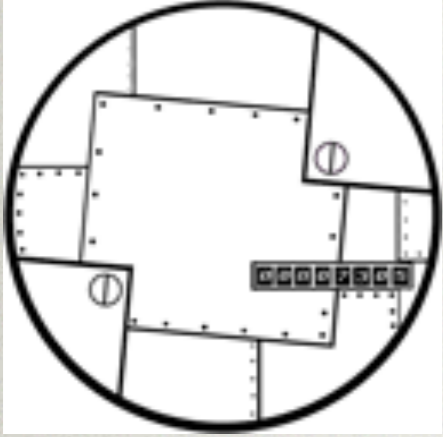
Luis Faria
lfaria@iantt.pt

What is RODA?

A digital repository created for archives, with the following features:

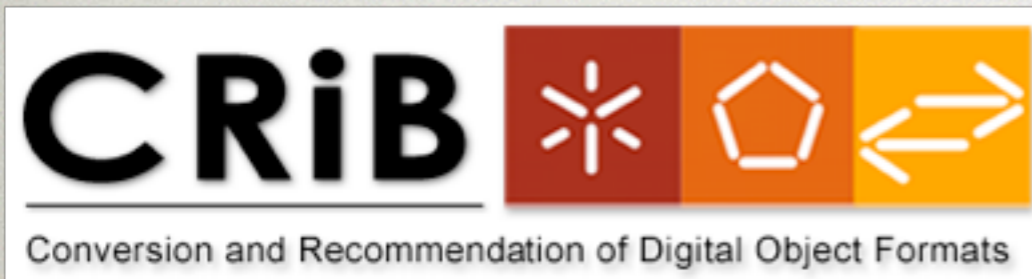
- Long term preservation and Autenticity
- Standards based (OAIS, EAD, PREMIS, METS, etc.)
- Following TRAC demands (Trustworthy Repositories Audit & Certification)
- Security policies
- Service Oriented Architecture (SOA)
- Nice and simple design
- Open source
- DGArc and University of Minho

CONTEXT



RODA (2006-2009)

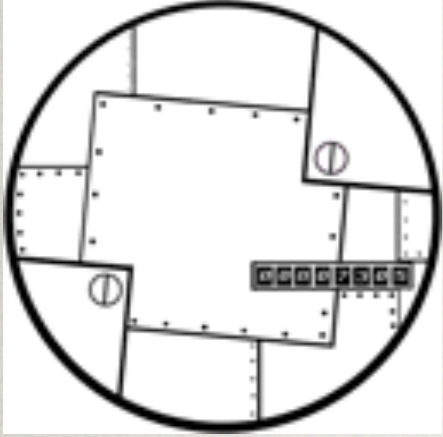
- Metadata management (EAD)
- Digital Object management (...)
- Digital Preservation Policies and protocols
- National Project (Archives National Board)



CRiB: Preservation Services Digital Repositories (2005-2008)

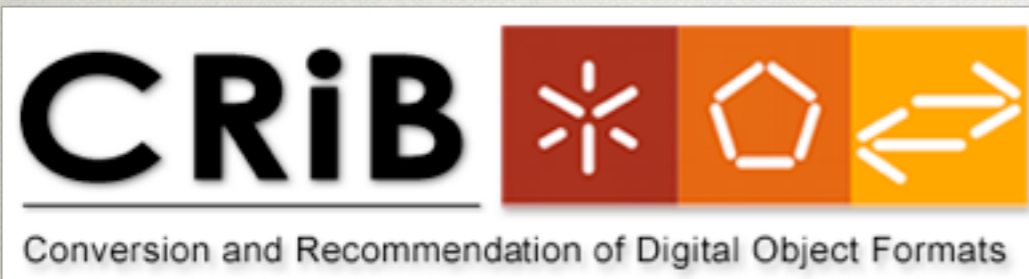
- Distributed Migration Service
- Migration Service supported by knowledge base
- phd Thesis / U. Minho (Miguel Ferreira)

CONTEXT



RODA (2006-2009)

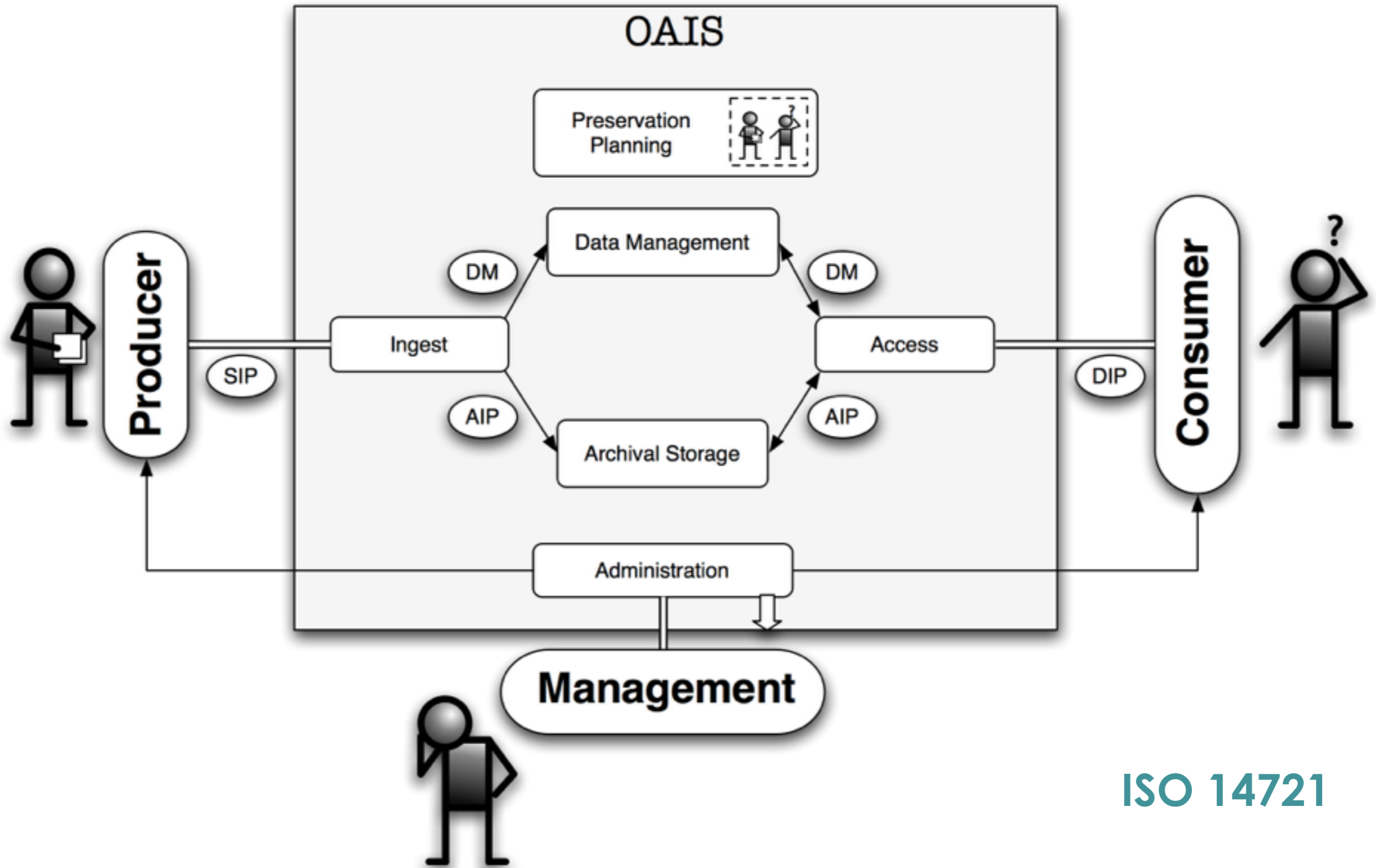
- Metadata management (EAD)
- Digital Object management (...)
- Digital Preservation Policies and protocols
- National Project (Archives National Board)



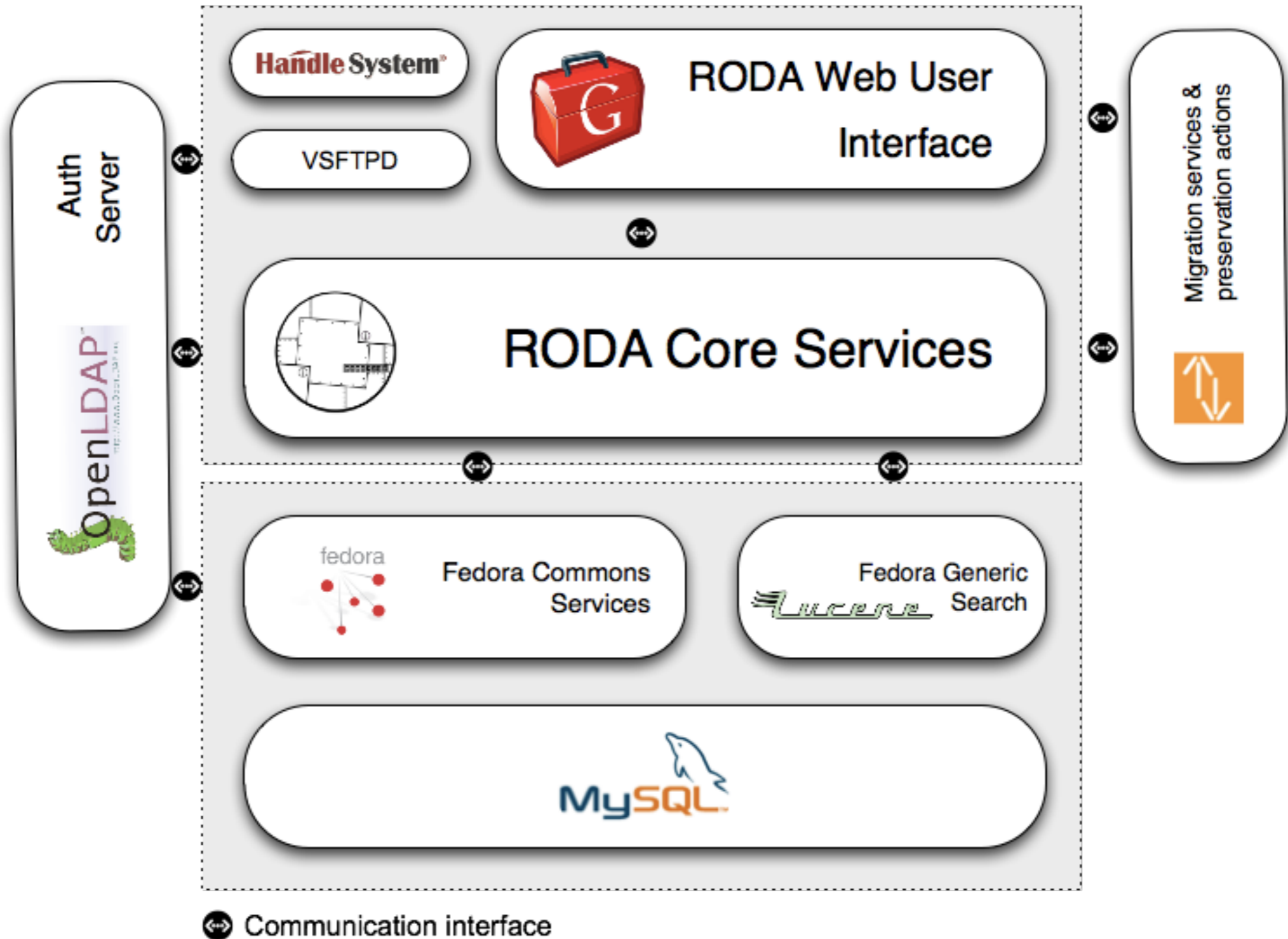
CRiB: Preservation Services Digital Repositories (2005-2008)

- Distributed Migration Service
- Migration Service supported by knowledge base
- phd Thesis / U. Minho (Miguel Ferreira)

