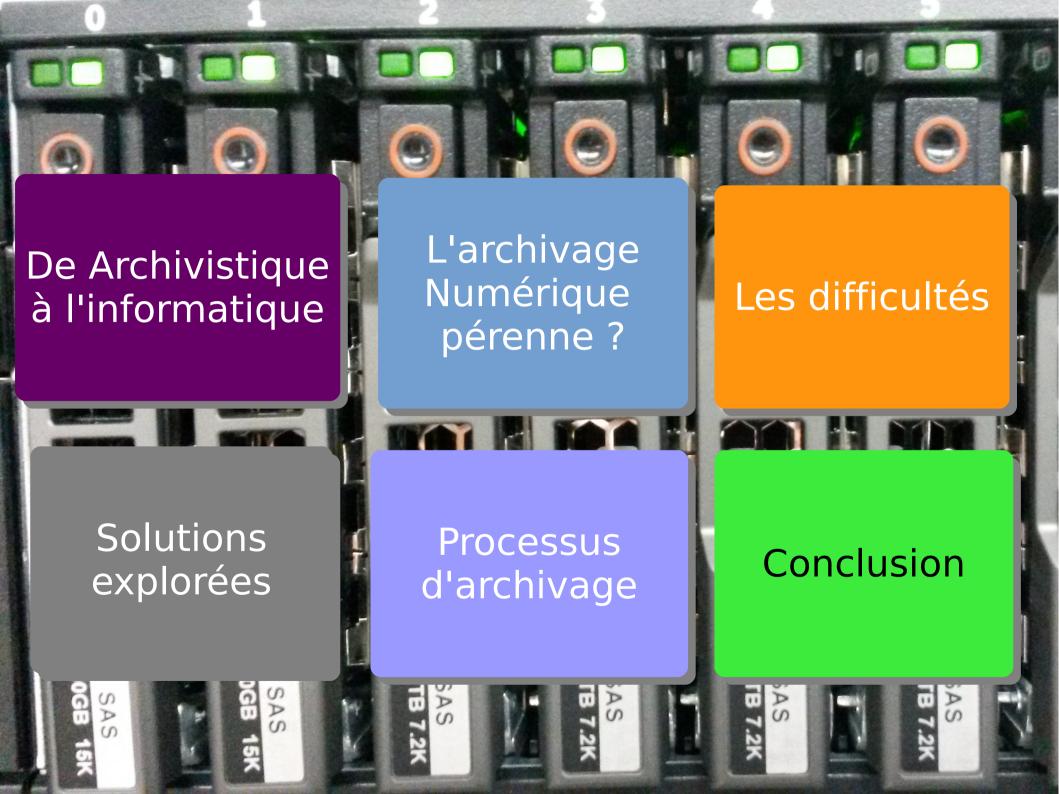
# Pourquoi est-il si difficile d'archiver?

Philippe Saby Capitoul 26 février 2015







# De l'archivistique à l'informatique

















# En Archivistique

#### Archives?

- L'ensemble des documents, quels que soient leur date, leur lieu de conservation, leur forme et leur support, produits ou reçus par toute personne physique ou morale et par tout service ou organisme public ou privé dans l'exercice de leur activité (Code du Patrimoine, art. L. 211-1),
- ·les services et institutions qui se chargent de leur gestion,
- les espaces de stockage de ces documents.

#### La notion d'archives est indépendante :

- de la date :
  - tout document est considéré comme archive dès son élaboration ou sa réception,
- de la forme :
  - document manuscrit, texte imprimé ou dactylographié, plan, affiche etc.
- du support :
  - papier, photographie, bande audiovisuelle, document numérique etc.

