

# RODA: repositório de objectos digitais autênticos

---

José Carlos Ramalho  
Dep. Informática  
Universidade do Minho  
[jcr@di.uminho.pt](mailto:jcr@di.uminho.pt)

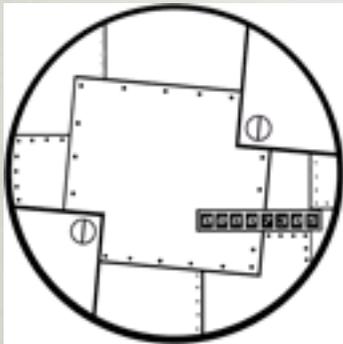
# CONTEXTO

---



## Digitarq (2003-2008)

- gestão de metainformação (EAD)
- gestão de ODs (NISO MIX)



## RODA (2006-2009)

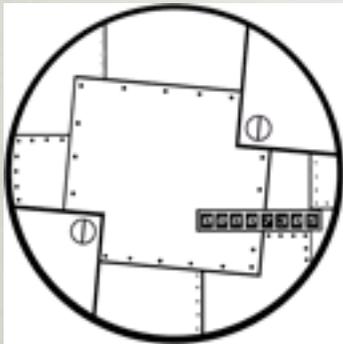
- gestão de metainformação (EAD)
- gestão de ODs (NISO MIX, METS, DBML, etc)
- políticas e protocolos de preservação digital

# CONTEXTO



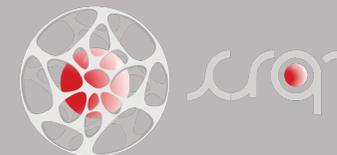
DigitArq (2003-2008)

- gestão de metainformação (EAD)
- gestão de ODs (NISO MIX)



RODA (2006-2009)

- gestão de metainformação (EAD)
- gestão de ODs (NISO MIX, METS, DBML, etc)
- políticas e protocolos de preservação digital



# RODA: MOTIVAÇÃO

---

- Hoje a história é digital;
- A produção de objectos digitais cresce a cada dia: desmaterialização na administração pública;
- Não há estruturas capazes de suportar a incorporação, a gestão e a preservação a longo prazo dos objectos digitais;
- É necessário preservar a memória digital, a herança e o testemunho das instituições.
  - Exemplo: SGU, processos judiciais, registos paroquiais, etc.

# REQUISITOS

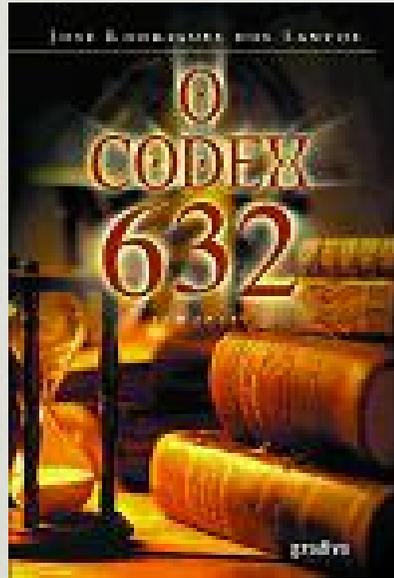
---

- Como garantir a autenticidade?
- Como descrever e classificar ODs?
- Como implementar a preservação digital?
- ...
- Como manter a solução depois de desenvolvida?



# AUTENTICIDADE

---



“O Codex 632” de José Rodrigues dos Santos

Assunto: Quem foi Cristovão Colombo?

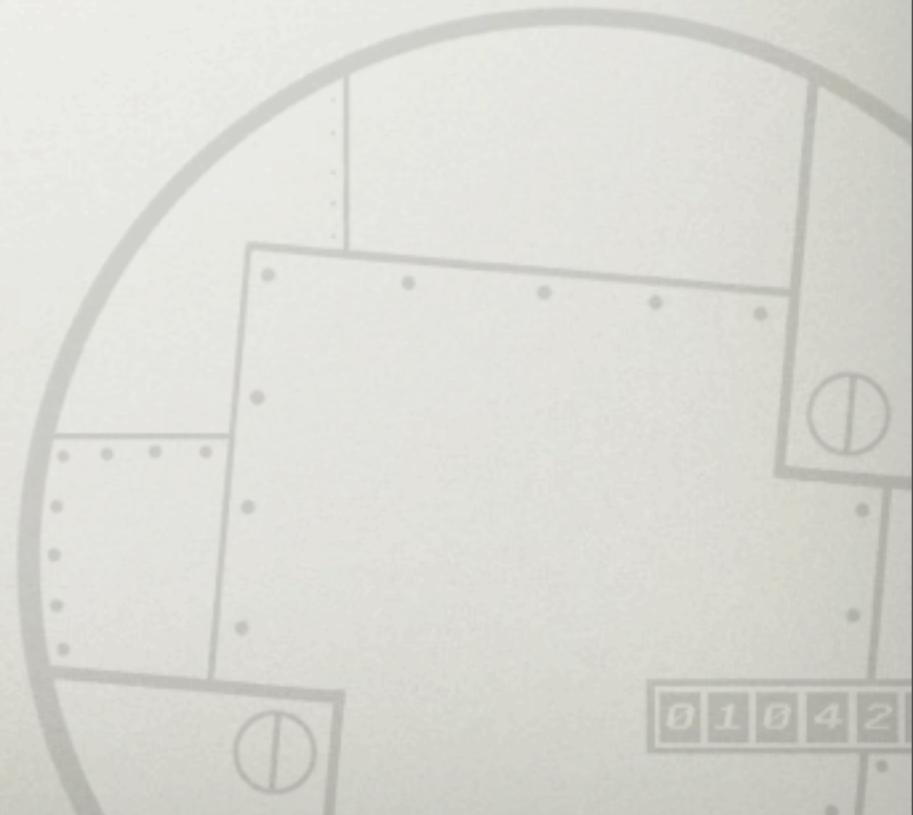
Um italiano? Espanhol? Ou um português de uma família judaica?

# AUTENTICIDADE

---

Temos de confiar nas nossas fontes: na História Antiga não há discurso directo ou evidências.

EX: a bíblia



# AUTENTICIDADE

---

Temos de confiar nas nossas fontes: na História Antiga não há discurso directo ou evidências.

EX: a bíblia

Como nos tornamos dignos de confiança?

# AUTENTICIDADE

---

Temos de confiar nas nossas fontes: na História Antiga não há discurso directo ou evidências.

EX: a bíblia

Como nos tornamos dignos de confiança?

- Reputação
- Documentando todo o ciclo de vida de um OD
- Seguindo o **TRAC** (Trustworthy Repositories Audit & Certification)

01042

# FORMATOS SUPORTADOS

---



# FORMATOS SUPORTADOS

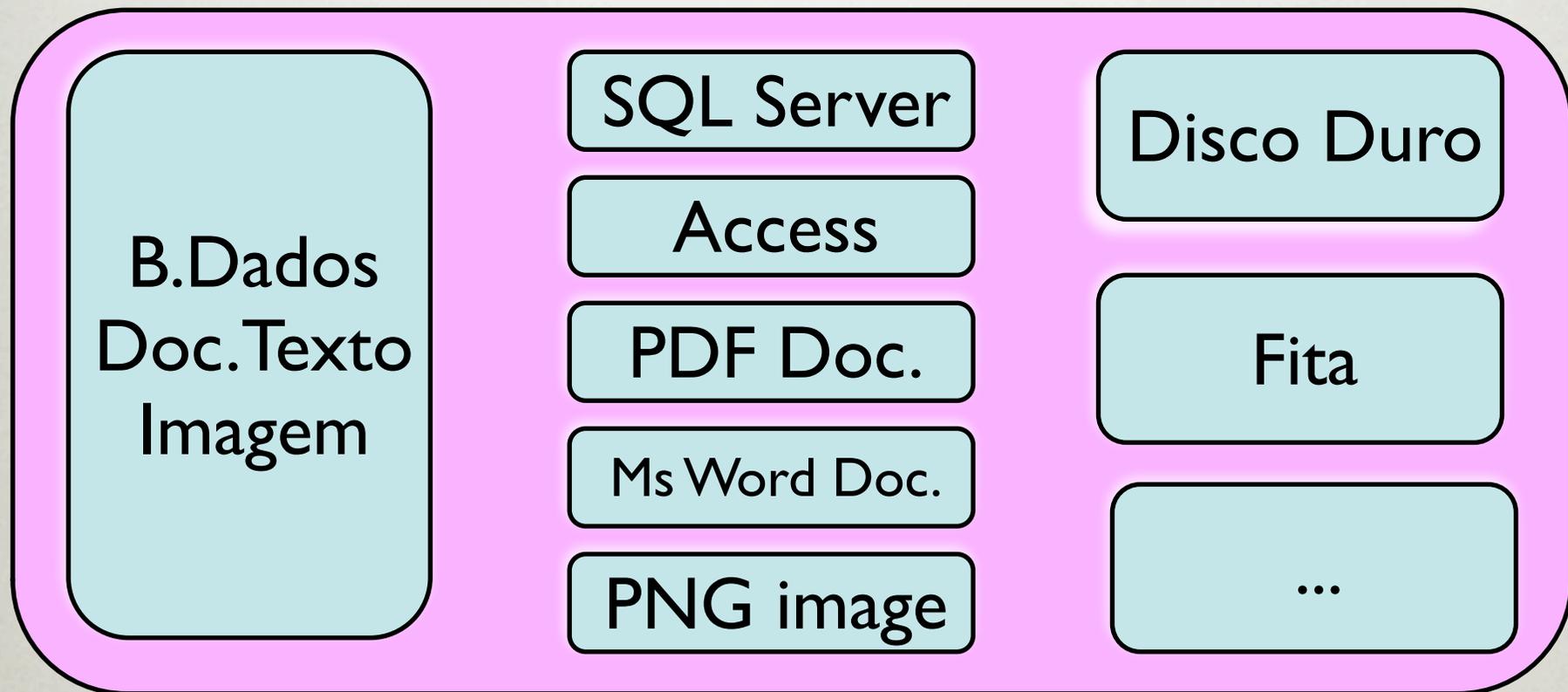


# Anatomia de um OD

Nível  
Conceptual

Nível  
Lógico

Nível  
Físico



# Anatomia de um OD

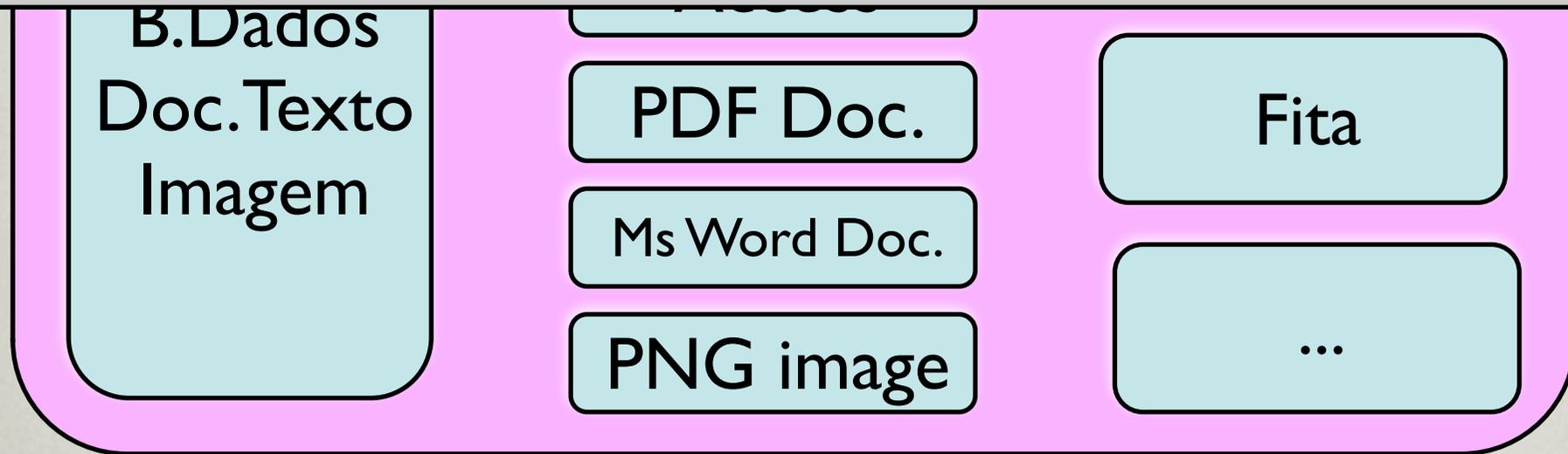
---

Nível  
Conceptual

Nível  
Lógico

Nível  
Físico

**Se um destes níveis se tornar obsoleto deixamos de ter acesso ao OD**



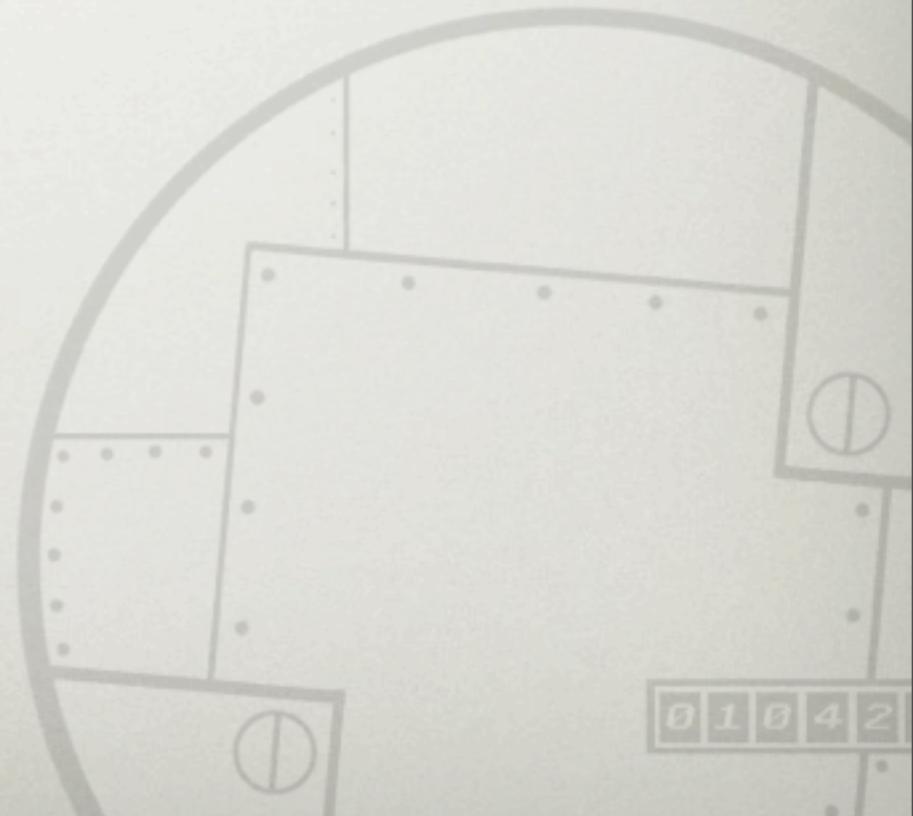
# Estratégias de Preservação

---

- Focada no **objecto físico**
  - o Centrada na preservação da informação no seu **formato lógico e/ou no suporte físico**;
  - o Usa a tecnologia original associada ao objecto para assegurar o acesso ao mesmo;
  - o **Preservação tecnológica.**
- Focada no **objecto conceptual**
  - o Centrada na preservação das propriedades do objectos de **forma independente** do hardware e software;
  - o **Preservação do objecto conceptual.**

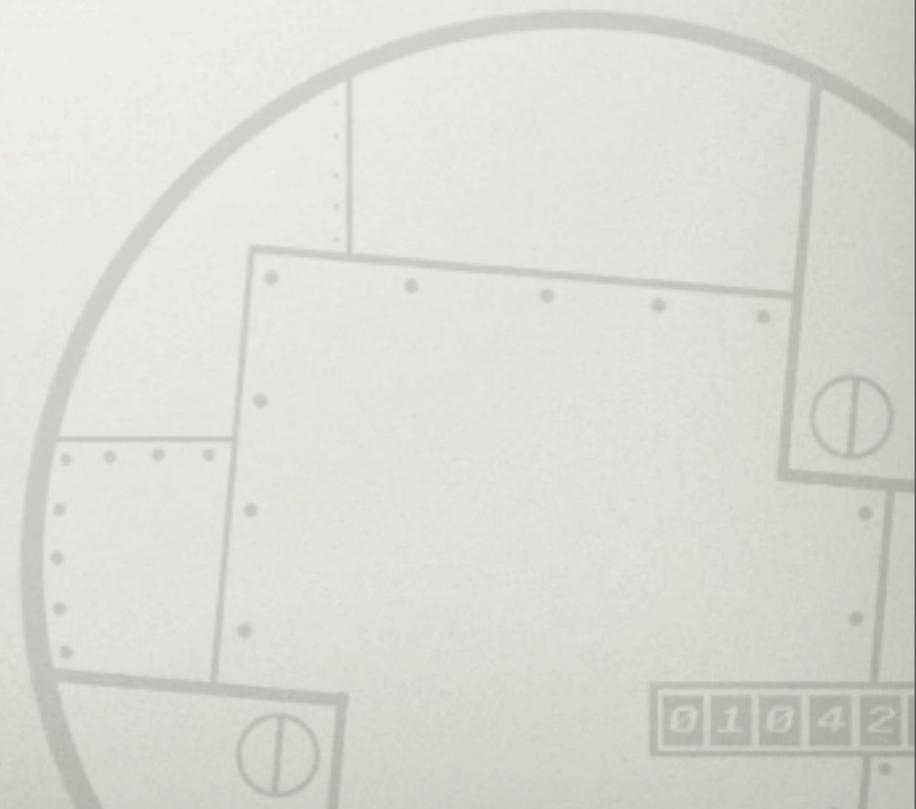
# Emulação

---



# Emulação

**Emulador: aplicação** capaz de reproduzir o comportamento de uma plataforma de hardware/software: ZX Spectrum, GBA, ...



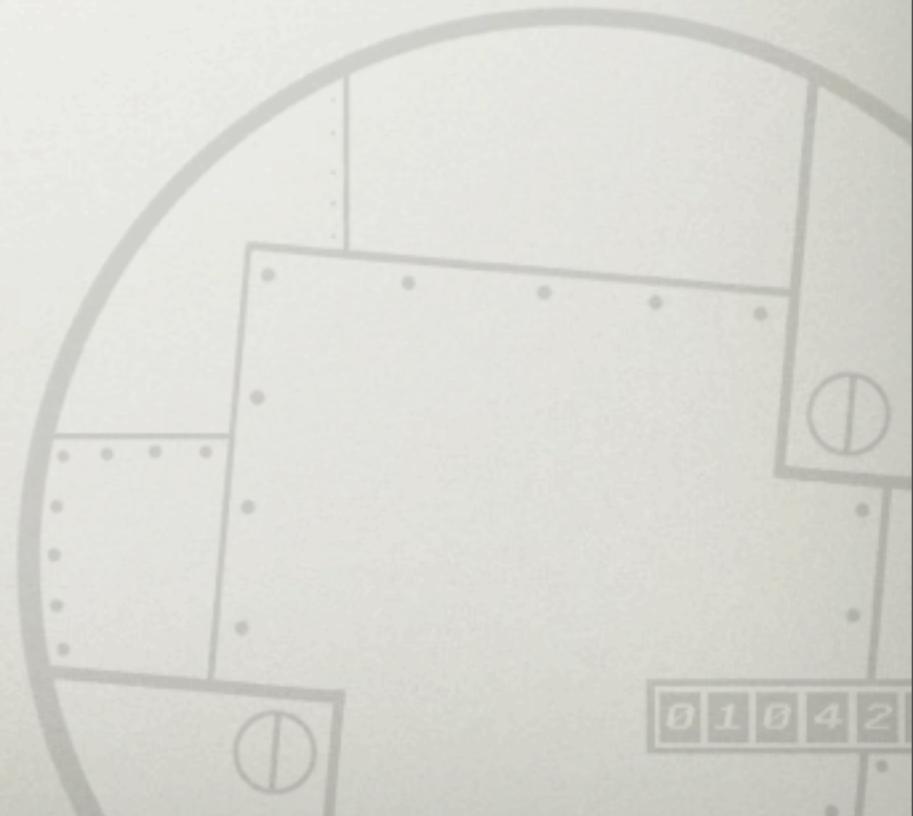
# Emulação

**Emulador: aplicação** capaz de reproduzir o comportamento de uma plataforma de hardware/software: ZX Spectrum, GBA, ...

- Vantagens
  - o Reprodução do contexto tecnológico original;
  - o Preservação do *look & feel* do objecto.
- Desvantagens
  - o Os emuladores também se tornam obsoletos;
  - o Os utilizadores têm de trabalhar com sistemas obsoletos;
  - o Criar um emulador é uma tarefa complexa;
  - o Problemas de Copyright;
  - o Preservar todo um sistema para poder visualizar um documento pode ser demasiado!
  - o A reutilização da informação não é garantida.

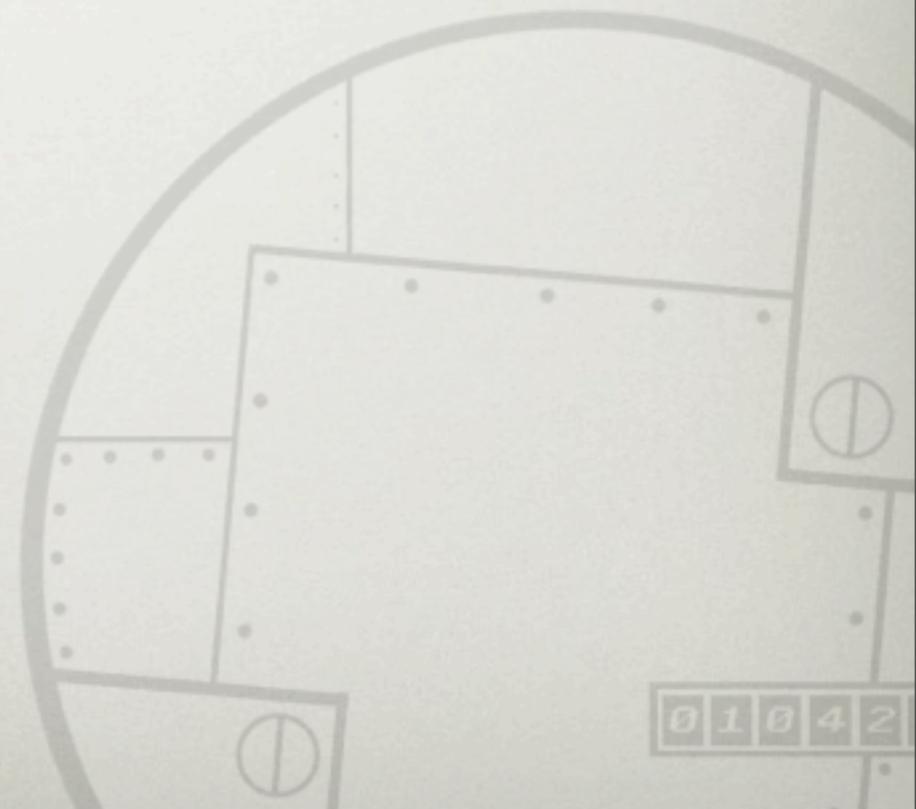
# Encapsulamento

---



# Encapsulamento

Preservação da **bit stream original** juntamente com metainformação suficiente para assegurar a sua interpretação e acesso no futuro



# Encapsulamento

Preservação da **bit stream original** juntamente com metainformação suficiente para assegurar a sua interpretação e acesso no futuro

- Vantagens
  - o Permite adiar as responsabilidades da preservação;
  - o Deve ser utilizada em objectos cujo acesso será desejado num futuro a longo prazo;
- Desvantagens
  - o **Objectos complexos têm especificações complexas;**
  - o Uma **especificação incompleta** pode ter efeitos

# Preservação do Objecto Conceptual

---

**Migração:** transferência periódica do OD de uma configuração de hw/sw para outra mais actualizada (preservando as propriedades significativas do objecto em lugar da sequência de bits original).