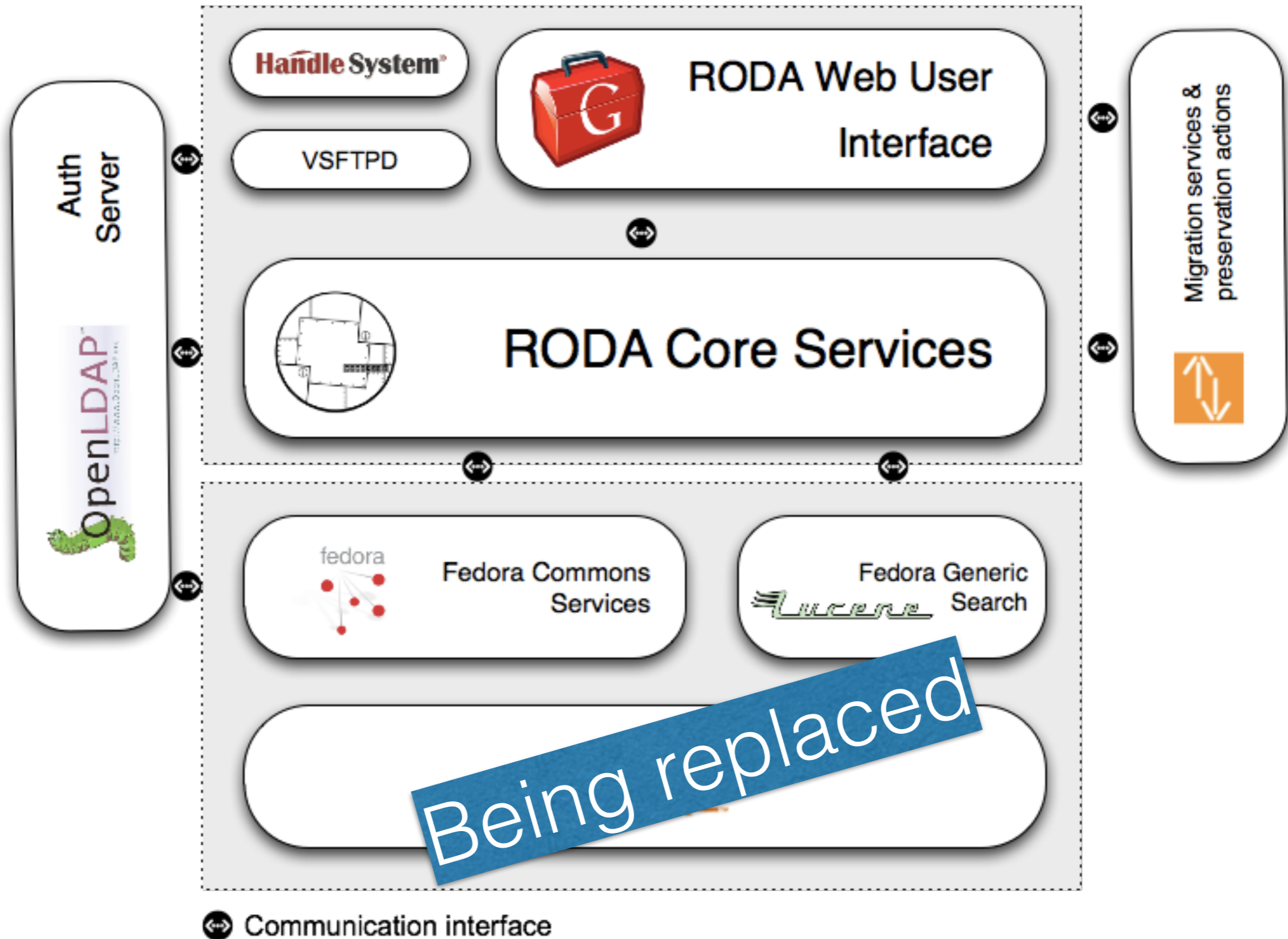
OAIS



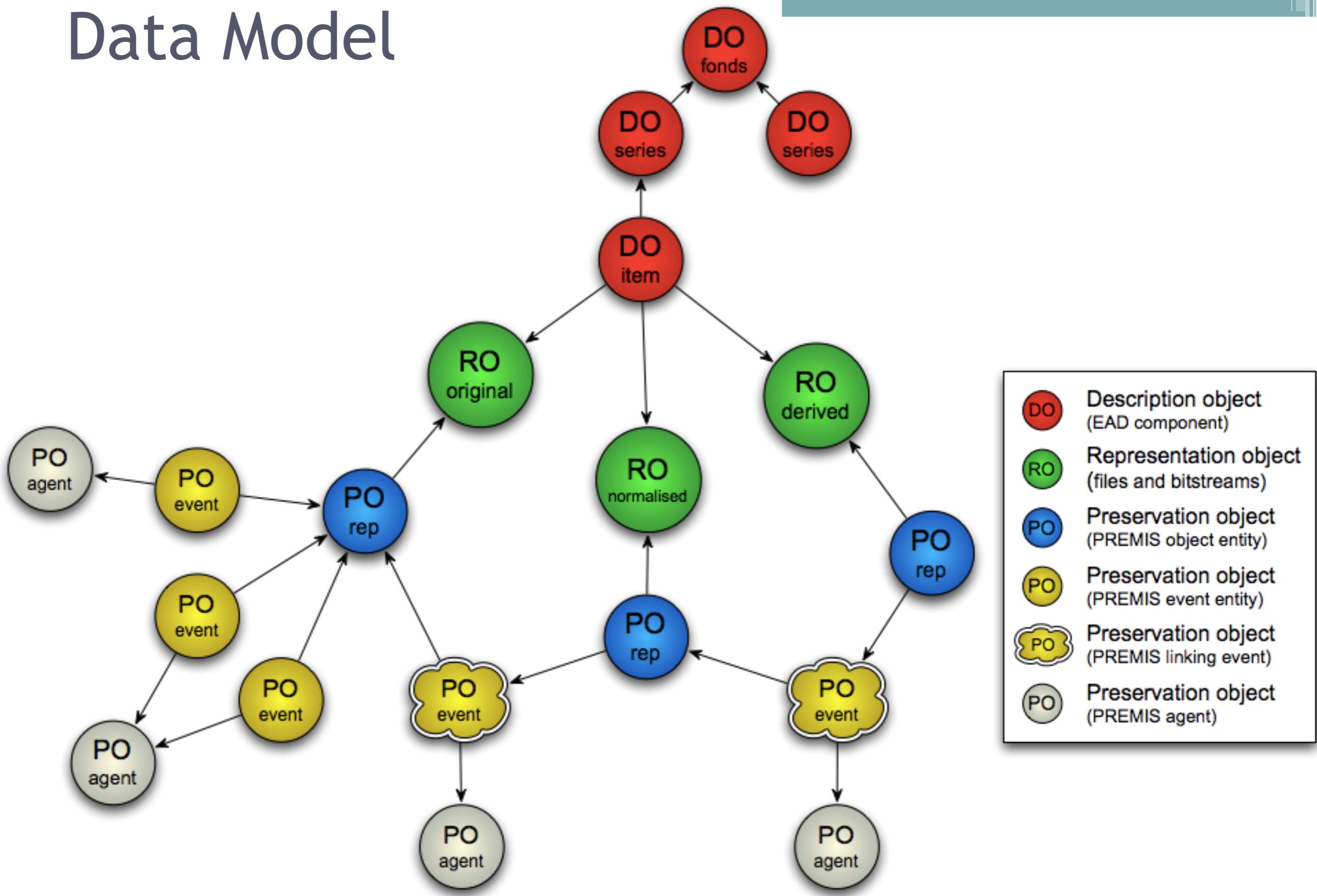
Architecture



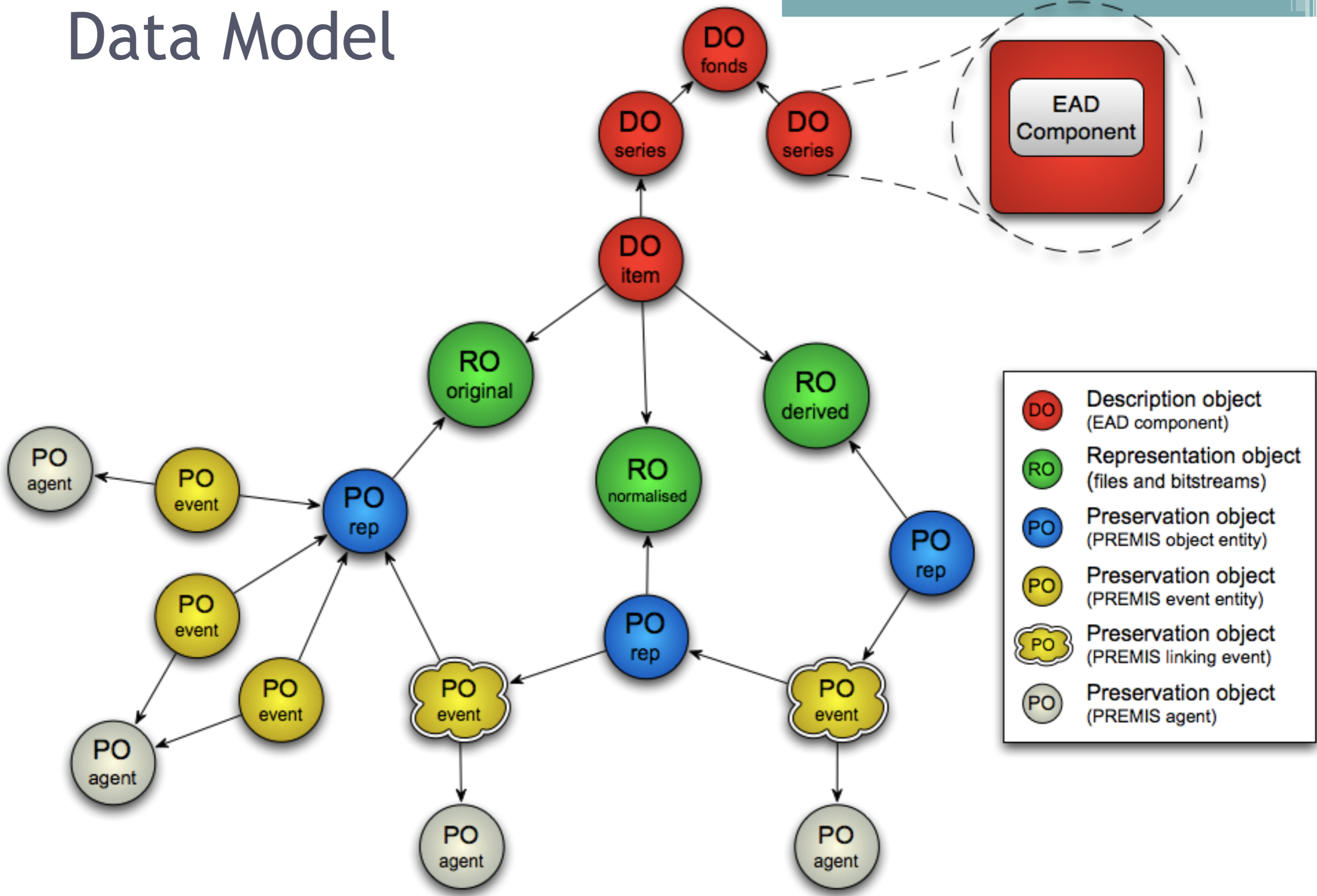
Architecture



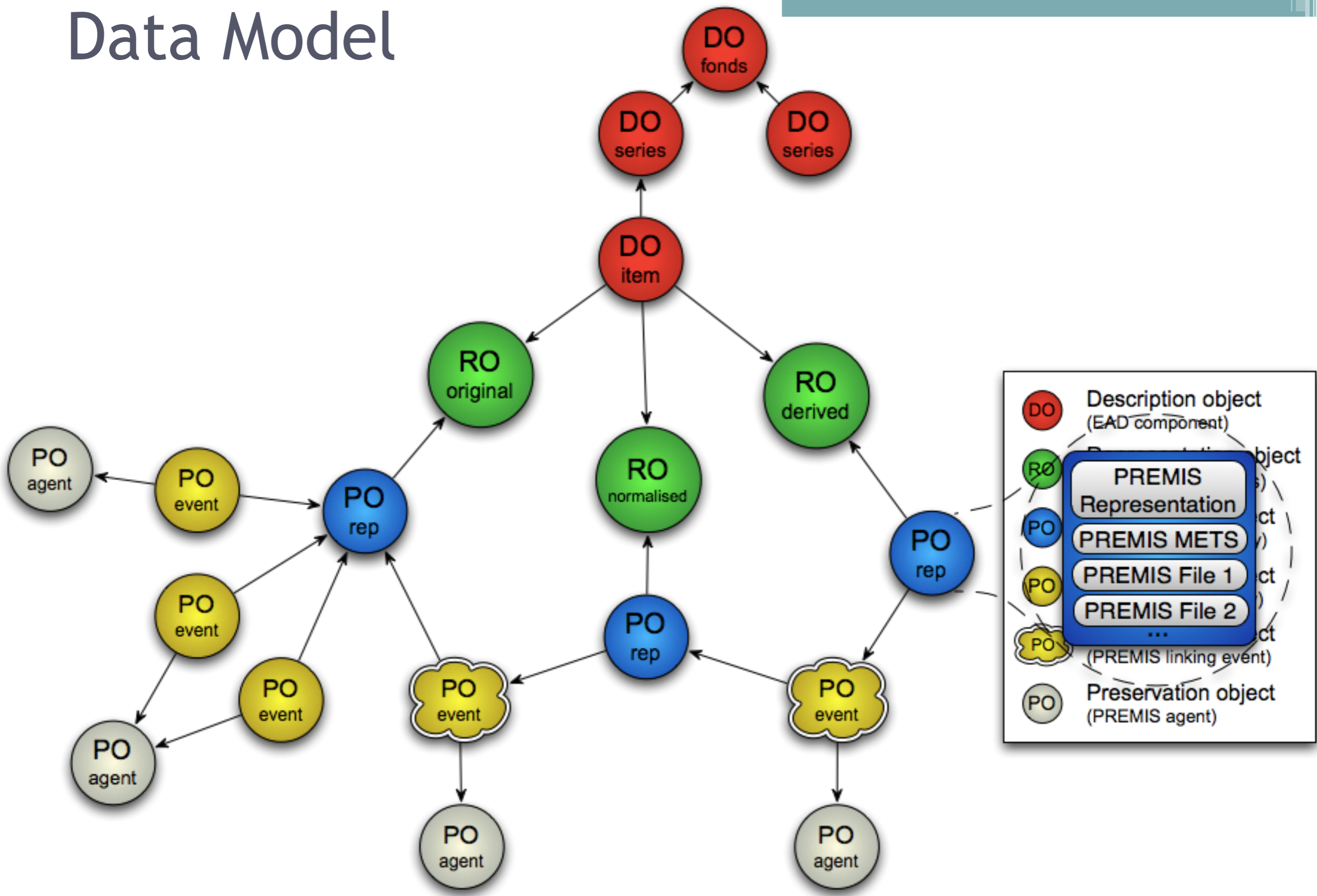
Data Model



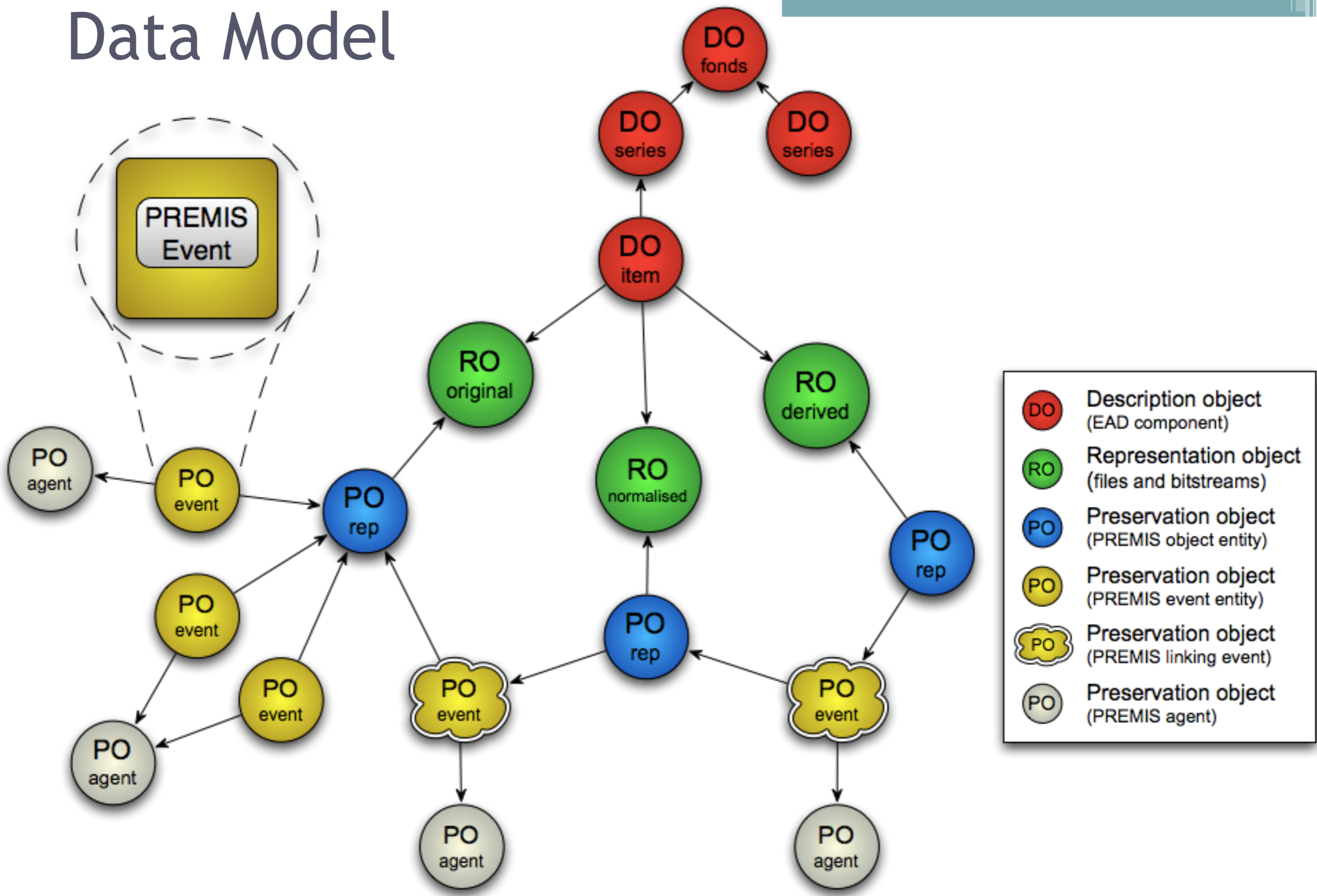
Data Model



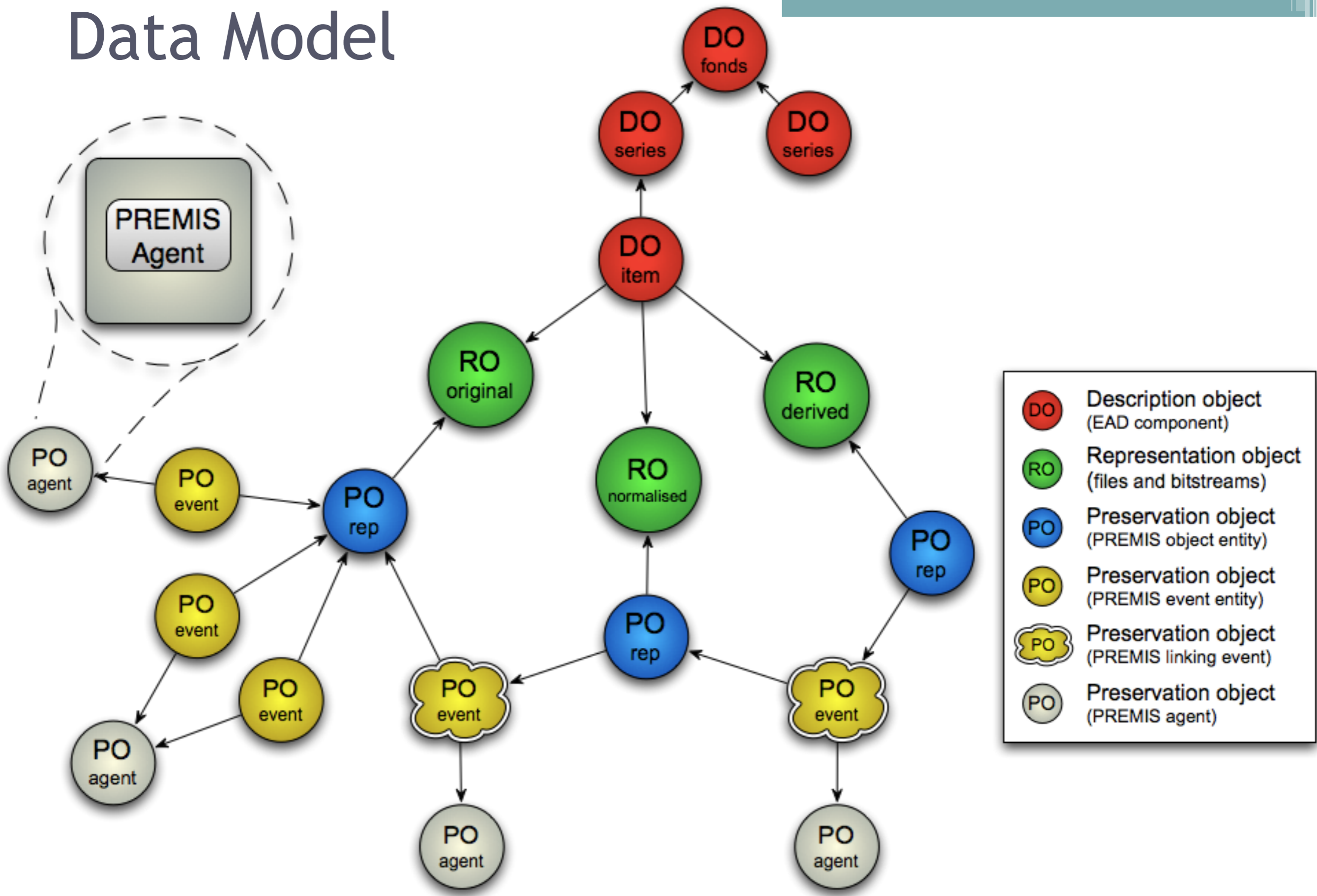
Data Model



Data Model



Data Model



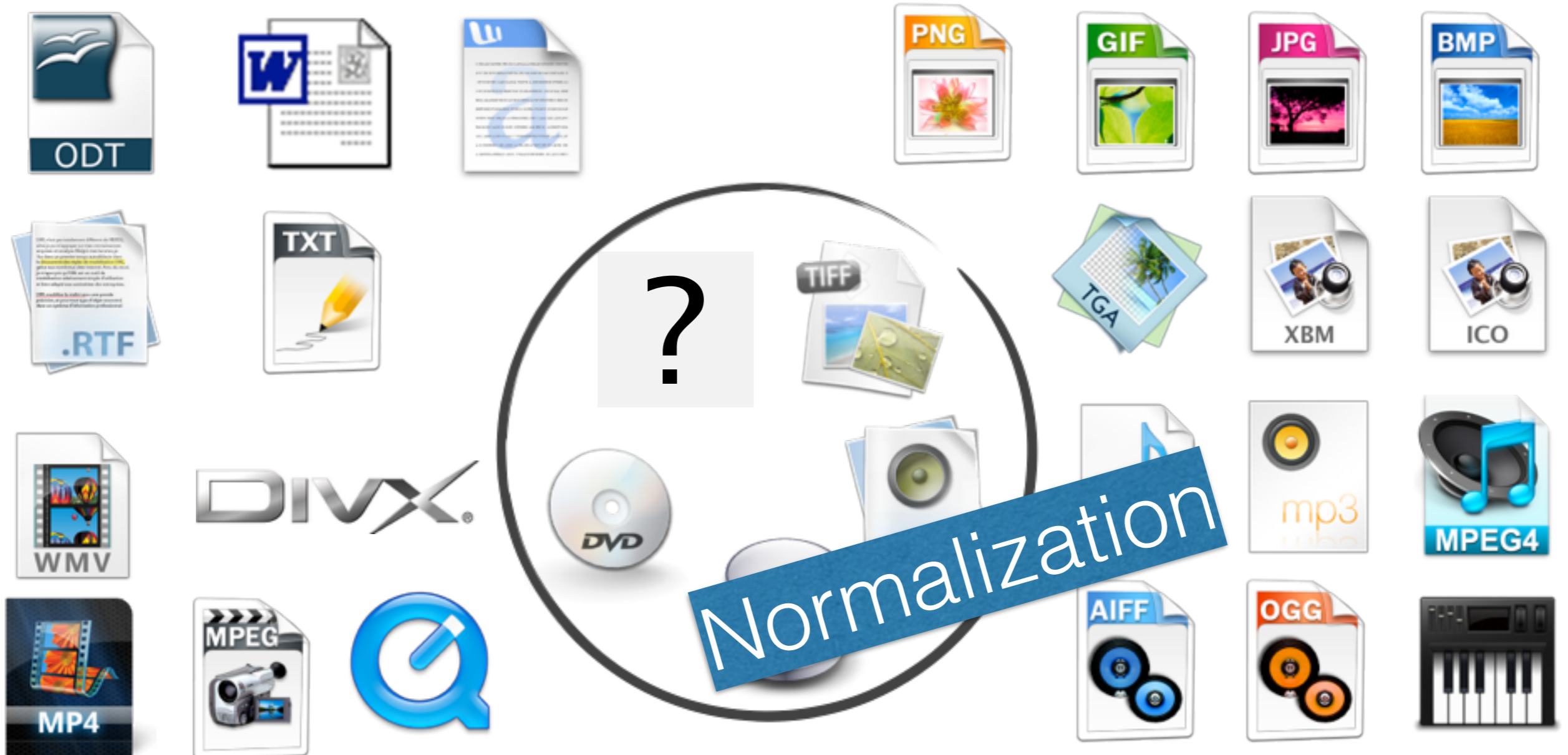
Digital Object Classes



Digital Object Classes



Digital Object Classes



Digital Object Classes

ODT

Word documents

PDF

PNG

GIF

JPG

BMP

.RTF

TXT

TIFF

TGA

XBM

ICO

WMV

DIVX

DVD

mp3

MPEG4

MP4

MPEG

AIFF

OGG

Keyboard icon

?