# En Archivistique

#### Cycle de vie des documents ?

- •<u>Dossiers actifs</u> (aussi appelés <u>archives courantes</u> ou archives administratives) qui regroupent les documents qui sont nécessaires à l'activité courantes des organismes qui les ont produits (créés ou reçus dans le cadre de leurs activités);
- •<u>Dossiers clôturés</u> (<u>archives intermédiaires</u>) à conserver de façon intègre tant qu'ils ont une valeur probante (utilité administrative ou légale) pour la traçabilité des activités de l'organisme producteur;
- Dossiers inactifs qui n'ont plus de valeur probante

-Tri: éliminés ou archivés

archivage numérique pérenne

30 ans

## De l'archivistique à ... l'informatique

#### OAIS

#### **Open Archival Information System**

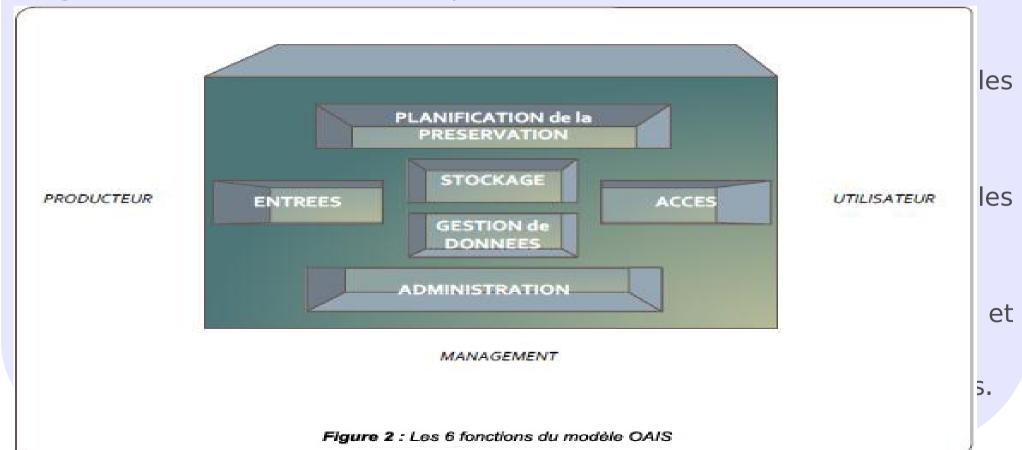
- Modèle conceptuel destiné à la gestion, à l'archivage et à la préservation à long terme de documents numériques
- Entités fonctionnelles du modèle
  - -<u>les entrées</u> : reçoivent les paquets d'information à verser et les transmettent au stockage,
  - -le stockage: stocke et sauvegarde les paquets d'information archivés,
  - -<u>la gestion des données</u> : met à disposition toutes les informations utiles au fonctionnement de l'archive,
  - -<u>l'administration</u> : pilote le système,
  - -<u>la planification de la préservation</u> : assure une veille technologique et propose les évolutions et les stratégies pour prévenir l'obsolescence,
  - -<u>l'accès</u> : communique les paquets d'information diffusés aux utilisateurs.

## De l'archivistique à ... l'informatique

#### **OAIS**

#### **Open Archival Information System**

• Modèle conceptuel destiné à la gestion, à l'archivage et à la préservation à long terme de documents numériques



# L'archivage numérique pérenne?

















## Stockage / Sauvegarde / Archivage

#### Stockage

- Données de travail
- Local (poste), lan (NFS, CIFS ...), WAN (cloud, objet ...)