## Vantagens

- Os ODs são disseminados em formatos familiares aos utilizadores;
- Não há necessidade de manter a plataforma original de hw/sw;
- A estratégia mais usada e a única que funcionou até hoje.

## Desvantagens

- Possível perda de informação nas conversões;
- Manutenção constante e continuada;
- **A longo prazo os custos poderão ser altos.**

# Preservação do Objecto Conceptual

---

**Migração:** transferência periódica do OD de uma configuração de hw/sw para outra mais actualizada (preservando as propriedades significativas do objecto em lugar da sequência de bits original).

## Vantagens

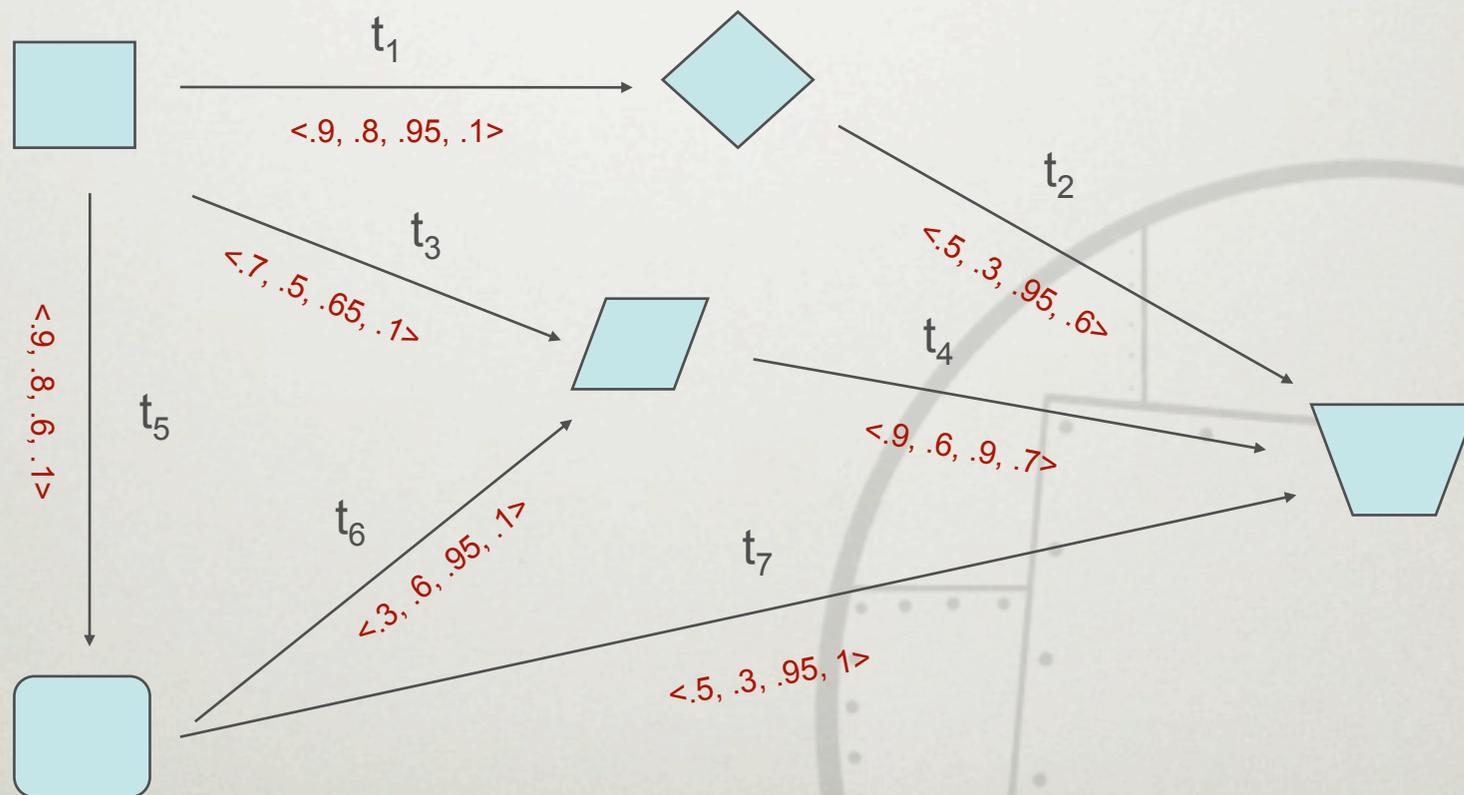
- Os ODs são disseminados em formatos familiares aos utilizadores;
- Não há necessidade de manter a plataforma original de hw/sw;
- A estratégia mais usada e a única que funcionou até hoje.

## Desvantagens

- Possível perda de informação nas conversões;
- Manutenção constante e continuada;
- **A longo prazo os custos poderão ser altos.**

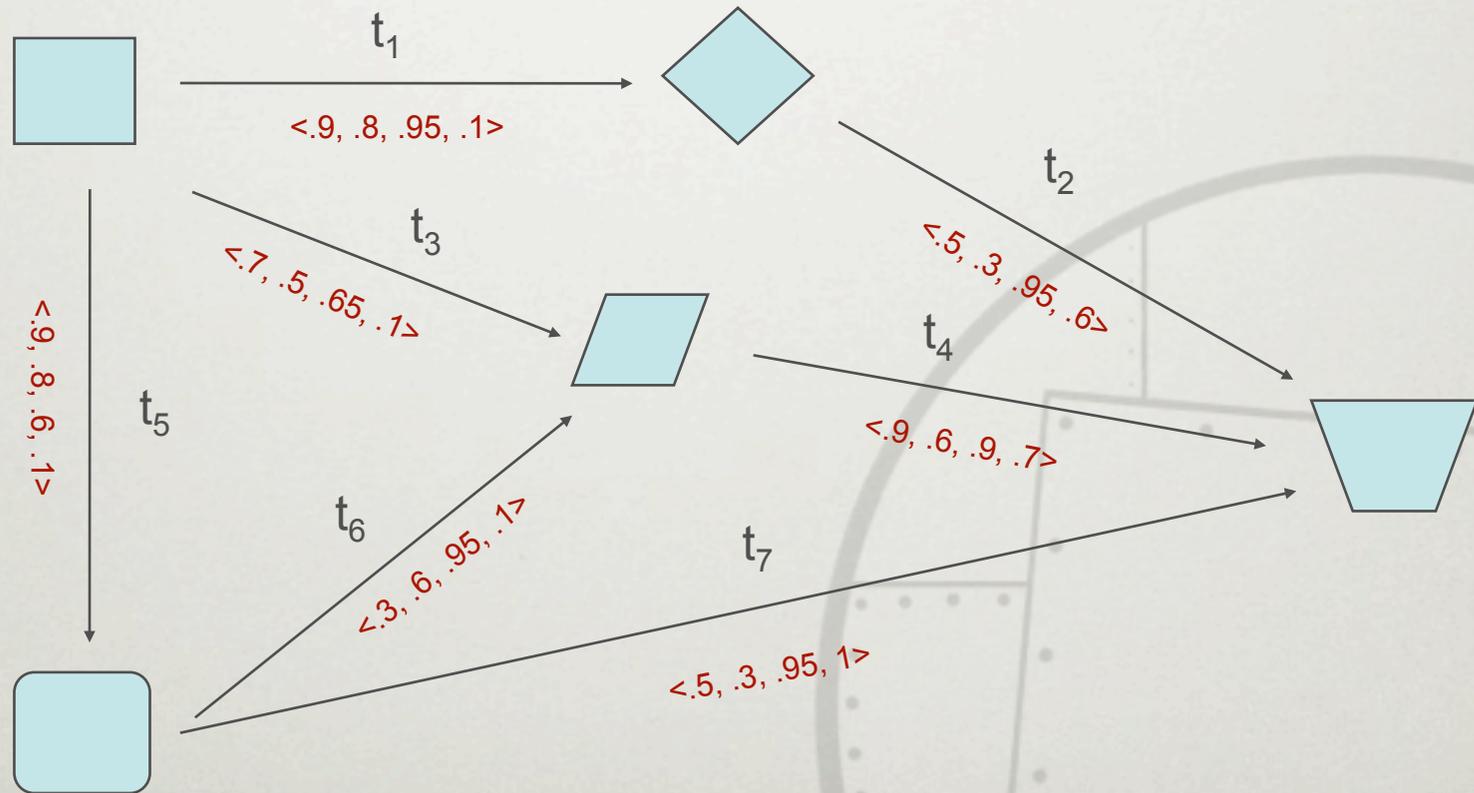
**O que são propriedades significativas?**

# Serviços de Preservação

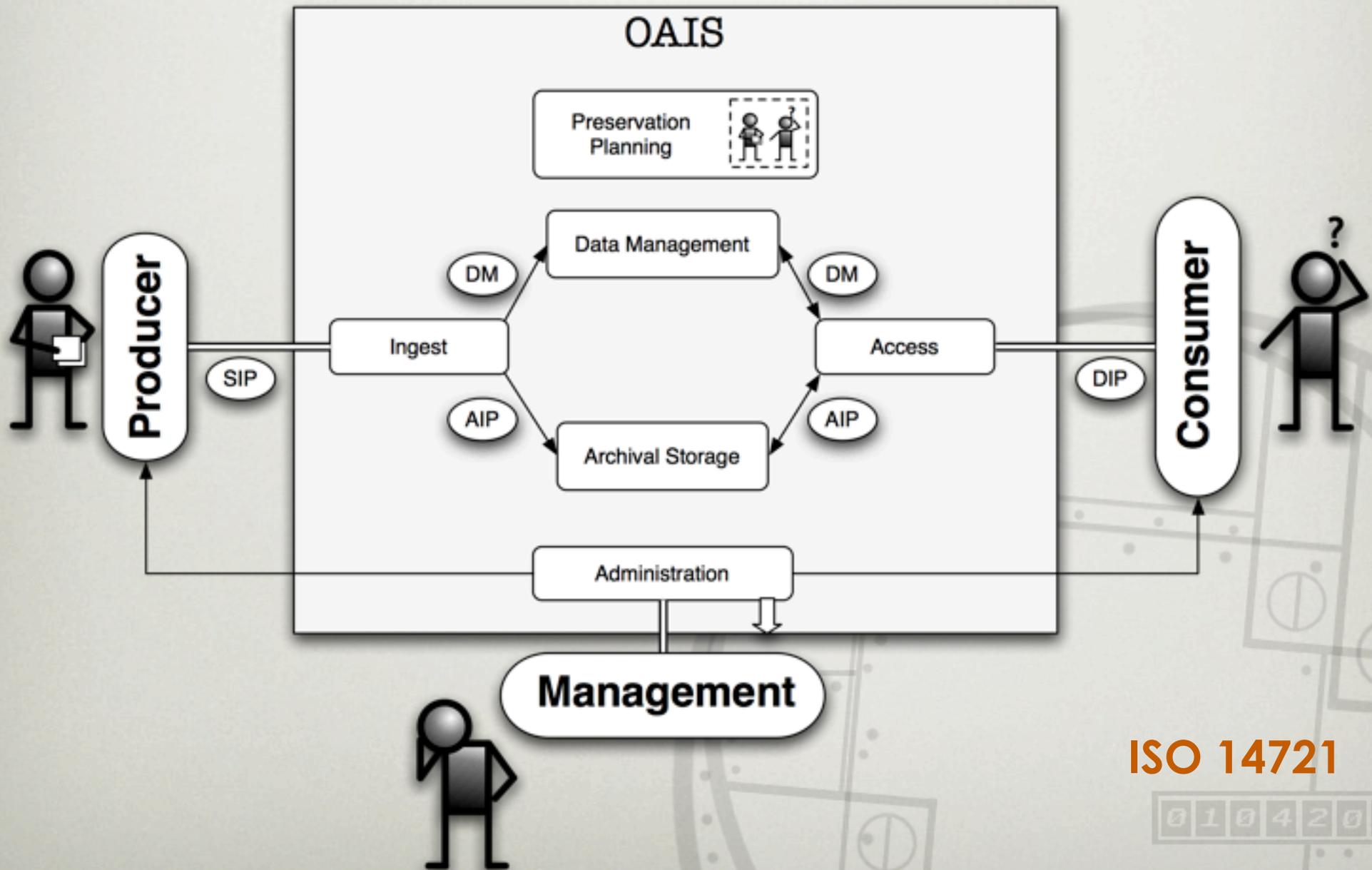


# Serviços de Preservação

**CRiB project:** <http://crib.dsi.uminho.pt>



# OPEN ARCHIVAL INFORMATION SYSTEM (OAIS)



ISO 14721

01042005

# OAIS (COMPONENTES FUNCIONAIS)

---

## ❖ Ingestão

- **Recepção, validação, transformação/normalização**, descrição do pacote enviado pelo produtor;

## ❖ Armazenamento

- Assegura a preservação da informação ao nível físico/lógico (e.g. refrescamento, migração, verificação de integridade, tolerância a falhas, etc.)

## ❖ Gestão de metainformação

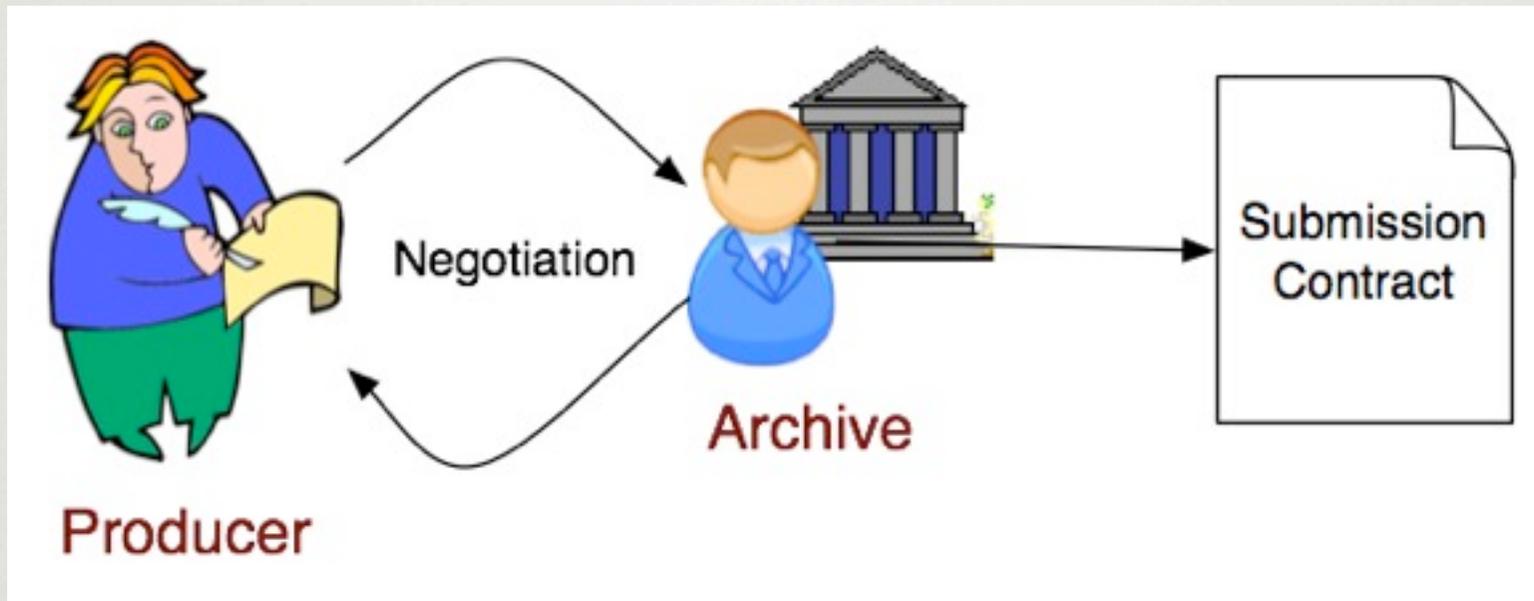
- Responsável pela gestão dos ODs armazenados

# OAIS (“PACOTES” DE INFORMAÇÃO)

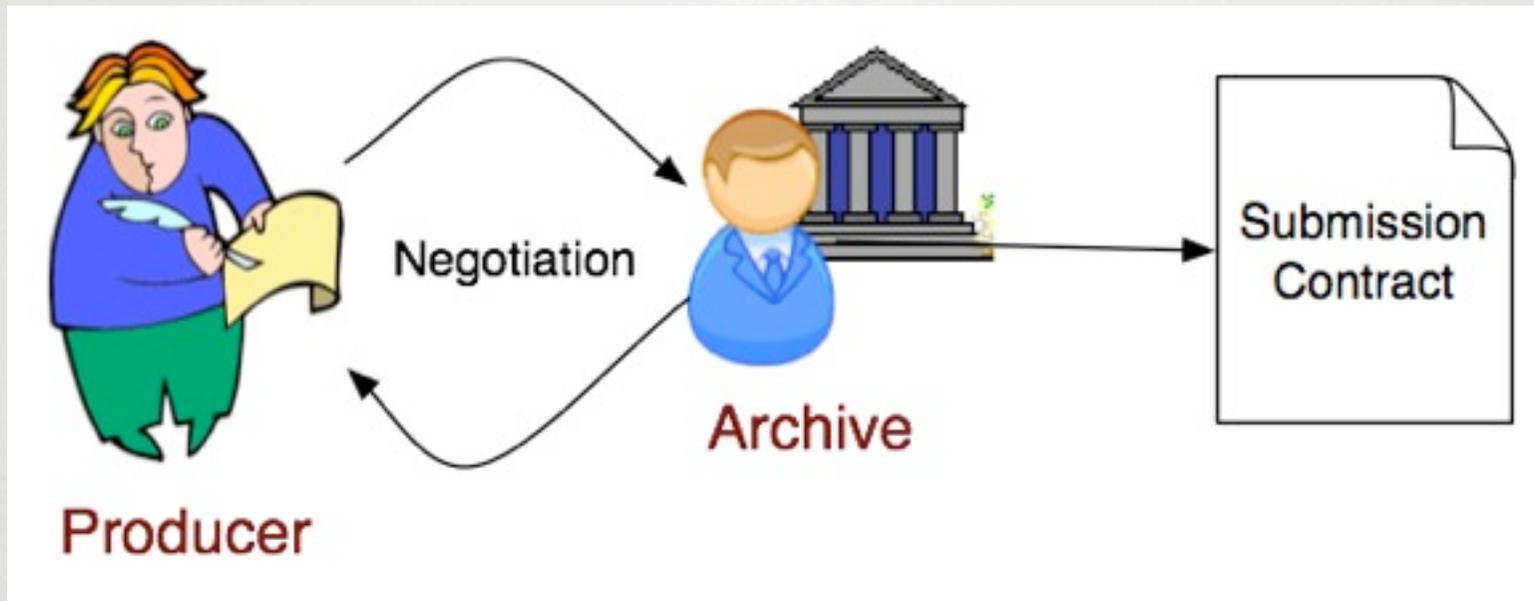
---

- Submission Information Package (SIP)
  - \* **Objecto Digital (representação ou representações);**
  - \* **Metainformação criada pelo produtor**
    - ▶ demasiado livre...
- Archival Information Package (AIP)
  - \* **Objecto digital arquivado (processado);**
  - \* **Metainformação:** a suficiente para assegurar a preservação e o acesso aos ODs;
    - ▶ modelo definido pelo PREMIS;
- Dissemination Information Package (DIP)
  - OD transformado no **formato** a ser **entregue ao consumidor;**
  - **Metainformação.**

# INGESTÃO



# INGESTÃO



## Contrato de Submissão/Envio

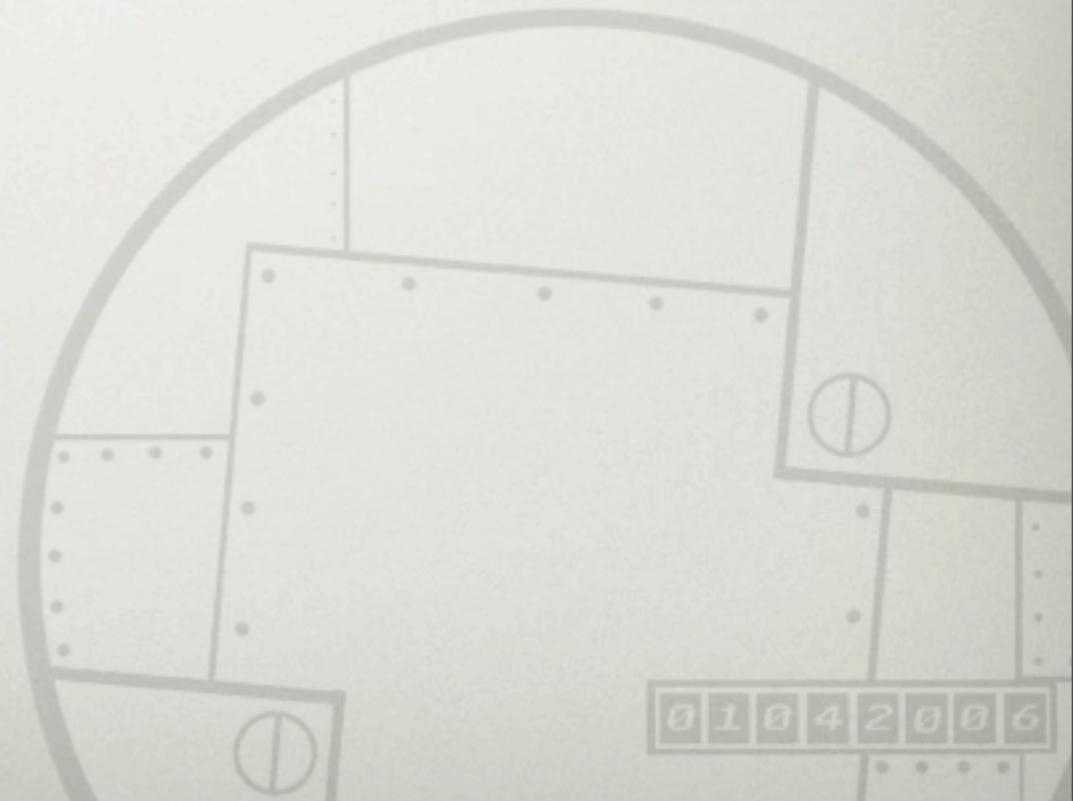
- Especificação do SIP;
- Especificação do workflow de ingestão.