DVD

Normalization

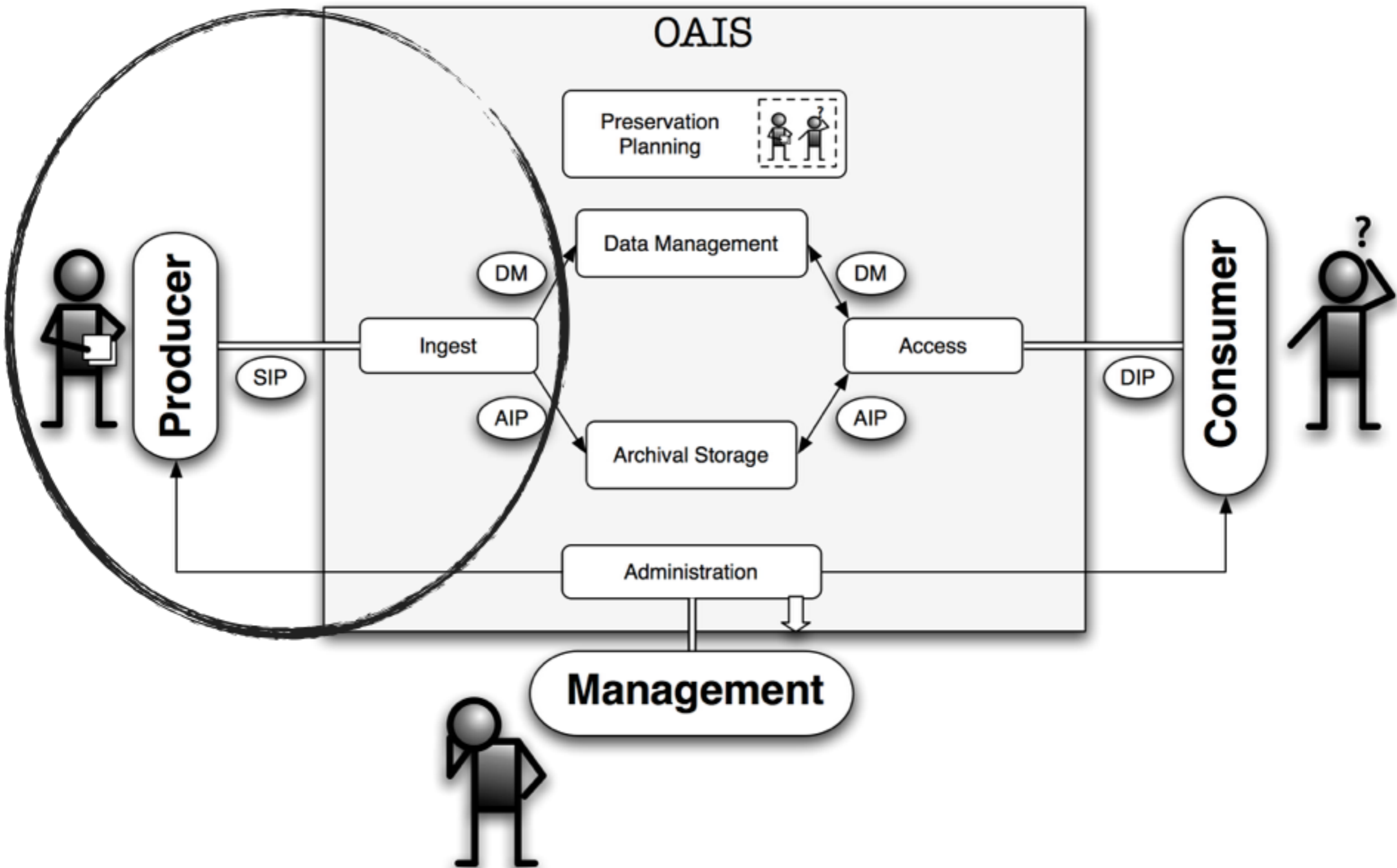
Significant properties

Microsoft SQL Server

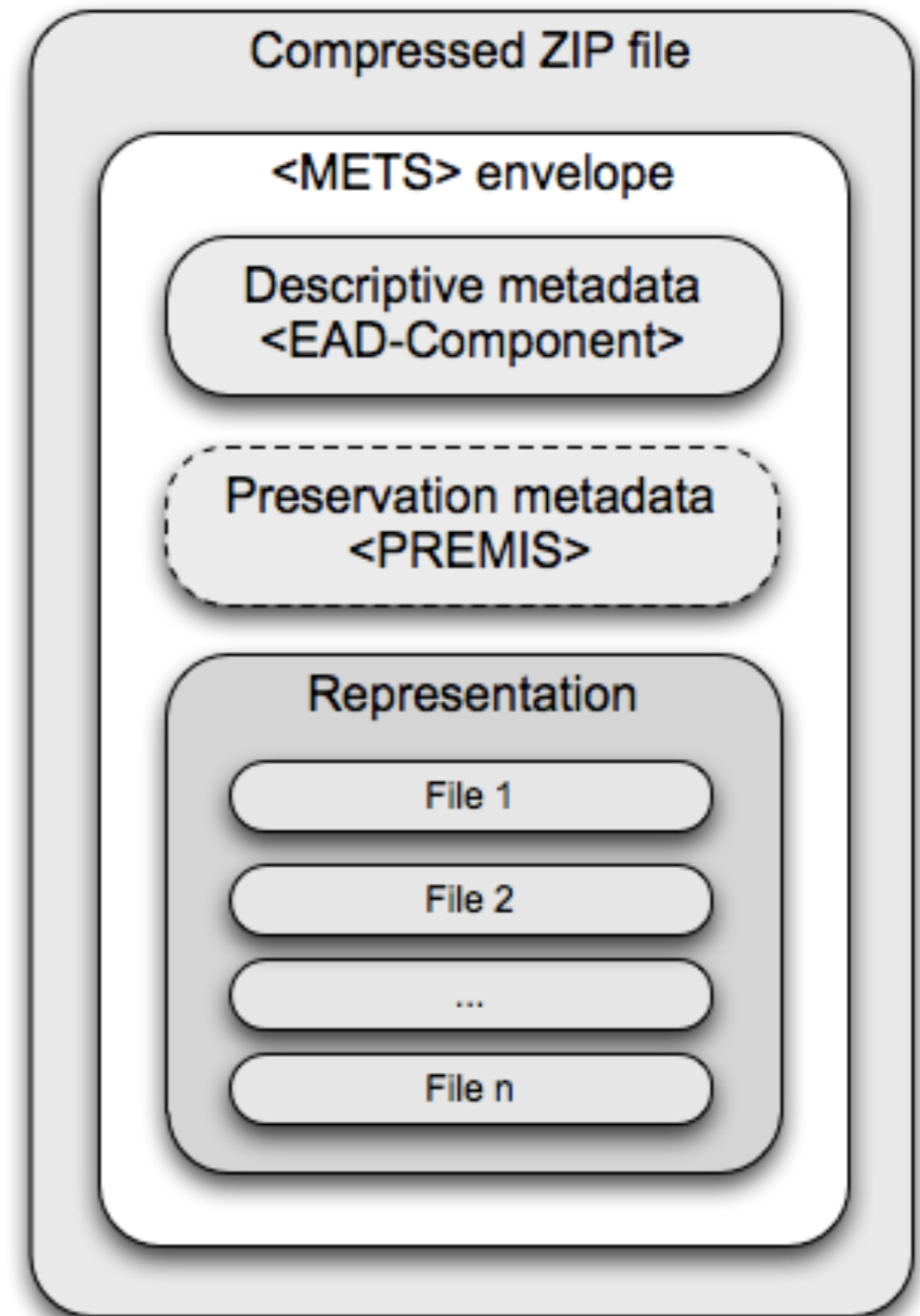
MySQL

Microsoft Office Access

Ingest

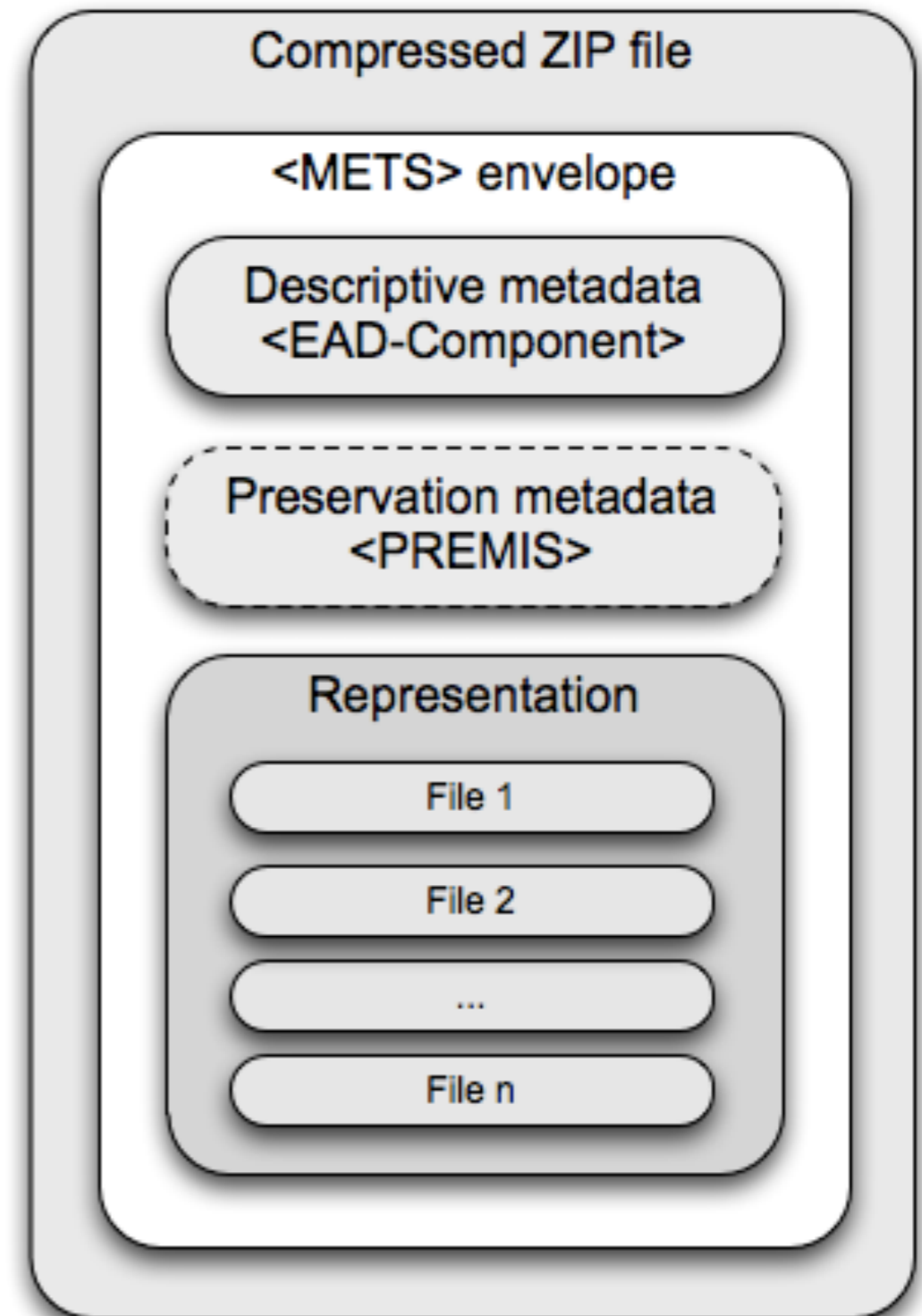


SIP structure



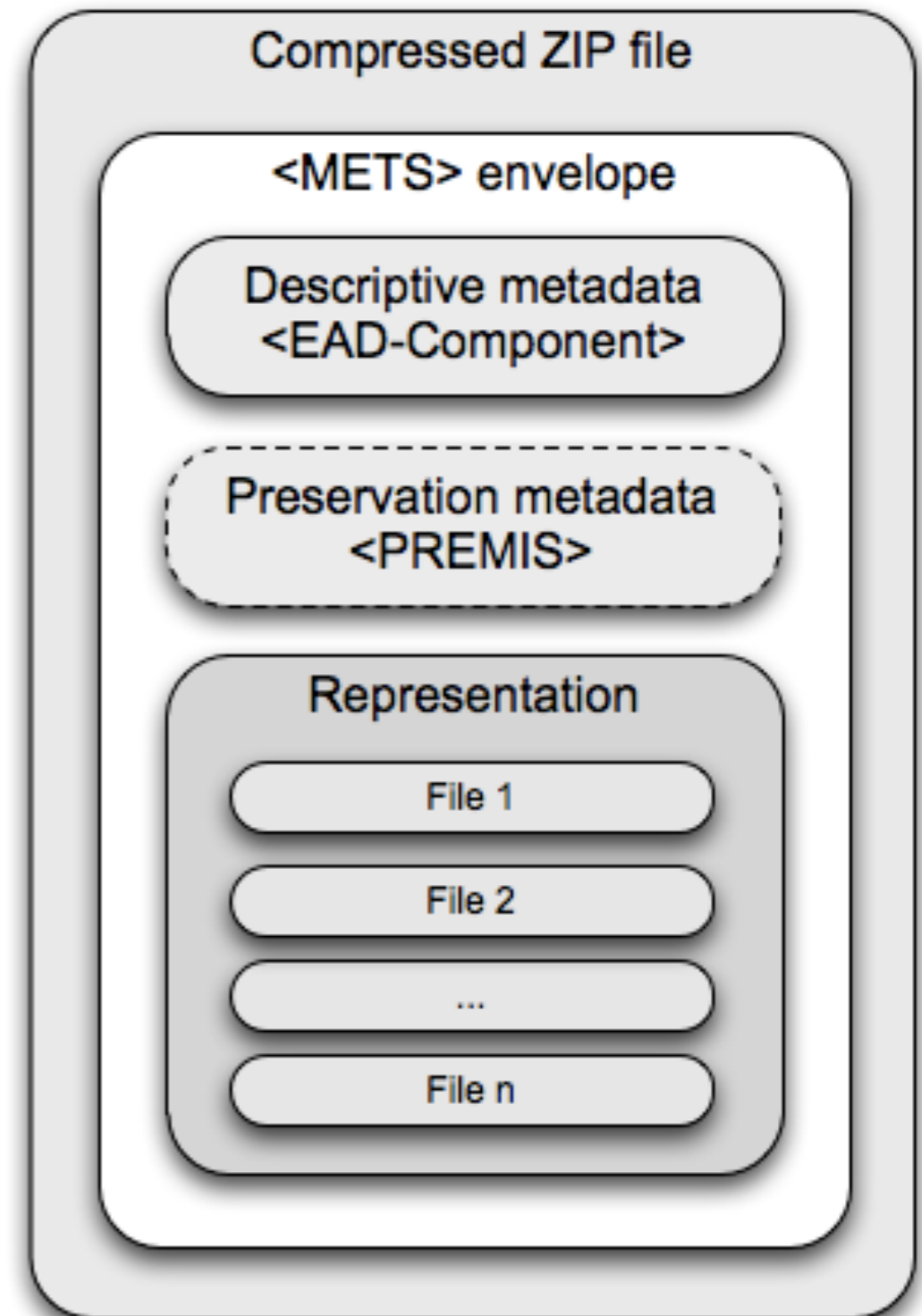
SIP structure

- Object class and format



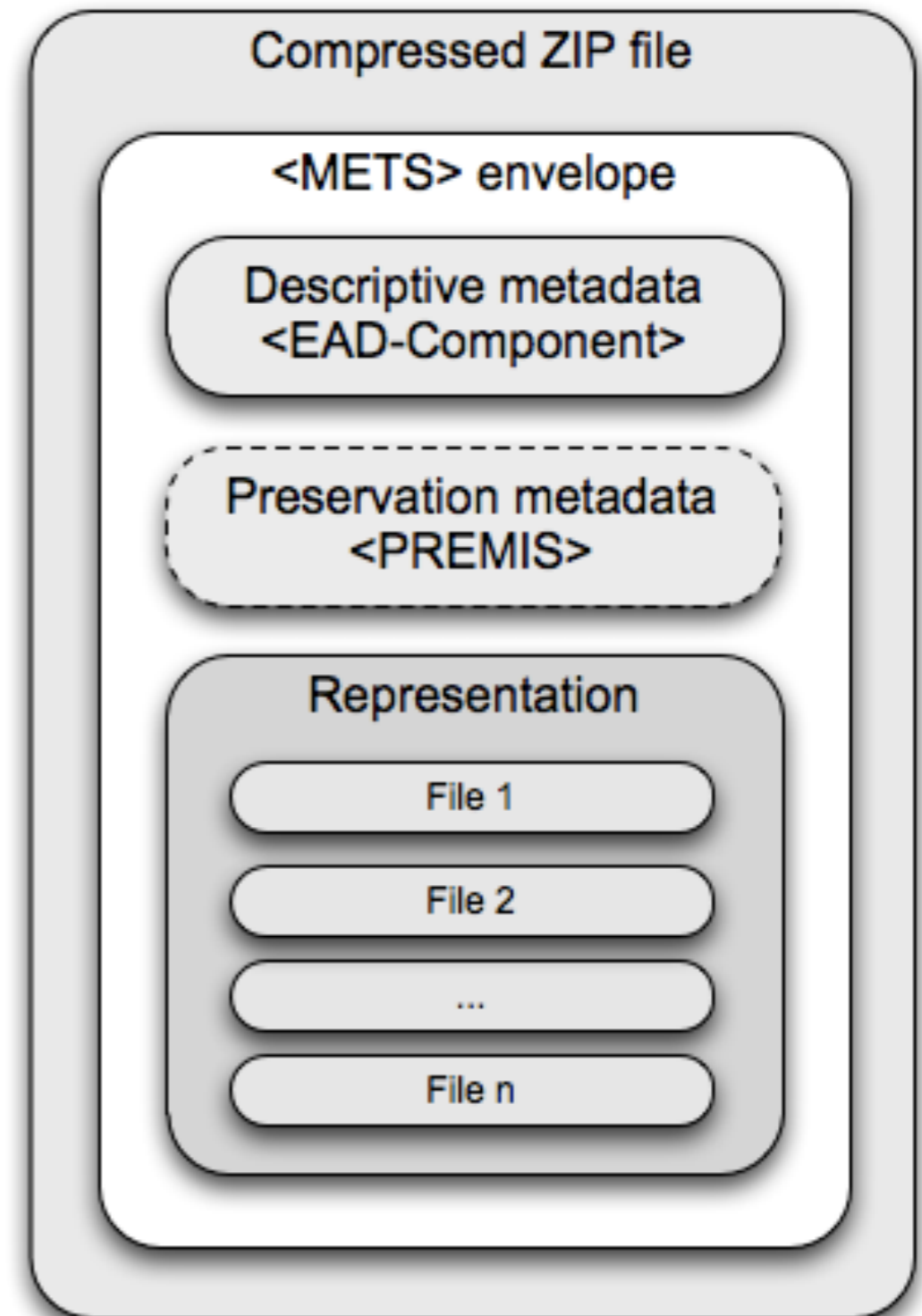
SIP structure

- Object class and format
- Place in classification plan



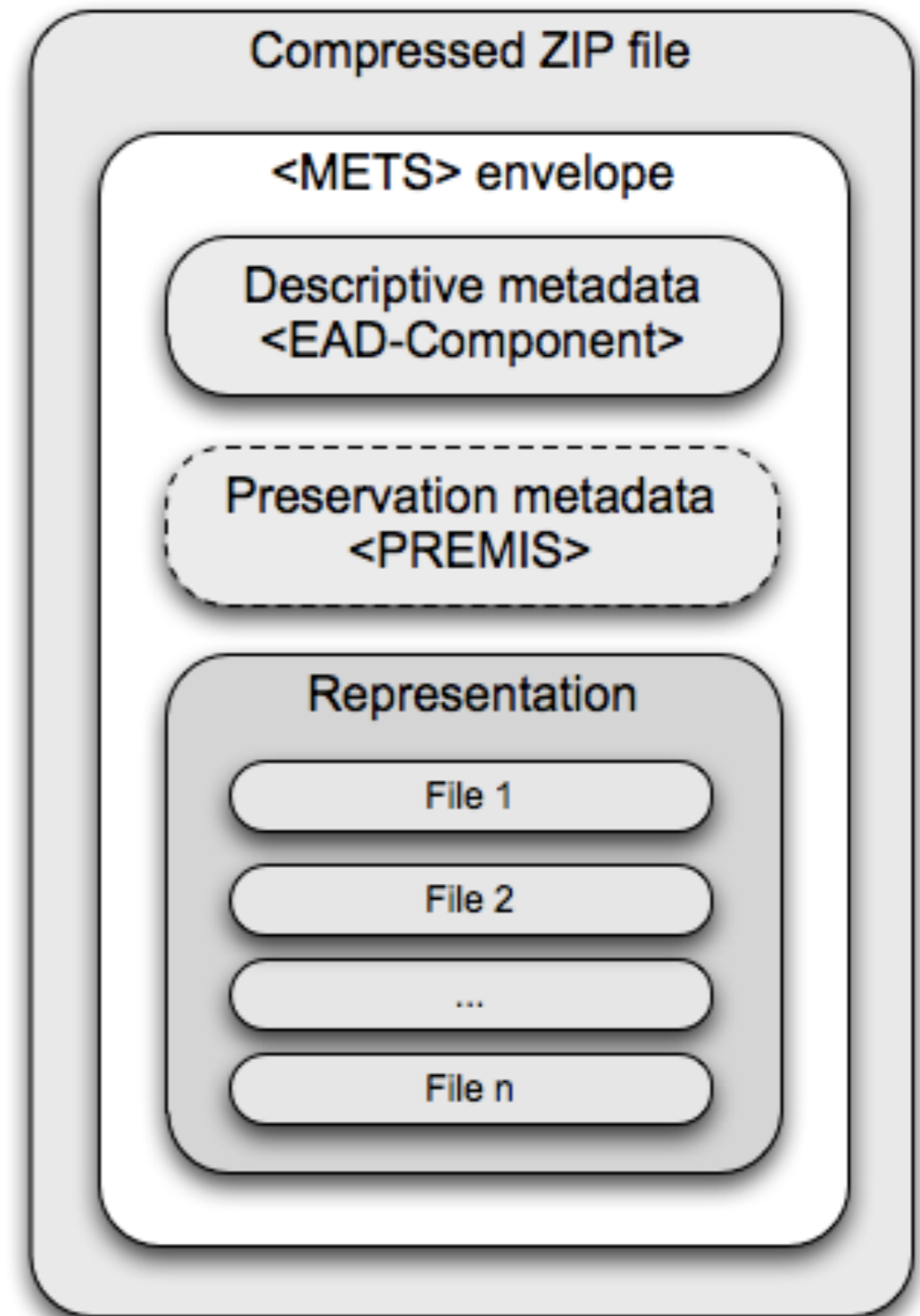
SIP structure

- Object class and format
- Place in classification plan
- Descriptive metadata



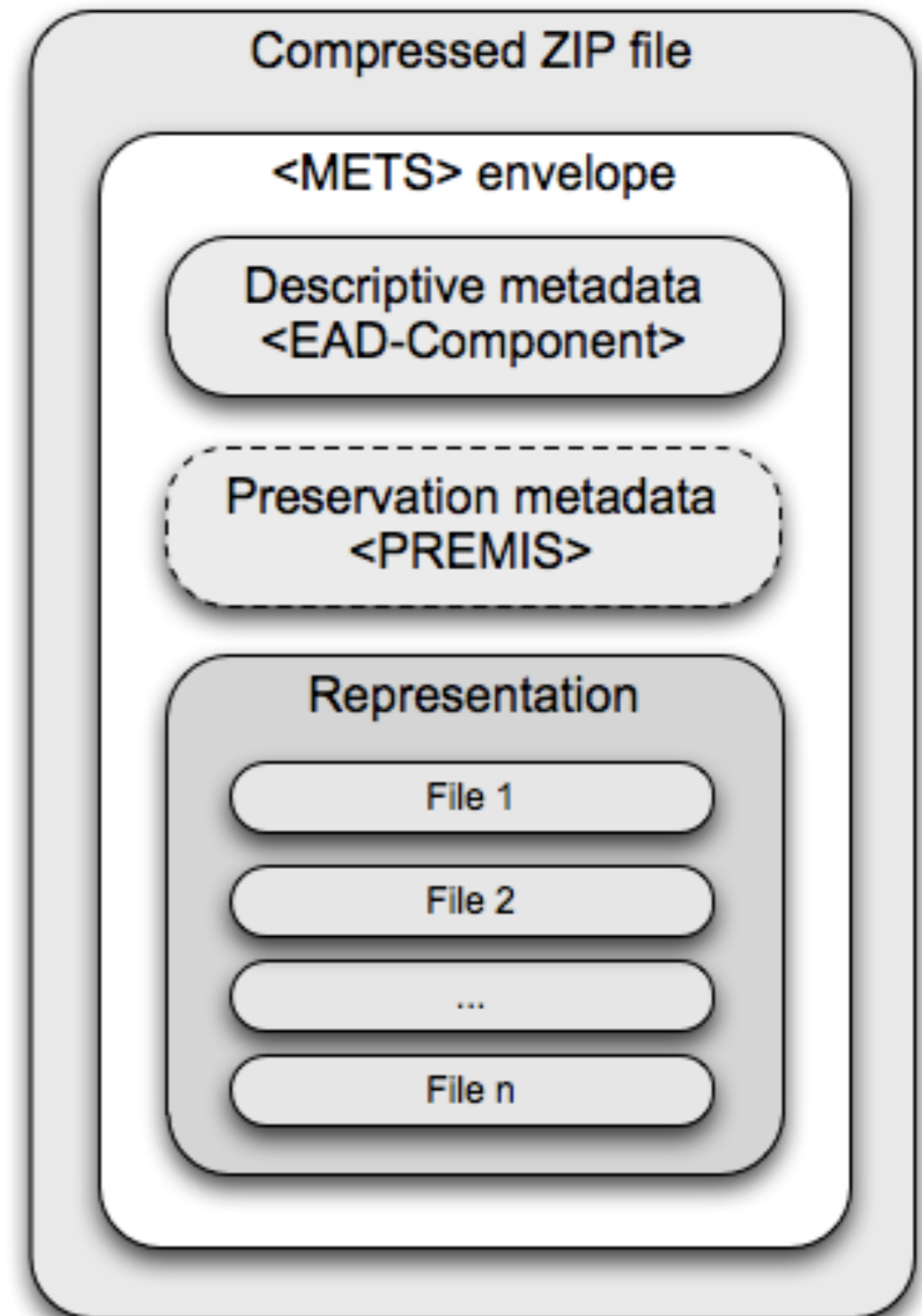
SIP structure

- Object class and format
- Place in classification plan
- Descriptive metadata
- Preservation metadata



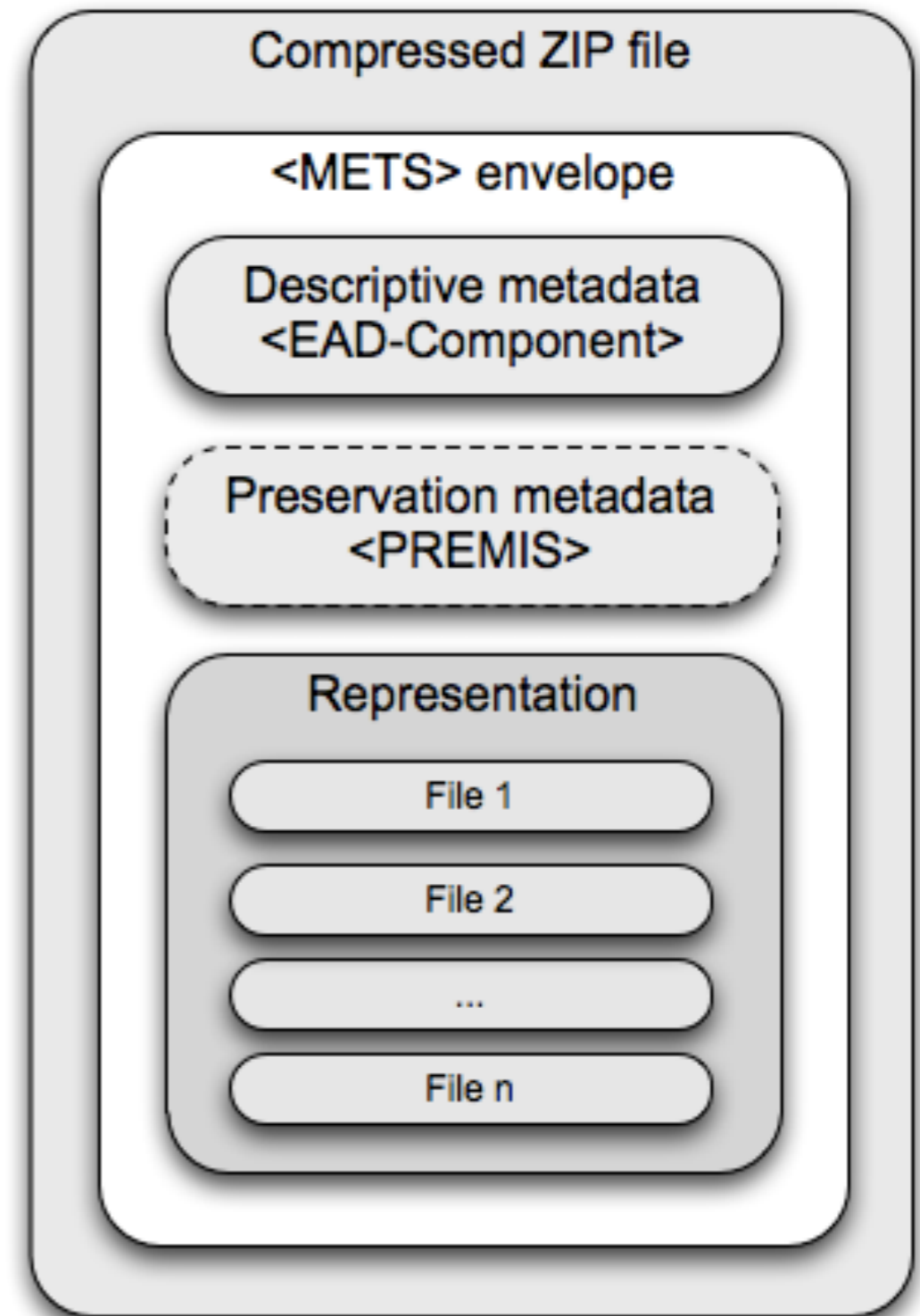
SIP structure

- Object class and format
- Place in classification plan
- Descriptive metadata
- Preservation metadata
- Technical metadata



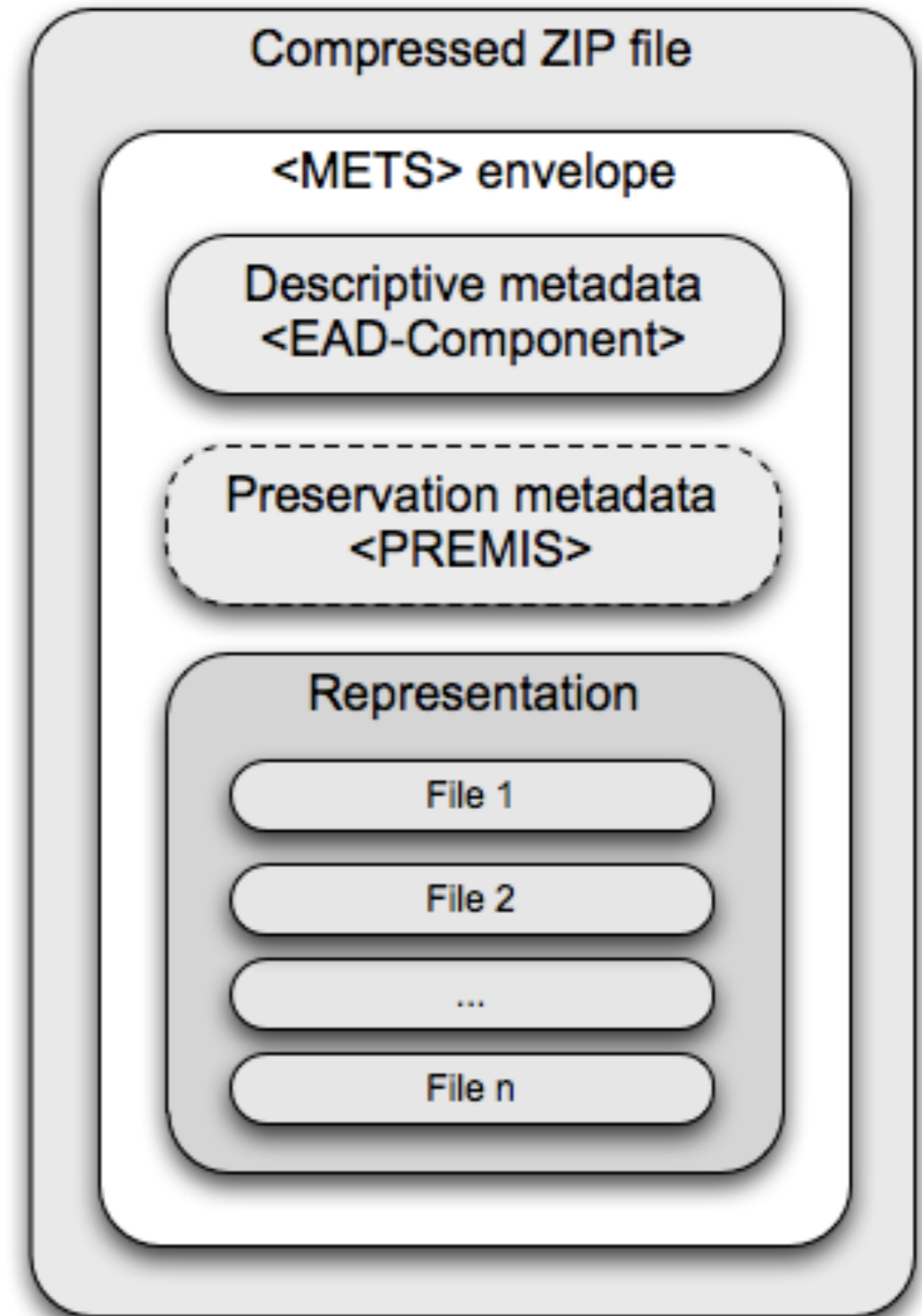
SIP structure

- Object class and format
- Place in classification plan
- Descriptive metadata
- Preservation metadata
- Technical metadata
- Representation files

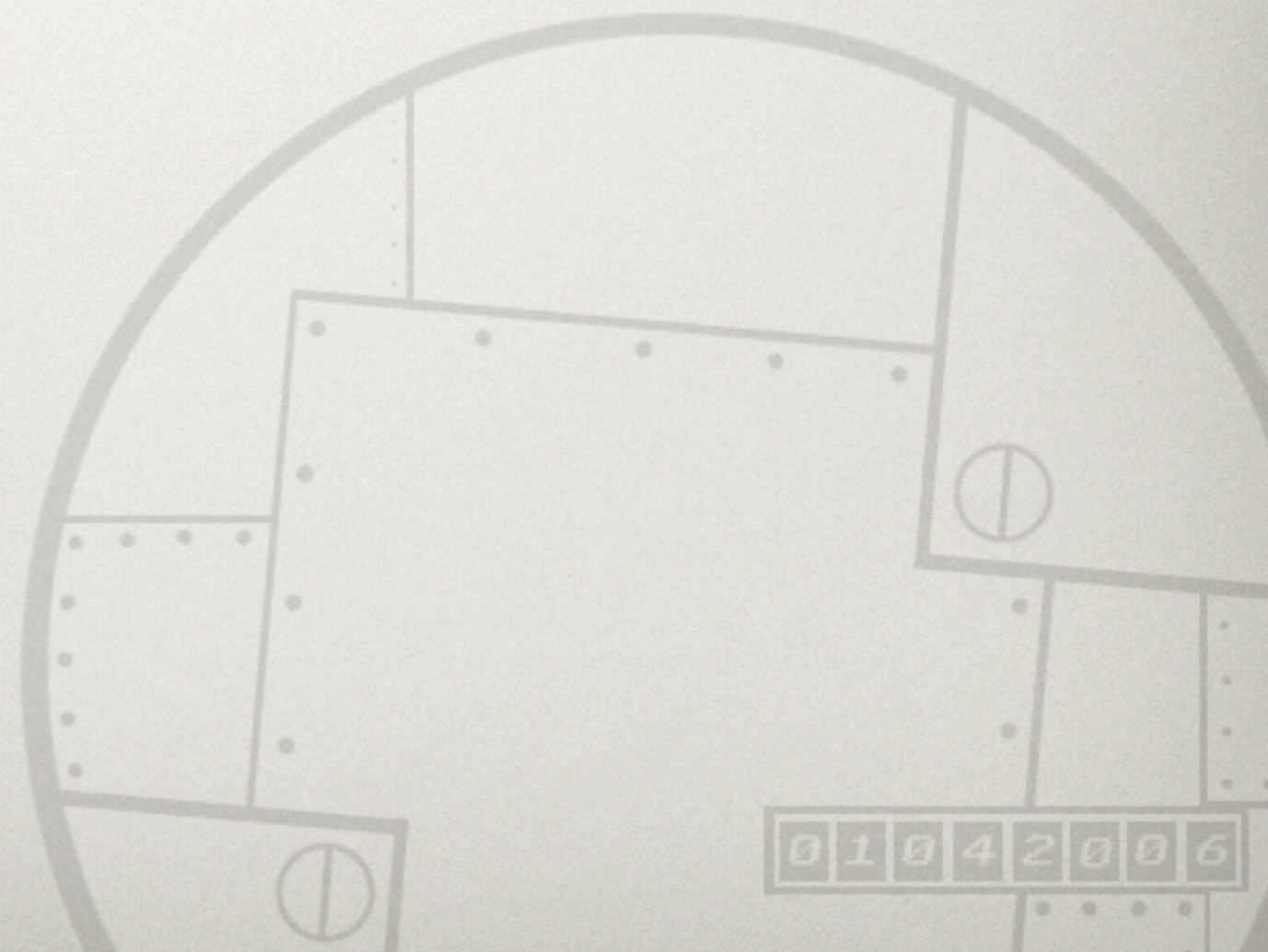


SIP structure

- Object class and format
- Place in classification plan
- Descriptive metadata
- Preservation metadata
- Technical metadata
- Representation files
- Identification of root file