• 1 copie

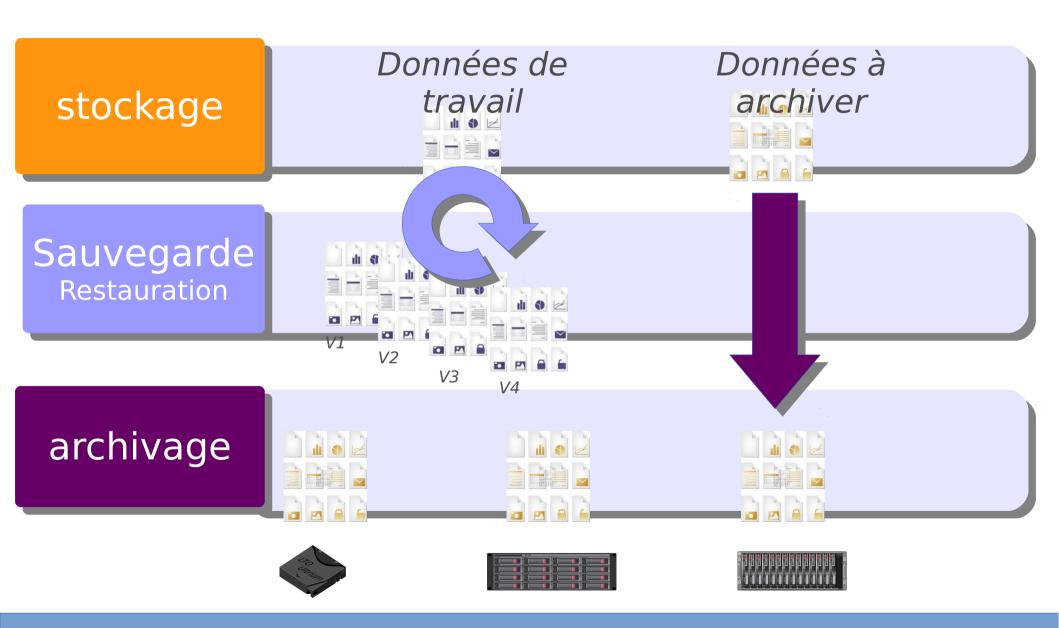
#### Sauvegarde Restauration

- Conserver plusieurs versions d'un fichier du stockage
- Processus automatique (système)
- Restauration sécurisée et autonome (utilisateur et/ou admin)

#### Archivage

- Mettre en sécurité (plusieurs copies) d'un fichier d'archive
- Processus manuel (utilisateur)
- Restitution sécurisée et autonome (utilisateur et/ou admin)

## Stockage / Sauvegarde / Archivage



# L'archivage numérique pérenne

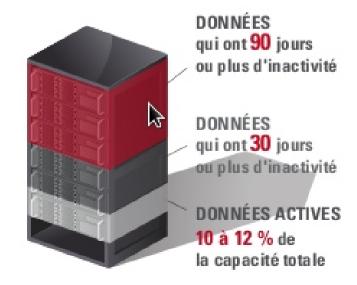
a pour but

- Conserver le document
- Le rendre accessible
- En préserver l'intelligibilité

30 ans

ce n'est pas

- Une sauvegarde
- Un service de HSM
- La dernière étape avant l'oubli ou la perte définitive des données



# Pourquoi est-il si difficile d'archiver?

















## Les difficultés

## les solutions

- l'obsolescence / panne matérielle,
  - Supports: Bandes, disques
  - Lecteurs : bandes (version LTO)
- l'obsolescence logicielle,
  - Changement version, arrêt maintenance
  - disparition société donc logiciel
- l'obsolescence du format de fichier,
  - Format propriétaire, local
  - Dépendance système disparu
- la perte de la signification du contenu,
  - Fichier seul, nom approximatif
- l'humain
  - J'ai pas l'temps de trier ...
  - Ah non j'peux pas effacer!

- Plusieurs copies, plusieurs sites
- Rafraichir supports et migration
- Veille technologique, choix standards
- Virtualisation, émulation
- Formats « durables » et normes
- conversion
- Décrire -> Métadonnées
- Information, formation, motivation
- Pas de solution ...

## Les difficultés

### les solutions

- l'obsolescence / panne matérielle,
  - Supports: Bandes, disques
  - Lecteurs : bandes (version LTO)

- Plusieurs copies, plusieurs sites
- Rafraichir supports et migration

• |'o'
• (

choix standards

Sur quels critères puis-je agir ?

» et normes

on

- Fichier seul, nom approximatif
- Décrire -> Métadonnées

• l'humain

• l'o

la

- J'ai pas l'temps de trier ...
- Ah non j'peux pas effacer!

- Information, formation, motivation
- Pas de solution ...