# SIP STRUCTURE (EXAMPLE)

---



one still image

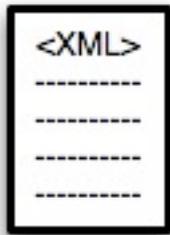


# SIP STRUCTURE (EXAMPLE)

---



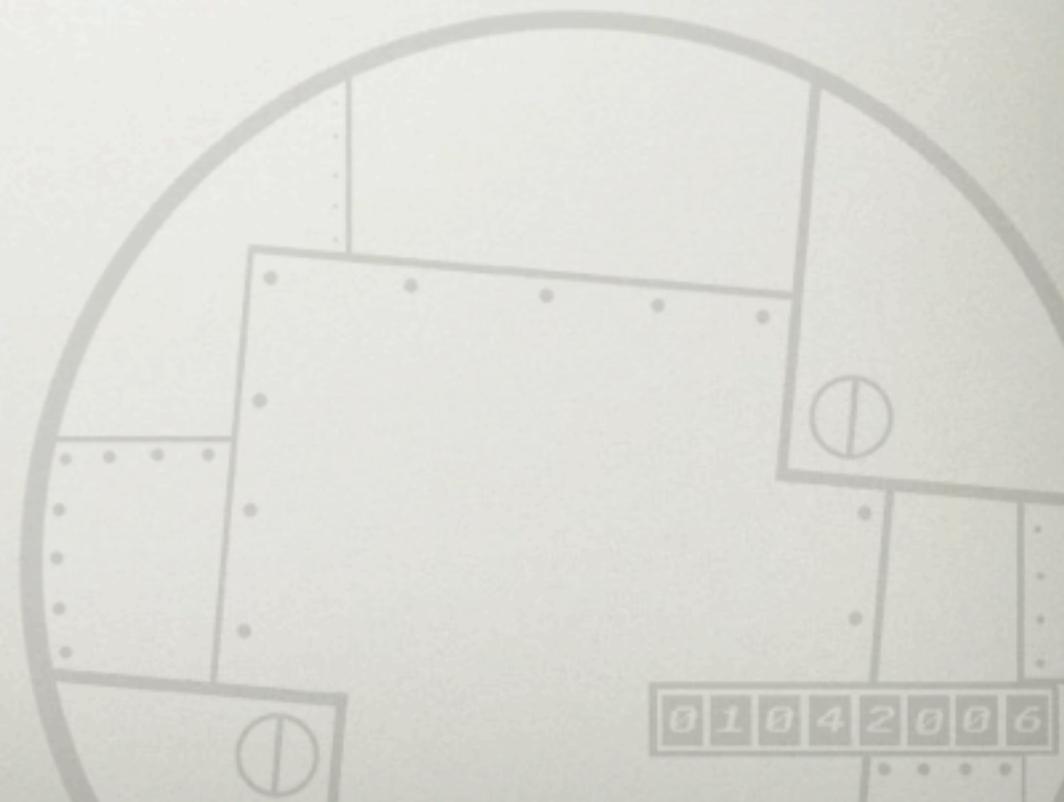
one still image



creation

properties:

- date
- hardware
- ...

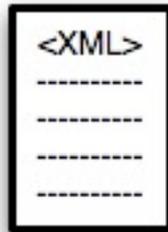


# SIP STRUCTURE (EXAMPLE)

---

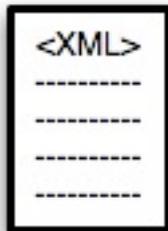


one still image



Technical Metadata:

- color
- dimensions
- ...



creation

properties:

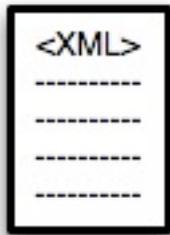
- date
- hardware
- ...



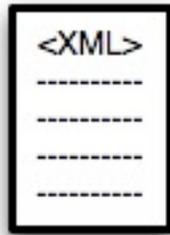
# SIP STRUCTURE (EXAMPLE)



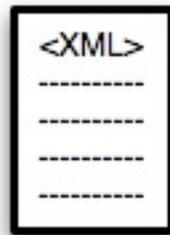
one still image



creation  
properties:  
- date  
- hardware  
- ...



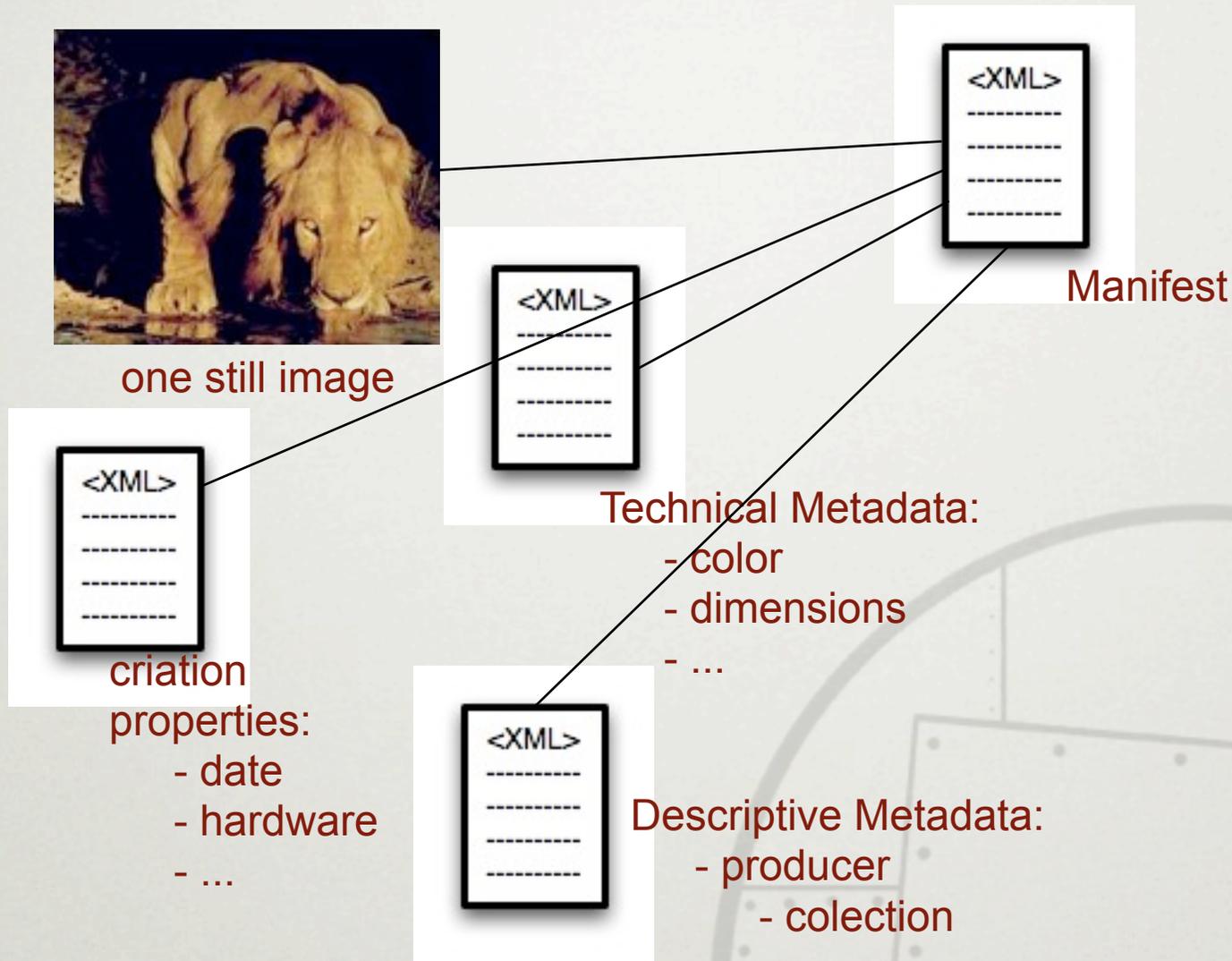
Technical Metadata:  
- color  
- dimensions  
- ...



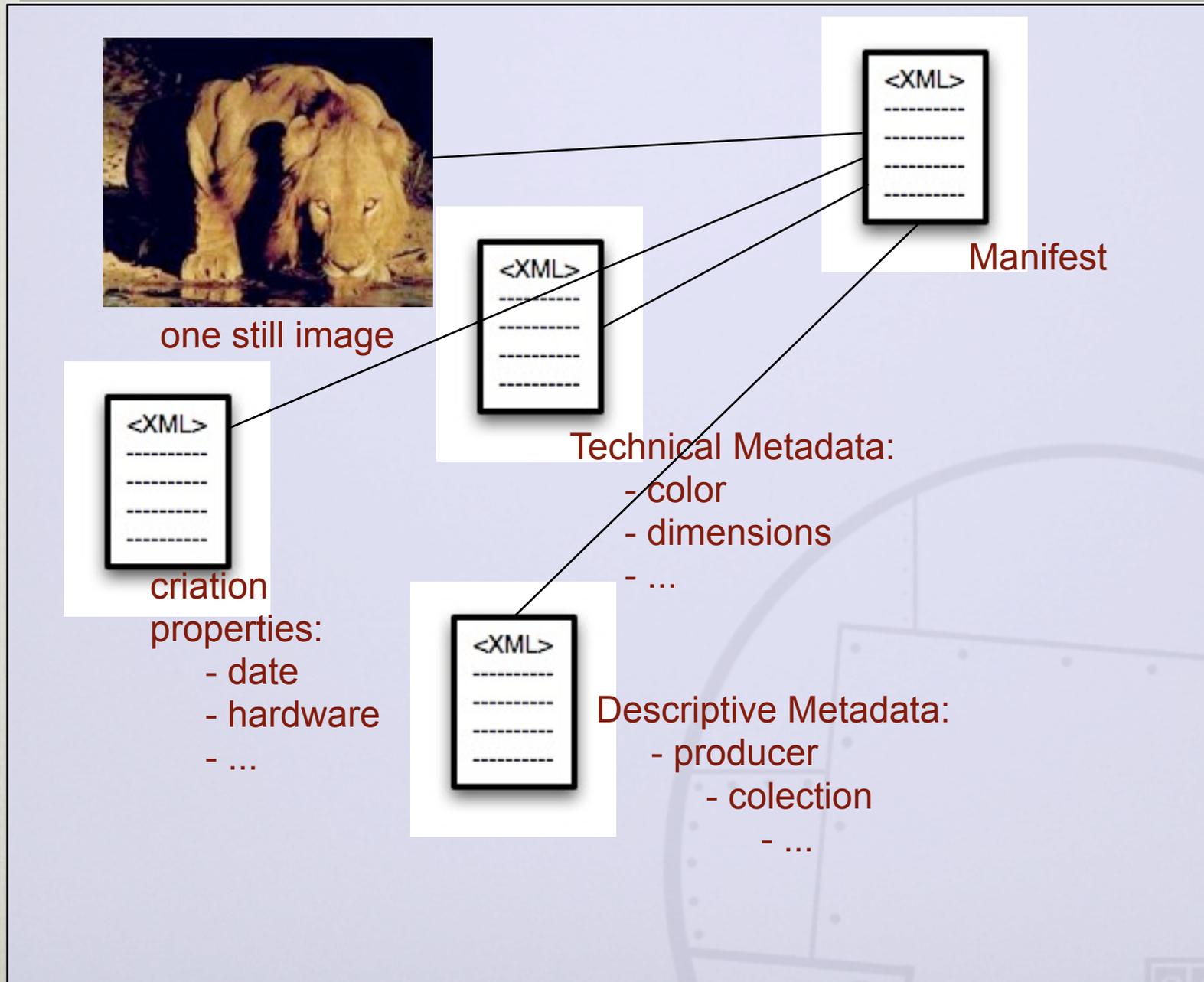
Descriptive Metadata:  
- producer  
- collection  
- ...

01042006

# SIP STRUCTURE (EXAMPLE)



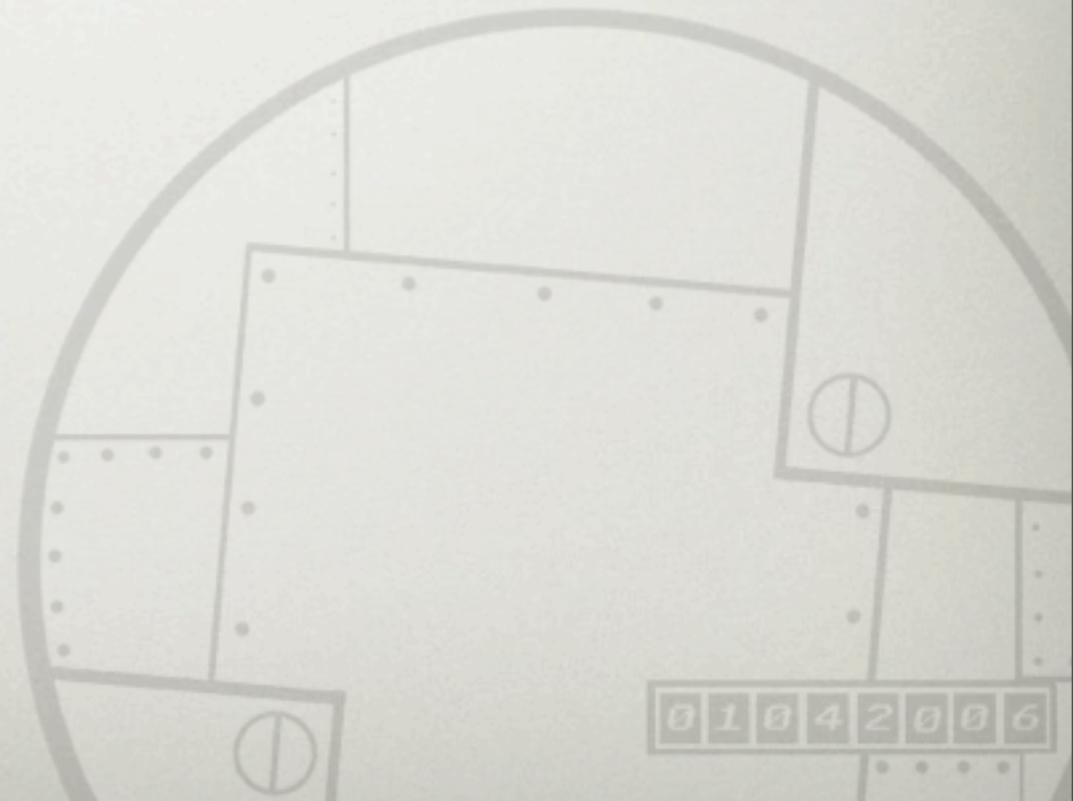
# SIP STRUCTURE (EXAMPLE)



**Compressed File**

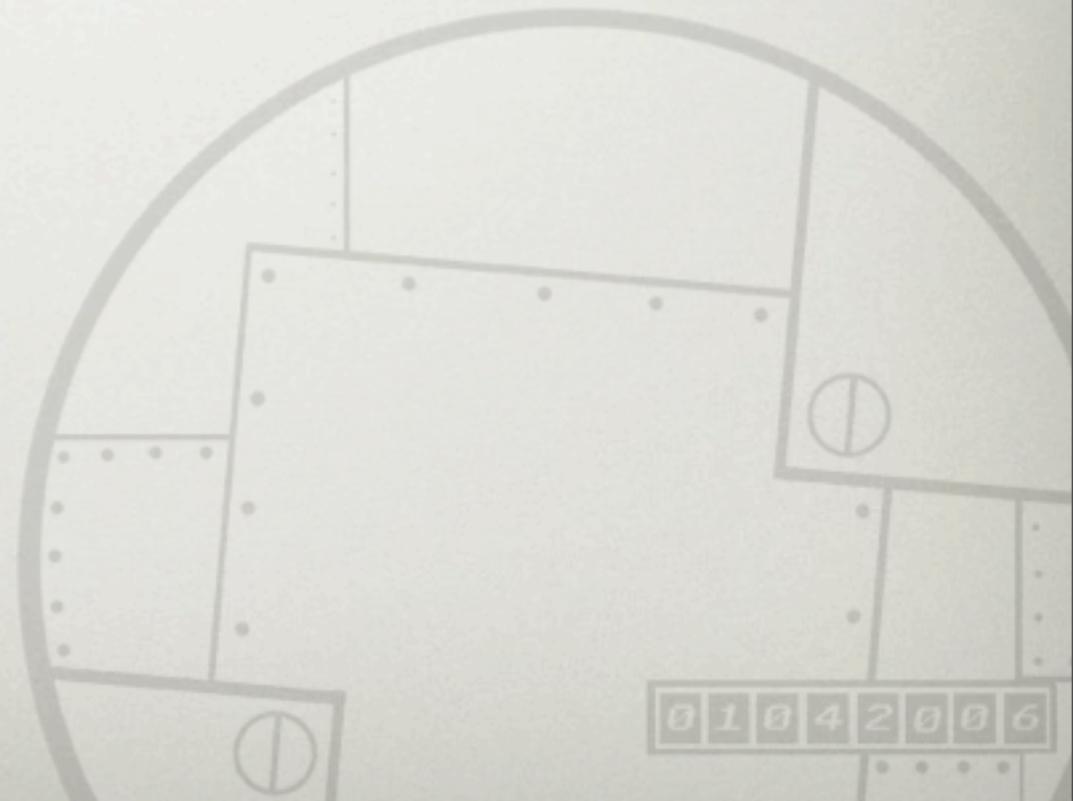
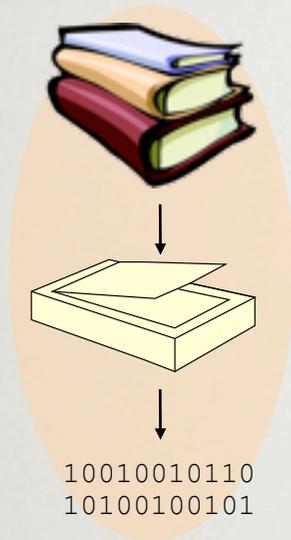
# SIP STRUCTURE (+COMPLEX)

---



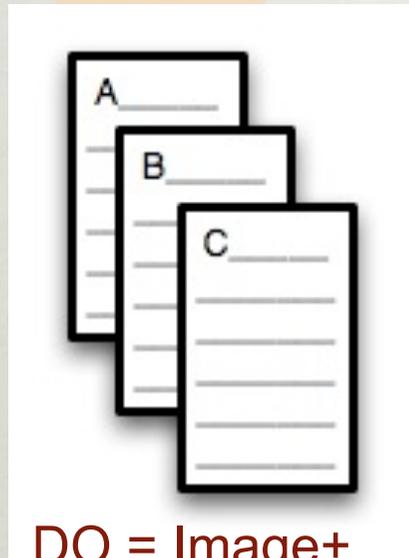
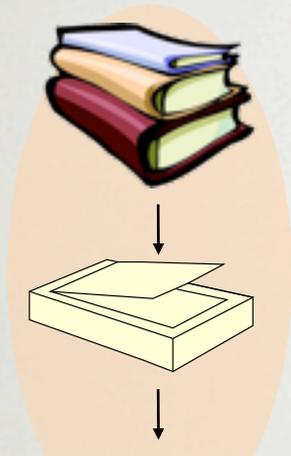
# SIP STRUCTURE (+COMPLEX)

---

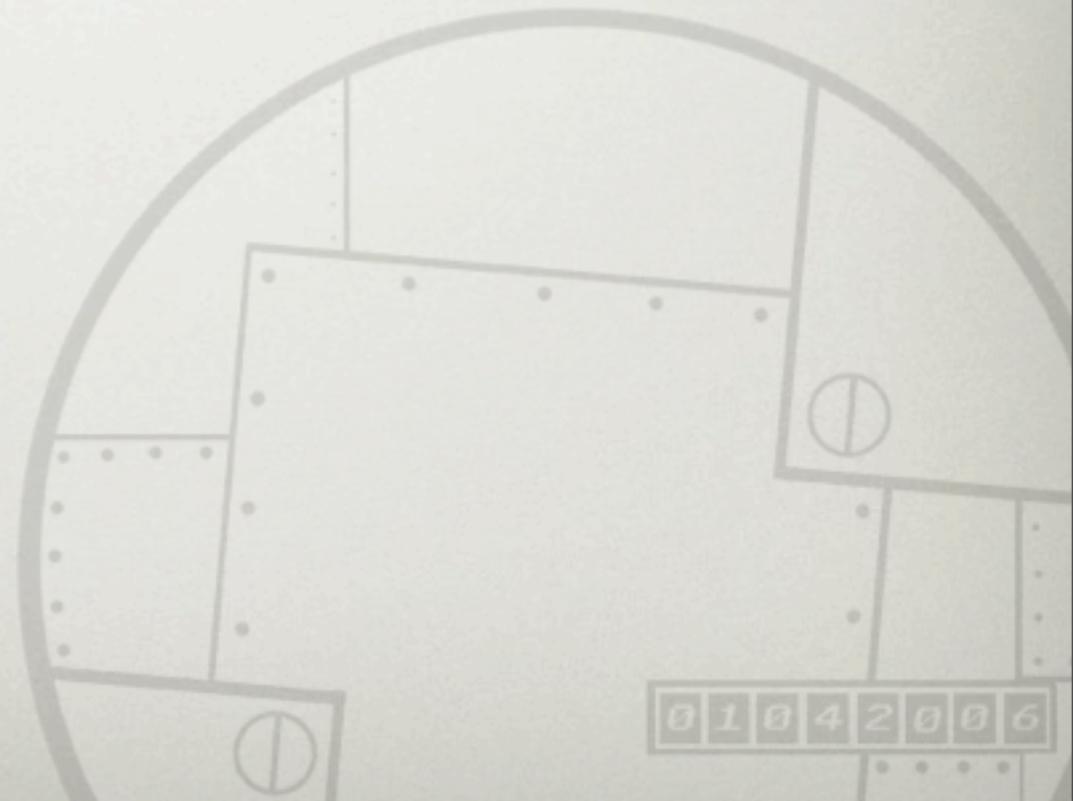


# SIP STRUCTURE (+COMPLEX)

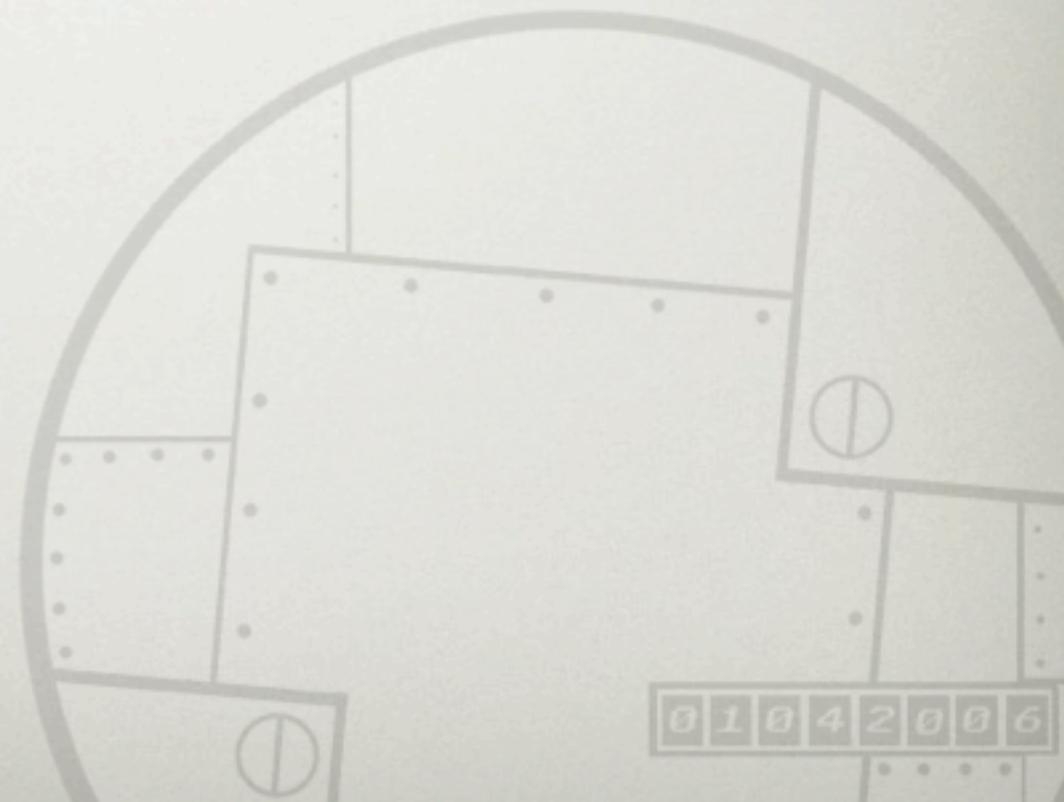
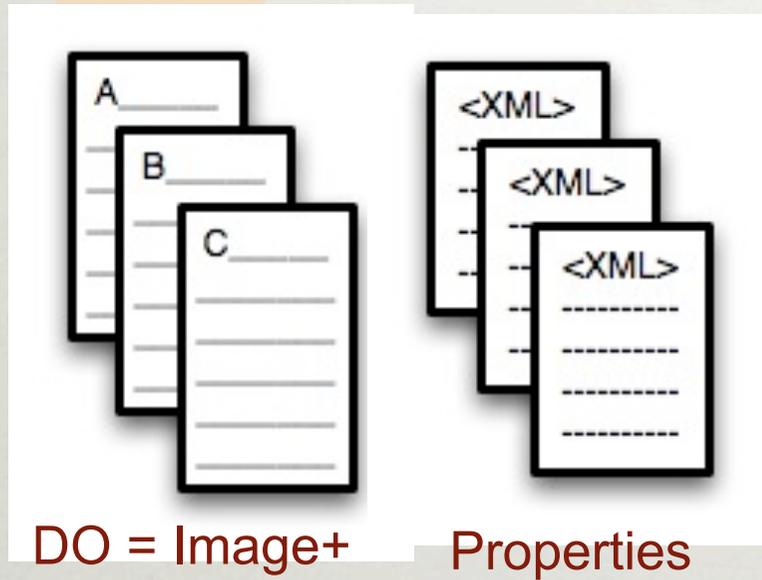
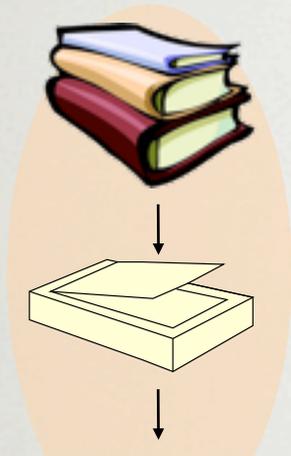
---