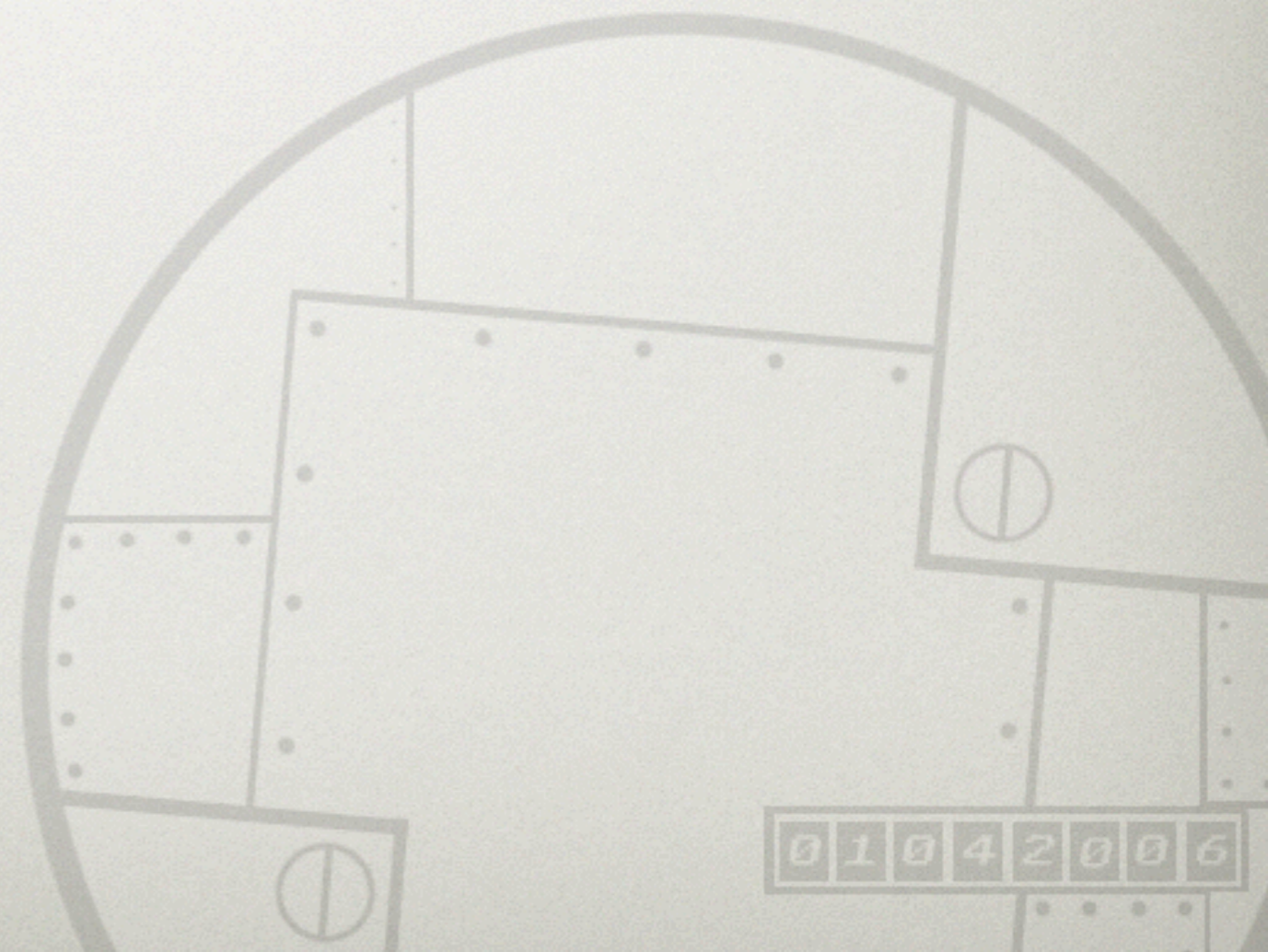


DATABASES



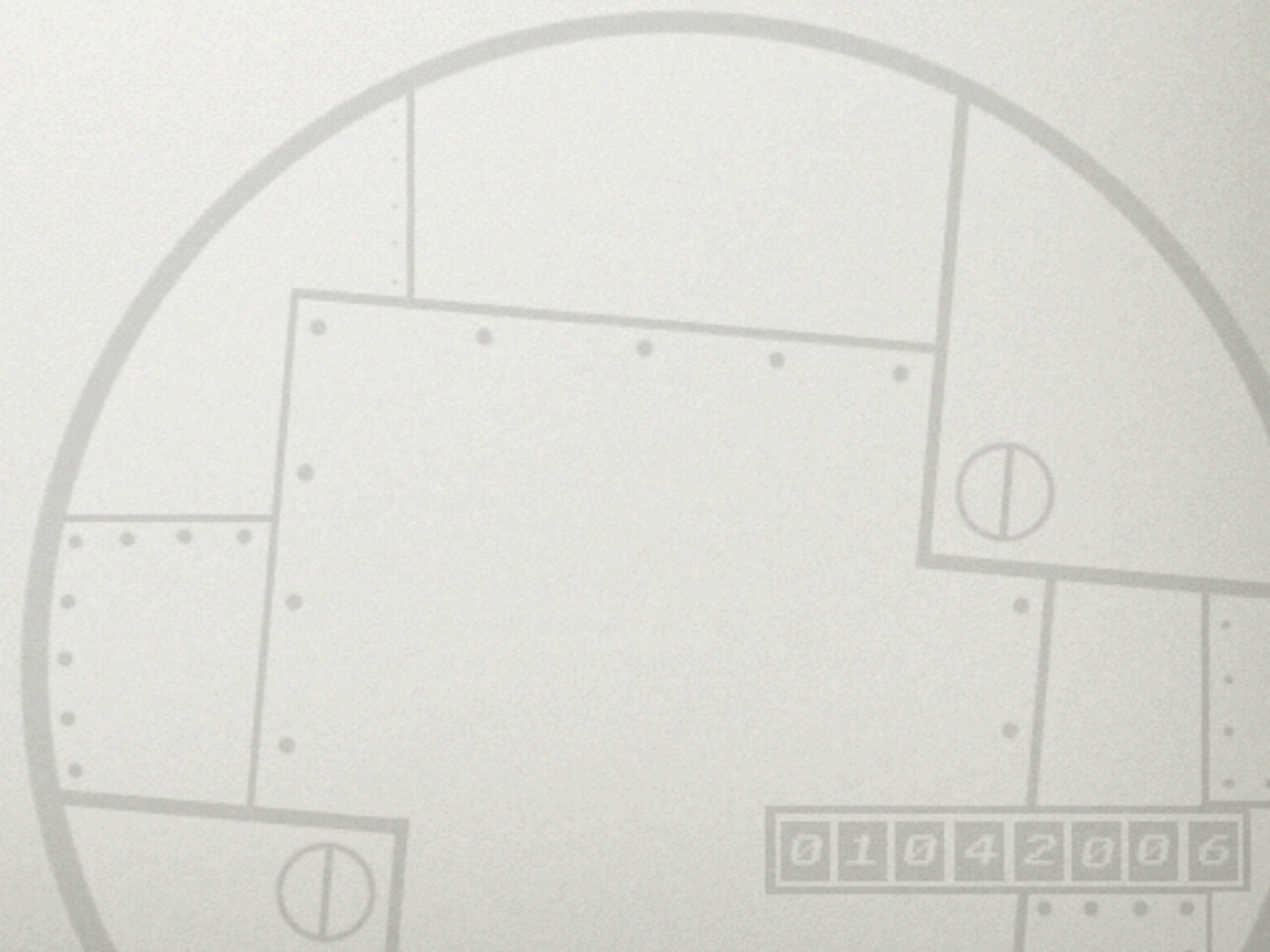
DATABASES

- **How do we keep archived databases readable and usable in the long term (at acceptable cost)?**



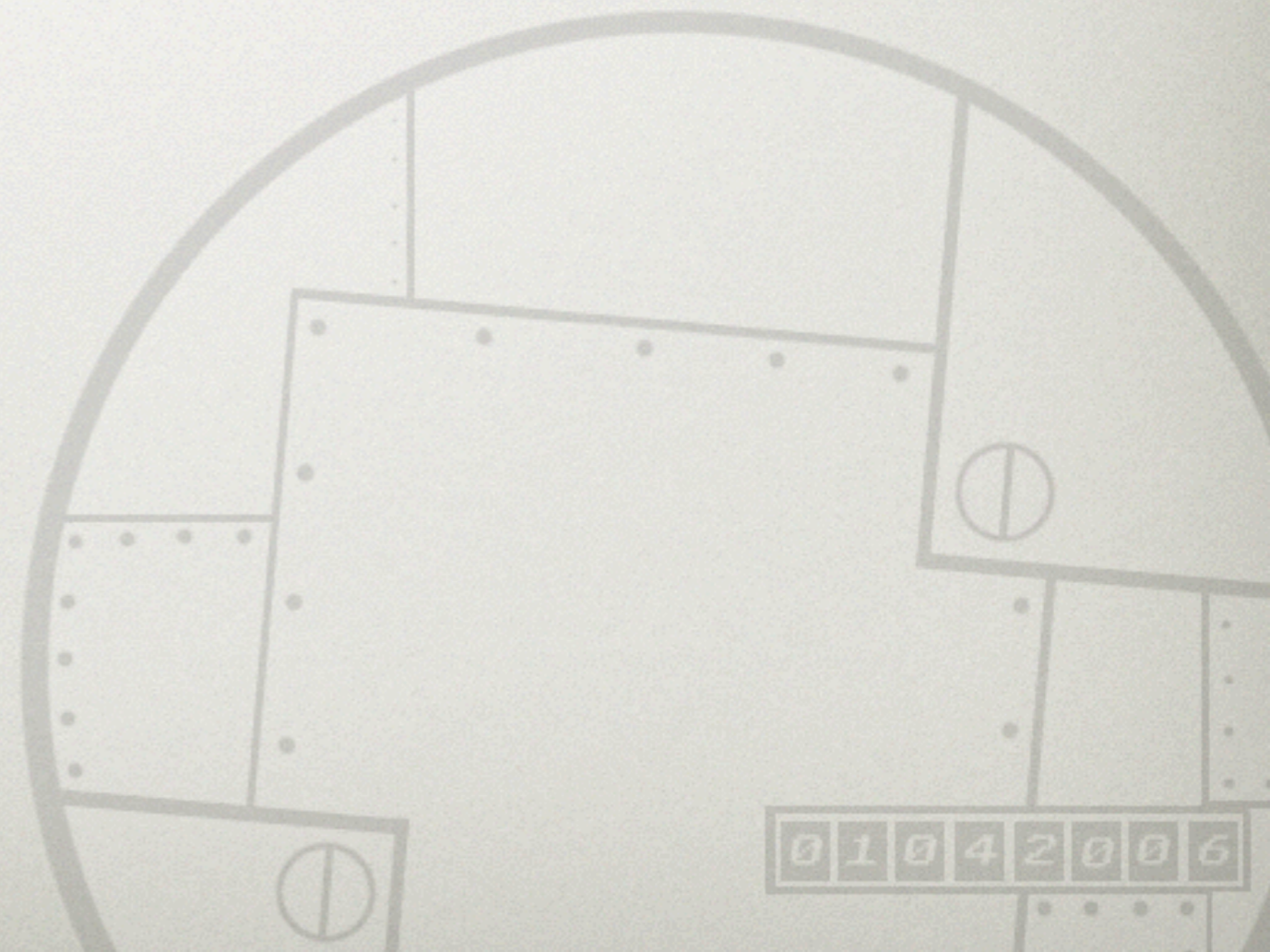
DATABASES

- How do we keep archived databases readable and usable in the long term (at acceptable cost)?
- How do we **separate the data** from a specific database management environment?



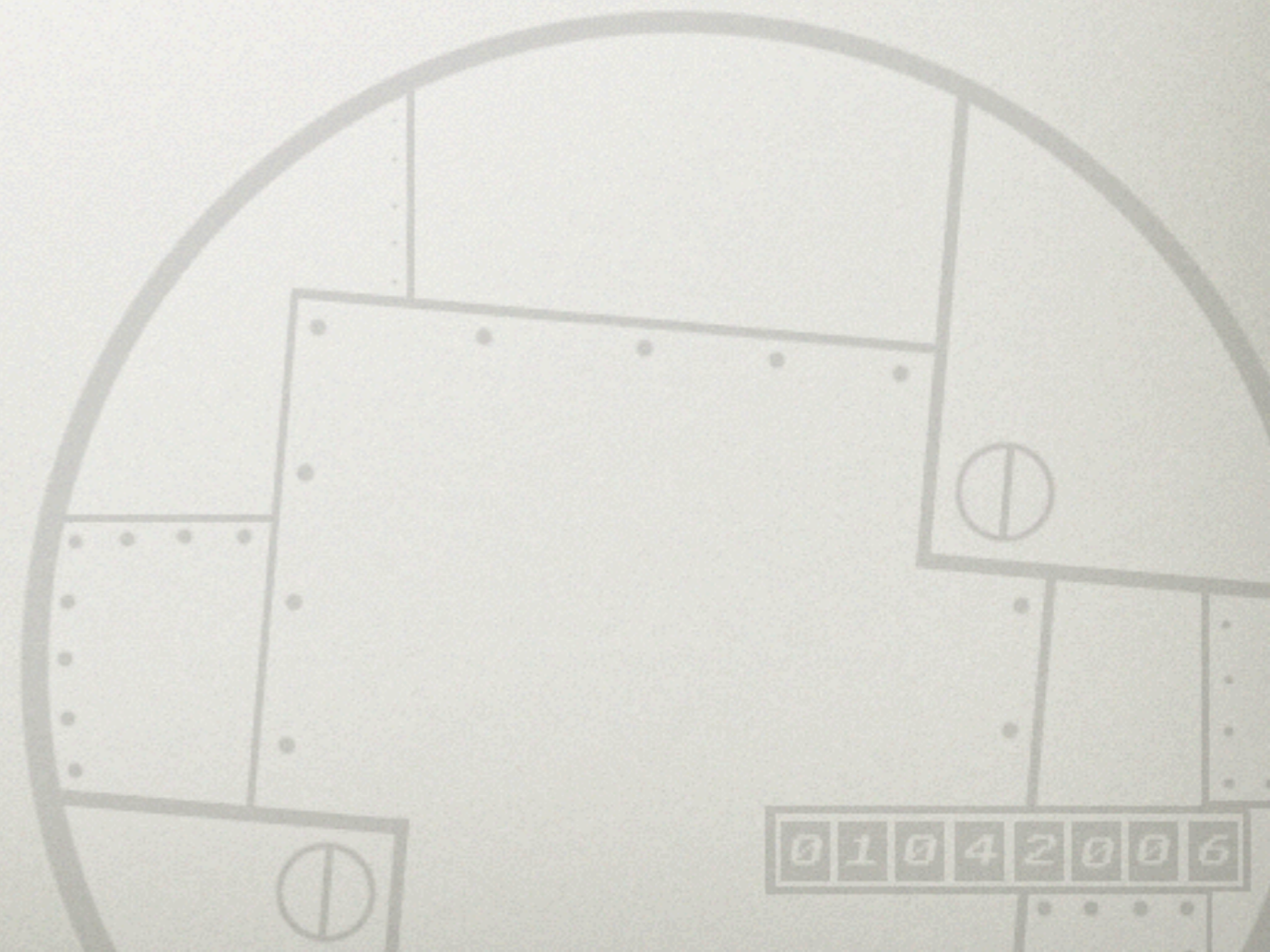
DATABASES

- How do we keep archived databases readable and usable in the long term (at acceptable cost)?
- How do we **separate the data** from a specific database management environment?
- How can we **preserve** the original **data semantics** and **structure**?



DATABASES

- How do we keep archived databases readable and usable in the long term (at acceptable cost)?
- How do we **separate the data** from a specific database management environment?
- How can we **preserve** the original **data semantics** and **structure**?
- How can we **preserve authenticity** and **provenance** of databases?



DATABASES

- How do we keep archived databases readable and usable in the long term (at acceptable cost)?
- How do we **separate the data** from a specific database management environment?
- How can we **preserve** the original **data semantics** and **structure**?
- How can we **preserve authenticity** and **provenance** of databases?
- How can we **preserve data** while it continues to **evolve**? (here we can split the problem in two: operational databases and frozen databases)

DATABASES

- How do we keep archived databases readable and usable in the long term (at acceptable cost)?
- How do we **separate the data** from a specific database management environment?
- How can we **preserve** the original **data semantics** and **structure**?
- How can we **preserve authenticity** and **provenance** of databases?
- How can we **preserve data** while it continues to **evolve**? (here we can split the problem in two: operational databases and frozen databases)
- How can we have efficient preservation frameworks, while retaining the ability to query different database versions?

DATABASES

- How do we keep archived databases readable and usable in the long term (at acceptable cost)?
- How do we **separate the data** from a specific database management environment?
- How can we **preserve** the original **data semantics** and **structure**?
- How can we **preserve authenticity** and **provenance** of databases?
- How can we **preserve data** while it continues to **evolve**? (here we can split the problem in two: operational databases and frozen databases)
- How can we have efficient preservation frameworks, while retaining the ability to query different database versions?
- How can multi-user online access be provided to hundreds of archived databases containing terabytes of data?

DATABASES

- How do we keep archived databases readable and usable in the long term (at acceptable cost)?
- How do we **separate the data** from a specific database management environment?
- How can we **preserve** the original **data semantics** and **structure**?
- How can we **preserve authenticity** and **provenance** of databases?
- How can we **preserve data** while it continues to **evolve**? (here we can split the problem in two: operational databases and frozen databases)
- How can we have efficient preservation frameworks, while retaining the ability to query different database versions?
- How can multi-user online access be provided to hundreds of archived databases containing terabytes of data?
- Can we move from a centralized model to a distributed, redundant model of database preservation?

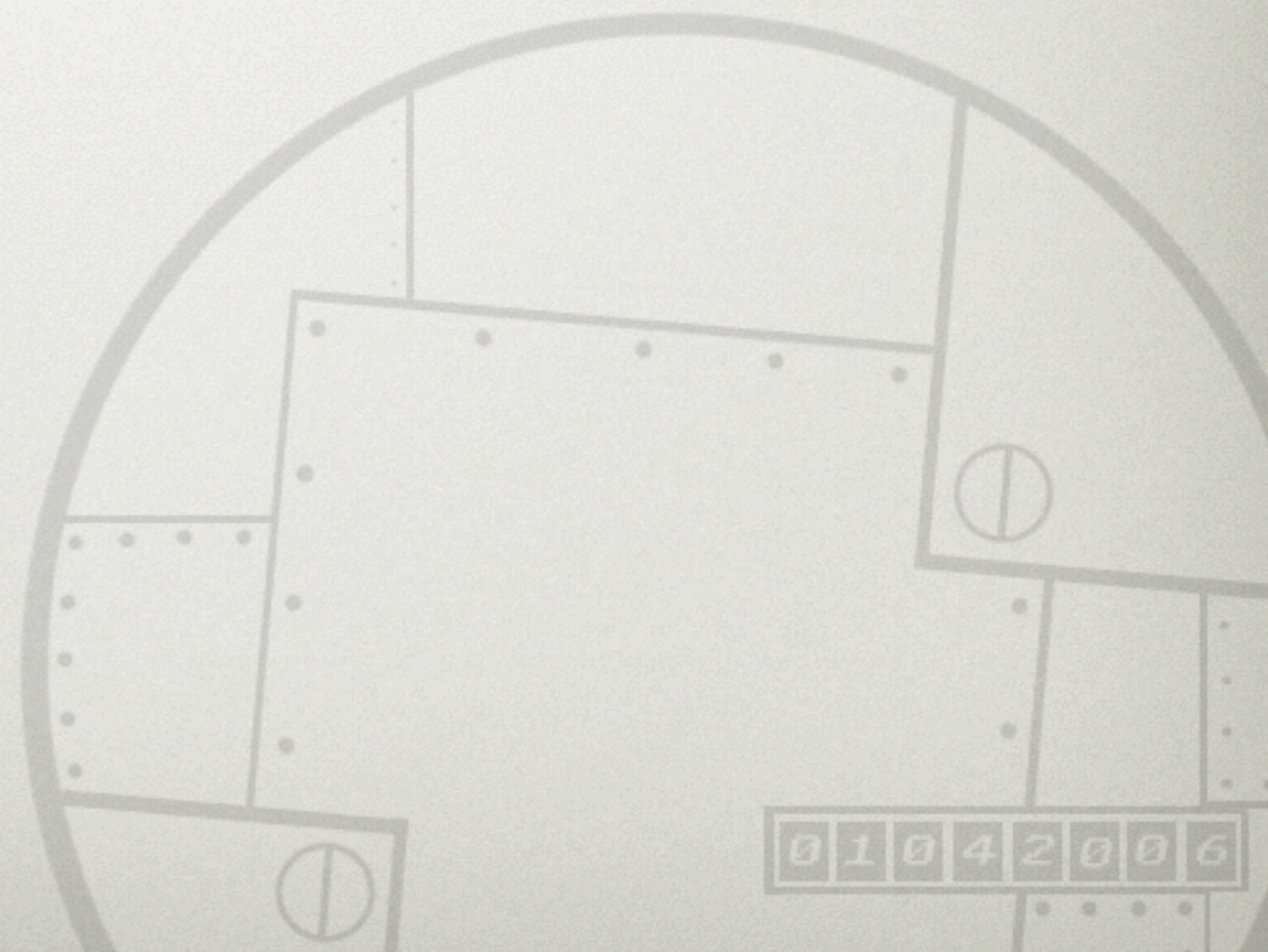
DATABASES

- How do we keep archived databases readable and usable in the long term (at acceptable cost)?
- How do we **separate the data** from a specific database management environment?
- How can we **preserve** the original **data semantics** and **structure**?
- How can we **preserve authenticity** and **provenance** of databases?
- How can we **preserve data** while it continues to **evolve**? (here we can split the problem in two: operational databases and frozen databases)
- How can we have efficient preservation frameworks, while retaining the ability to query different database versions?
- How can multi-user online access be provided to hundreds of archived databases containing terabytes of data?
- Can we move from a centralized model to a distributed, redundant model of database preservation?

General big questions...

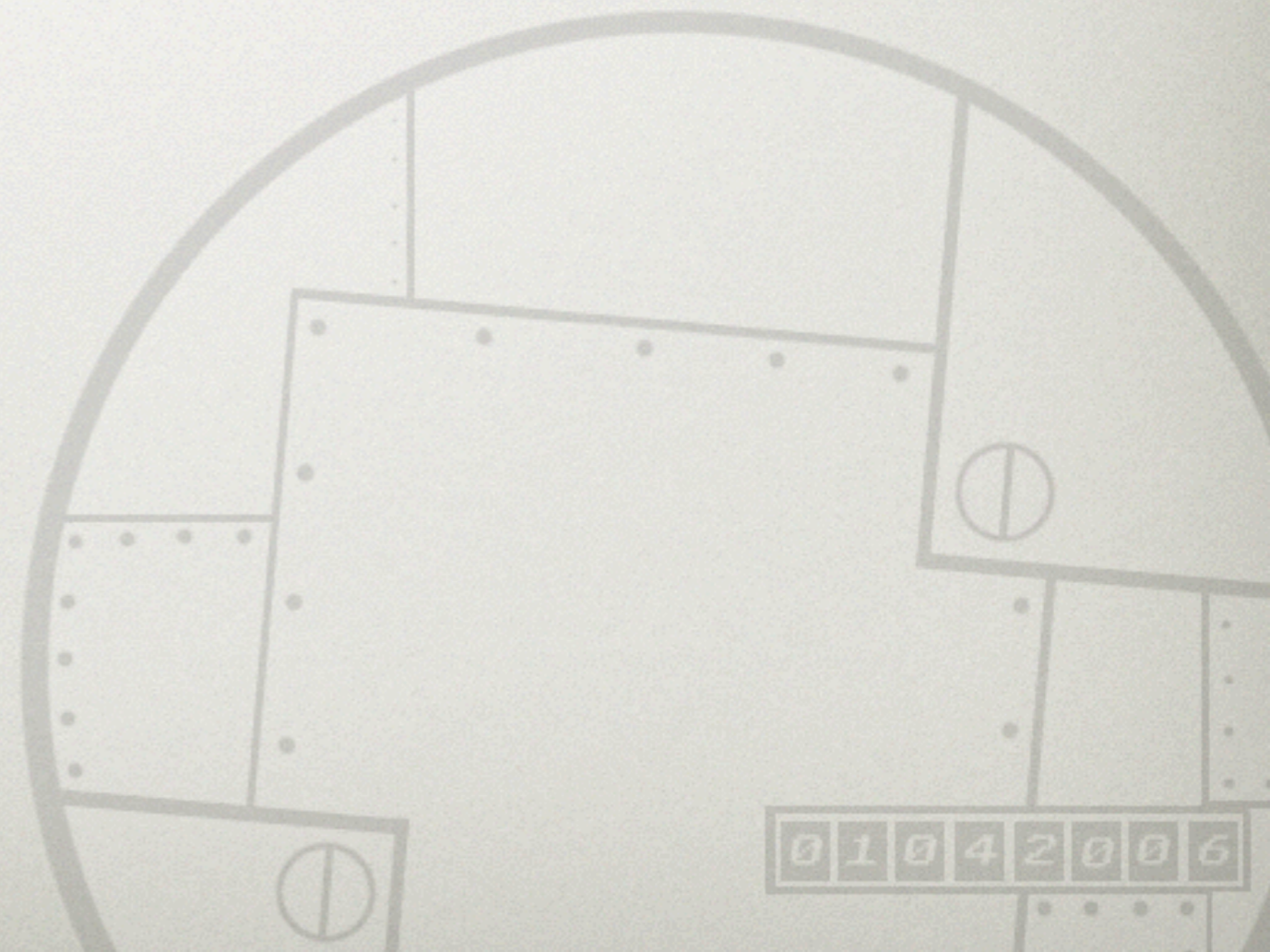
01042006

DATABASES



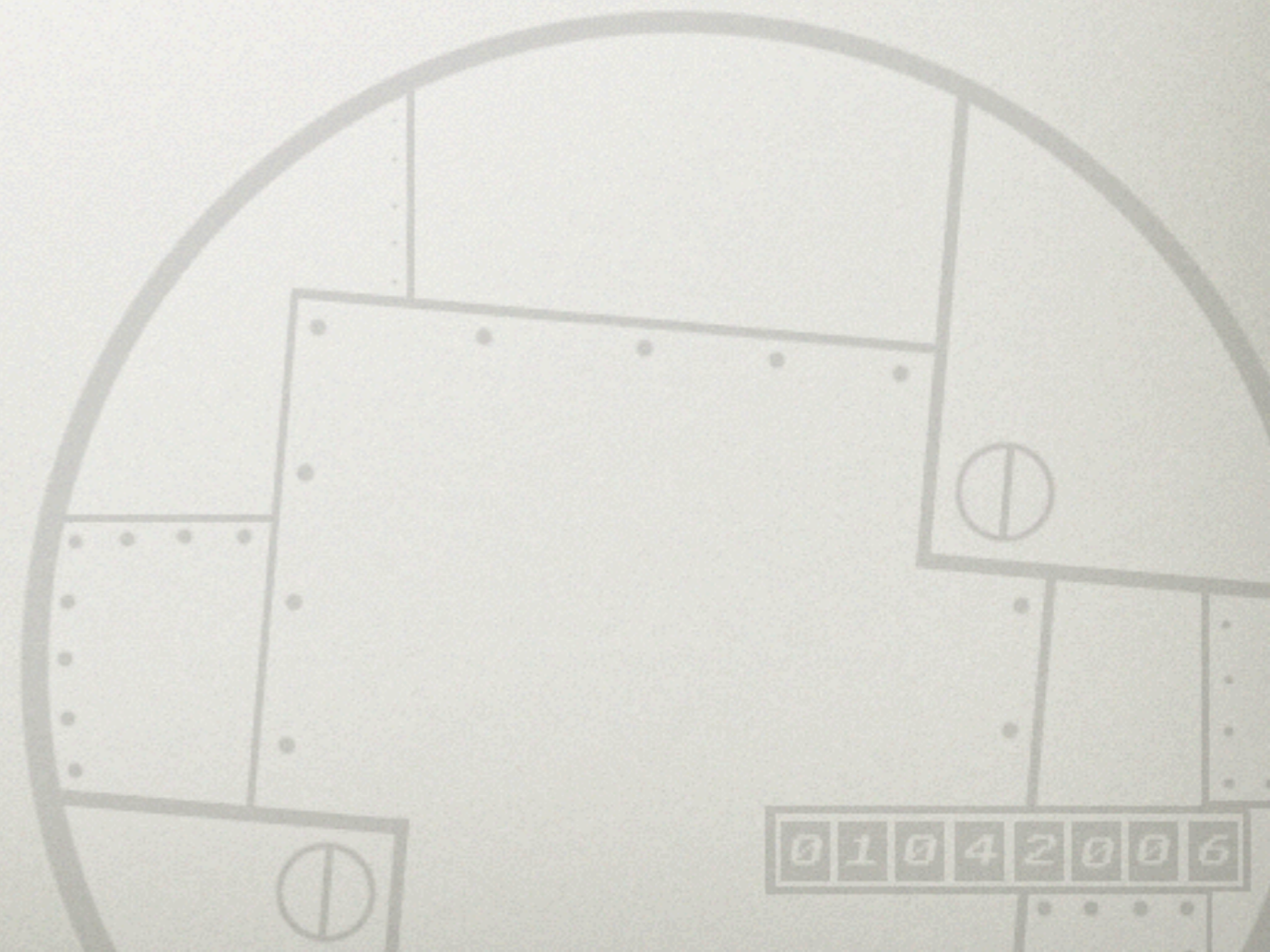
DATABASES

- What are the salient features of a database that should be preserved? (*significant properties...*)