## Les difficultés

### les solutions

- l'obsolescence / panne matérielle,
  - Supports: Bandes, disques
  - Lecteurs : bandes (version LTO)



- Changement version, arrêt maintenance
- disparition société donc logiciel
- l'obsolescence du format de fichier,
  - Format propriétaire, local
  - Dépendance système disparu



• Rafraichir supports et migration Progiciels ou outils standards

Veille technologique, choix standards
Virtualisation, émulation

Travail « métier » en amont rormats « aurapies » et normes conversion



- la perte de la signification du contenu,
  - Fichier seul, nom approximatif



Décrire -> Métadonnées

Moyens coercitifs pour la saisie

- l'humain
  - J'ai pas l'temps de trier ...
  - Ah non j'peux pas effacer!



Information, formation, motivation

Best effort!

















# Cahier charges (fin 2013)

- Proposer un service d'archivage pour des données scientifiques à l'OMP
  - Types : divers et variés
  - Impact minimum sur les machines
  - Plusieurs dizaines utilisateurs
  - 100 et 200 To + 30 To /an

Travail réalisé en collaboration avec le LEGOS (B. Buisson)

# Cahier charges (fin 2013)

- Proposer un service d'archivage pour des données scientifiques à l'OMP
  - Types : divers et variés
  - Impact minimum sur les machines
  - Plusieurs dizaines utilisateurs
  - 100 et 200 To + 30 To /an

Travail réalisé en collaboration avec le LEGOS (B. Buisson)

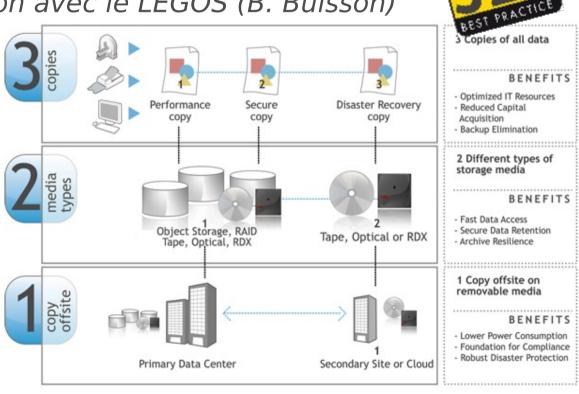
- Pistes explorées :
  - Externalisation
  - Solution interne

# Cahier charges (fin 2013)

- Proposer un service d'archivage pour des données scientifiques à l'OMP
  - Types : divers et variés
  - Impact minimum sur les machines
  - Plusieurs dizaines utilisateurs
  - 100 et 200 To + 30 To /an

Travail réalisé en collaboration avec le LEGOS (B. Buisson)

- Pistes explorées :
  - Externalisation
  - Solution interne



### Externalisation

Le C.I.N.E.S. (Centre Informatique National de l'Enseignement Supérieur)

- trois missions stratégiques nationales :
- le calcul numérique intensif,
- l'archivage pérenne de données électroniques,







# Externalisation Liste des formats archivables

Liste des formats archivables par la plateforme PAC :

Format	Version	Encodage	Identifiant Pronom	Туре МІМ
MPEG-4	ALL	AAC	fmt/199	-
AIFF	ALL	PCM	-	audio/x-aif
FLAC	ALL	FLAC	fmt/279	-
GIF	87a	-	fmt/3	image/gif
GIF	89a	-	fmt/4	image/gif
GEOTIFF	-	-	fmt/155	image/tiff
HTML	3.2	-	fmt/98	text/html
HTML	4.0	-	ftm/99	text/html
HTML	4.01	-	ftm/100	text/html
JPEG	Raw JPEG	-	fmt/41	image/jpeg
JPEG	1.0	-	fmt/42	image/jpeg
JPEG	1.01	-	fmt/43	image/jpeg
JPEG	1.02	-	fmt/44	image/jpeg
JPEG2000	-	-	fmt/151	-
MPEG-4	ALL	AVC/AAC	fmt/199	-
MPEG-4	ALL	AVC	fmt/199	-
MKV	ALL	AVC/FLAC	fmt/569	-