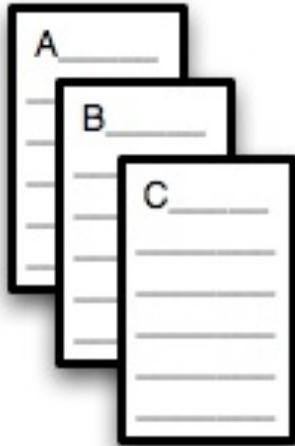
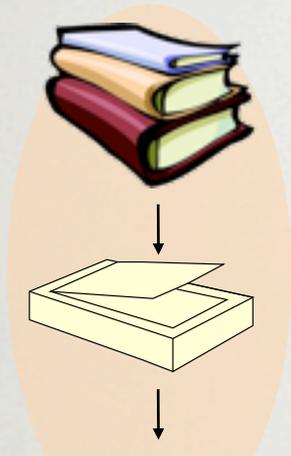
DO = Image+



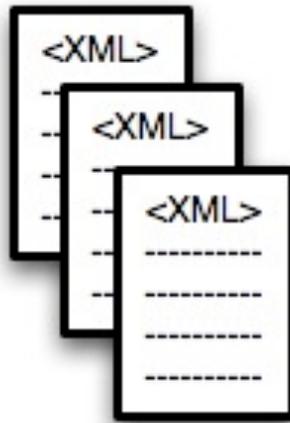
# SIP STRUCTURE (+COMPLEX)



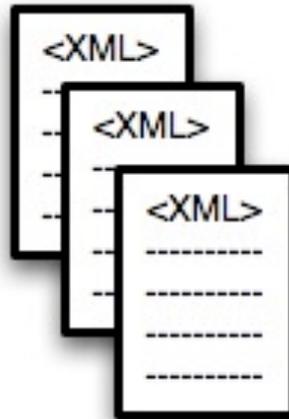
# SIP STRUCTURE (+COMPLEX)



DO = Image+

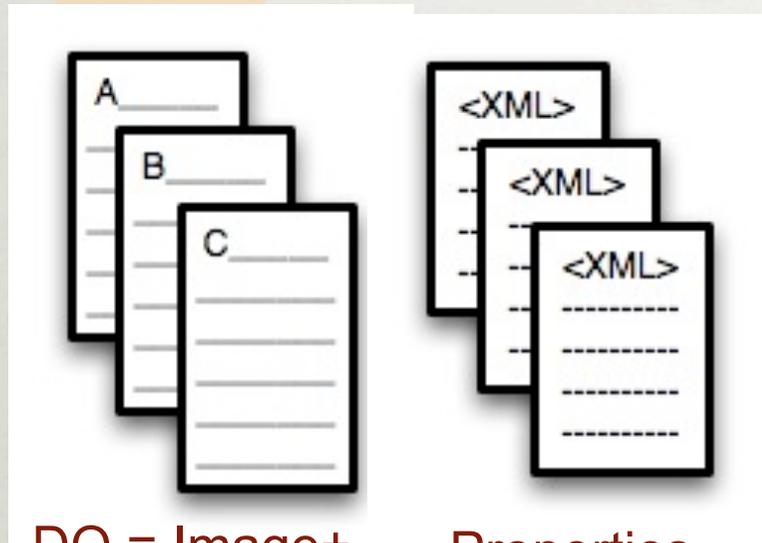
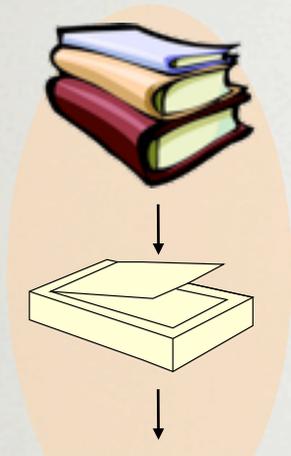


Properties



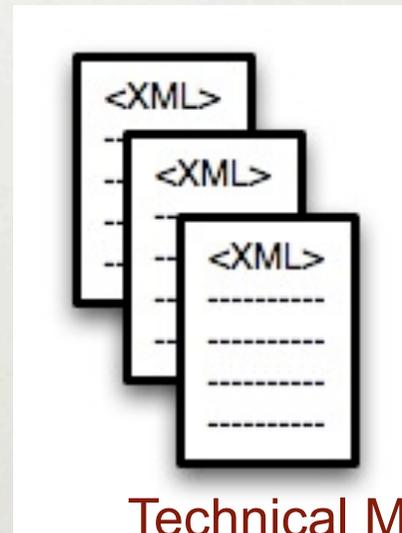
Technical Metadata

# SIP STRUCTURE (+COMPLEX)

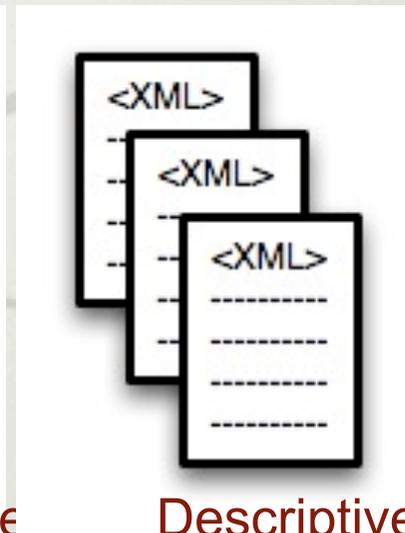


DO = Image+

Properties

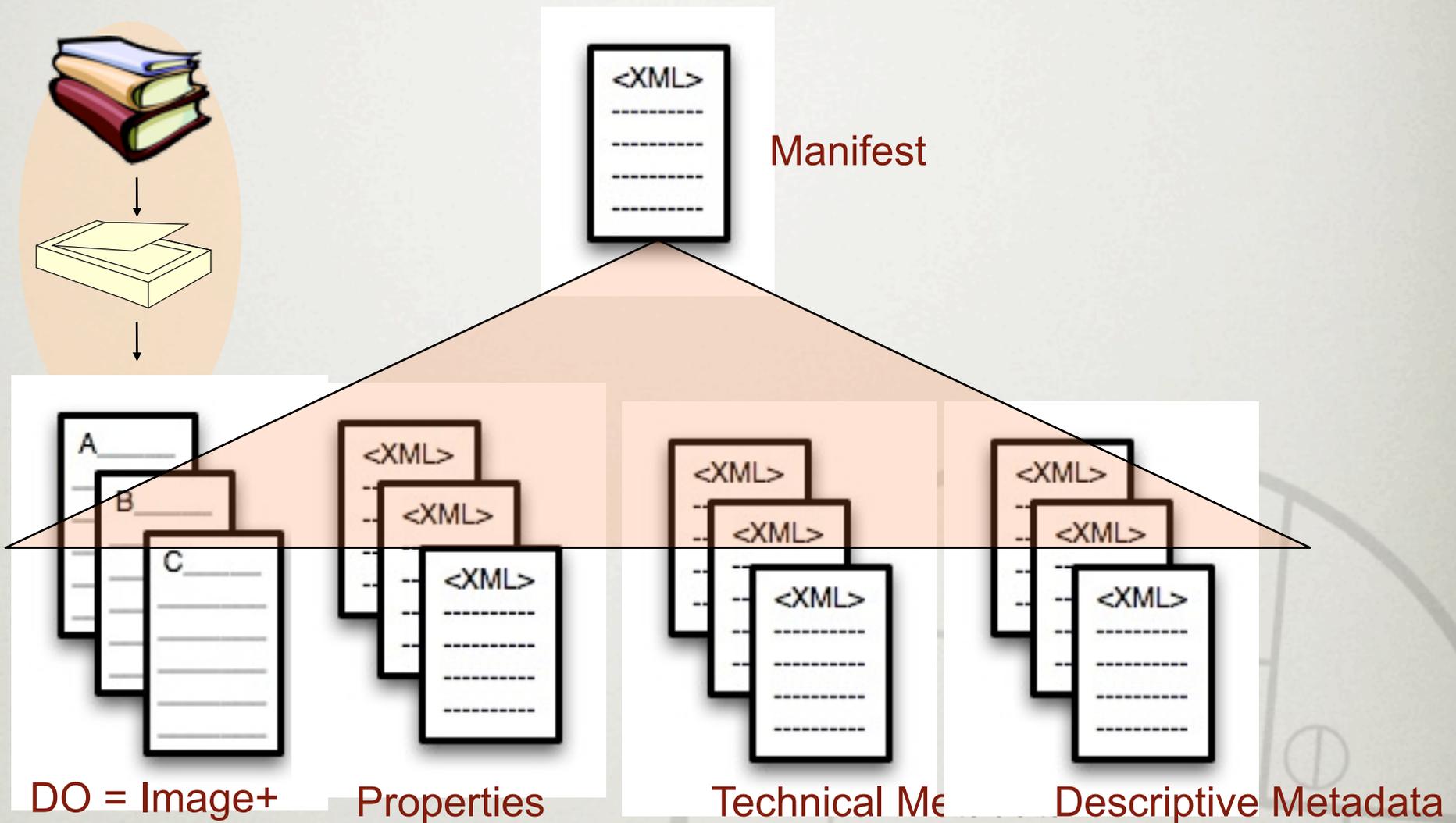


Technical Me

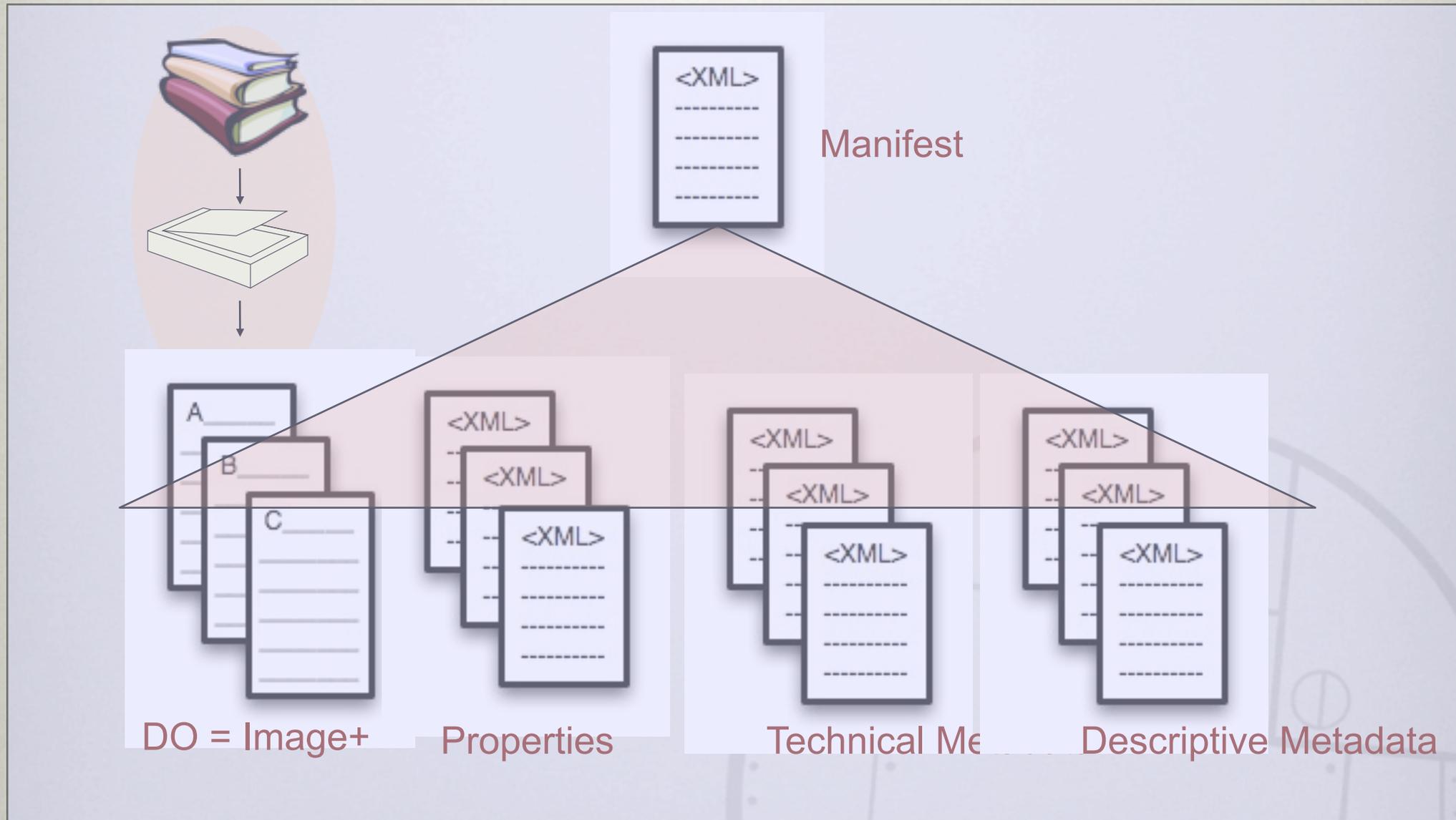


Descriptive Metadata

# SIP STRUCTURE (+COMPLEX)



# SIP STRUCTURE (+COMPLEX)



**Compressed File**

01042006

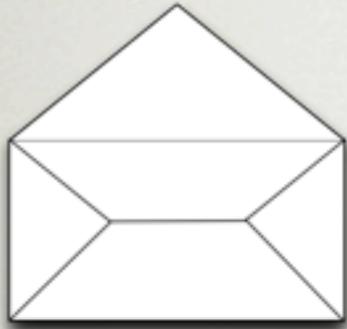
# WORKFLOW DE INGESTÃO

---

**SIP**

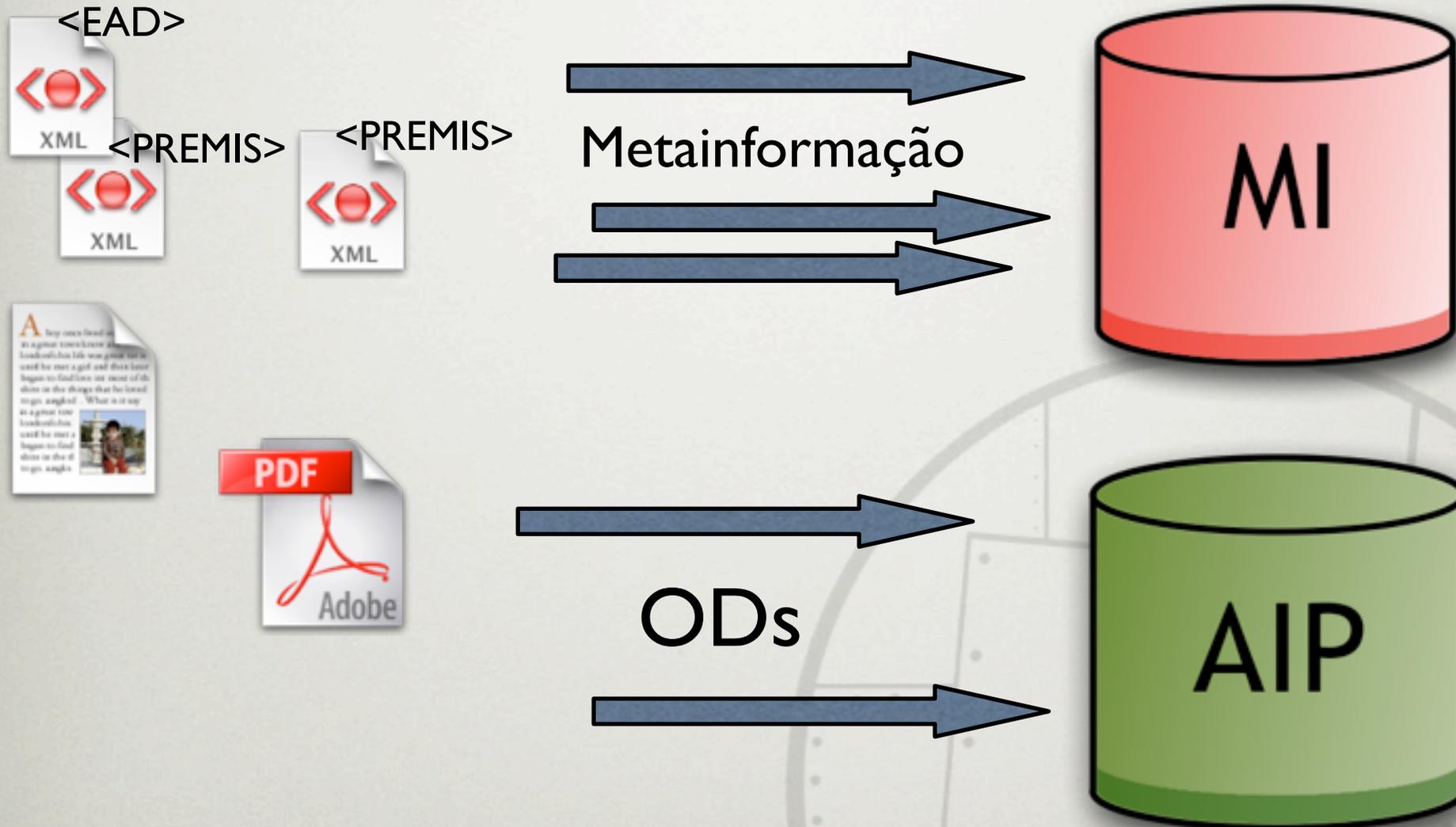


**AIP**

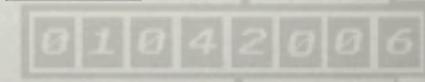
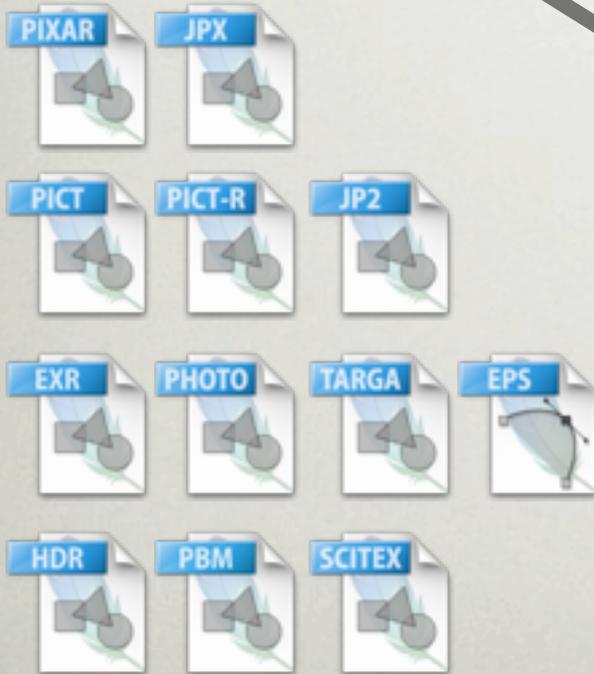


- verificação da integridade do pacote;
- verificação de vírus;
- geração da metainformação de preservação inicial;
- conversão num formato normalizado;
- geração de metainformação técnica;
- geração de metainformação de preservação (após normalização).