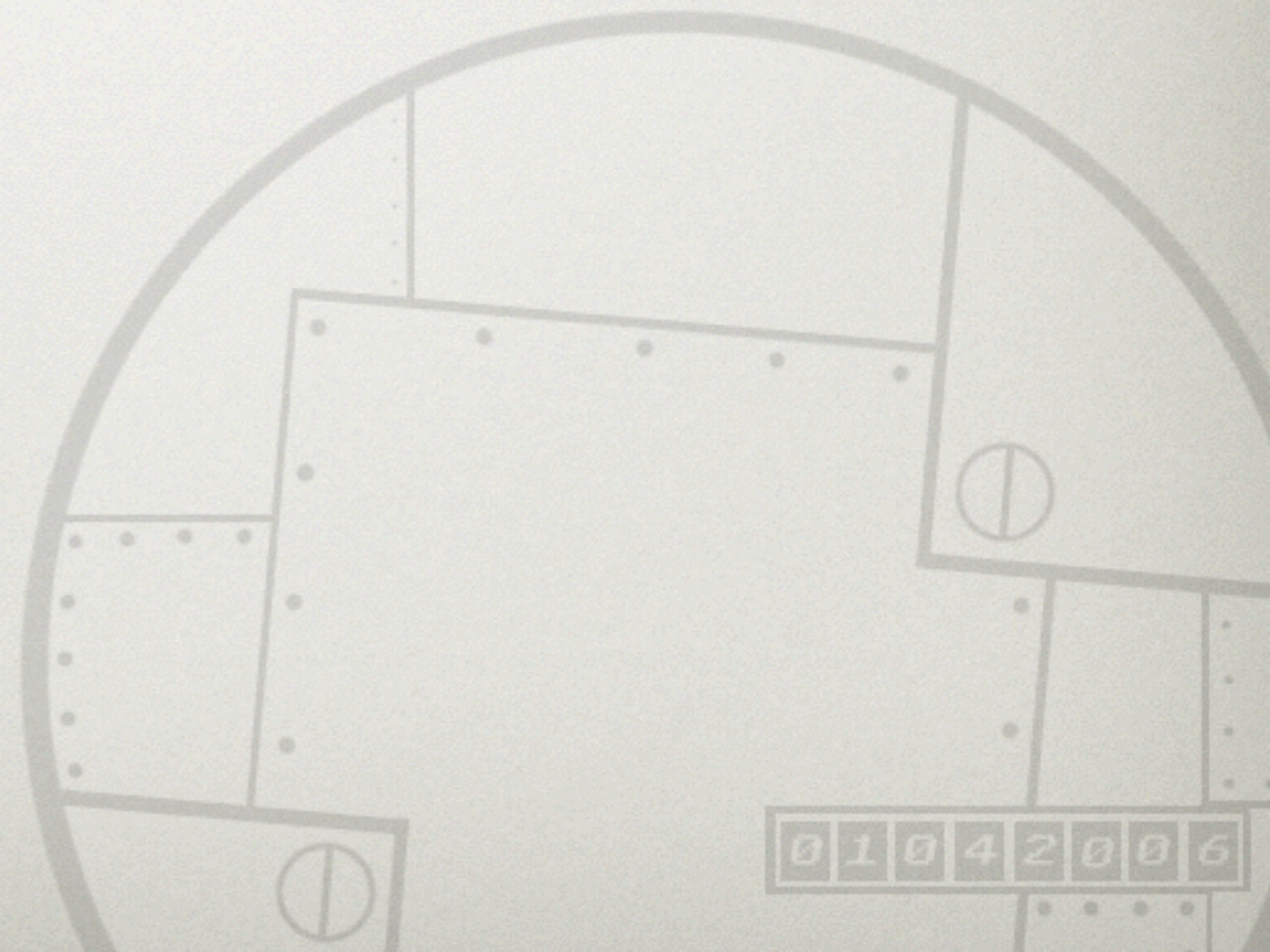
DATABASES

- What are the salient features of a database that should be preserved? (*significant properties...*)
- What are the different **stages** in the database preservation's life cycle?



DATABASES

- What are the salient features of a database that should be preserved? (*significant properties...*)
- What are the different **stages** in the database preservation's life cycle?
- What **documentation** is preserved together with a database, and in what **format**?



DATABASES

- What are the salient features of a database that should be preserved? (*significant properties...*)
- What are the different **stages** in the database preservation's life cycle?
- What **documentation** is preserved together with a database, and in what **format**?
- What are the **legal encumbrances** on database preservation?

DATABASES

- What are the **salient features** of a database that should be **preserved**? (*significant properties...*)
- What are the different **stages** in the database preservation's life cycle?
- What **documentation** is preserved together with a database, and in what **format**?
- What are the **legal encumbrances** on database preservation?
- What can be learned from traditional archival appraisal for the selection of databases for preservation?

DATABASES

- What are the **salient features** of a database that should be **preserved**? (*significant properties...*)
- What are the different **stages** in the database preservation's life cycle?
- What **documentation** is preserved together with a database, and in what **format**?
- What are the **legal encumbrances** on database preservation?
- What can be learned from traditional archival appraisal for the selection of databases for preservation?
- To what extent can the preservation strategies, and procedural policies developed by archivists be adapted for databases?

DATABASES

- What are the salient features of a database that should be preserved? (*significant properties...*)
- What are the different **stages** in the database preservation's life cycle?
- What **documentation** is preserved together with a database, and in what **format**?
- What are the **legal encumbrances** on database preservation?
- What can be learned from traditional archival appraisal for the selection of databases for preservation?
- To what extent can the preservation strategies, and procedural policies developed by archivists be adapted for databases?
- How can we measure the quality of preservation strategies when they are applied to data-bases? What are **DB significant properties**? (*quality assurance...*)

DATABASES

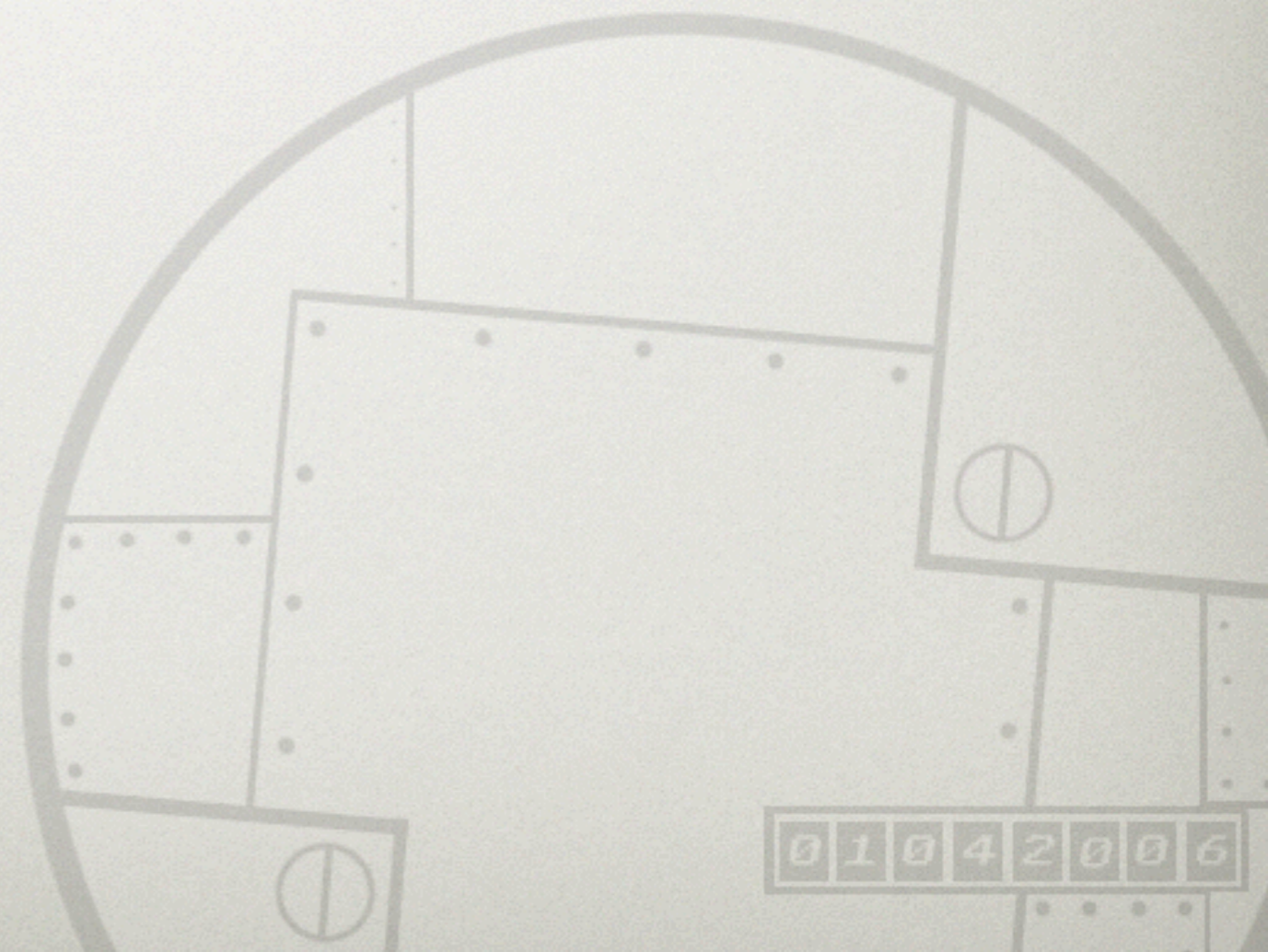
- What are the salient features of a database that should be preserved? (*significant properties...*)
- What are the different **stages** in the database preservation's life cycle?
- What **documentation** is preserved together with a database, and in what **format**?
- What are the **legal encumbrances** on database preservation?
- What can be learned from traditional archival appraisal for the selection of databases for preservation?
- To what extent can the preservation strategies, and procedural policies developed by archivists be adapted for databases?
- How can we measure the quality of preservation strategies when they are applied to data-bases? What are **DB significant properties**? (*quality assurance...*)

Technical questions...

01042006

DATABASES: GOALS

- How do we store them?
- How do we access them?



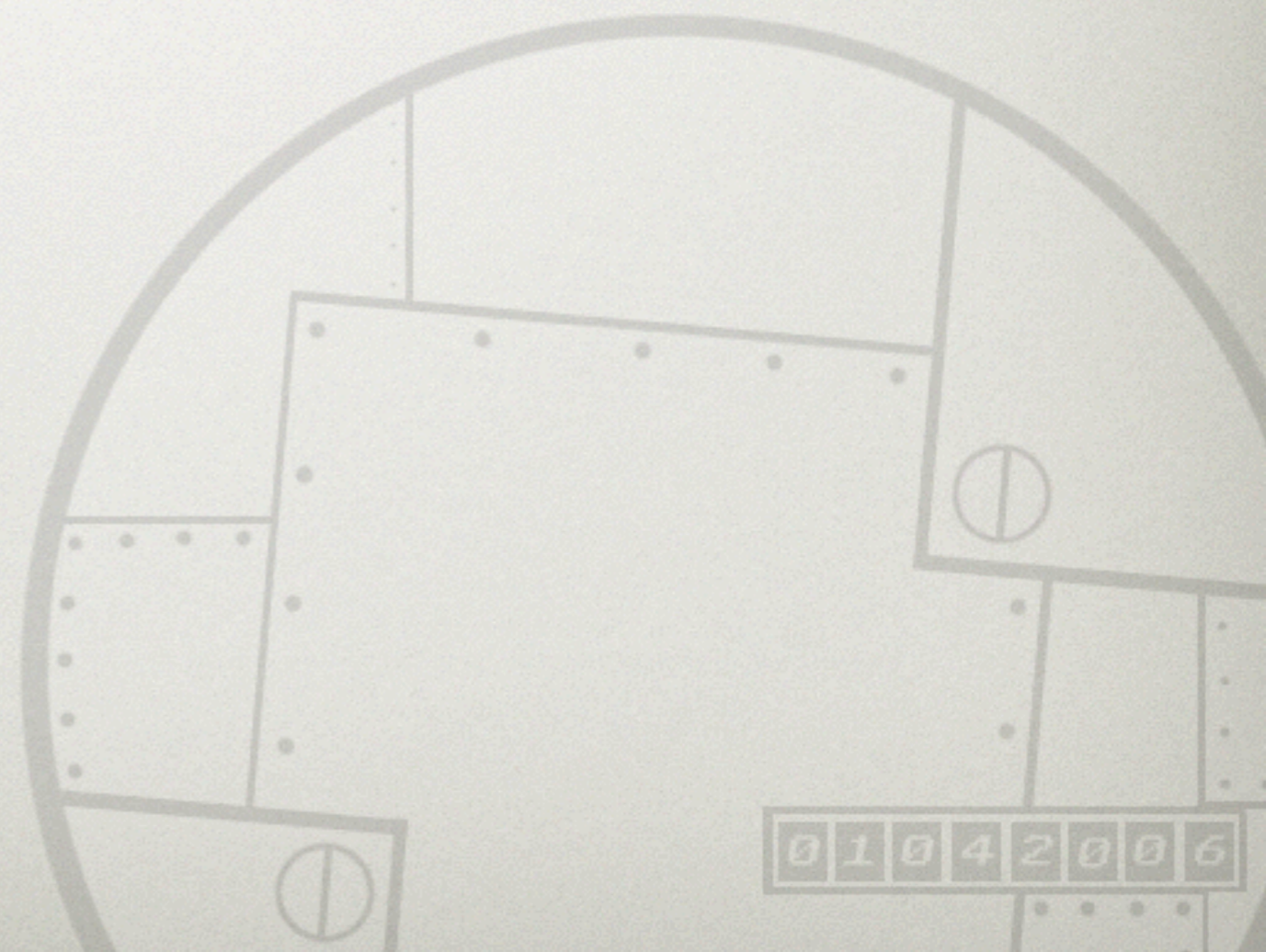
DATABASES: GOALS

- How do we store them?
- How do we access them?

RODA questions...

01042006

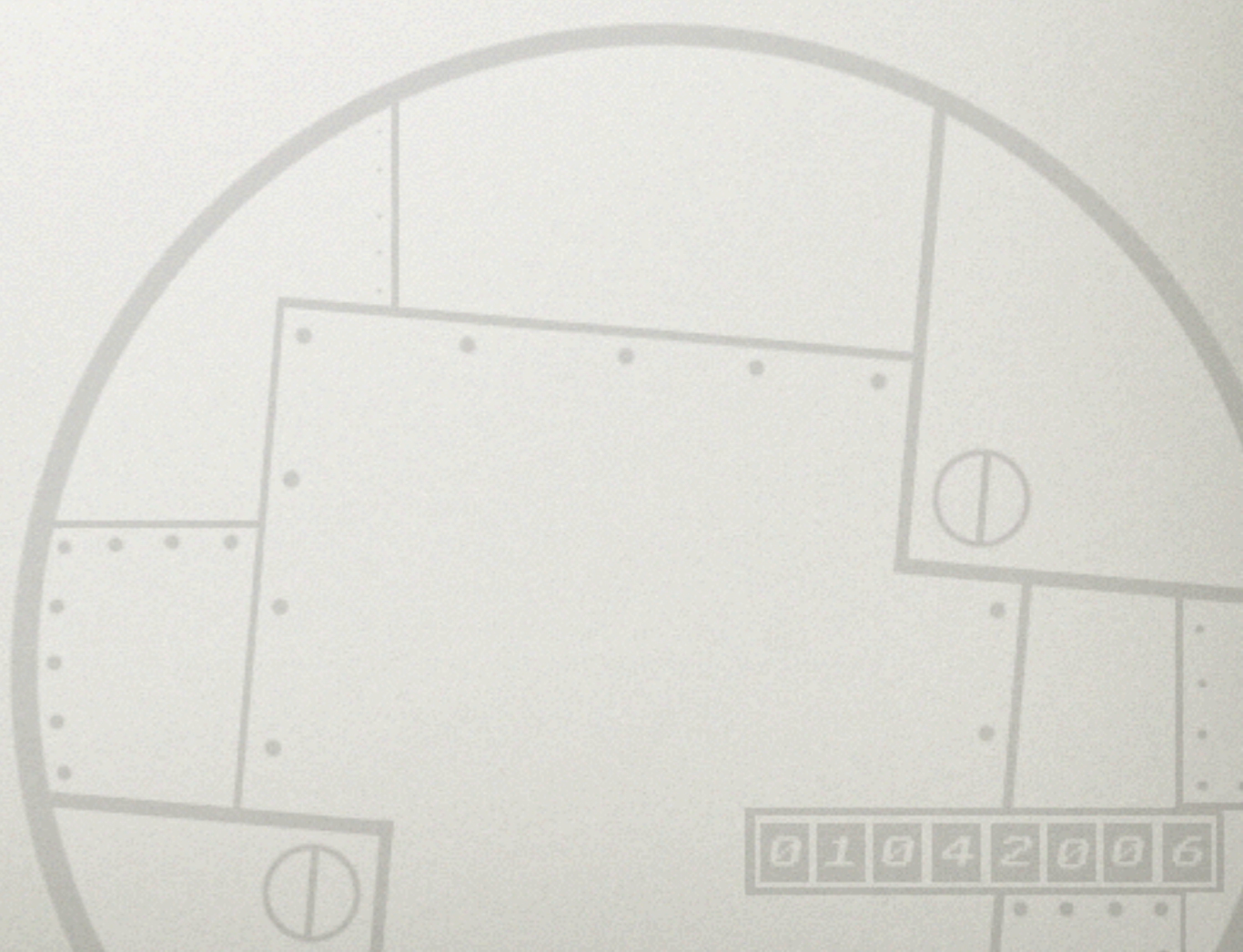
DATABASES



DATABASES

Normal evolution path:

Data => Structure => Semantics



DATABASES

Normal evolution path:

Data => Structure => Semantics

First prototype:

- Data
- Structure
- Only “frozen” DBs

DATABASES

- Data?

Normal evolution path:

Data => Structure => Semantics

First prototype:

- Data
- Structure
- Only “frozen” DBs

DATABASES

Normal evolution path:

Data => Structure => Semantics

- Data?
- Structure?

First prototype:

- Data
- Structure
- Only “frozen” DBs

DATABASES

Normal evolution path:

Data => Structure => Semantics

- Data?
- Structure?
- Views?

First prototype:

- Data
- Structure
- Only “frozen” DBs

DATABASES

Normal evolution path:

Data => Structure => Semantics

- Data?
- Structure?
- Views?
- Reports?

First prototype:

- Data
- Structure
- Only “frozen” DBs

DATABASES

Normal evolution path:

Data => Structure => Semantics

- Data?
- Structure?
- Views?
- Reports?
- Stored Procedure

First prototype:

- Data
- Structure
- Only “frozen” DBs

DATABASES

Normal evolution path:

Data => Structure => Semantics

- Data?
- Structure?
- Views?
- Reports?
- Stored Procedure
- ...

First prototype:

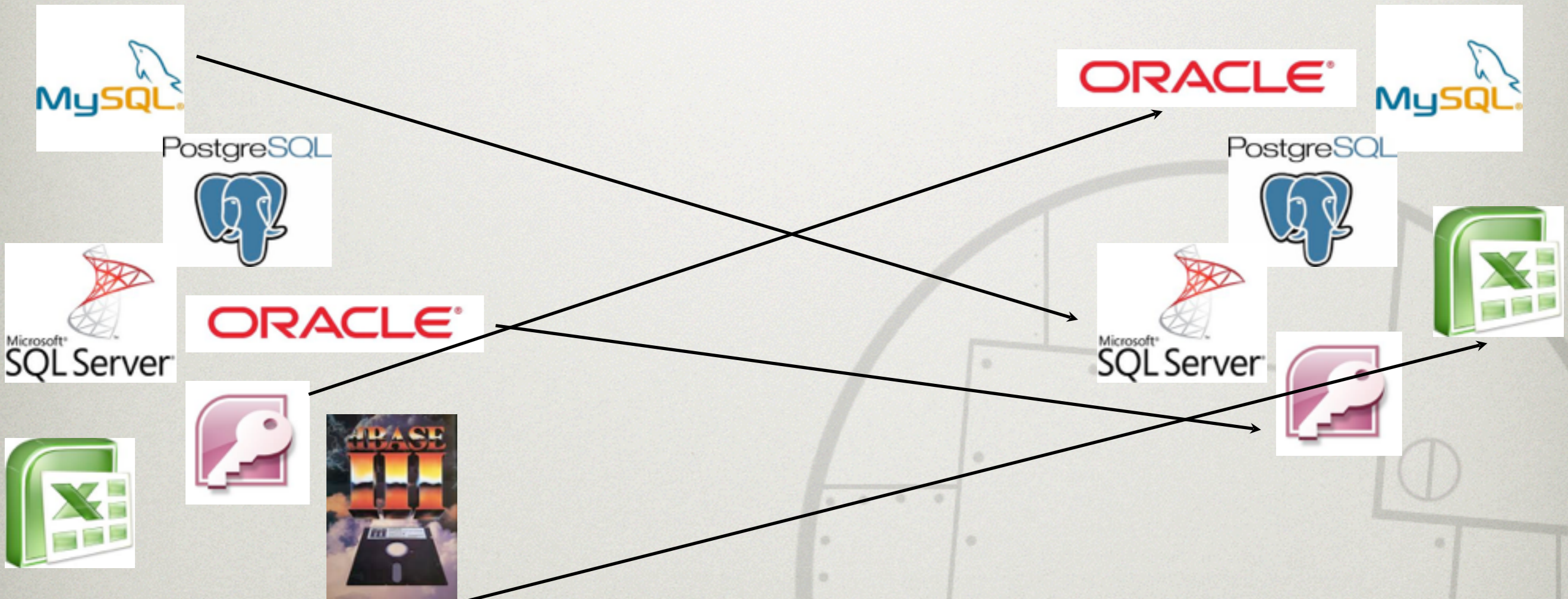
- Data
- Structure
- Only “frozen” DBs

HOW?