OGG         ALL         THEORA/ VORBIS           ODT         1.0         -           ODT         1.1         -           ODT         1.2         -         fmt/291           PDF         1.4         -         fmt/18           PDF         1.5         -         fmt/19           PDF         1.6         -         fmt/20           PDF/A         1b         -         fmt/354           PDF/A         1a         -         fmt/95           PDF         1.7         -         fmt/276           PNG         1.0         -         fmt/11           PNG         1.0         -         fmt/12           PNG         1.1         -         fmt/13           SVG         1.0         -         fmt/91           SVG         1.1         -         fmt/92           TIFF         -         -         fmt/353           TXT         ALL         UTF-8         -           WAVE         ALL         Wave/PCM,data         fmt/6, fmt/141           XML         1.1         -         fmt/101			. and	
ODT 1.1 - fmt/291  PDF 1.4 - fmt/18  PDF 1.5 - fmt/19  PDF 1.6 - fmt/20  PDF/A 1b - fmt/354  PDF/A 1a - fmt/95  PDF 1.7 - fmt/276  PNG 1.0 - fmt/11  PNG 1.1 - fmt/12  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF - fmt/353  TXT ALL UTF-8  VMAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	ogg	ALL	E383	
ODT 1.2 - fmt/291  PDF 1.4 - fmt/18  PDF 1.5 - fmt/19  PDF 1.6 - fmt/20  PDF/A 1b - fmt/354  PDF/A 1a - fmt/95  PDF 1.7 - fmt/276  PNG 1.0 - fmt/11  PNG 1.1 - fmt/12  PNG 1.2 - fmt/91  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF fmt/353  TXT ALL UTF-8 -  WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	ODT	1.0	-	
PDF 1.4 - fmt/18  PDF 1.5 - fmt/19  PDF 1.6 - fmt/20  PDF/A 1b - fmt/354  PDF/A 1a - fmt/95  PDF 1.7 - fmt/276  PNG 1.0 - fmt/11  PNG 1.1 - fmt/12  PNG 1.2 - fmt/13  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF - fmt/353  TXT ALL UTF-8 -  WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	ODT	1.1	-	
PDF 1.5 - fmt/19  PDF 1.6 - fmt/20  PDF/A 1b - fmt/354  PDF/A 1a - fmt/95  PDF 1.7 - fmt/276  PNG 1.0 - fmt/11  PNG 1.1 - fmt/12  PNG 1.2 - fmt/13  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF fmt/353  TXT ALL UTF-8 - WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	ODT	1.2	-	fmt/291
PDF 1.6 - fmt/20  PDF/A 1b - fmt/354  PDF/A 1a - fmt/95  PDF 1.7 - fmt/276  PNG 1.0 - fmt/11  PNG 1.1 - fmt/12  PNG 1.2 - fmt/13  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF fmt/353  TXT ALL UTF-8 -  WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	PDF	1.4	-	fmt/18
PDF/A 1b - fmt/354  PDF/A 1a - fmt/95  PDF 1.7 - fmt/276  PNG 1.0 - fmt/11  PNG 1.1 - fmt/12  PNG 1.2 - fmt/13  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF fmt/353  TXT ALL UTF-8 - WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	PDF	1.5	-	fmt/19
PDF/A 1a - fmt/95  PDF 1.7 - fmt/276  PNG 1.0 - fmt/11  PNG 1.1 - fmt/12  PNG 1.2 - fmt/13  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF fmt/353  TXT ALL UTF-8 - WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	PDF	1.6	-	fmt/20
PDF 1.7 - fmt/276  PNG 1.0 - fmt/11  PNG 1.1 - fmt/12  PNG 1.2 - fmt/13  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF fmt/353  TXT ALL UTF-8 - WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	PDF/A	1b	-	fmt/354
PNG 1.0 - fmt/11  PNG 1.1 - fmt/12  PNG 1.2 - fmt/13  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF fmt/353  TXT ALL UTF-8 - WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	PDF/A	1a	-	fmt/95
PNG 1.1 - fmt/12  PNG 1.2 - fmt/13  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF fmt/353  TXT ALL UTF-8 - WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	PDF	1.7	-	fmt/276
PNG 1.2 - fmt/13  SVG 1.0 - fmt/91  SVG 1.1 - fmt/92  TIFF fmt/353  TXT ALL UTF-8 - WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	PNG	1.0	-	fmt/11
SVG       1.0       -       fmt/91         SVG       1.1       -       fmt/92         TIFF       -       -       fmt/353         TXT       ALL       UTF-8       -         WAVE       ALL       Wave/PCM,data       fmt/6, fmt/141         XML       1.0       -       fmt/101	PNG	1.1	-	fmt/12
SVG         1.1         -         fmt/92           TIFF         -         -         fmt/353           TXT         ALL         UTF-8         -           WAVE         ALL         Wave/PCM,data         fmt/6, fmt/141           XML         1.0         -         fmt/101	PNG	1.2	-	fmt/13
TIFF fmt/353  TXT ALL UTF-8	svg	1.0	-	fmt/91
TXT ALL UTF-8 -  WAVE ALL Wave/PCM,data fmt/6, fmt/141  XML 1.0 - fmt/101	SVG	1.1	-	fmt/92
WAVE         ALL         Wave/PCM,data         fmt/6, fmt/141           XML         1.0         -         fmt/101	TIFF	-	-	fmt/353
XML 1.0 - fmt/101	TXT	ALL	UTF-8	-
	WAVE	ALL	Wave/PCM,data	fmt/6, fmt/141
XML 1.1 - fmt/101	XML	1.0	-	fmt/101
	XML	1.1	-	fmt/101