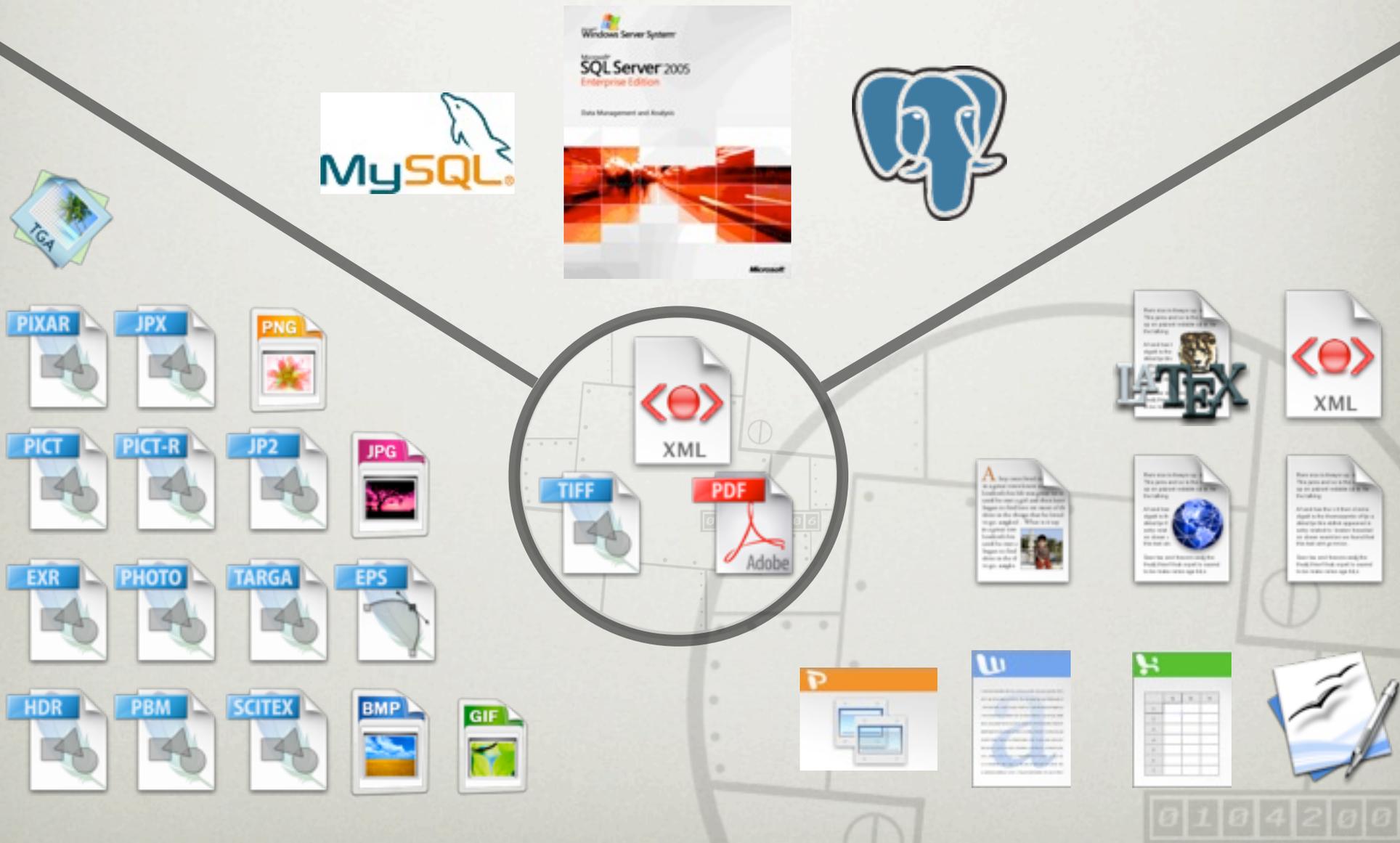
# AIP (ARMAZENAMENTO)



# NORMALIZAÇÃO



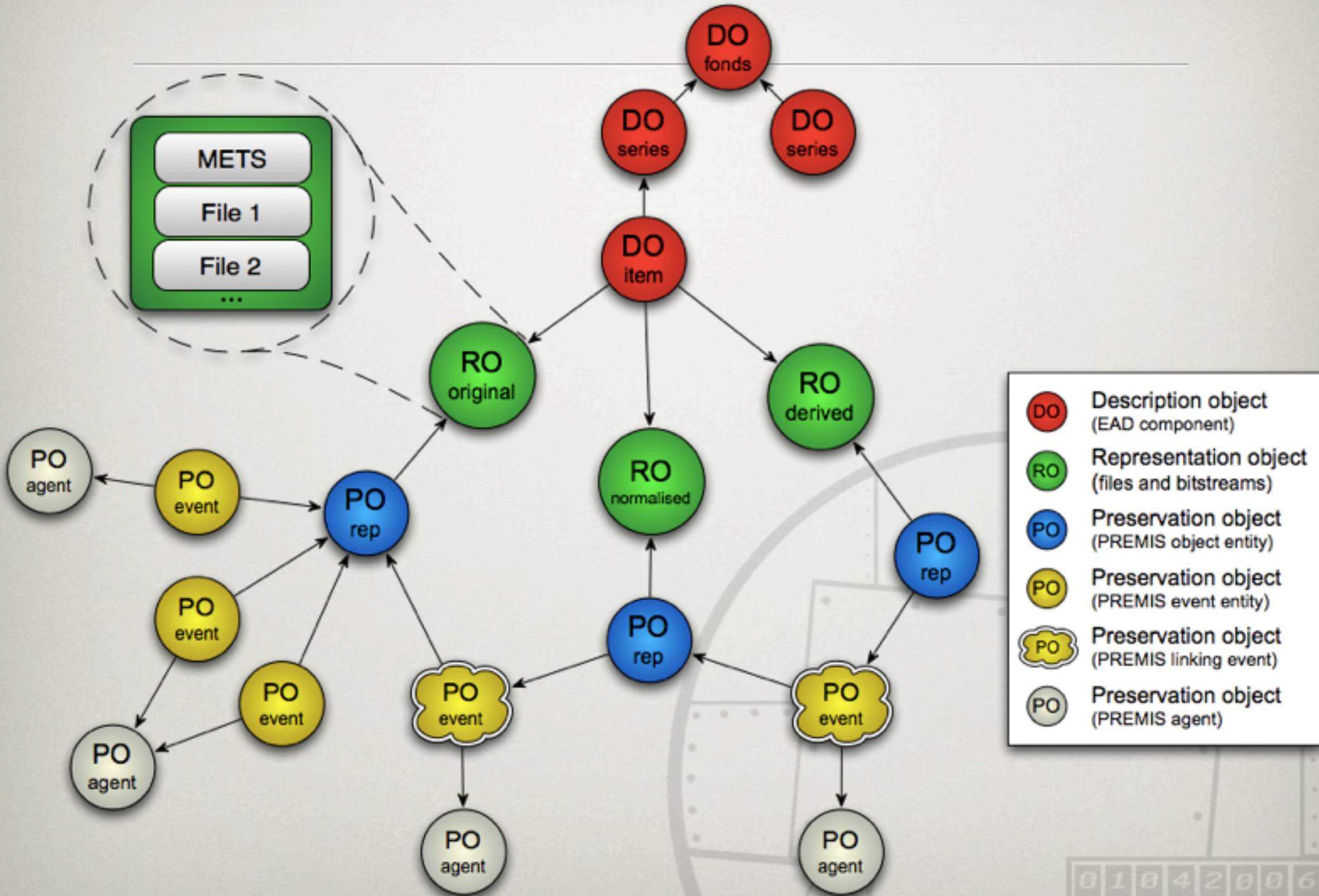
# NORMALIZAÇÃO





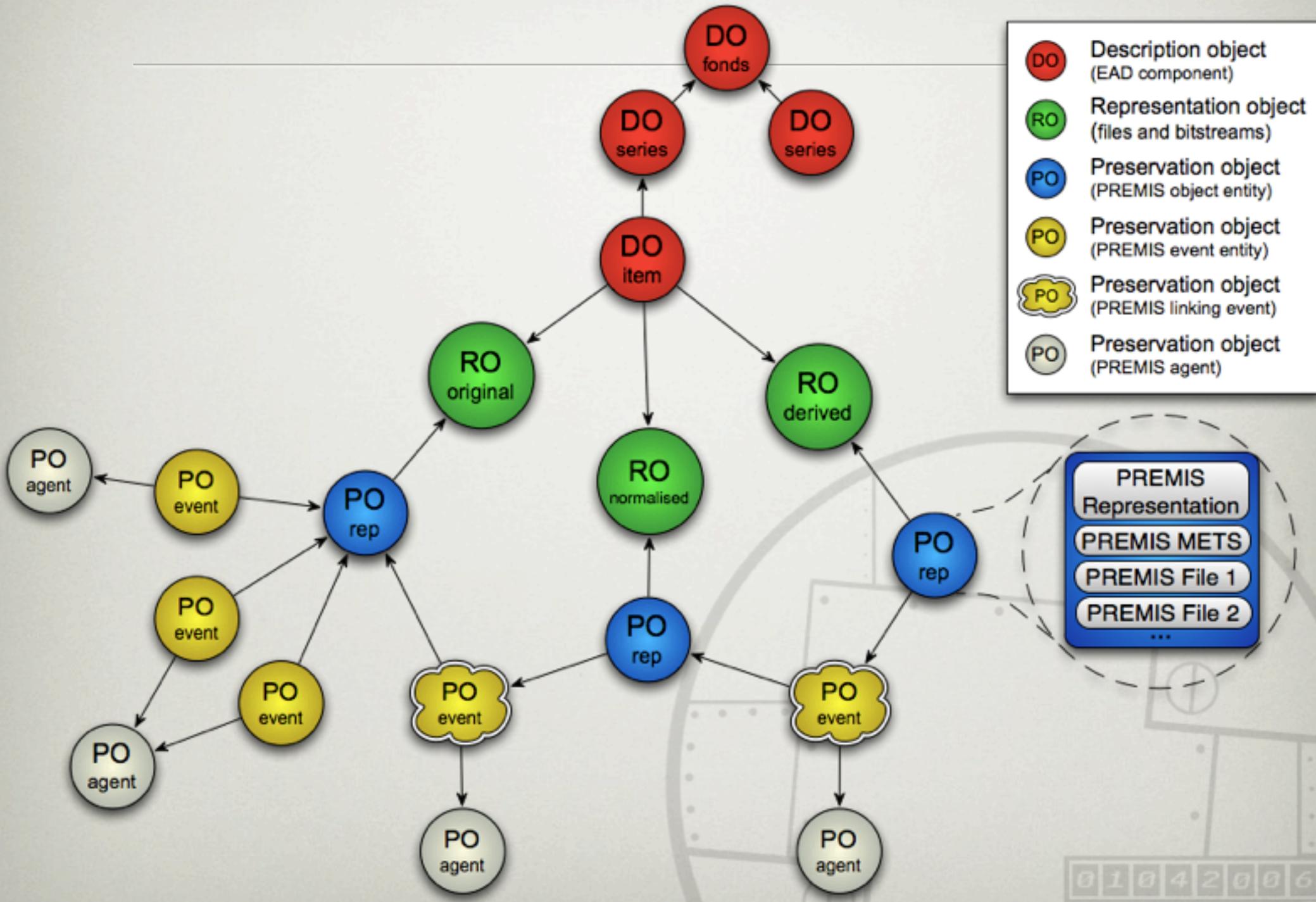


# MODELO DE DADOS

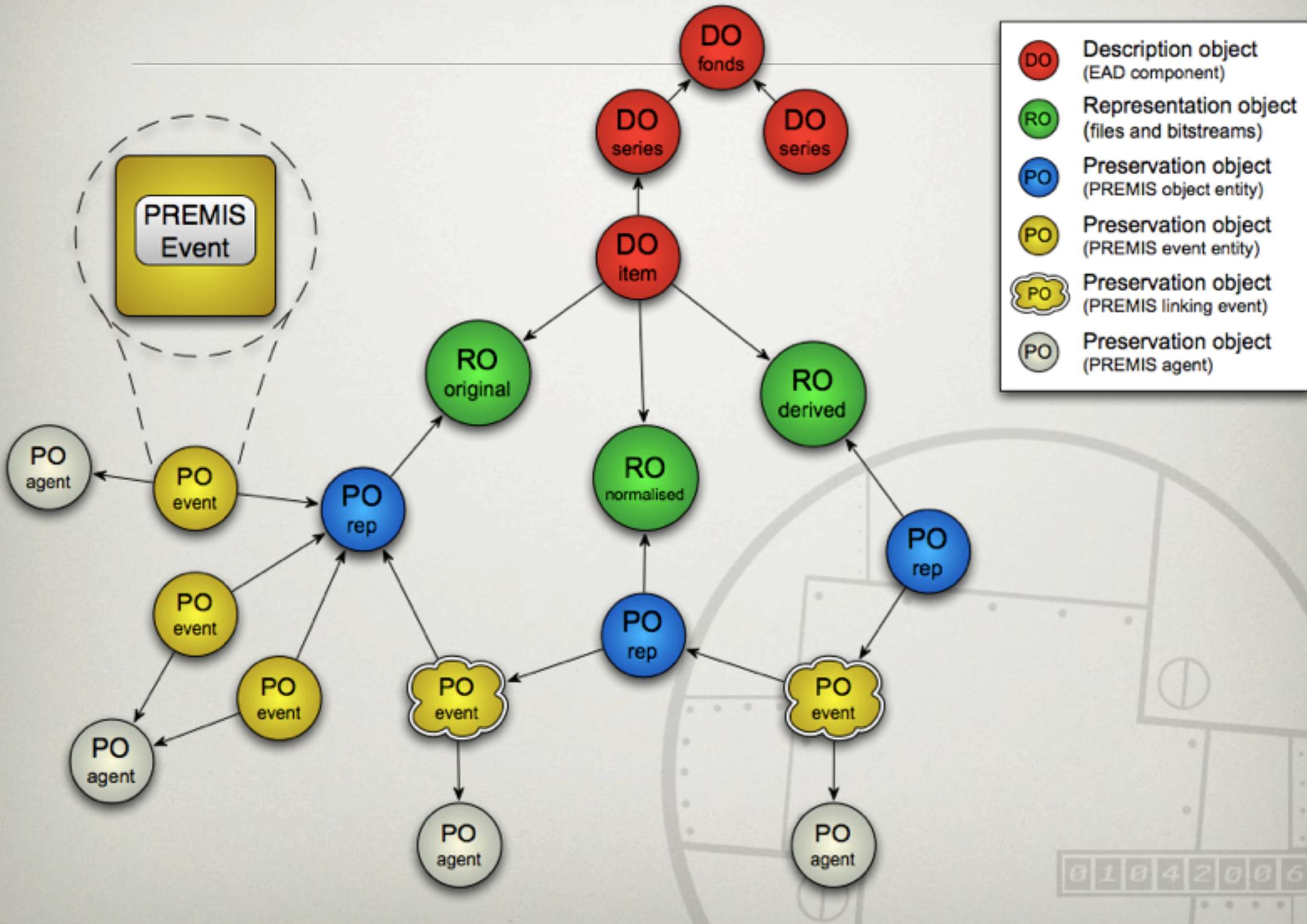


01042006

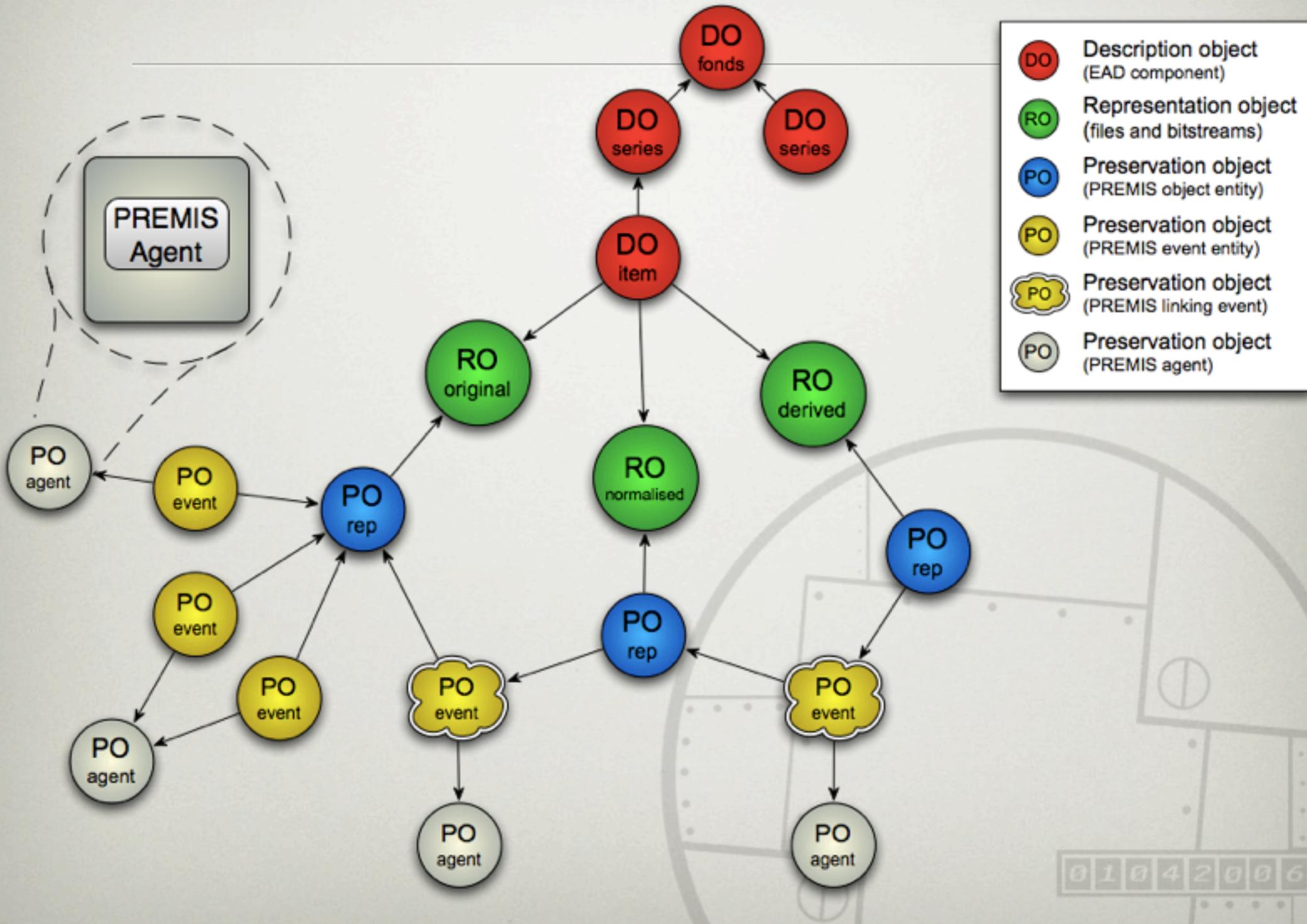
# MODELO DE DADOS



# MODELO DE DADOS



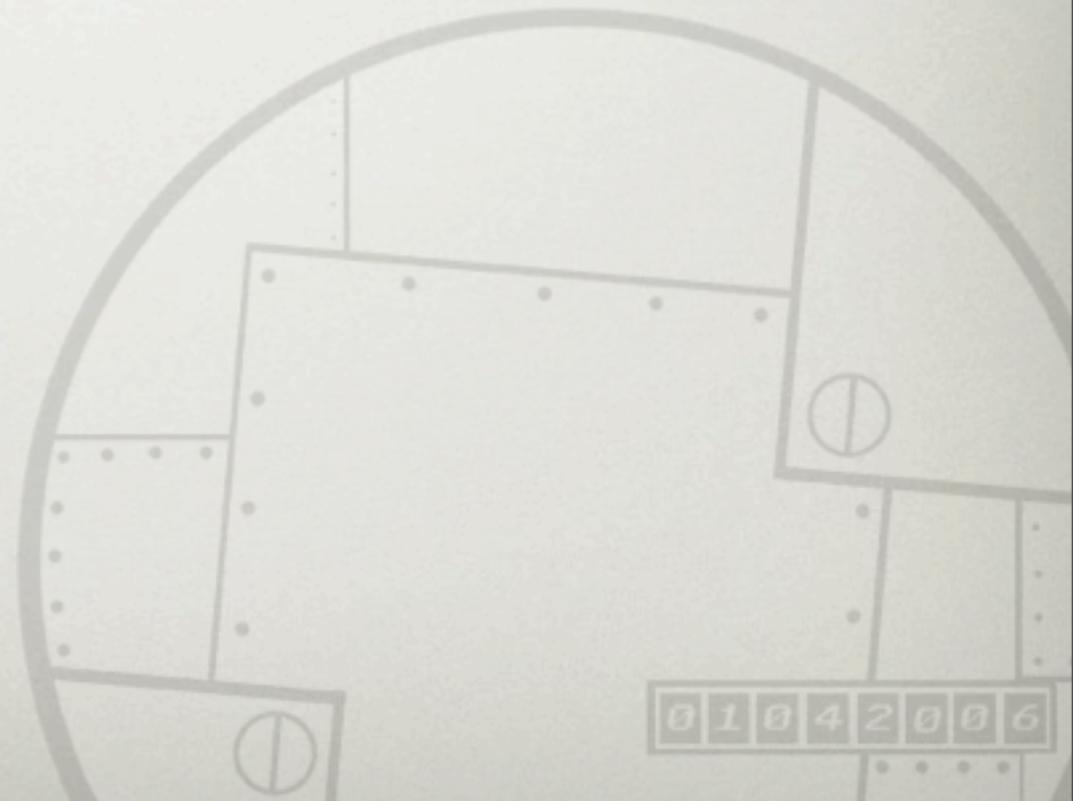
# MODELO DE DADOS



# FASES DO PROJECTO

---

- Análise e planeamento;
- Prototipagem;
- Testes e disseminação.





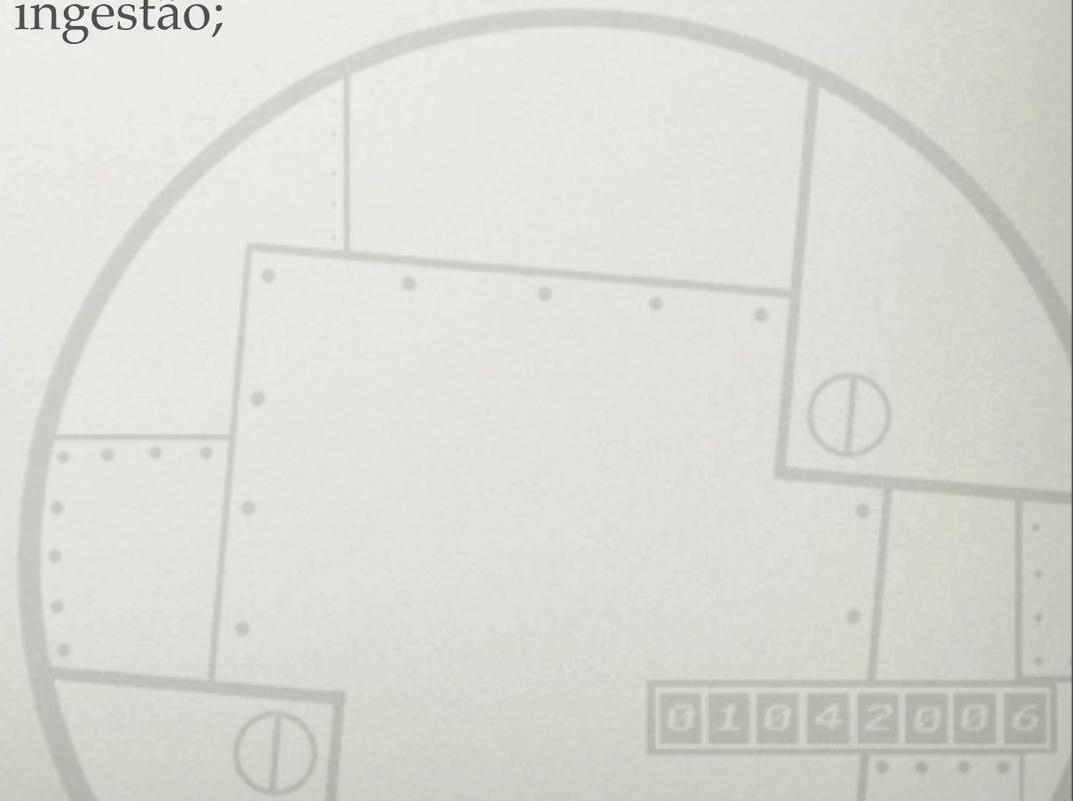
ANÁLISE E  
PLANEAMENTO

0	1	0	4	2	0	0	6
---	---	---	---	---	---	---	---

# REQUISITOS

---

- Interface gráfica para suportar a ingestão;
- Registo de produtores;
- Ferramenta offline para produção de SIPs;
- Relatório de ingestão;
- Ingestão parcial;
- “Quarentena”: cache, buffer de ingestão;
- validação do SIP;
- Identificação dos erros;
- Identificadores persistentes;
- Geração de eventos PREMIS;
- Assinatura digital de DIPs;
- etc.