The need for an intermediate representation

Input Formats (M)

Output Formats (N)

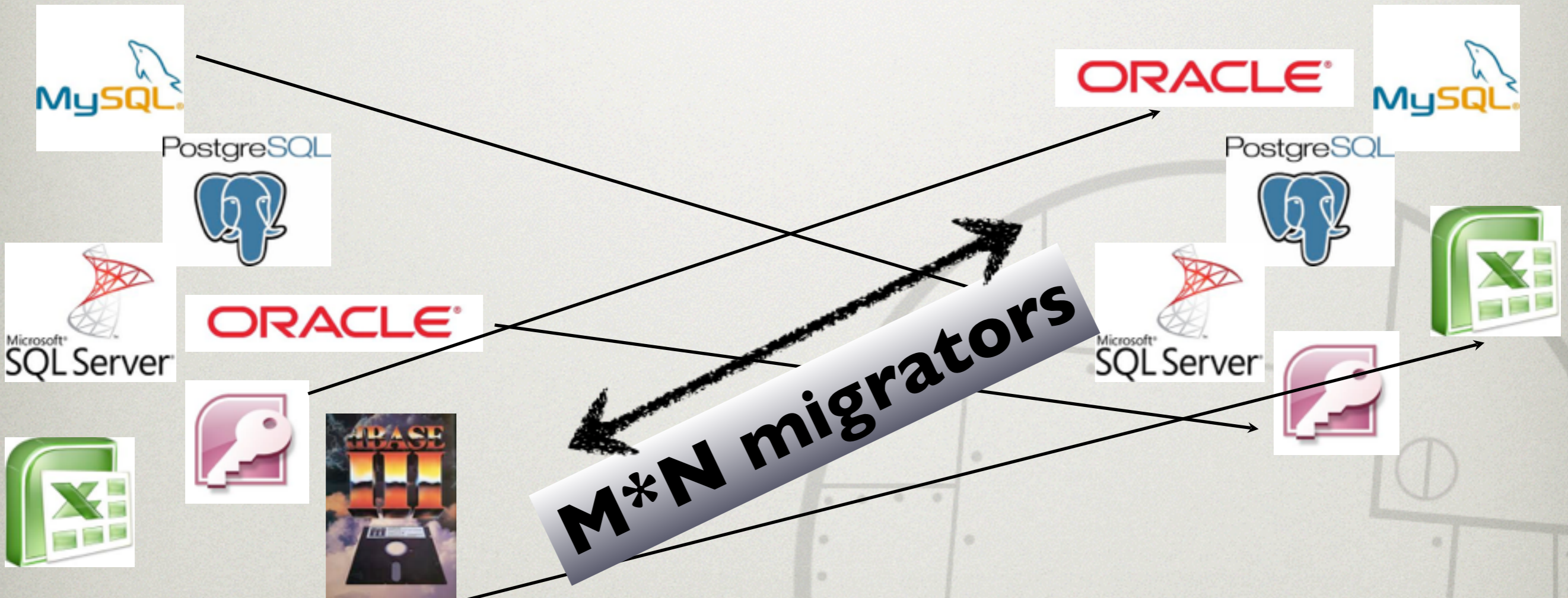


HOW?

The need for an intermediate representation

Input Formats (M)

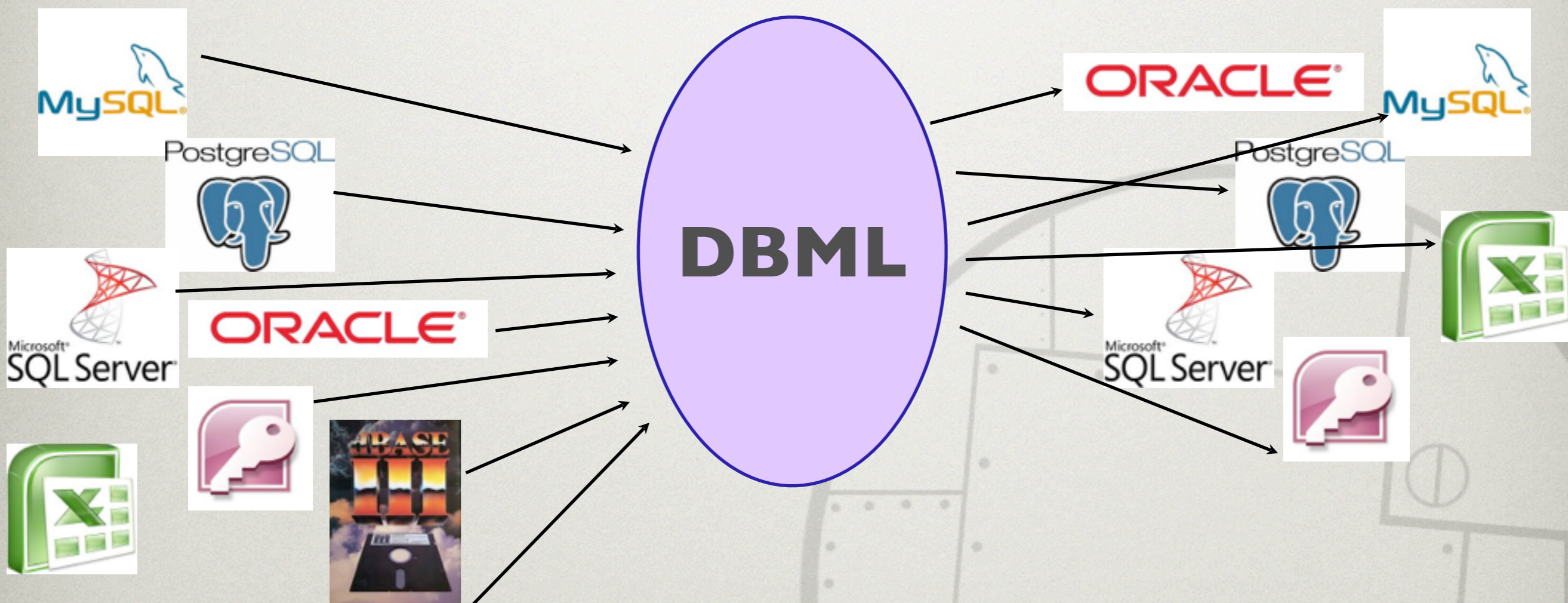
Output Formats (N)



IR: DBML

Input Formats (M)

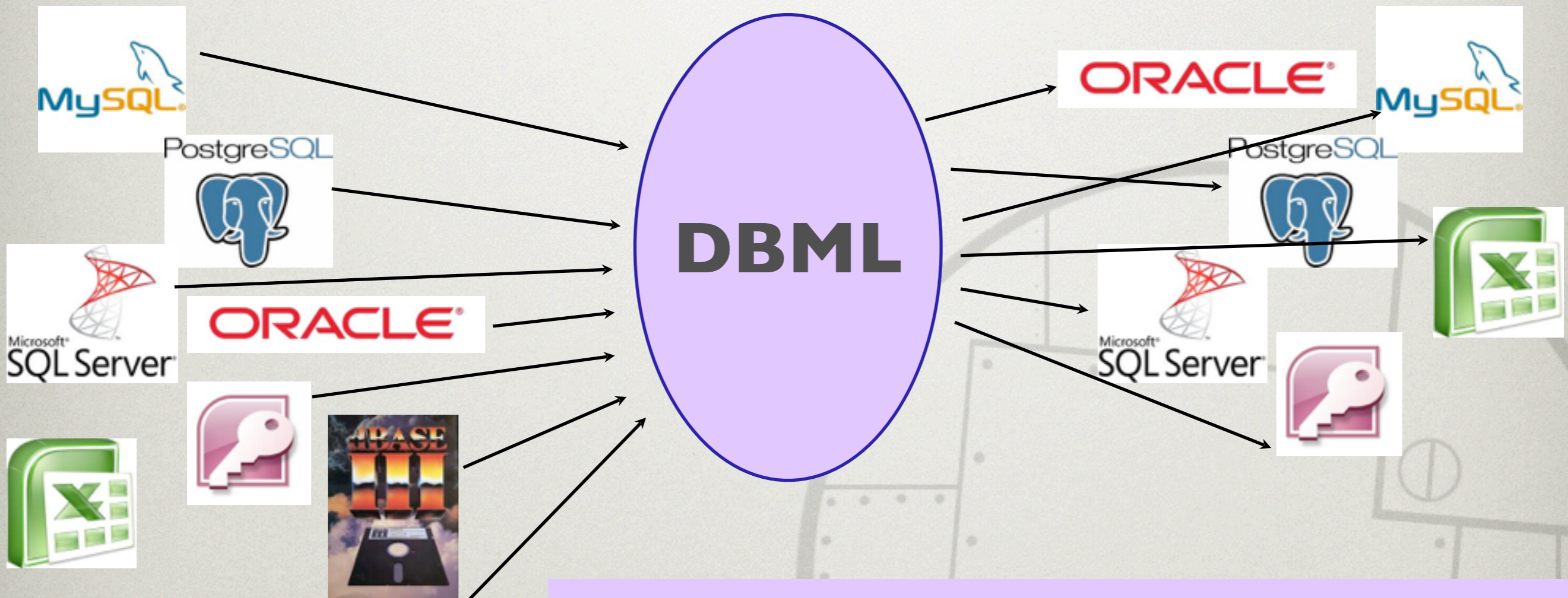
Output Formats (N)



IR: DBML

Input Formats (M)

Output Formats (N)



For each new output or input format you only need to code one new migrator.



DBML DESIGN PRINCIPLES

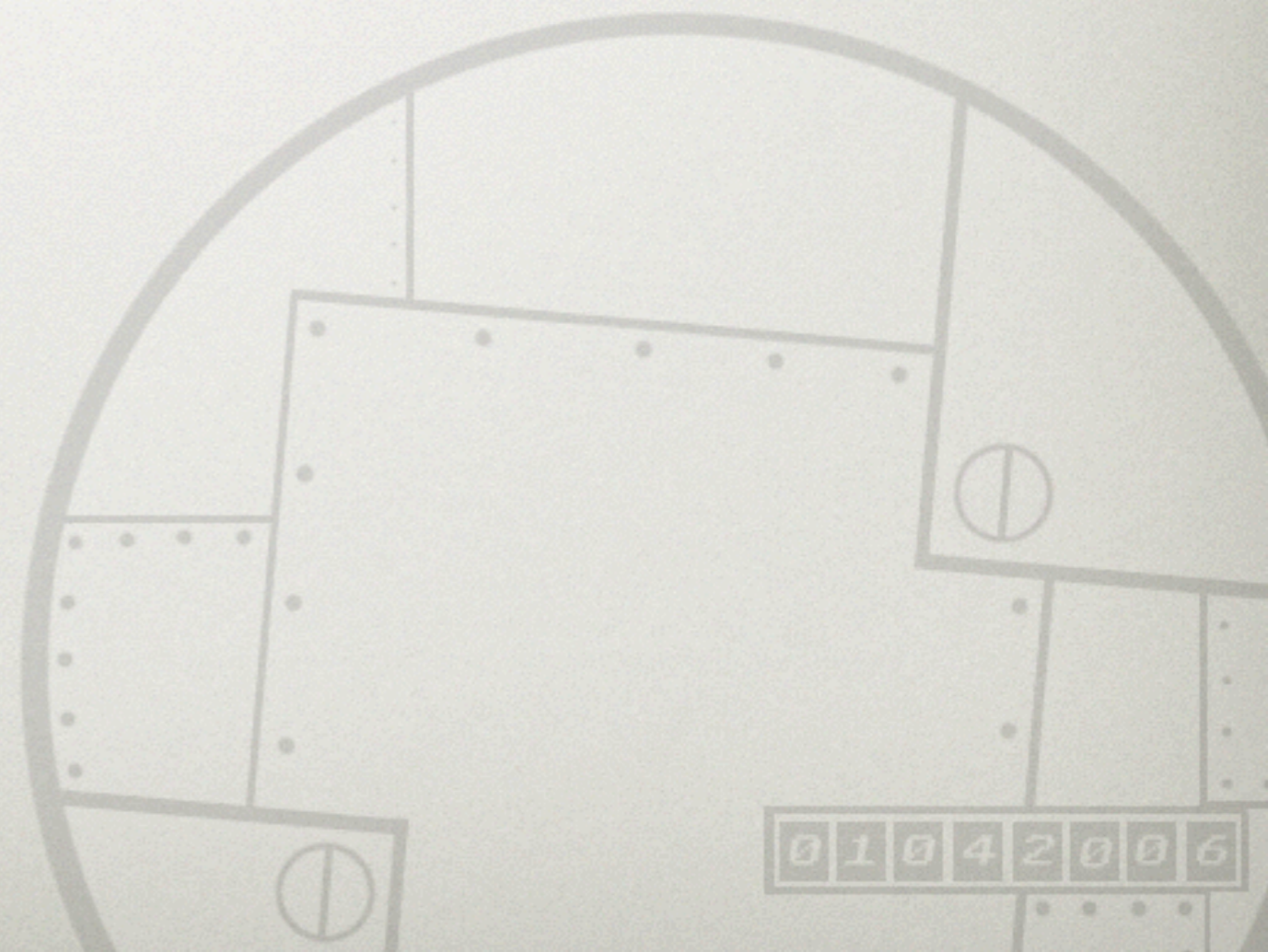
- Hardware independent;
- Software independent;
- Easy to process;
- Descriptive;
- It should be possible to add metadata;
- It should be possible to add semantics;

DBML DESIGN PRINCIPLES

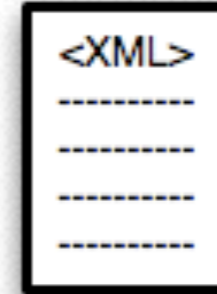
- Hardware independent;
- Software independent;
- Easy to process;
- Descriptive;
- It should be possible to add metadata;
- It should be possible to add semantics;

XML was the obvious choice

SIP STRUCTURE (DB EXAMPLE)

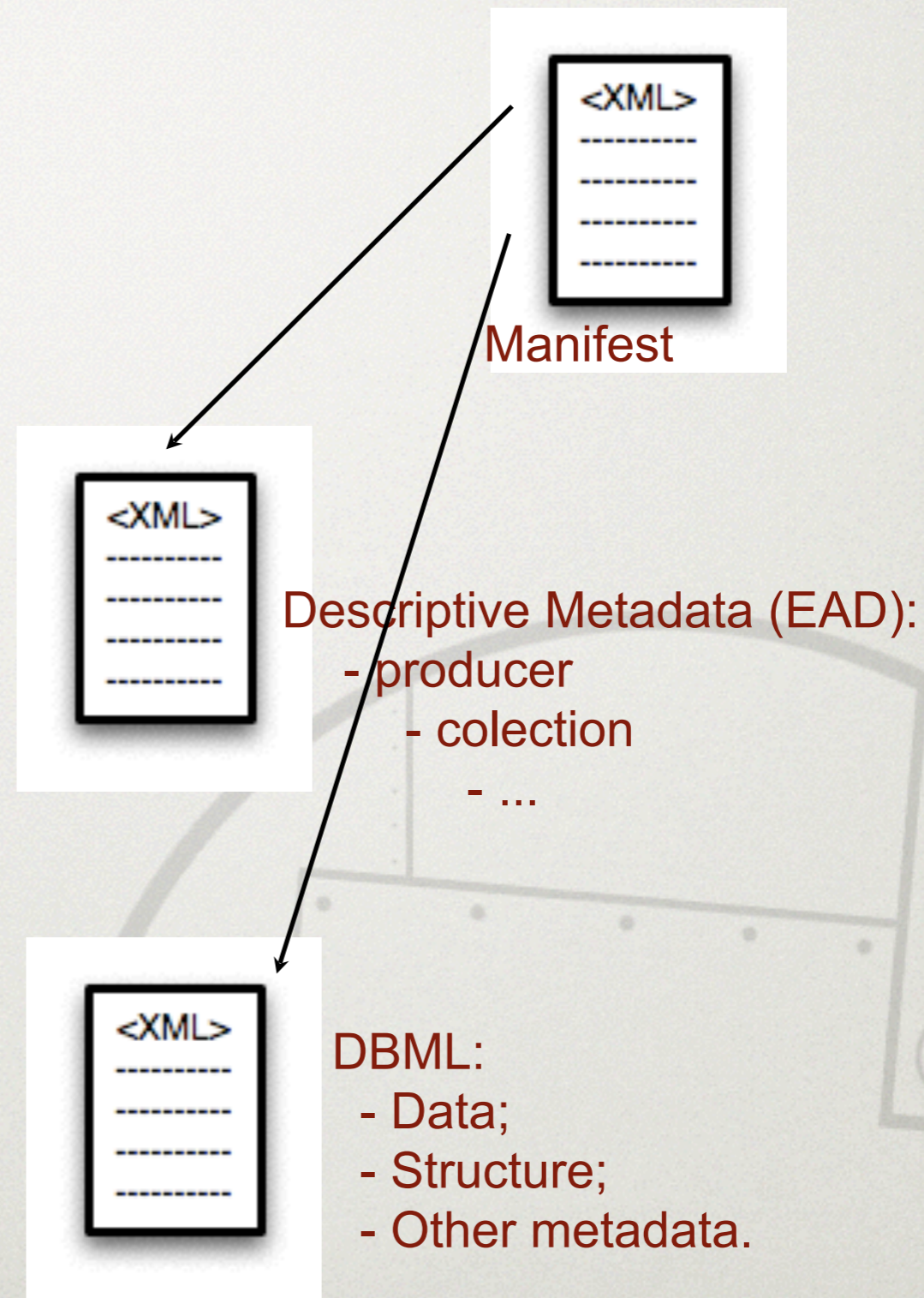


SIP STRUCTURE (DB EXAMPLE)



Manifest

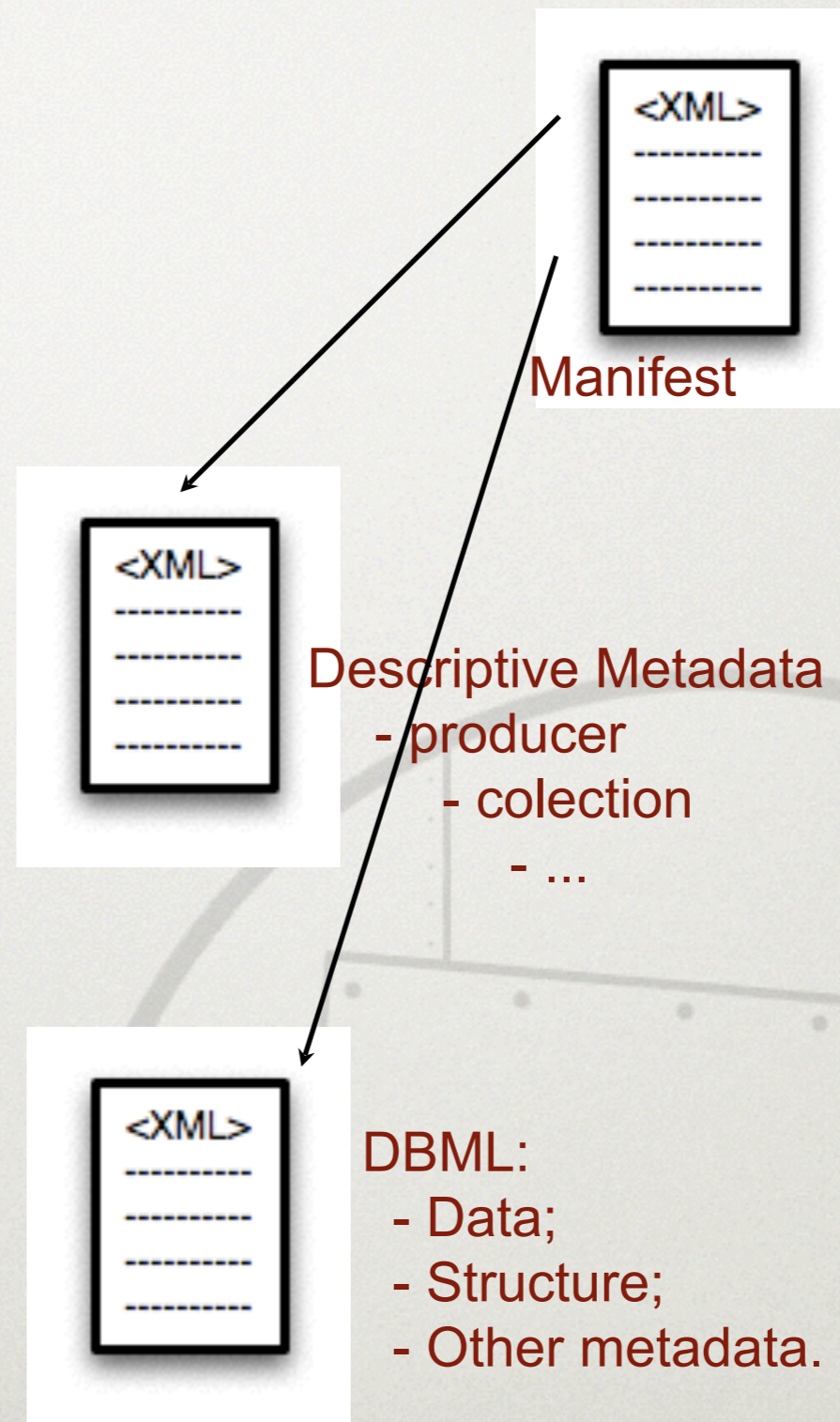
SIP STRUCTURE (DB EXAMPLE)



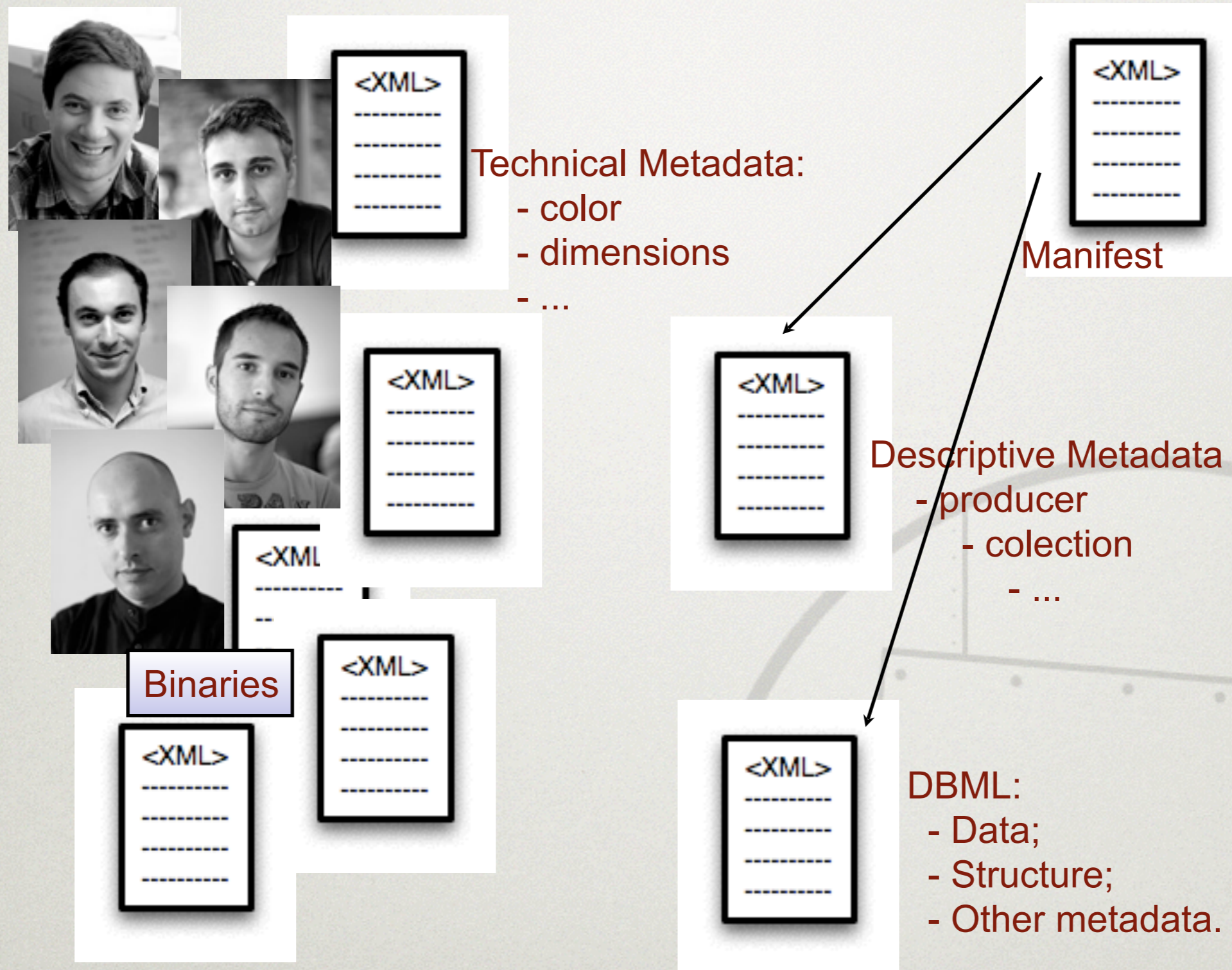
SIP STRUCTURE (DB EXAMPLE)