### Externalisation

Le C.I.N.E.S. (Centre Informatique National de l'Enseignement Supérieur)

- trois missions stratégiques nationales :
- le calcul numérique intensif,
- · l'archivage pérenne de données électroniques,



• l'hébergement de plates-formes informatiques d'envergure nationale

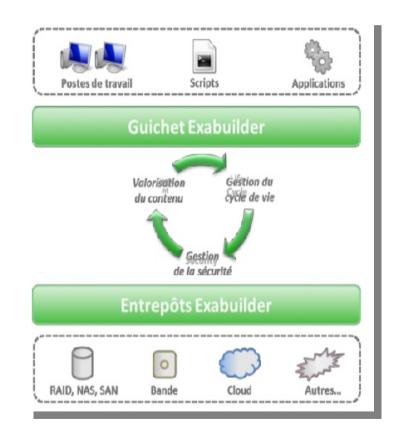
Une très bonne solution si nos formats sont supportés ...

Un exemple : FITS (Flexible Image Transport System) est le format de fichiers le plus communément utilisé en astronomie.

-> Une piste à explorer et un travail à faire sur les formats

- Exabuilder
- ADA (ASG),
- TINA module archivage (ASG),
- Serveur FTP sécurisé.

- Exabuilder
- Dédié archivage, entrepôts sur plusieurs supports avec migration automatique
- Clients : légers (java), lourds (windows LINUX, Mac OS)
- Erreur stratégie marketing -> liquidation judiciaire depuis juillet 2014 (reprise en cours ...)
- ADA (ASG),
- TINA module archivage (ASG),
- Serveur FTP sécurisé.

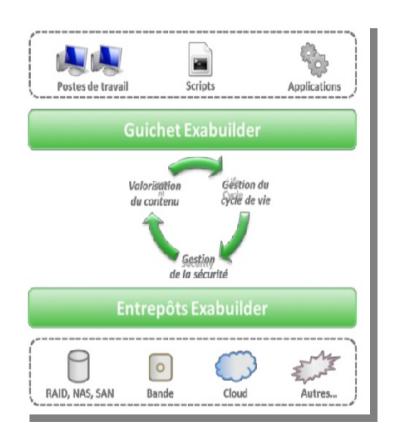


- Exabuilder
- Dédié archivage, entrepôts sur plusieurs supports avec migration automatique
- Clients : légers (java), lourds (windows LINUX, Mac OS)
- Erreur stratégie marketing -> liquidation judiciaire depuis juillet 2014
- ADA (ASG),
- TINA module archivage (ASG),
- Serveur FTP sécurisé.