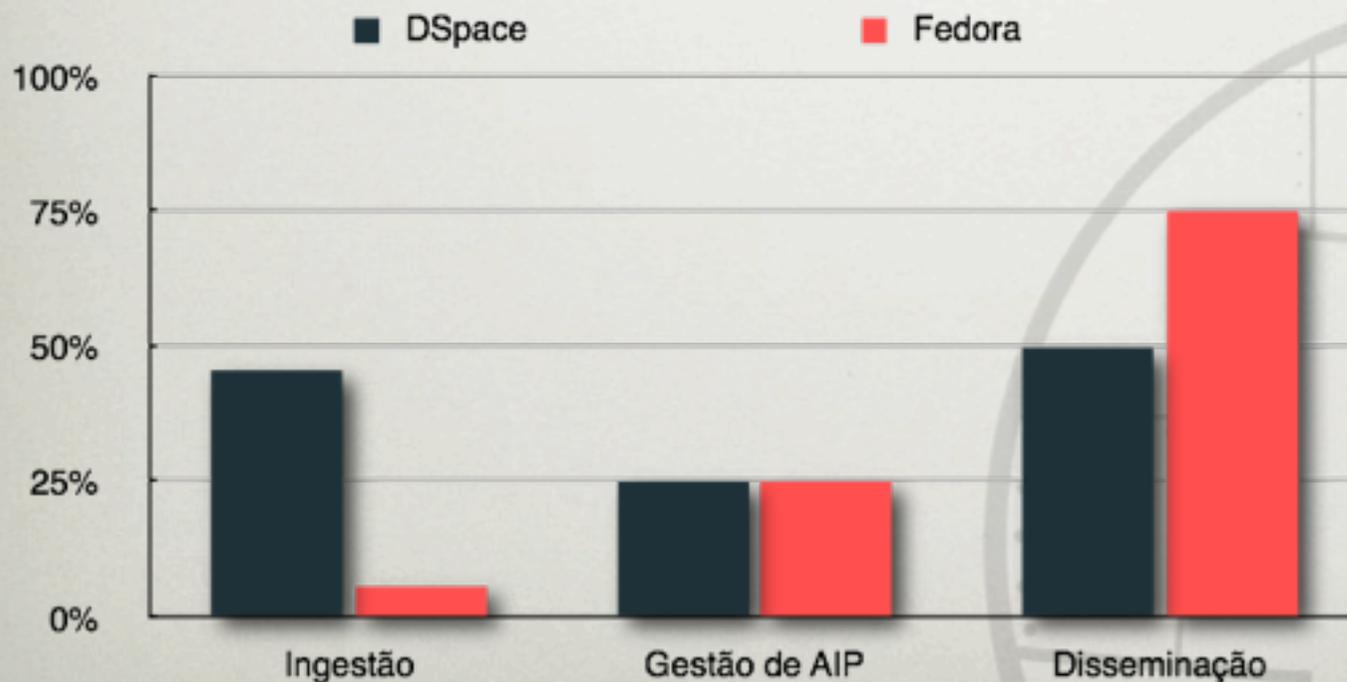
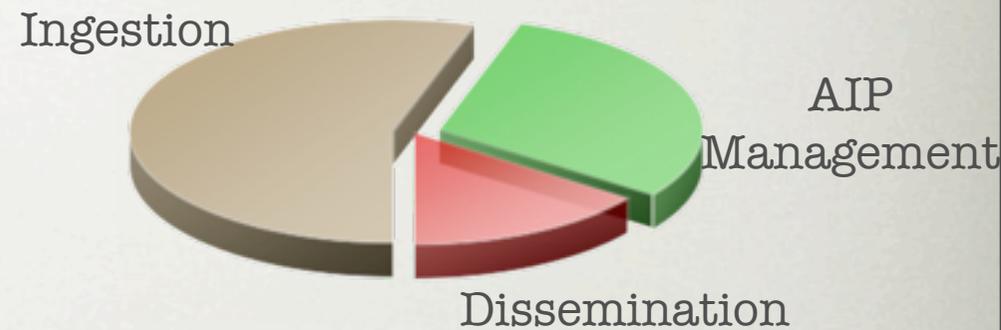


# FRAMEWORK

---

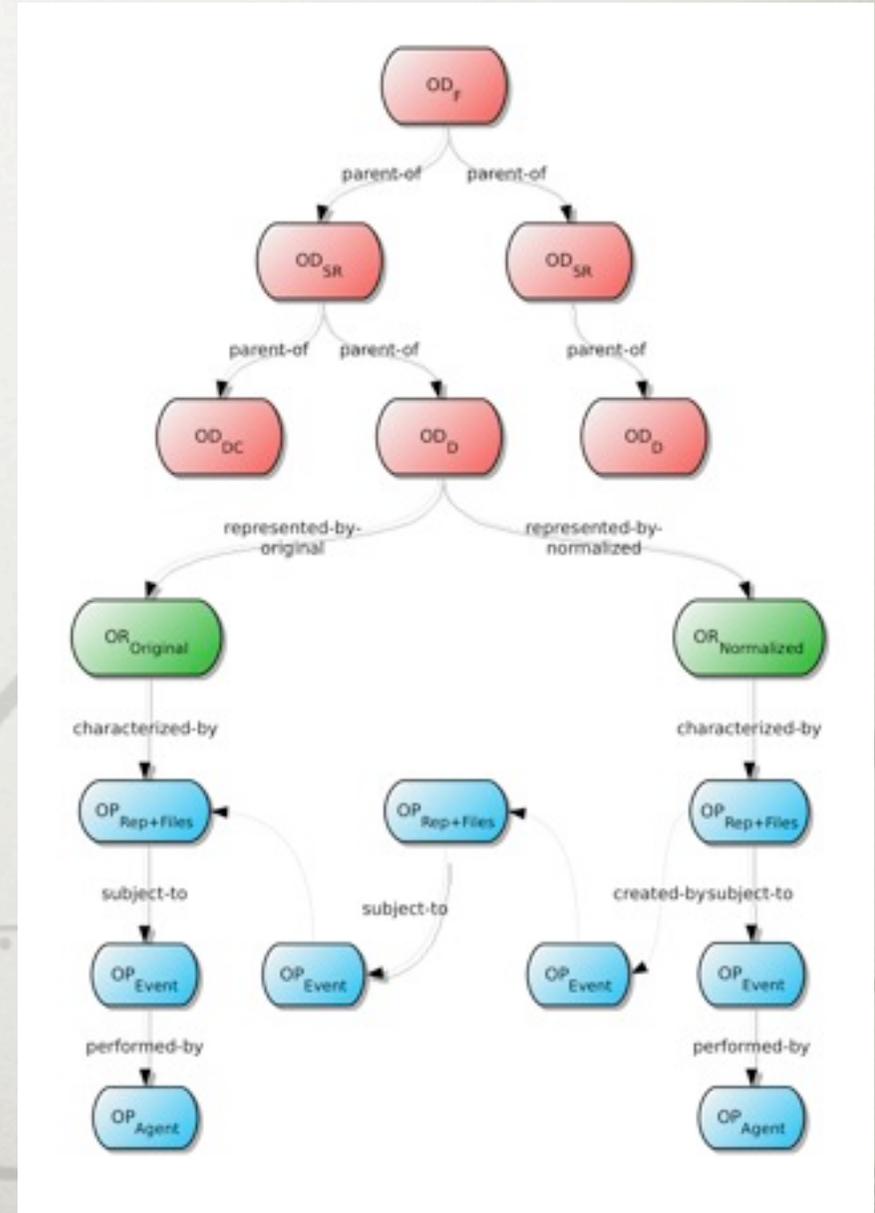
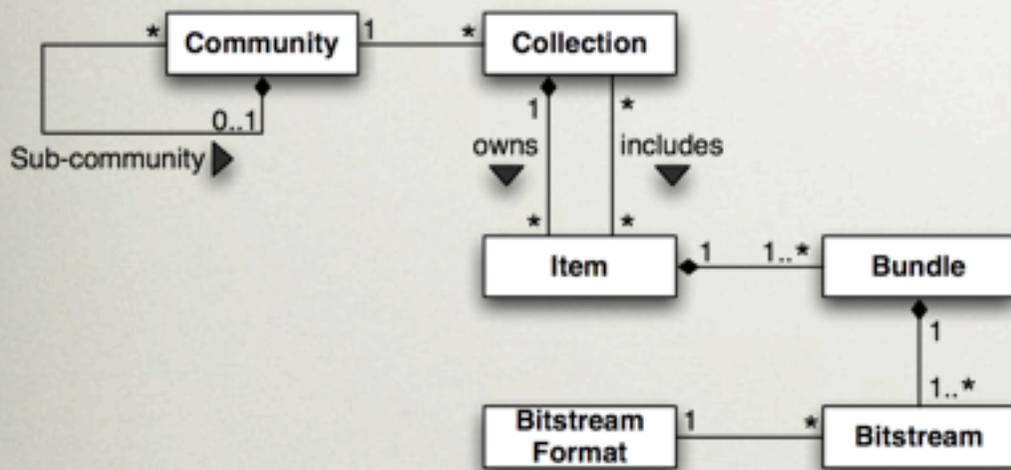


# COMPARAÇÃO BASEADA NOS REQUISITOS



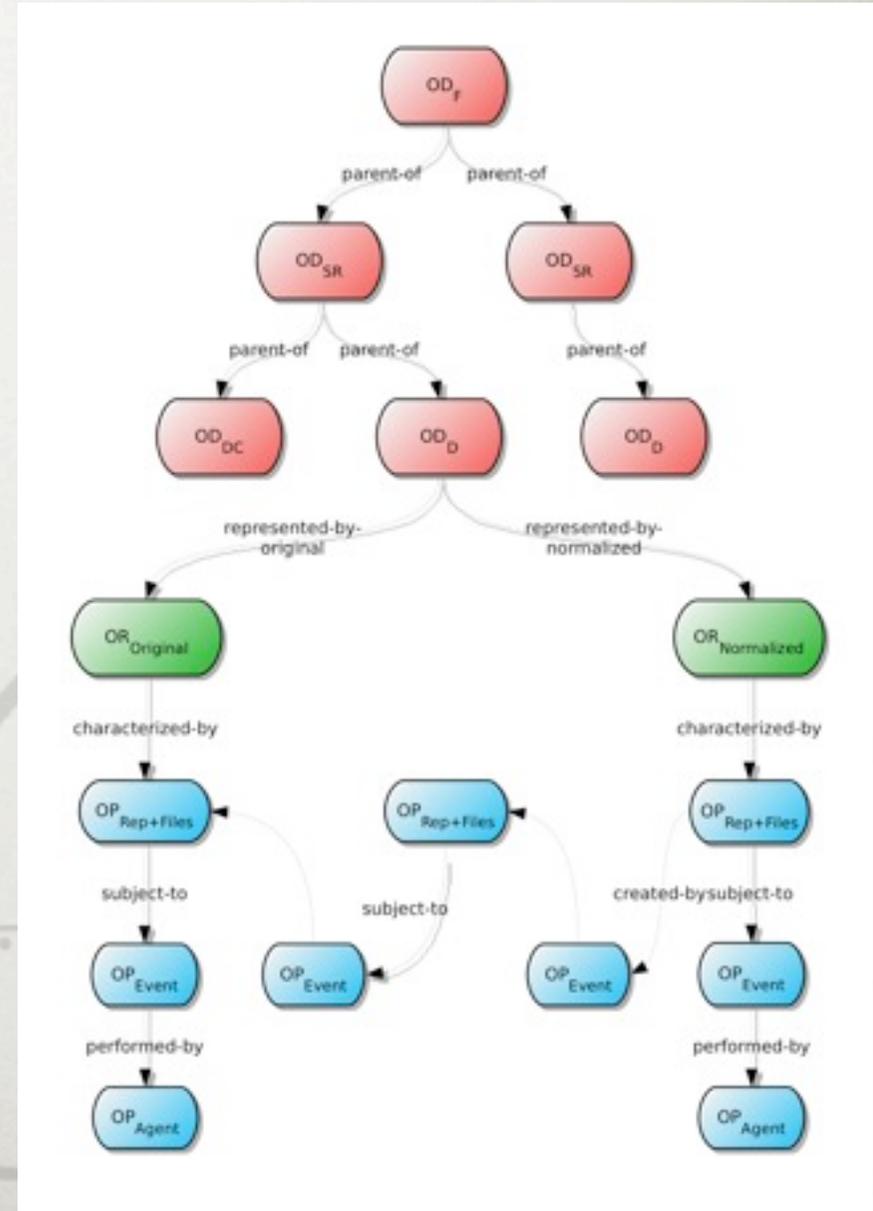
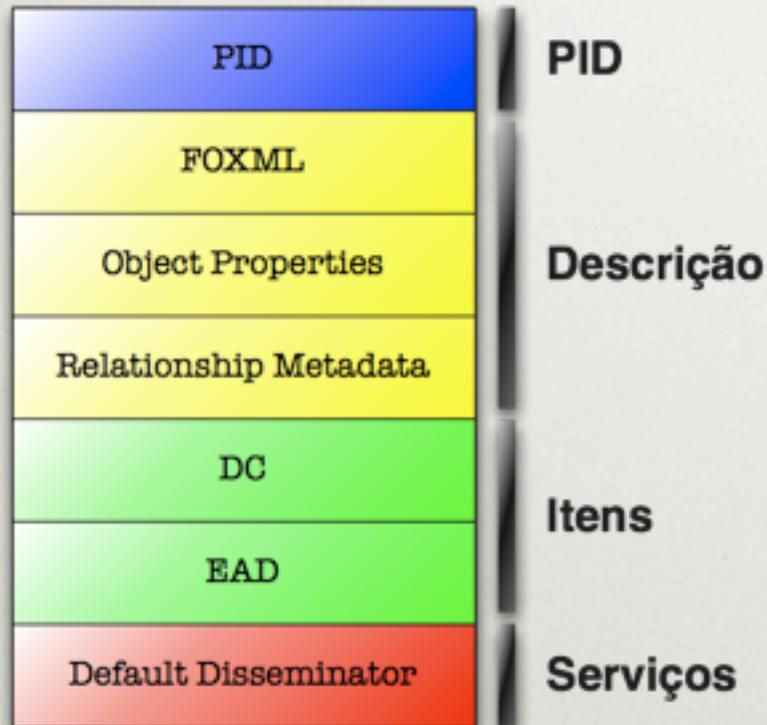
# COMPARAÇÃO DOS MODELOS

DSpace

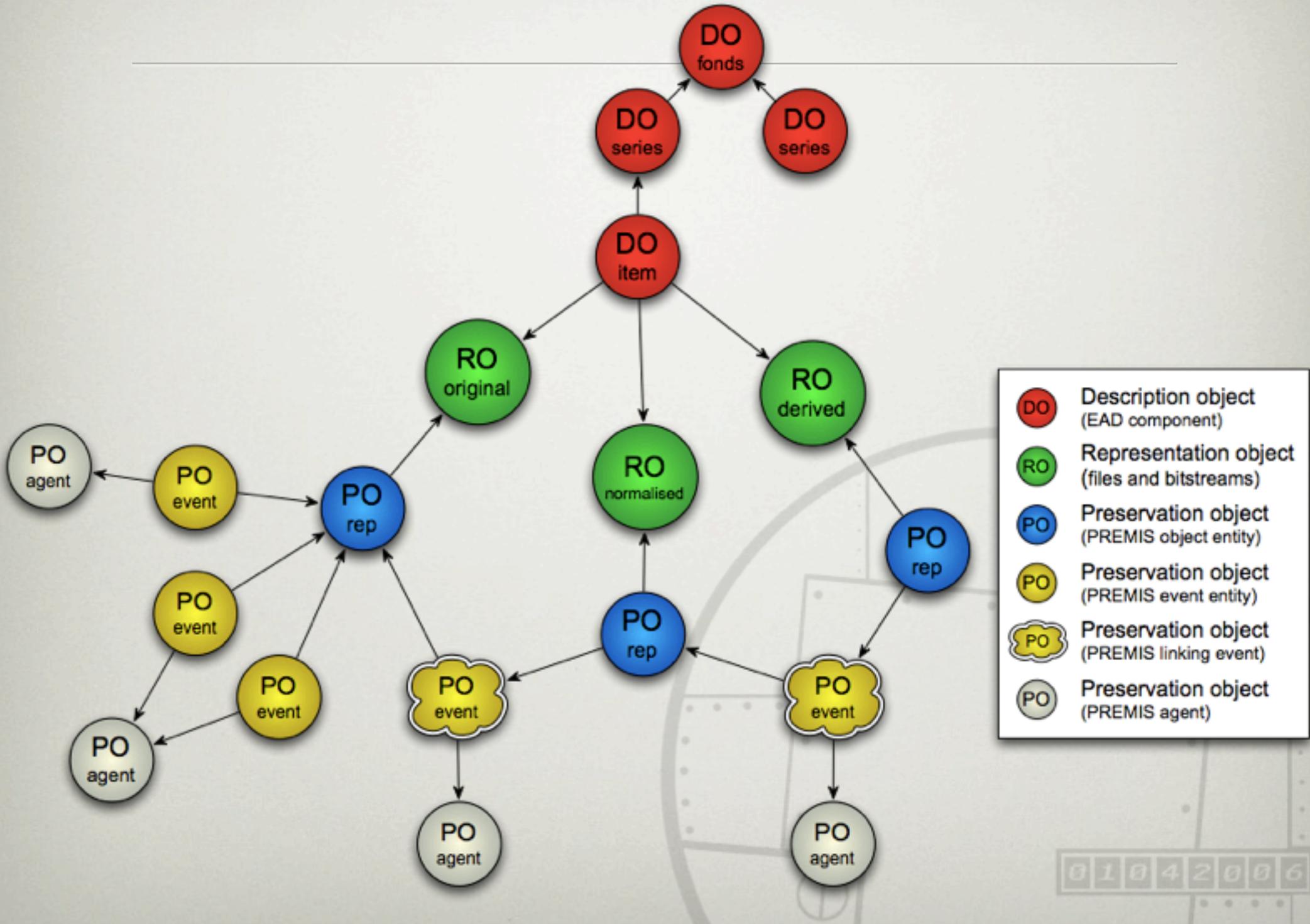


# COMPARAÇÃO DOS MODELOS

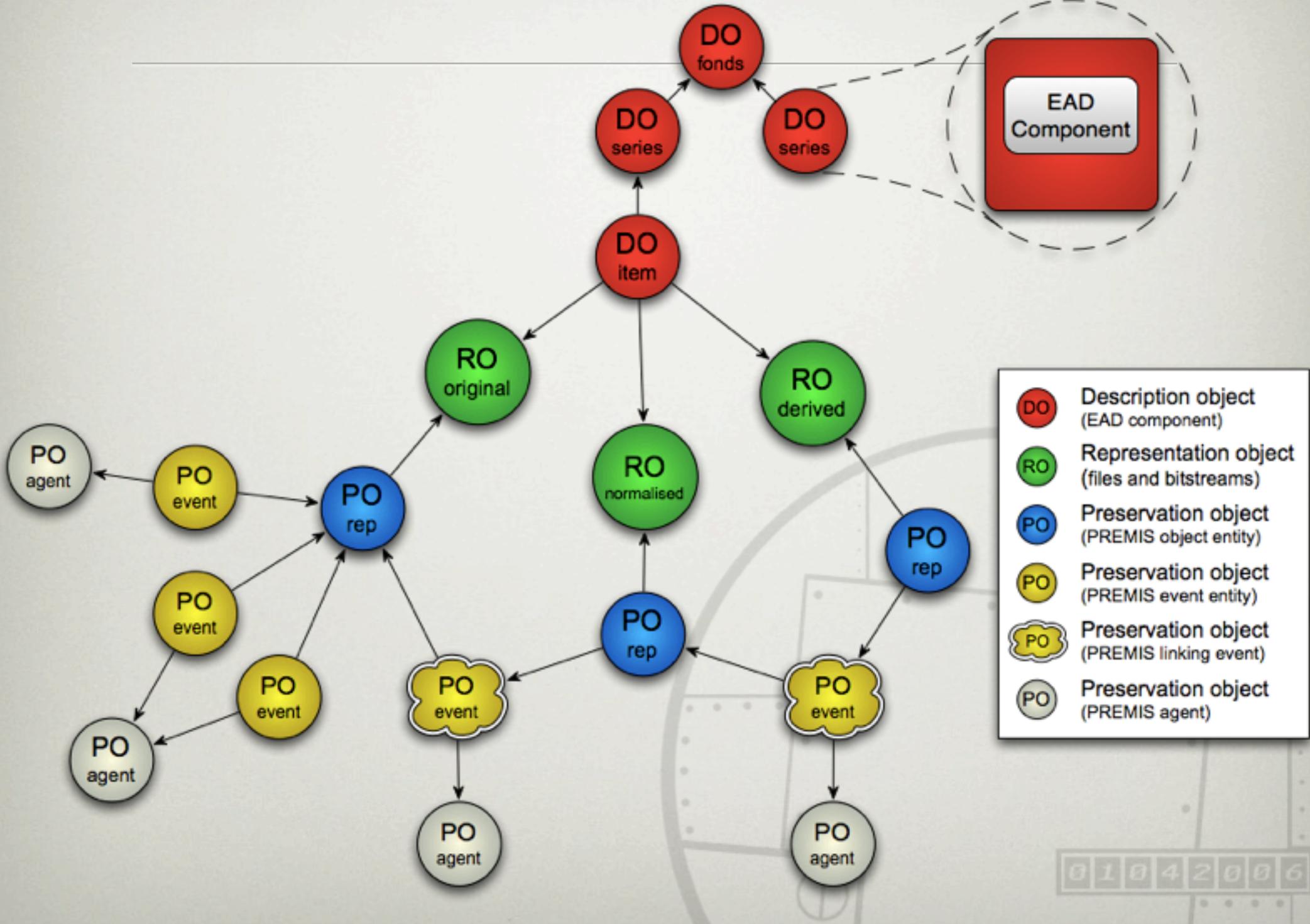
Fedora



# MODELO DE DADOS DO RODA



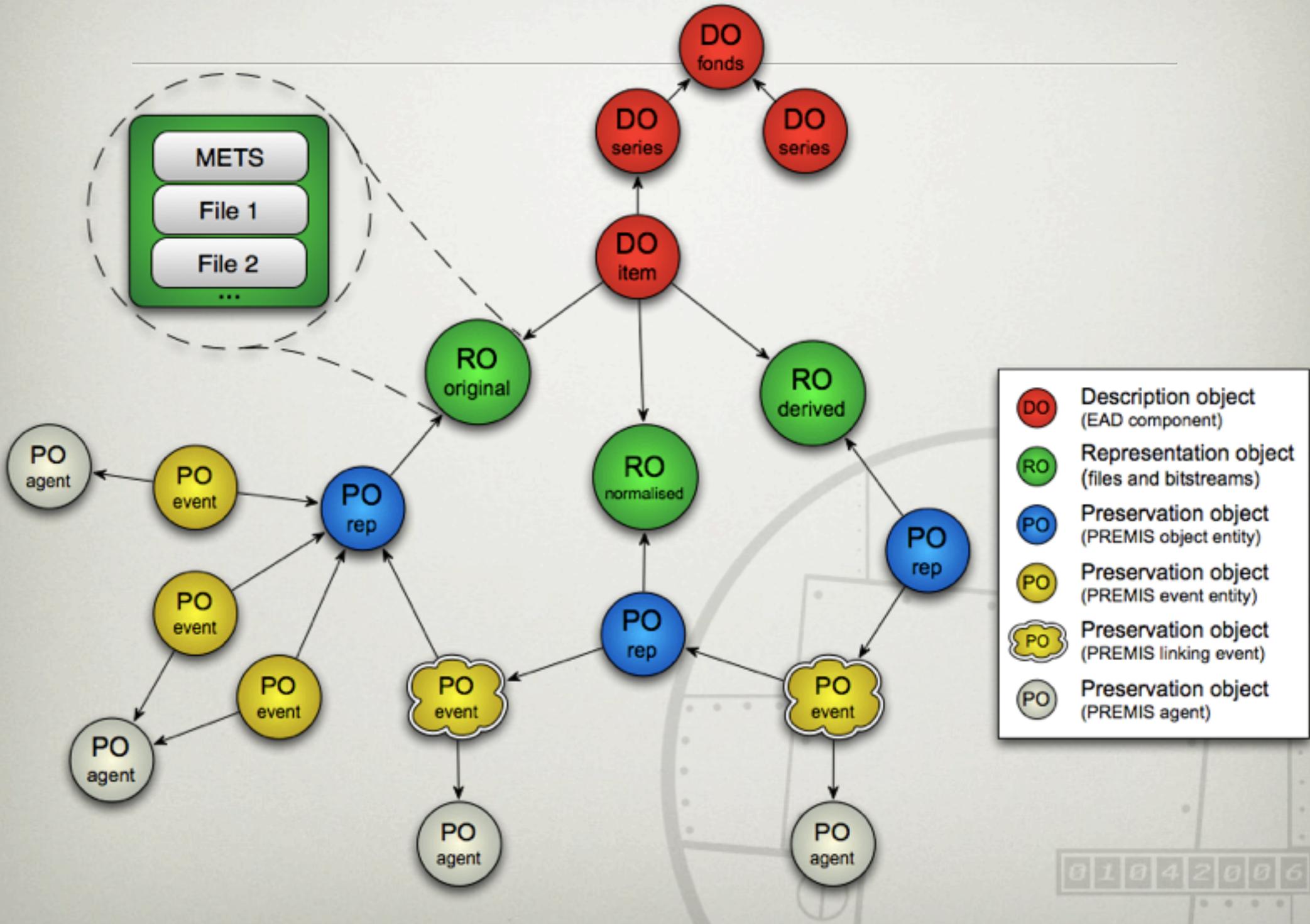
# MODELO DE DADOS DO RODA



- DO Description object (EAD component)
- RO Representation object (files and bitstreams)
- PO rep Preservation object (PREMIS object entity)
- PO event Preservation object (PREMIS event entity)
- PO Preservation object (PREMIS linking event)
- PO agent Preservation object (PREMIS agent)