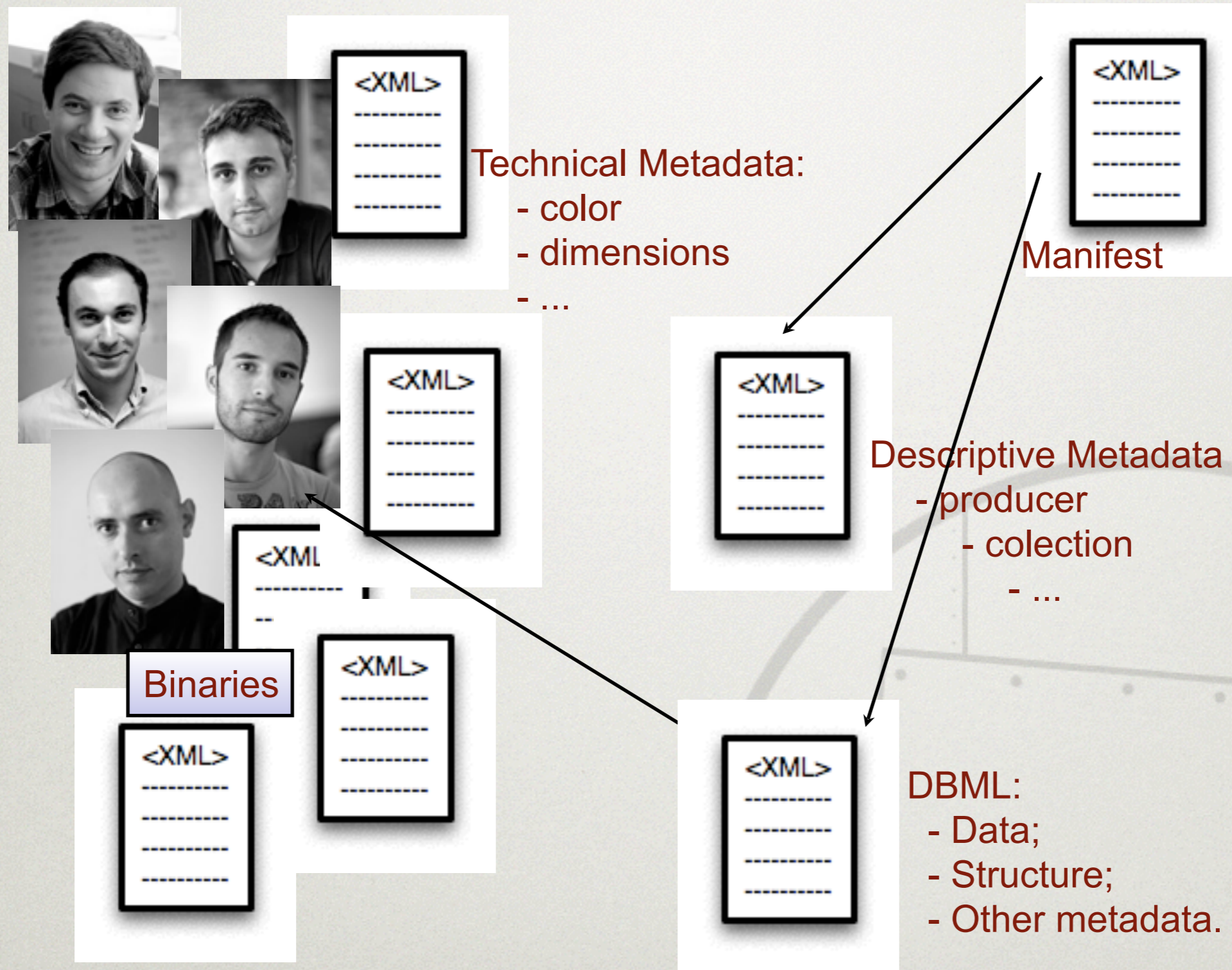
Binaries



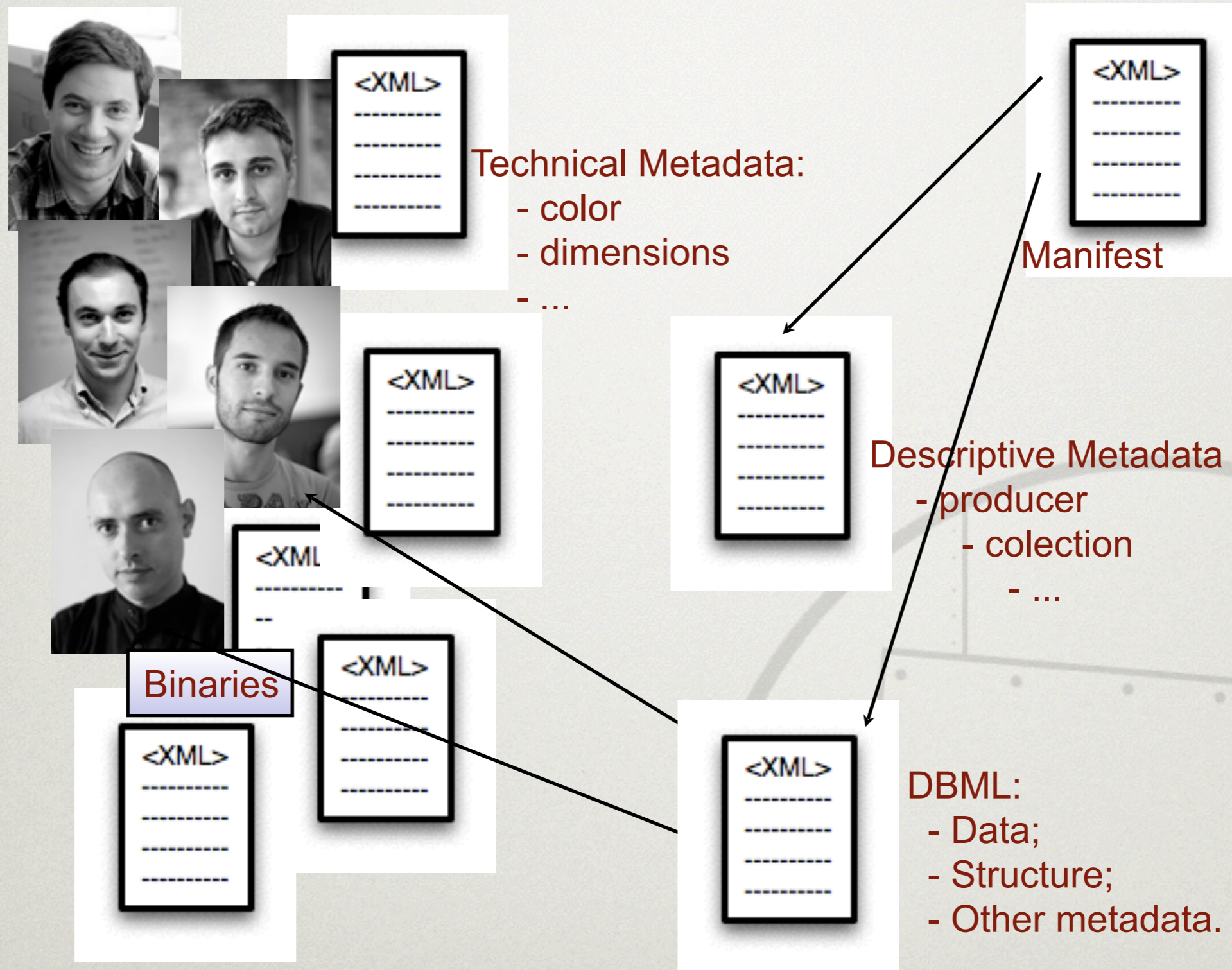
SIP STRUCTURE (DB EXAMPLE)



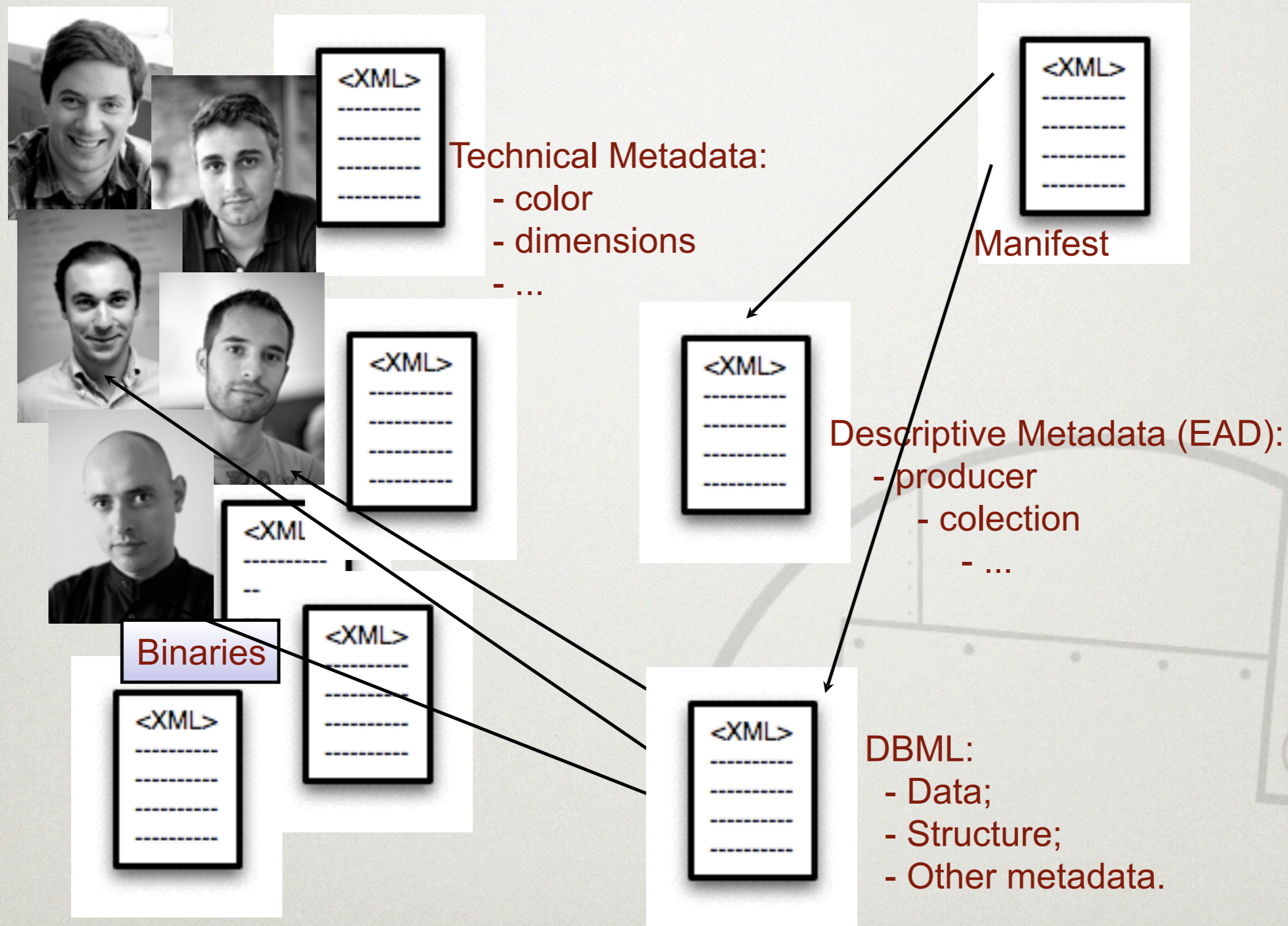
SIP STRUCTURE (DB EXAMPLE)



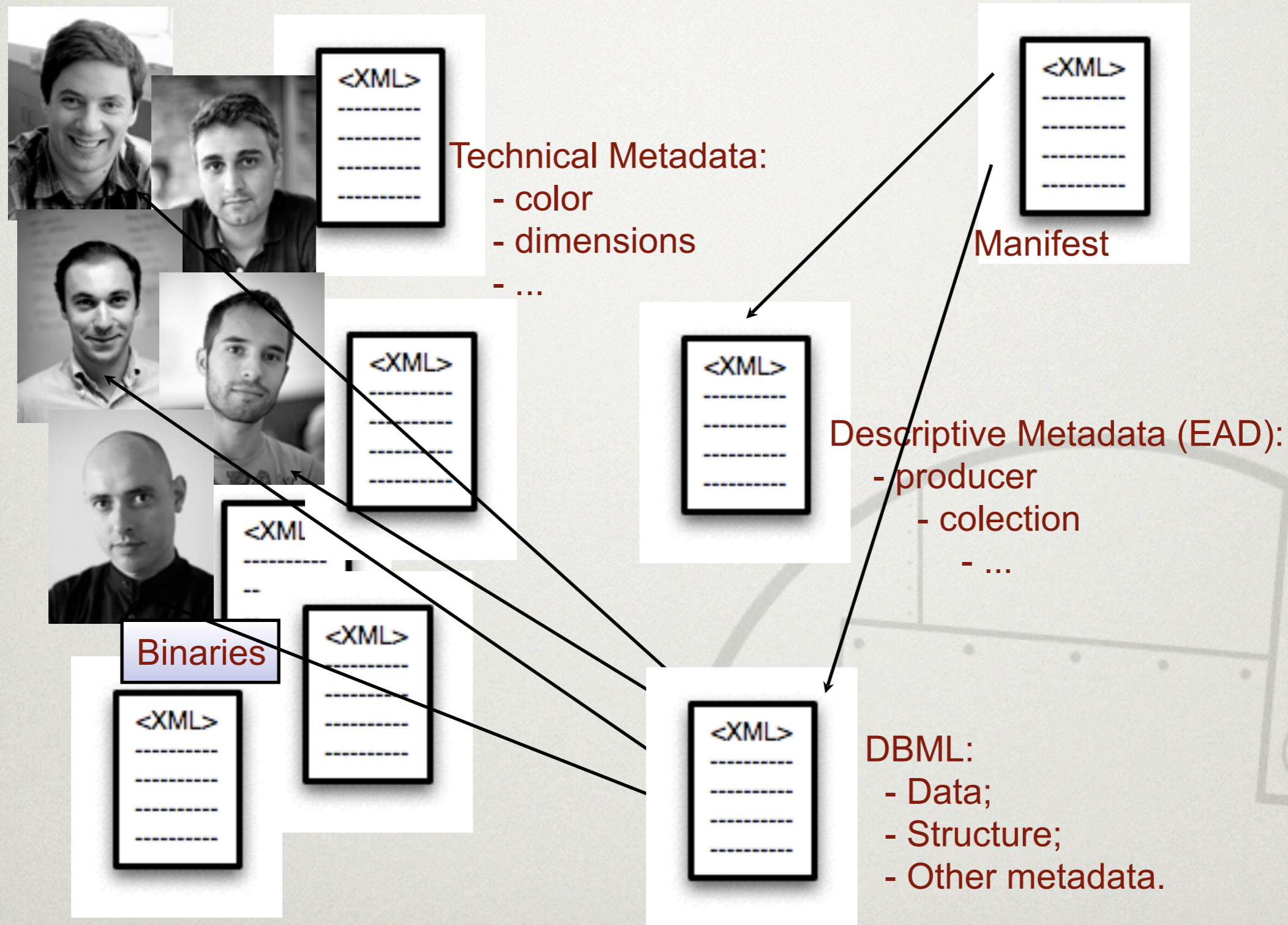
SIP STRUCTURE (DB EXAMPLE)



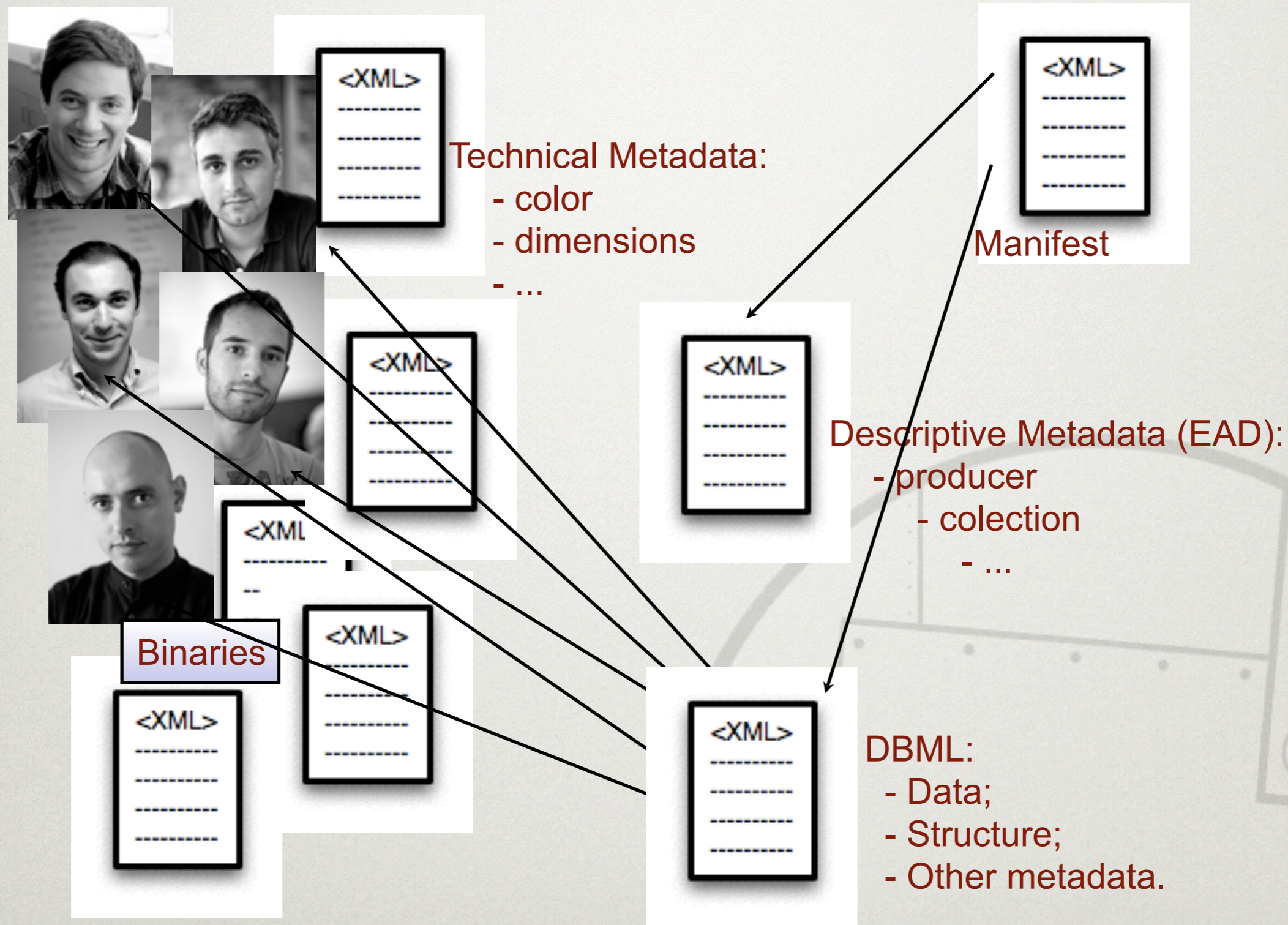
SIP STRUCTURE (DB EXAMPLE)



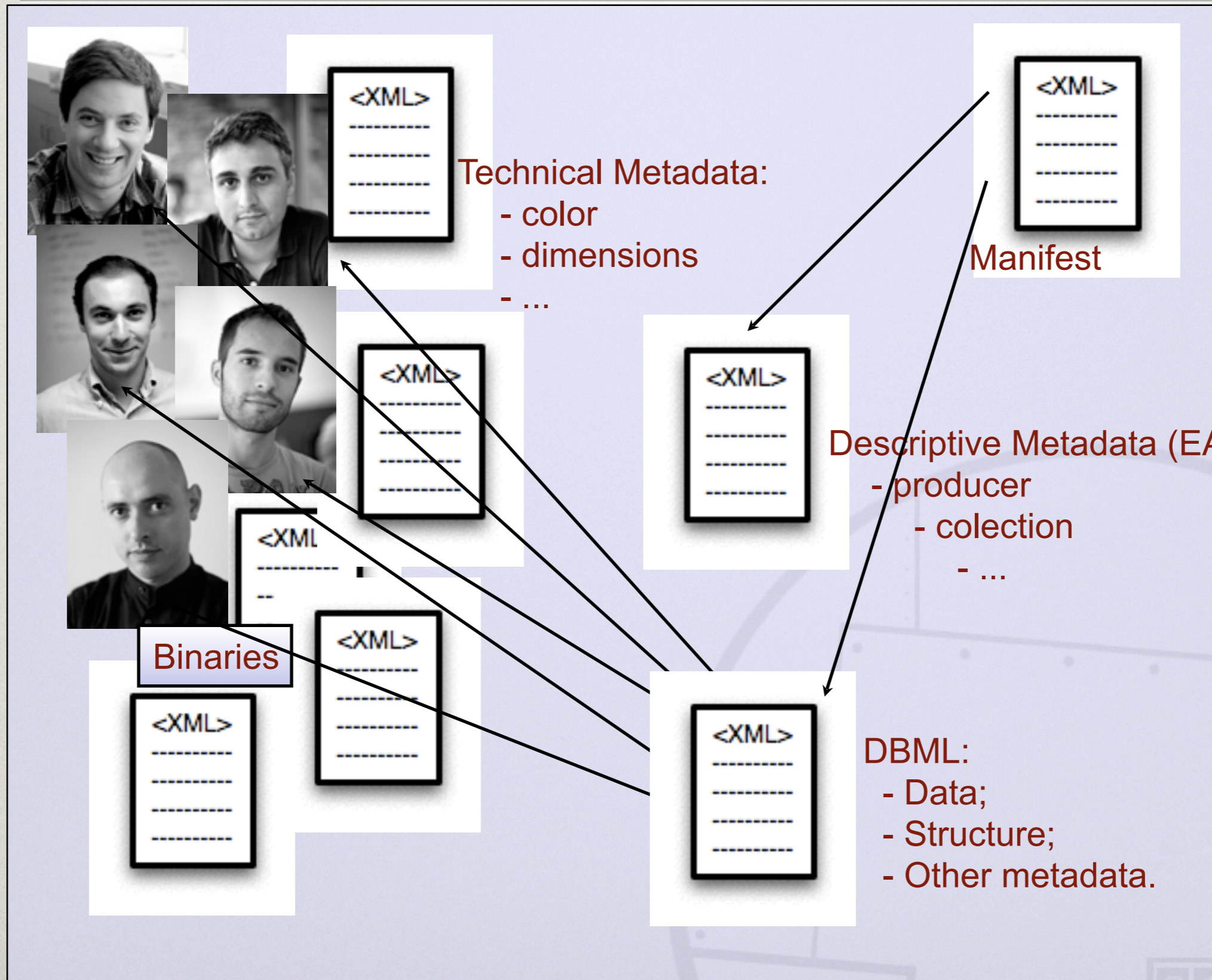
SIP STRUCTURE (DB EXAMPLE)



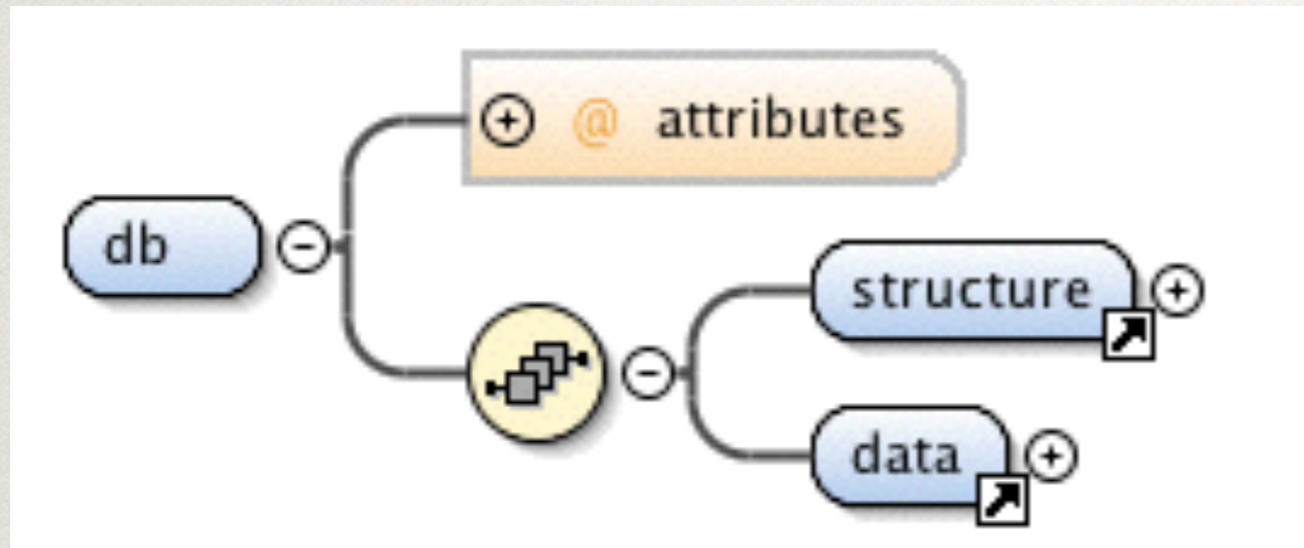
SIP STRUCTURE (DB EXAMPLE)