Levée de fonds en cours

- 400 et 500 k€/an pour Exabuilder
- D. Vinay
  - Investisseurs pour reprise

Renaissance 2016?



- Exabuilder
- Dédié archivage, entrepôts sur plusieurs supports avec migration automatique
- Clients : légers (java), lourds (windows LINUX, Mac OS)
- Erreur stratégie marketing -> liquidation judiciaire depuis juillet 2014 (reprise
- ADA (ASG Digital Archive),
- Archivage et HSM + « stubbing »,
- Performant mais lié aux serveurs de fichiers du LAN
- Clients : légers (java) ou lourds (windows LINUX, Mac OS)
- Pas adapté besoin et prix
- TINA module archivage (ASG),
- Serveur FTP sécurisé.

#### ADA SERVEUR SYSTEMES D'EXPLOITATION :

 Microsoft Windows, Linux, UNIX

#### INTERFACE UTILISATEUR SYSTEMES D'EXPLOITATION :

 Microsoft Windows, Mac OS, Java pour UNIX et platesformes Linux

#### STOCKAGES ET MEDIAS SUPPORTES :

Disque, Bande, Objet, Cloud, incluant :

Serveur de fichiers Windows, Mac OS, Linux & UNIX, HDS HCP, EMC Centera, NetApp, Amazon S3, EMC Atmos, DataDirect Networks WOS, Caringo Swarm, Scality RING, ObjectMatrix, et autres

#### NAS SUPPORTANT LE STUBBING :

 NetApp FAS series, EMC Celerra, HDS FAP 3000 et 4000 series

Exabuilder

Pour 70 To 20 K€ + 10 %

- Dédié archivage, entrepôts sur plusieurs supports avec migration automatique
- Clients : légers (java), lourds (windows LINUX, Mac OS)
- Erreur stratégie marketing -> liquidation judiciaire depuis juillet 2014 (reprise en cours ...)
- ADA (ASG Digital Archive),
- Archivage et HSM + « stubbing »,
- Performant mais lié aux serveurs de fichiers du LAN
- Clients: légers (java) ou lourds (windows LINUX, Mac OS)
- Pas de connexion avec le LAN et prix
- •TINA module archivage (ASG),
- Archivage, mots clés, dossiers d'archives
- Ancien, pas facile à utiliser
- Clients: lourds (windows LINUX, Mac OS)
- Inclus dans la licence TINA POC en cours de déploiement
- •Serveur FTP sécurisé.
  - Permettre à chacun d'archiver ses données sans agent
- Client ftps + authentification LDAP + pureftpd
- Couplé avec archivage TINA