01042006

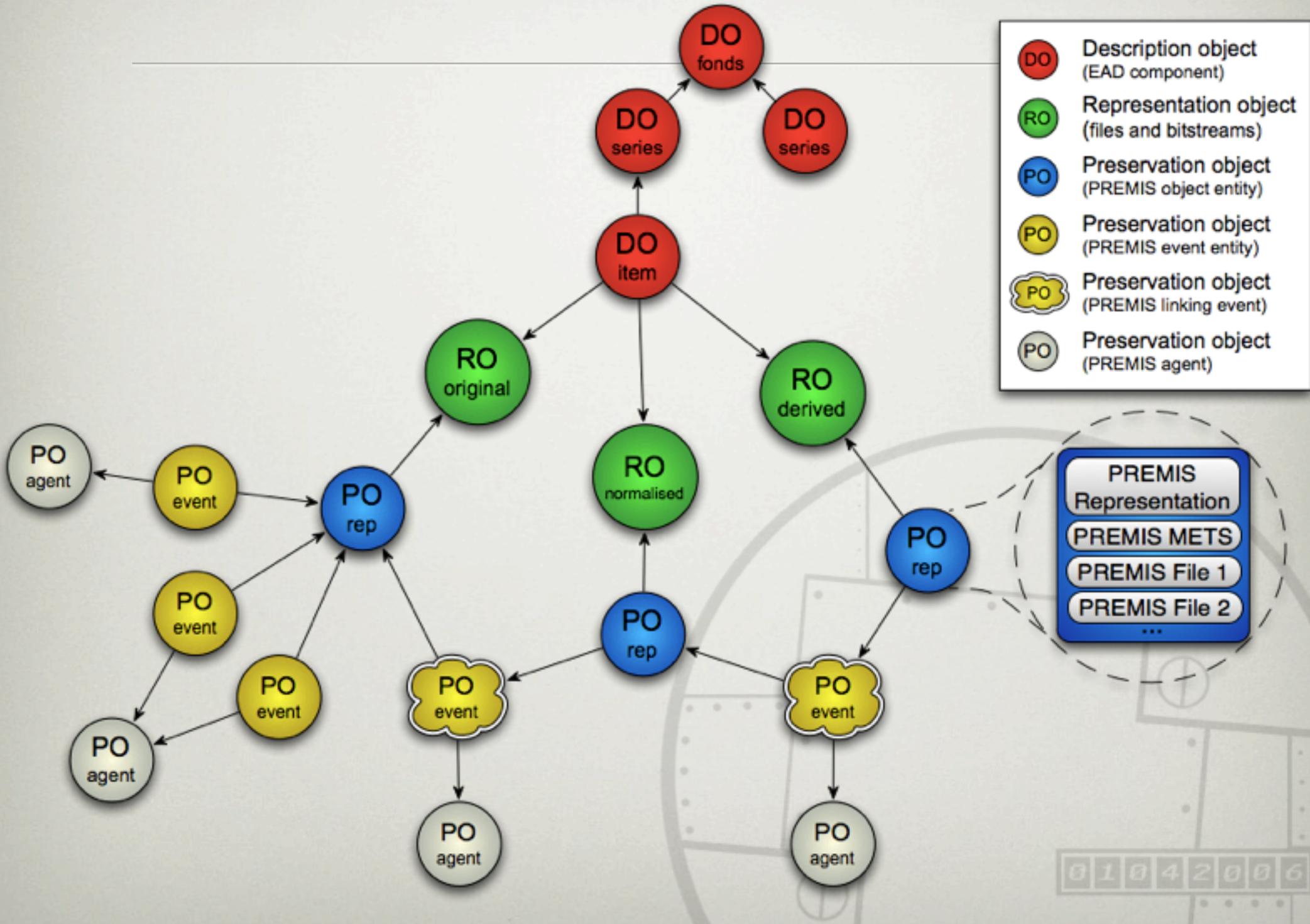
# MODELO DE DADOS DO RODA



- DO Description object (EAD component)
- RO Representation object (files and bitstreams)
- PO Preservation object (PREMIS object entity)
- PO Preservation object (PREMIS event entity)
- PO Preservation object (PREMIS linking event)
- PO Preservation object (PREMIS agent)

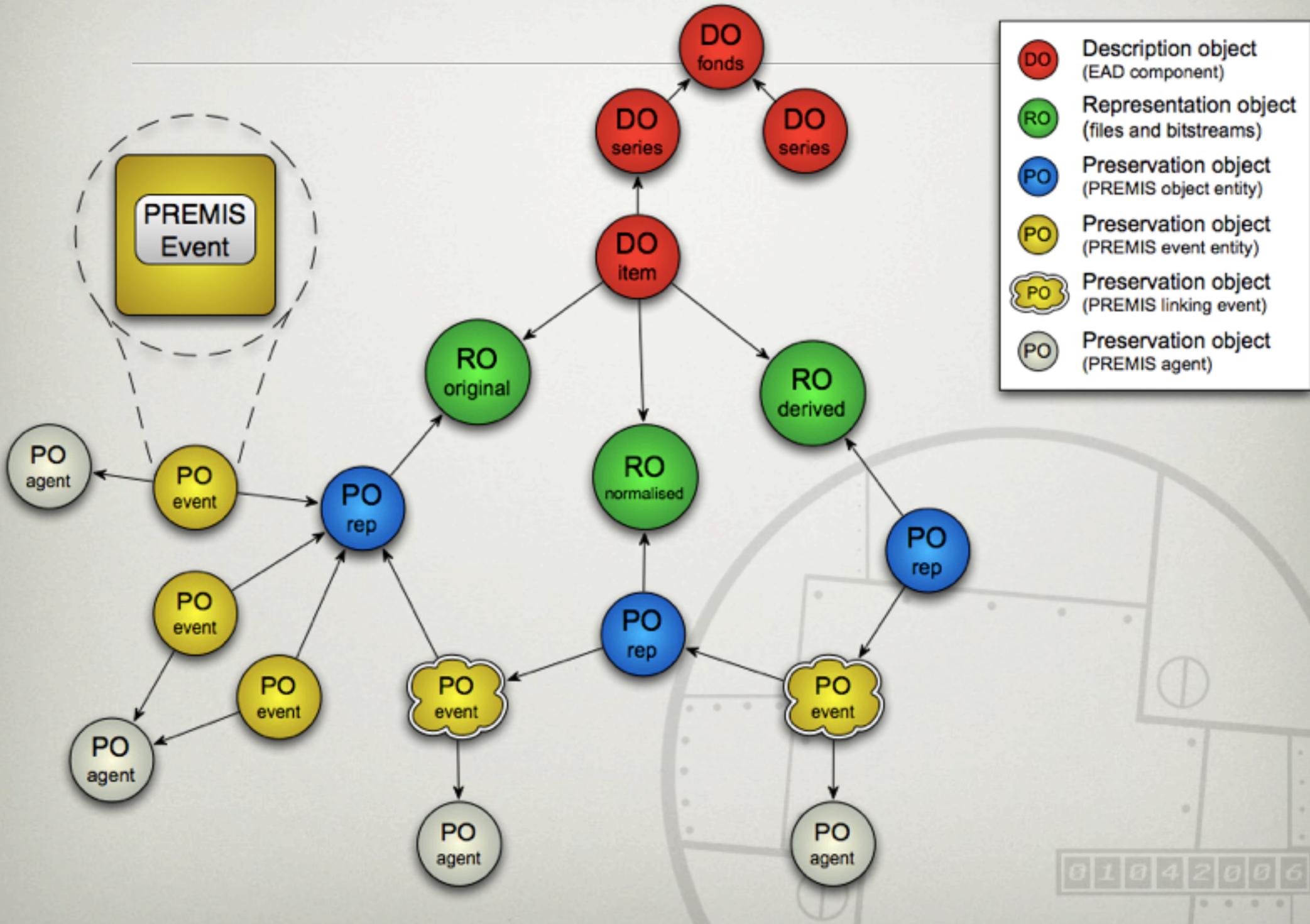
01042006

# MODELO DE DADOS DO RODA

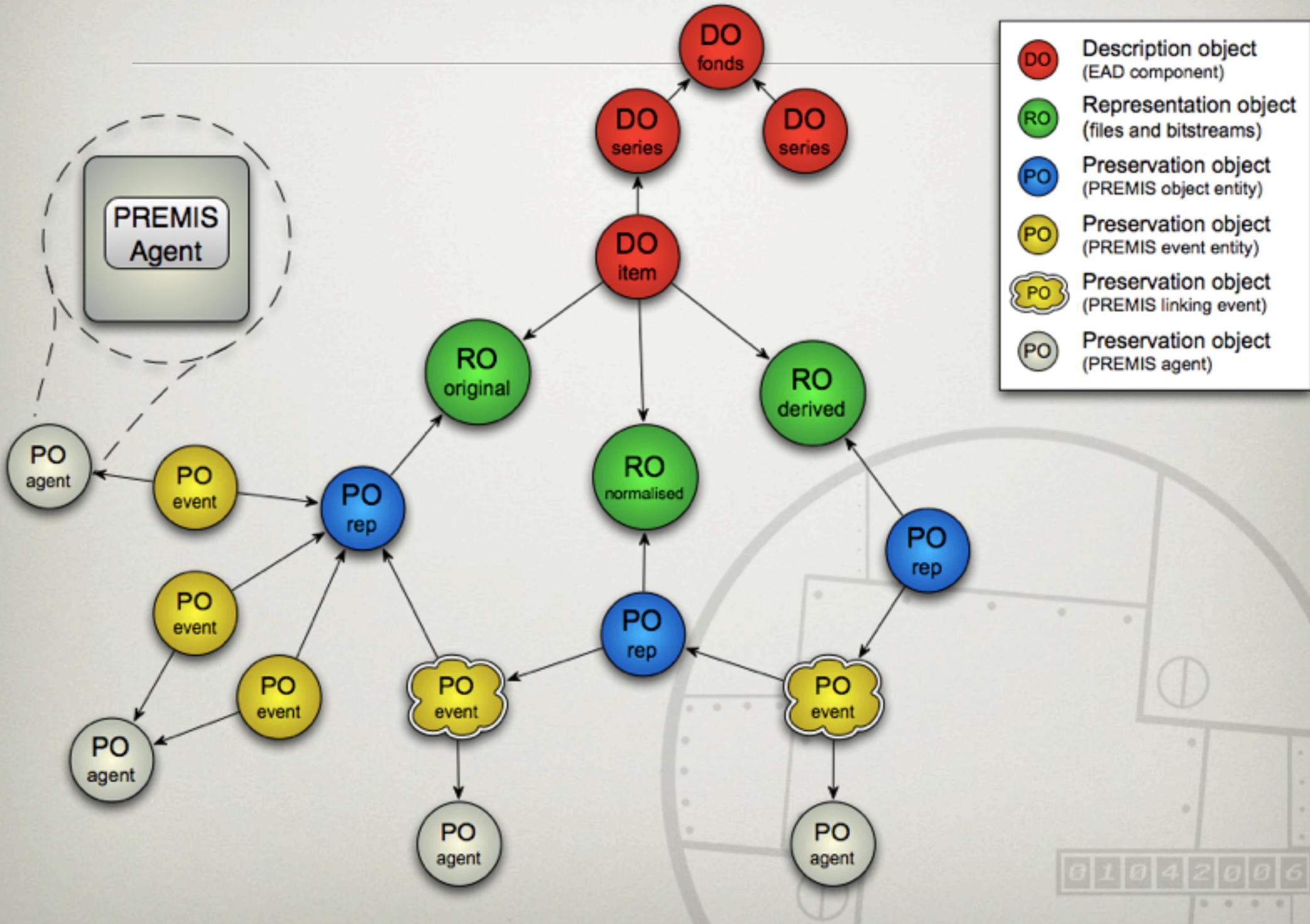


- DO Description object (EAD component)
- RO Representation object (files and bitstreams)
- PO Preservation object (PREMIS object entity)
- PO Preservation object (PREMIS event entity)
- PO Preservation object (PREMIS linking event)
- PO Preservation object (PREMIS agent)

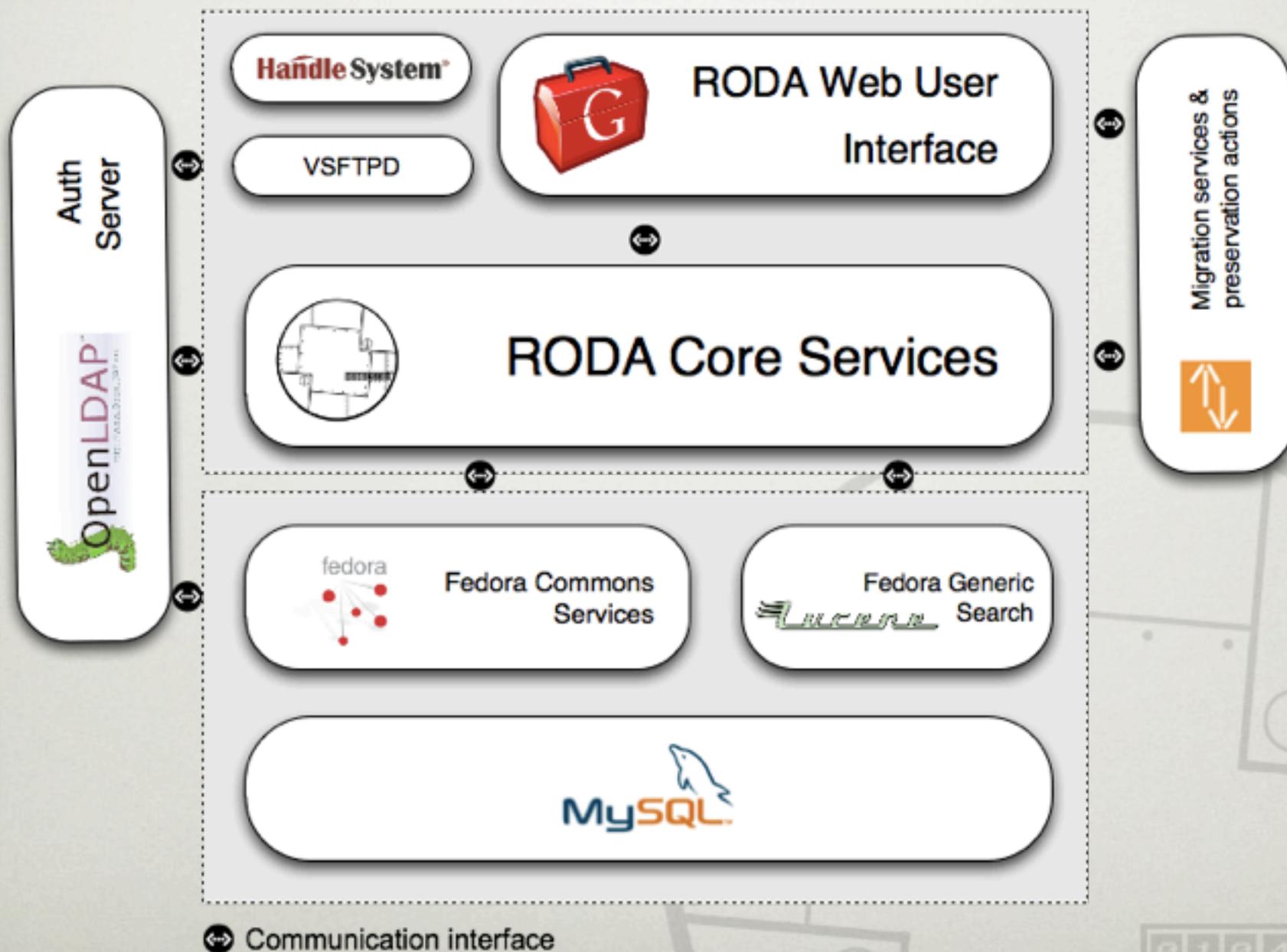
# MODELO DE DADOS DO RODA



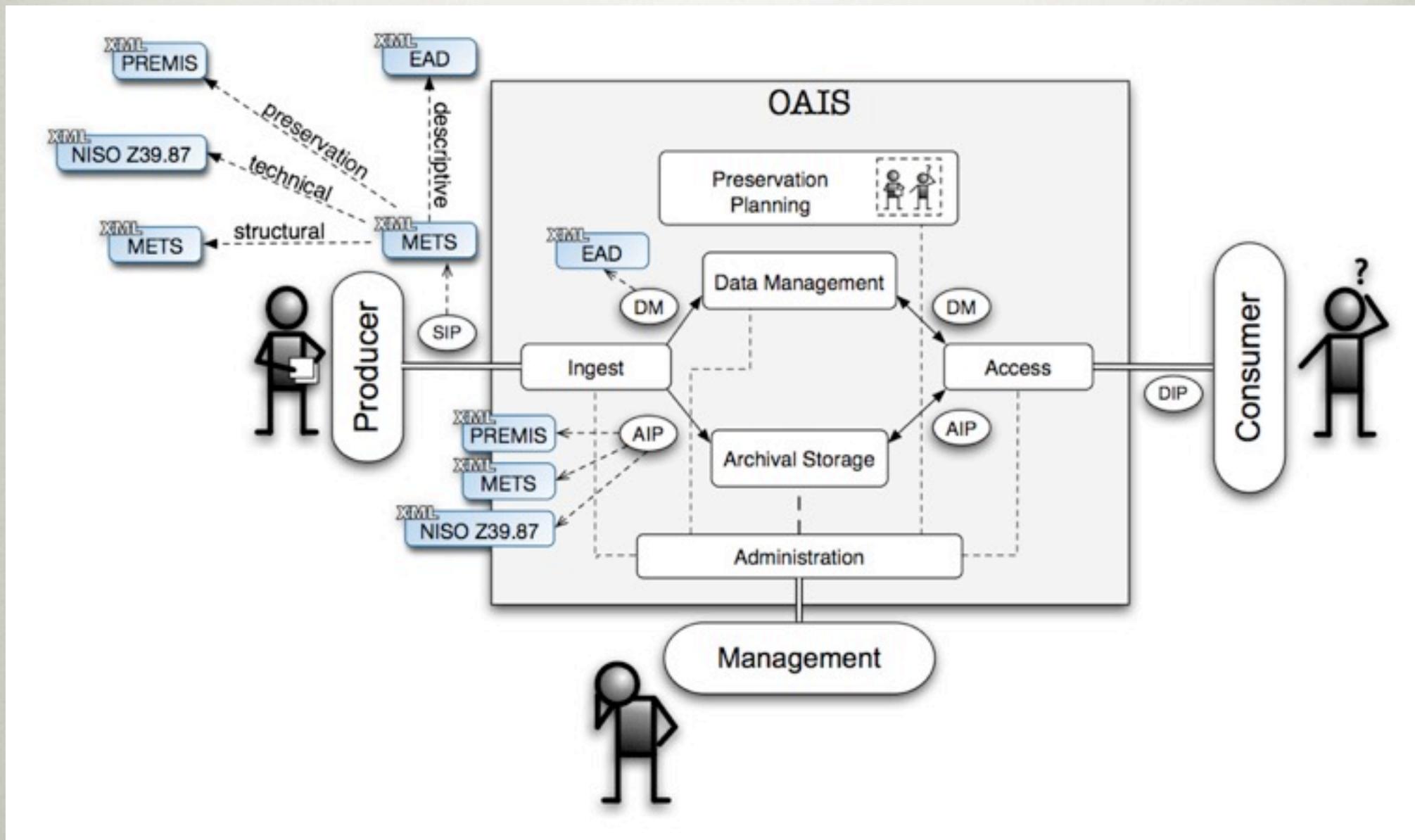
# MODELO DE DADOS DO RODA

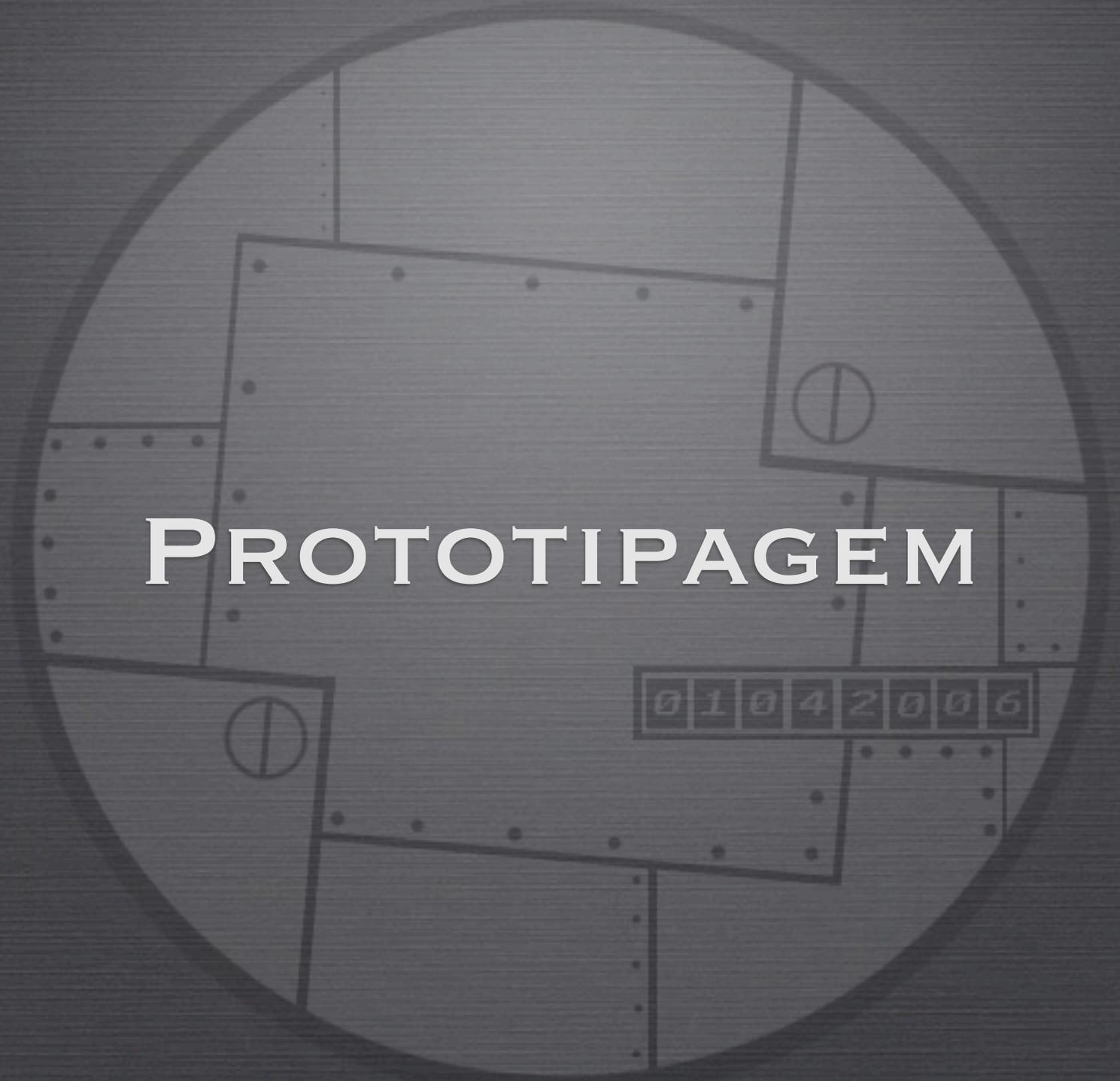


# ARQUITECTURA



# SCHEMAS DO RODA





# PROTOTIPAGEM

# Preservação do Objecto Conceptual

Nível  
Conceptual

Database  
Text Doc.  
Still Image

Nível  
Lógico

SQL Server

Access

PDF Doc.