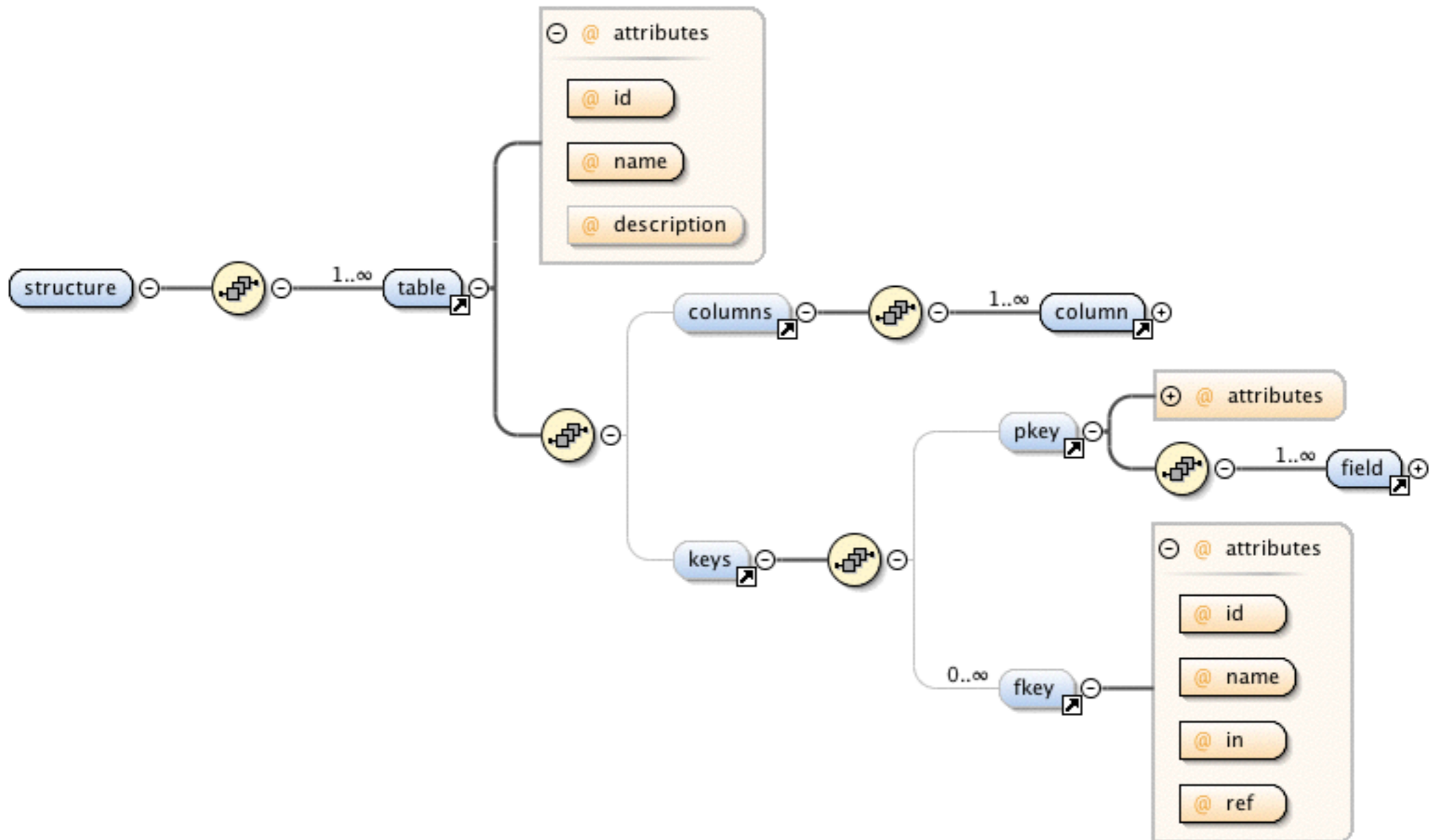
SIP STRUCTURE (DB EXAMPLE)



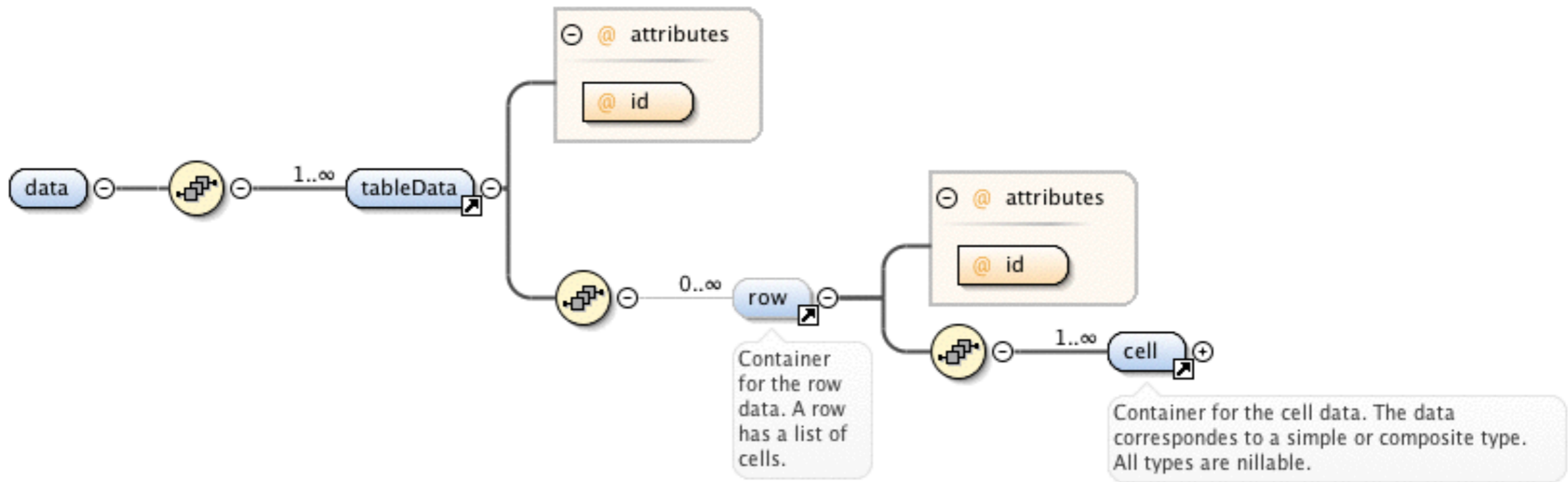
DBML



DBML: STRUCTURE



DBML: DATA



DATABASES: RODA ANSWERS

- How do we store them?
 - ★ DBML + binaries + technical metadata
- How do we access them?
 - ★ PhpMyAdmin (hacked version)

DATABASES: RODA ANSWERS

- How do we store them?
 - ★ DBML + binaries + technical metadata
- How do we access them?
 - ★ PhpMyAdmin (hacked version)

RODA answers...

01042006

SOME PROBLEMS

- Extracting data is easy:
 - ◆ `SELECT * FROM ...`
- Extracting the structure is not:
 - ◆ DBMS protect this information;
 - ◆ Each DBMS stores it differently;
 - ◆ Different versions of the same DBMS can also act differently;
 - ◆ We have to “prepare/hack” the DBMS.

Open Planets: “Database Archiving” (Copenhagen 2012)

Two approaches emerged:

- **Preserve the database and the environment allowing authentic access to the information and any accompanying applications**

- **Extract / migrate the raw data and table structure from the original database**



Open Planets: “Database Archiving” (Copenhagen 2012)

Two approaches emerged:

- Preserve the database and the environment allowing authentic access to the information and any accompanying applications

Relies on emulation strategy

- Extract / migrate the raw data and table structure from the original database



Open Planets: “Database Archiving” (Copenhagen 2012)

Two approaches emerged:

- Preserve the database and the environment all together, allowing authentic access to the information and any accompanying applications

Relies on emulation strategy

- Extract / migrate the raw data and table structure from the original database

Relies on migration strategy



Open Planets: “Database Archiving” (Copenhagen 2012)

Two approaches emerged:

- Preserve the database and the environment all together, allowing authentic access to the information and any accompanying applications

Relies on emulation strategy

- Extract / migrate the raw data and table structure from the original database

Relies on migration strategy

RODA+DBML



Open Planets: “Database Archiving” (Copenhagen 2012)

Two approaches emerged:

- Preserve the database and the environment all together, allowing authentic access to the information and any accompanying applications

Relies on emulation strategy

- Extract / migrate the raw data and table structure from the original database

Relies on migration strategy

RODA+DBML

SIARD



Open Planets: “Database Archiving” (Copenhagen 2012)

Two approaches emerged:

- Preserve the database and the environment all together, allowing authentic access to the information and any accompanying applications

Relies on emulation strategy

- Extract / migrate the raw data and table structure from the original database

Relies on migration strategy

RODA+DBML

SIARD

SIARD-DK



Open Planets: “Database Archiving” (Copenhagen 2012)

Two approaches emerged:

- Preserve the database and the environment all together, allowing authentic access to the information and any accompanying applications

Relies on emulation strategy

- Extract / migrate the raw data and table structure from the original database

Relies on migration strategy

RODA+DBML

SIARD

SIARD-DK

ADDML



DBML vs SIARD

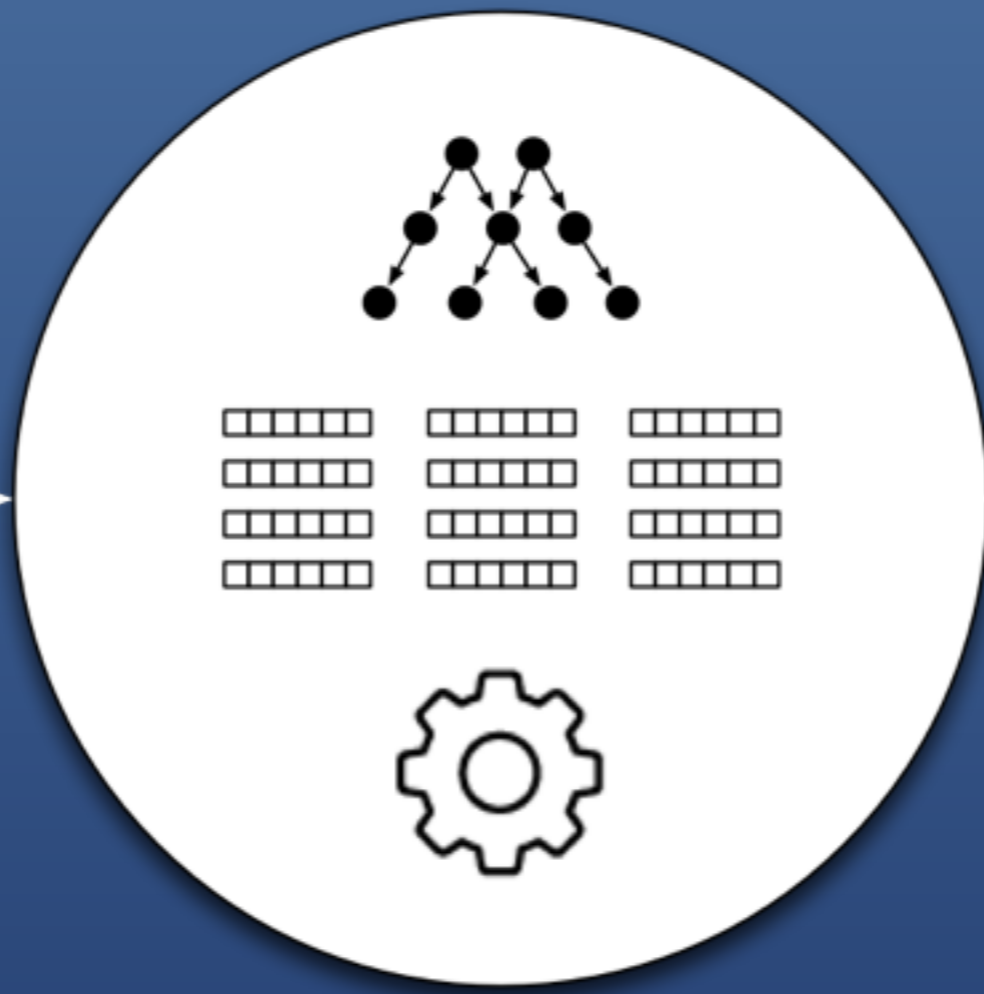
	DBML	SIARD
Data	Yes (no segmentation)	Yes (with segmentation)
Structure	Yes (XML)	Yes (XSD)
Stored Procedures	No	Yes (XML)
Triggers	No	Yes (XML)
Views	No	Yes (XML)

db-preservation-toolkit architecture

Import modules

- MySQL
- Oracle'12
- SQL Server
- PostgreSQL
- DB2
- MS Access
- ODBC
- DBML
- SIARD

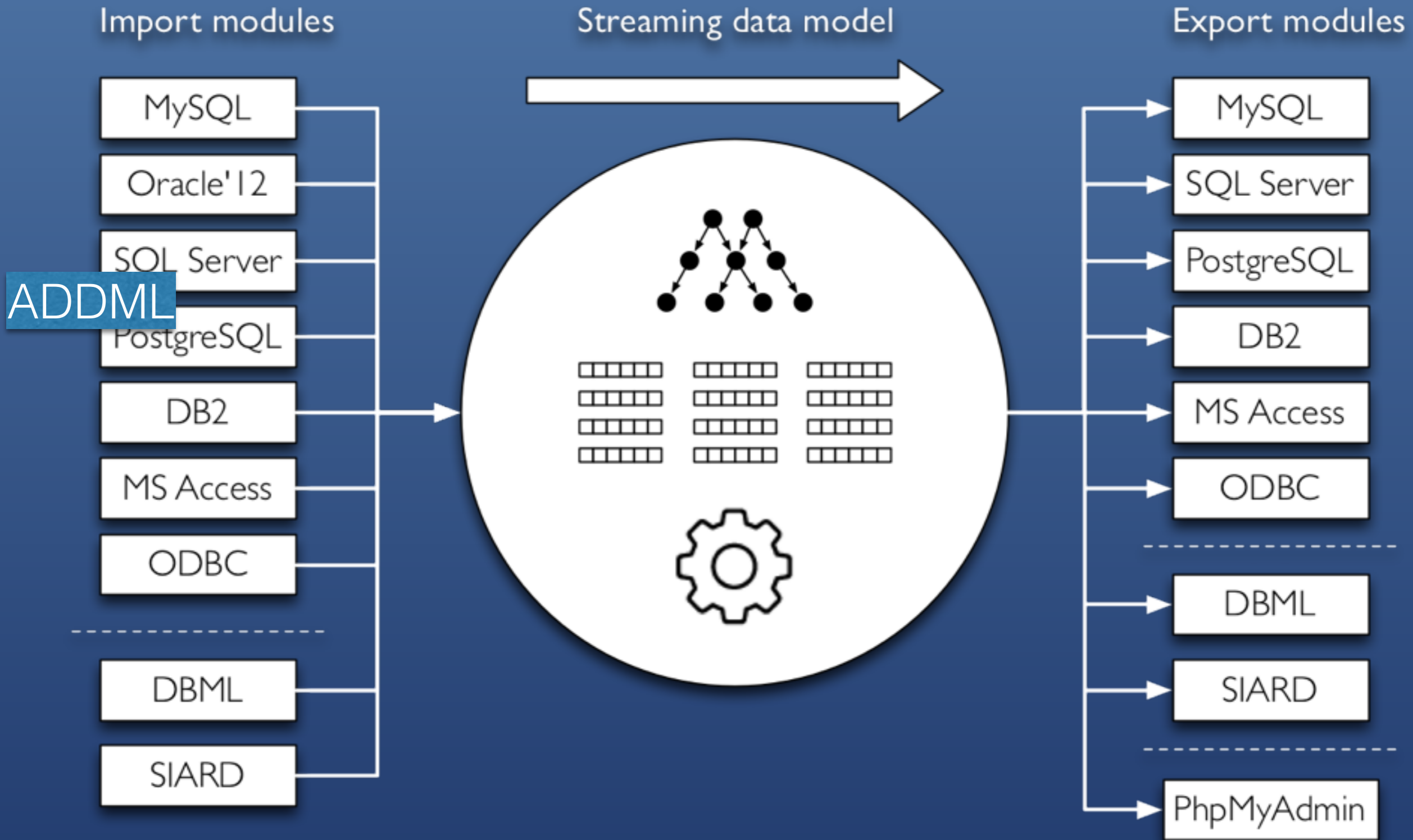
Streaming data model



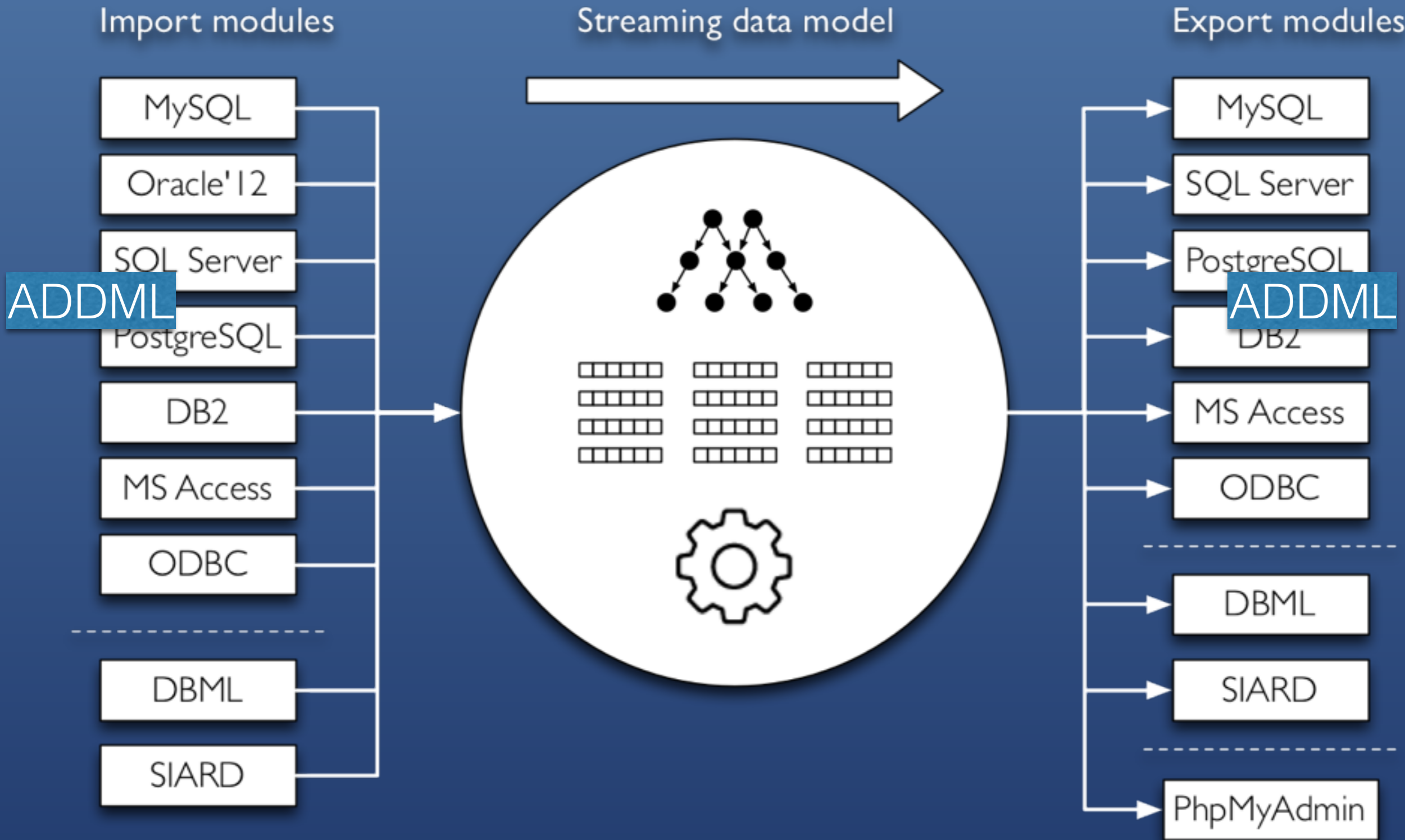
Export modules

- MySQL
- SQL Server
- PostgreSQL
- DB2
- MS Access
- ODBC
- DBML
- SIARD
- PhpMyAdmin

db-preservation-toolkit architecture

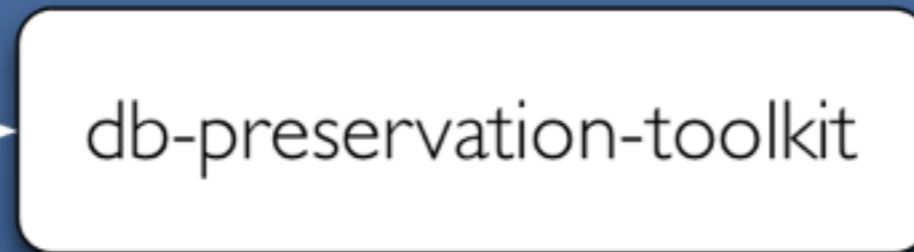
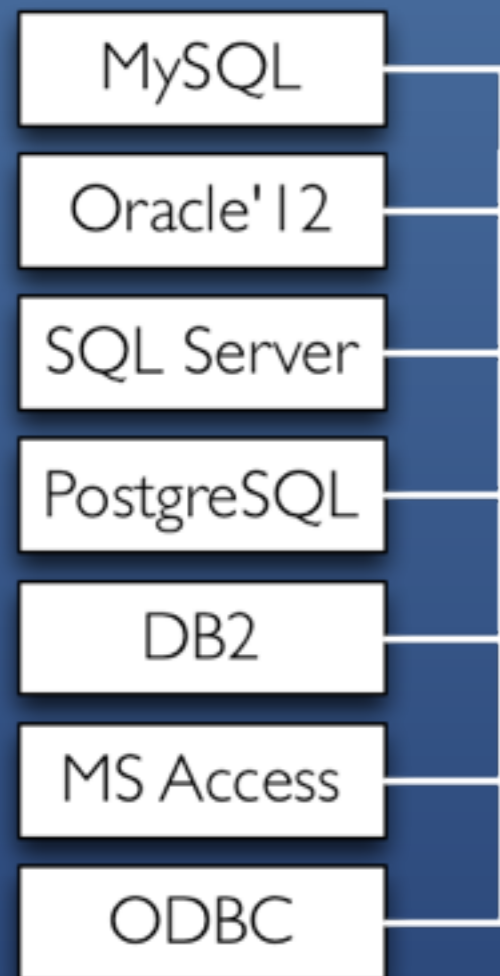


db-preservation-toolkit architecture



Pre-ingest / Production / SIP creation

Origin formats



Preservation formats

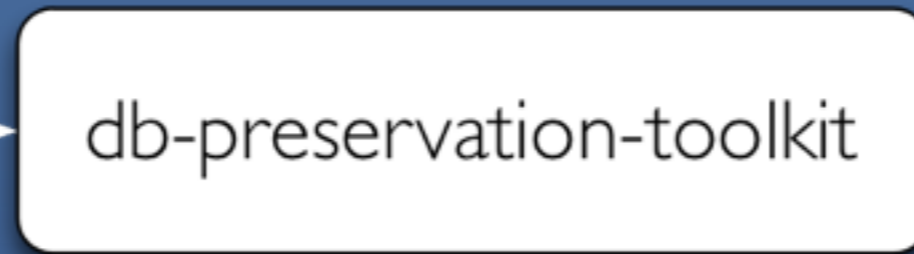


Pre-ingest / Production / SIP creation

Origin formats



new

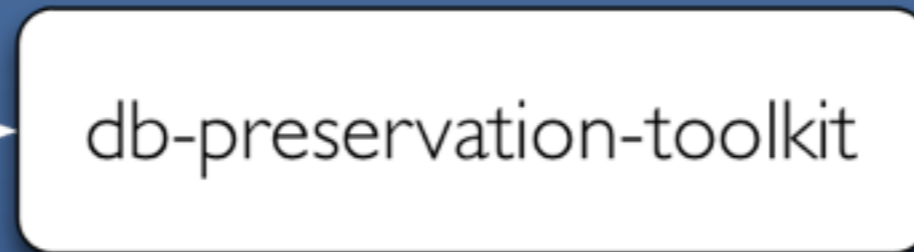


Preservation formats



Pre-ingest / Production / SIP creation

Origin formats



Preservation formats

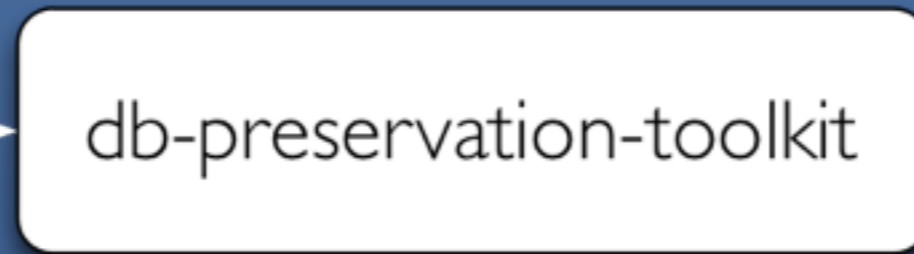


Pre-ingest / Production / SIP creation

Origin formats



new



Preservation formats



new

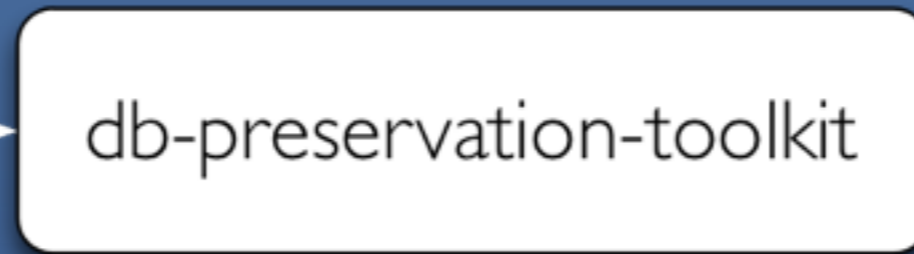
ADDML

Pre-ingest / Production / SIP creation

Origin formats



new



Preservation formats



new

ADDML to be implemented

Access

Preservation formats



Access

Preservation formats



Working on RDF viewer, we don't need to **stay stuck** in the relational model

New prototype

Msc thesis: *“Database Convert Tool: exploring SIARD”*

OWL model development

Triple Store Implementation

RDF endpoint implementation

SPARQL based access

Preservation formats

DBML

SIARD

db-preservation-toolkit

Fast viewer resources

Lucene / Solar index

Fast viewer application

Web interface & REST API

Future plan:
Fast database viewer



Recent work

Phd thesis: Graph view dissemination, relational model reverse engineering, algorithm development (from relational to graph model) [<http://repositorium.sdum.uminho.pt/handle/1822/25655>]

Msc thesis: SIARD support (DBML replacement)

Msc thesis: RDF representation for databases

Questions?



José Carlos Ramalho
Engineer / Researcher / Teacher
jcr@keep.pt / jcr@di.uminho.pt

KEEPSOLUTIONS
University of Minho SPIN-OFF

ARQUIVOS | BIBLIOTECAS | MUSEUS

www.keep.pt