30 K€ + 10 %

Licence TINA 20 K€ + 10 %

0 K€

# Finalement, un serveur FTP et un processus d'archivage...













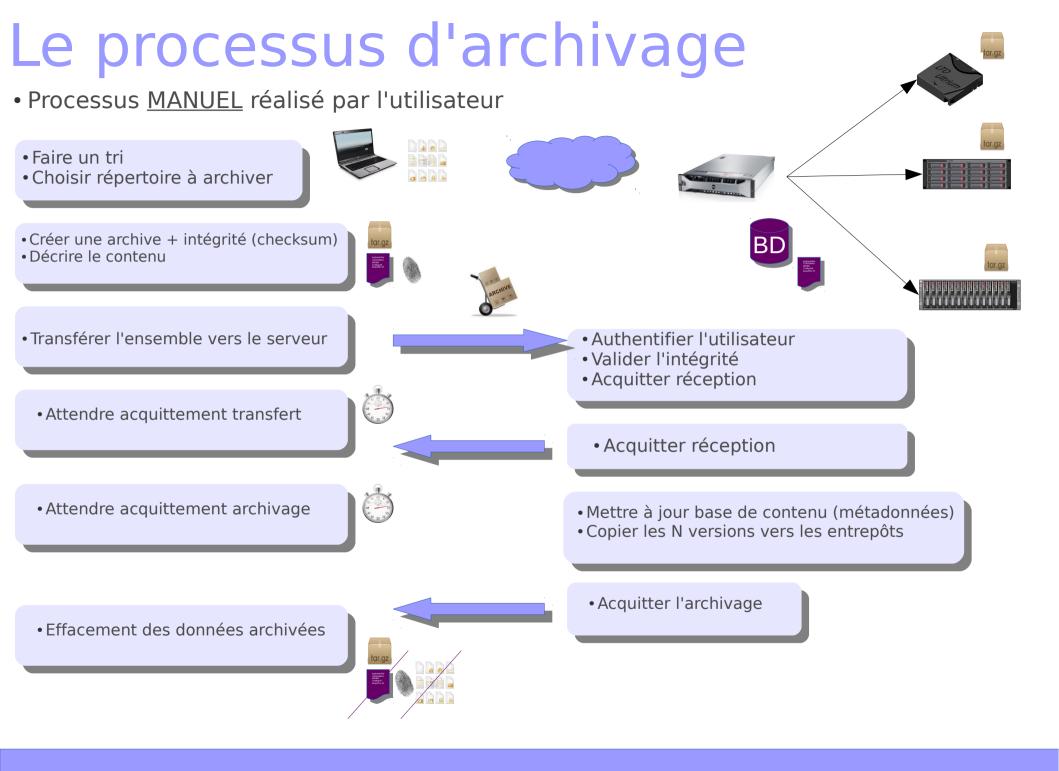


## Les métadonnées

- carte d'identité d'un document
  - identifier, décrire, expliquer l'origine de sa création, son utilité et ses destinataires
- Trois types
  - gestion, pour accéder au document,
  - description, pour en comprendre le contenu,
  - préservation, pour garantir la pérennité de l'accès et de la compréhension du document.
- 15 éléments de description fondamentaux (Dublin Core) sont :

title - format http://dublincore.org/
creator - source
subject - language
description - relation
publisher - coverage
contributor - rights
type - type

+ descripteurs spécifiques « métiers »



## Conclusion















### Pour archiver ...

- Nous devons
  - parler « l'archivistique », comprendre l'OAIS et les métadonnées,
  - conserver le document,
  - le rendre accessible,
  - en préserver l'intelligibilité.
- Lutter contre
  - l'obsolescence (matériel, logiciel, format),
  - la perte de la signification du contenu
  - nos utilisateurs (parfois)



comme l'archivage n'est pas « vendeur » nous n'avons pas ou peu de moyen

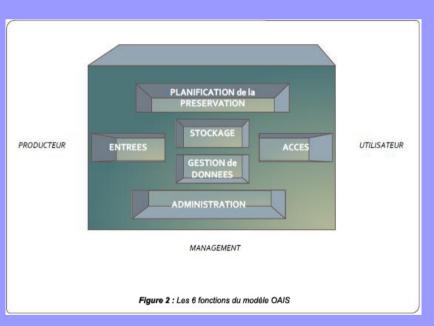


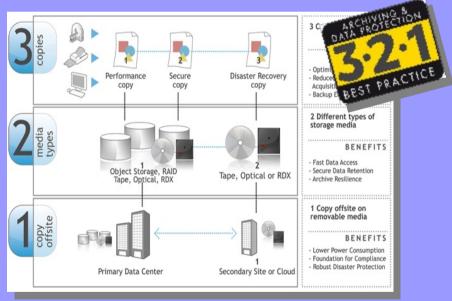






# Pourquoi est-il si difficile d'archiver?





Philippe Saby Capitoul 26 février 2015

de Archivistique à l'informatique l'archivage Numérique pérenne

Les difficultés Solutions explorées

Processus d'archivage

Conclusion