Ms Word Doc.

PNG image

Nível  
Físico

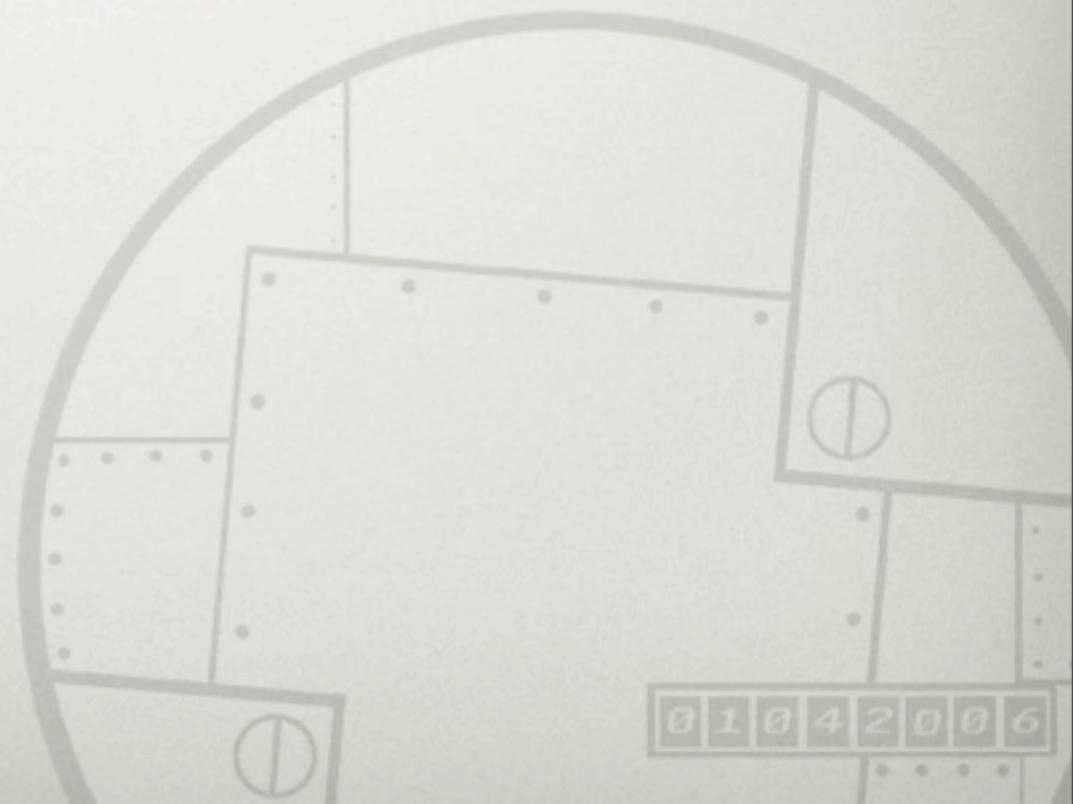
Hard Disc

Tape

...

# DOCUMENTOS DE TEXTO E IMAGENS

---

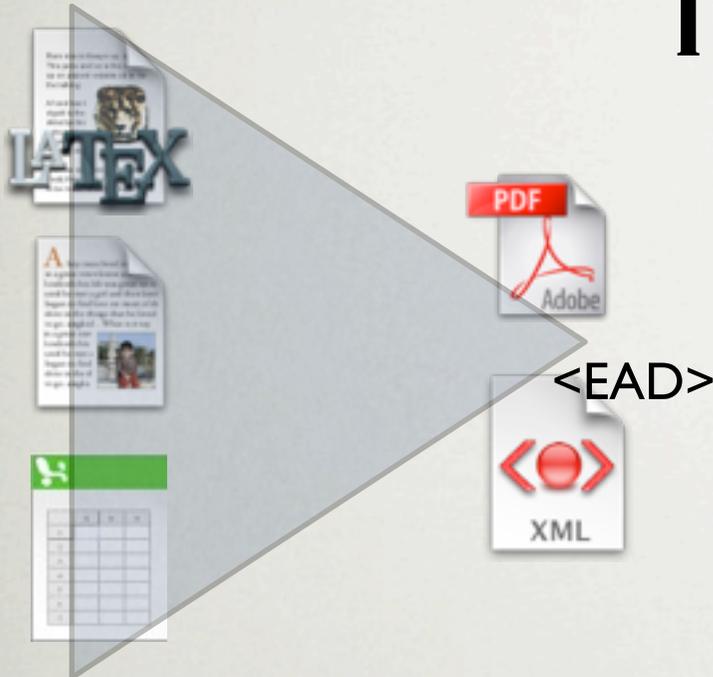


# DOCUMENTOS DE TEXTO E IMAGENS

---

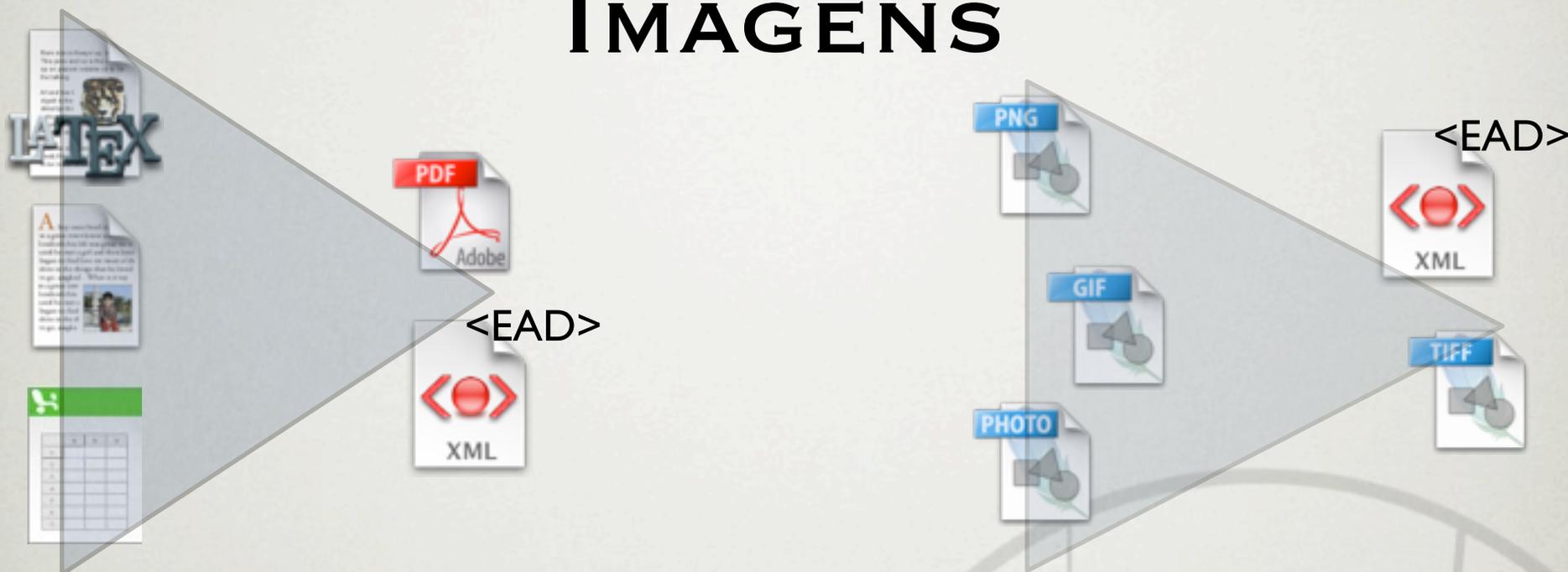
- O EAD retém a maior parte das propriedades significativas: proveniência, história, contexto, etc;
- O conteúdo é mantido num formato normalizado: PDF e TIFF não comprimido.

# DOCUMENTOS DE TEXTO E IMAGENS



- O EAD retém a maior parte das propriedades significativas: proveniência, história, contexto, etc;
- O conteúdo é mantido num formato normalizado: PDF e TIFF não comprimido.

# DOCUMENTOS DE TEXTO E IMAGENS

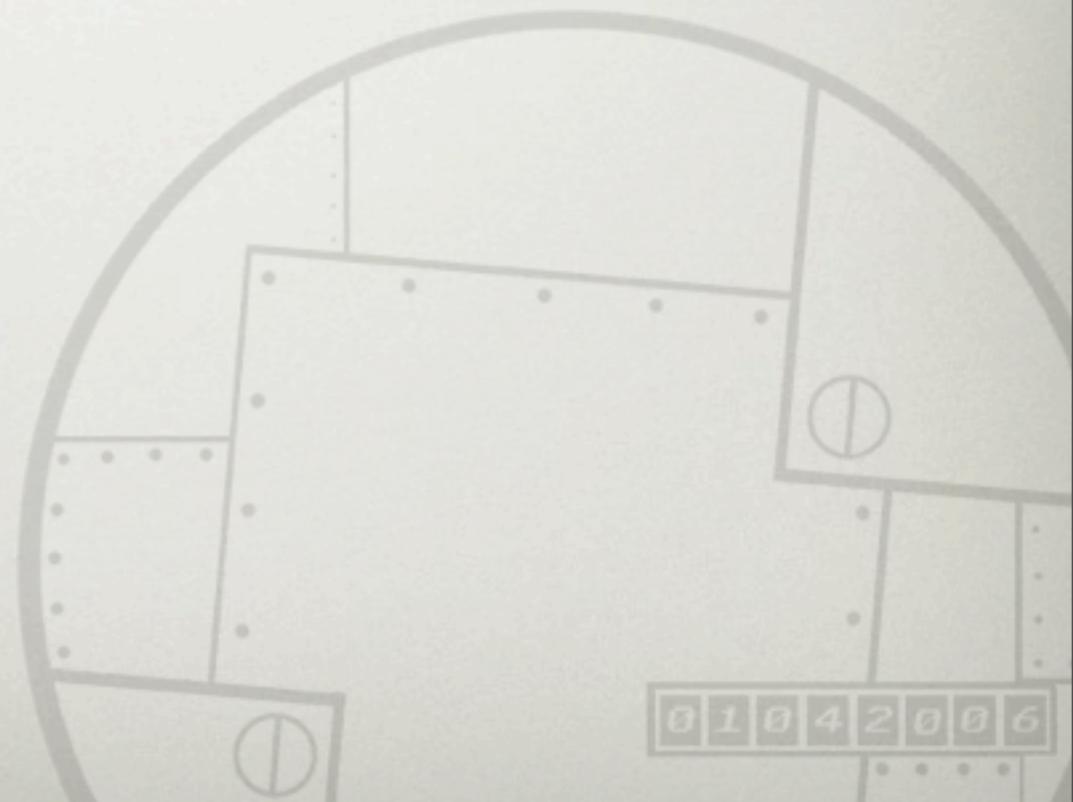


- O EAD retém a maior parte das propriedades significativas: proveniência, história, contexto, etc;
- O conteúdo é mantido num formato normalizado: PDF e TIFF não comprimido.

# DATABASES

---

- Dados?
- Estrutura?
- Views?
- Reports?
- Stored Procedures?
- ...



# DATABASES

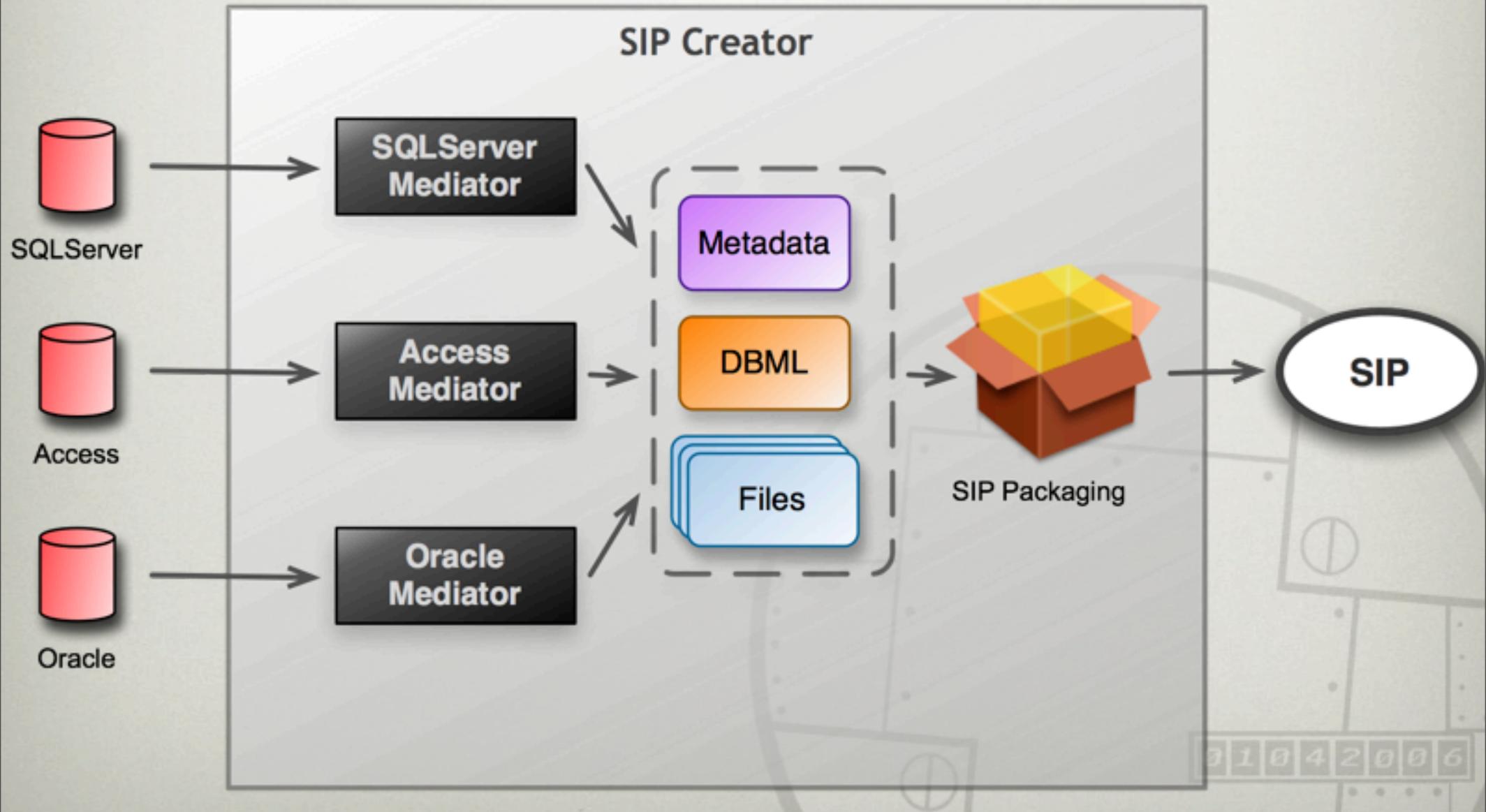
---

- Dados?
- Estrutura?
- Views?
- Reports?
- Stored Procedures?
- ...

## First prototype:

- Data
- Structure
- Modelo Conceptual

# CONSTRUTOR DE SIPs



# DBML

---

- Platform and RDBMS independent
- Stores the DB structure and information
- BLOBs are exported and preserved as standalone files in the representation
- Transformations to SQL and back are defined

# DBML

- Platform and RDBMS independent
- Stores
- BLOBs
- standard
- Transf
- define

```
<TABLE NAME="Districts">
  <COLUMNS>
    <COLUMN NAME="code" TYPE="int" NULL="no"/>
    ...
  </COLUMNS>
  <KEYS>
    <PKEY TYPE="simple">
      <FIELD NAME=""/>
    </PKEY>
    <PKEY TYPE="compound">
      <FIELD NAME=""/>
      <FIELD NAME=""/>
    </PKEY>
    <KEY NAME="" REF=""/>
    ...
  </KEYS>
</TABLE>
```

# DBML

- P
- S
- B
- S
- Transf
- define

```
...  
<DATA>  
  <products>  
    <products-REG>  
      <code> a122 </code>  
      <description> milk </description>  
    ...  
  </products-REG>  
</products-REG>  
...  
</products-REG>  
</products>  
...  
</DATA>  
...
```

MS independent

```
<COLUMN NAME="code" TYPE="int" NULL="no"/>  
...  
<UMNS>  
...  
TYPE="simple">  
  <FIELD NAME=""/>  
</>  
TYPE="compound">  
  <FIELD NAME=""/>  
  <FIELD NAME=""/>  
</PKEY>  
<KEY NAME="" REF=""/>  
...  
</KEYS>  
</TABLE>
```



# DBML

---

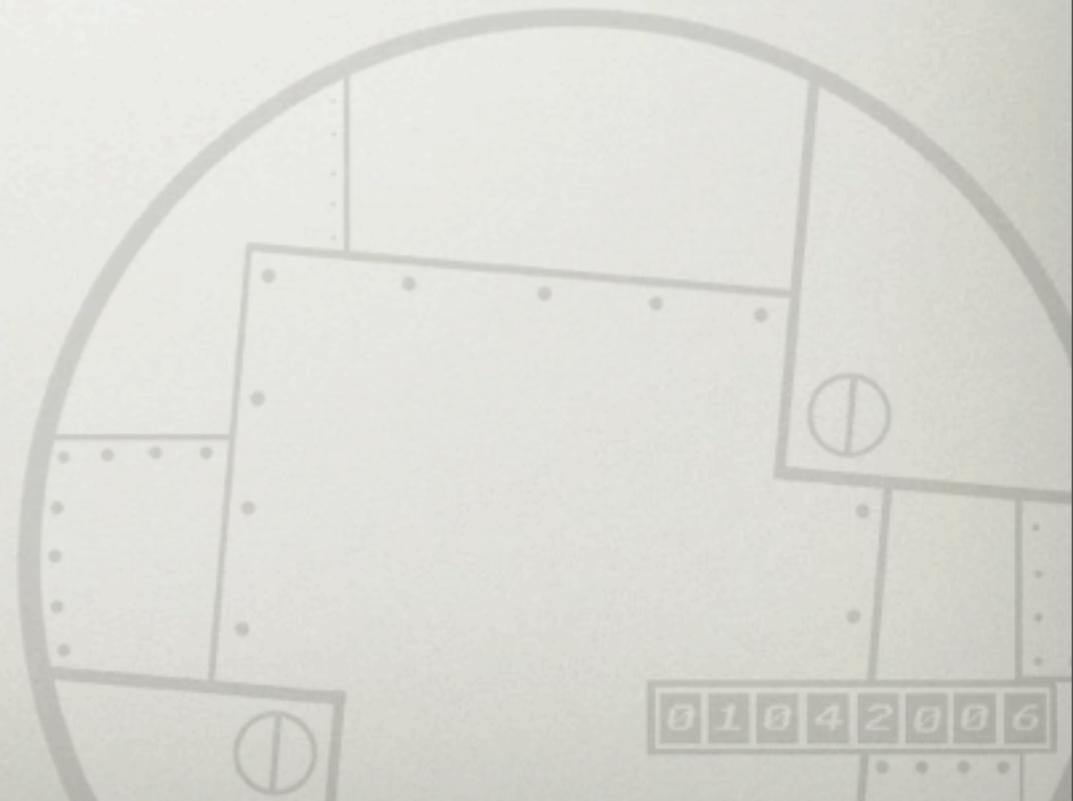
## Composição do SIP:

- METS como manifesto;
- EAD na descrição das propriedades intelectuais;
- ficheiro(s) DBML
- um OD para cada BLOB
- ficheiro METS + MIX para cada OD

# FINAL THOUGHTS

---

*“Data Preservation is a people problem”  
Michael Lesk*



# FINAL THOUGHTS

---

*“Data Preservation is a people problem”  
Michael Lesk*

- As pessoas precisam de ser treinadas para guardar os dados de modo apropriado;
- Preservar o quê? Dados, Estrutura, Semântica...
- A preservação é para utilizadores futuros mas são os utilizadores de hoje que votam o orçamento;
- Temos de fazer com que as pessoas responsáveis por colectar dados tenham preocupações de preservação;
- A preservação tem falhas. Todos os sistemas são imperfeitos.

# TRABALHO FUTURO

---

- 2 teses de doutoramento;
- 1 tese de mestrado;
- 1 projecto europeu:
  - \* [www.scape-project.eu](http://www.scape-project.eu)



- Criar uma comunidade;
- Promover o desenvolvimento Open Source;



LET'S PRESERVE TOMORROW'S HISTORY...

QUESTIONS?