

# La cellule Formats

## Axes de travail et premiers résultats

#PINFormats

# Composition de la cellule et programme de travail 2019-2020

Sous-groupe « Identification des expertises »

Présentation Dominique Naud (Siaf)

- Constitution d'un annuaire des expertises en France et à l'étranger

Sous-groupe « Connaissance des formats et critères d'obsolescence et de pérennisation »

Présentation Emeline Levasseur (AN)

- Recensement des formats existants, notamment les formats spécifiques
- Définition de critères d'obsolescence et de pérennité pour les formats

Sous-groupe « Outils et Corpus »

Présentation Thomas Ledoux (BNF)

- Recensement des corpus de fichiers librement réutilisables
- Utilisation de Wikidata pour recenser les outils

Sous-groupe « Traduction »

Présentation Edouard Vasseur (ENC)

- Traduction de la grille d'évaluation de la NDSA
- Traduction du Rapid Assessment Model du DPC
- Traduction du Handbook sur la préservation numérique élaboré par le DPC

# Identification des expertises

Base de  
connaissance

Outils

Stratégie de  
conservation

Coopération

Composé du  
Ministère des  
Armées  
(MINARM), du  
Commissariat à  
l'énergie  
atomique (CEA)  
et du Service  
interministériel  
des Archives de  
France (SIAF)

- Constitution d'un annuaire des expertises en France et à l'étranger

# Evolution 2019-2020

Base de  
connaissance

Outils

Stratégie de  
conservation

Coopération



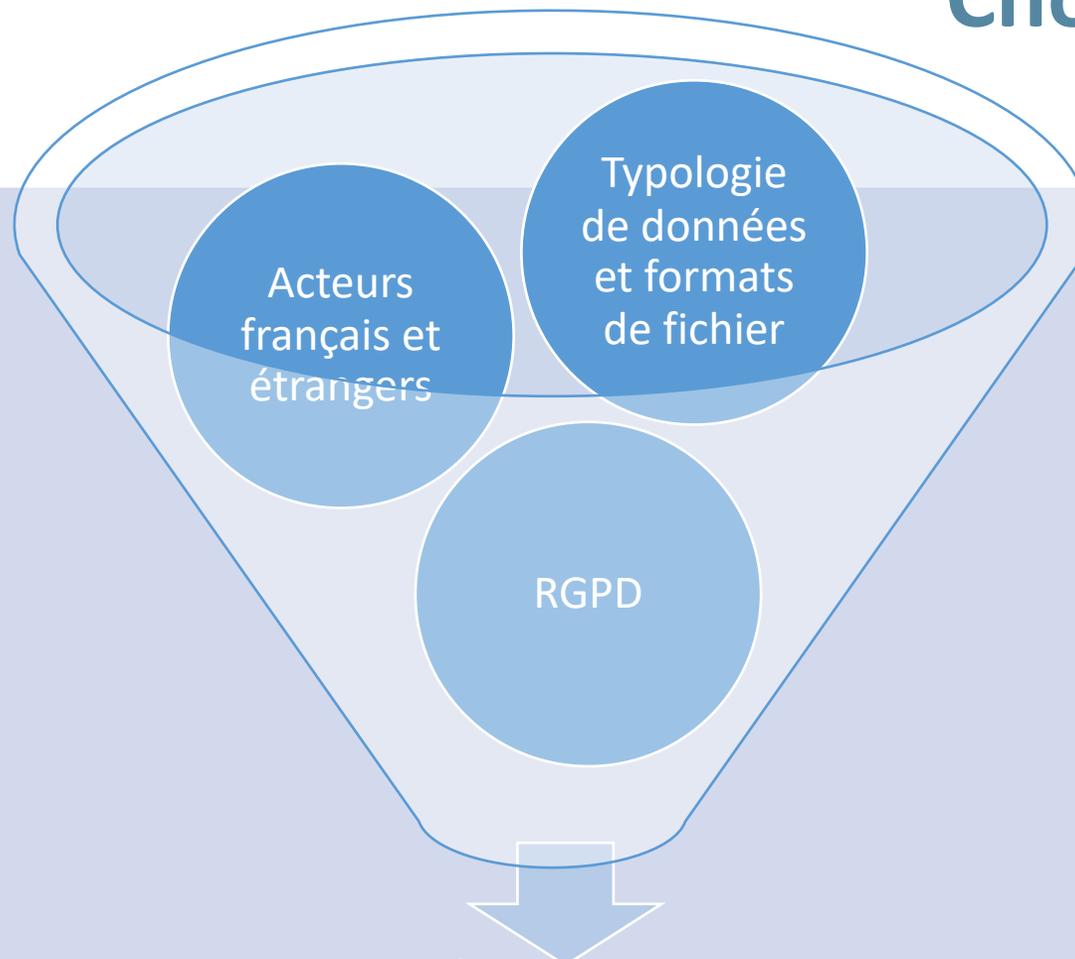
# Choix de publication

Base de  
connaissance

Outils

Stratégie de  
conservation

Coopération



**1<sup>ère</sup> version**

**Acteurs français du secteur public**  
**Sujet d'expertise = type de contenu**

Base de  
connaissance

Outils

Stratégie de  
conservation

Coopération

## GLOSSAIRE

*Pour faciliter votre circulation dans cet annuaire et vous permettre d'affiner votre recherche, les sujets d'expertises ont été regroupés selon la classification COPTR, dont voici ci-dessous leur équivalent en français.*

Pour toute contribution, et lorsque cela est possible, nous vous invitons à vous servir des listes déroulantes pour améliorer l'uniformité de nos données.

Content type	Type de contenu				
Audio	Audio				
Binary Data	Données binaires				
Container	Conteneur				
Database	Base de données				
Disk Image	Image disque				
Document	Document				
EBook	Livres électroniques				
Email	Courriel				
Geospatial	Géospatial				
Image	Image				
Project Management Data	Gestion de projet				
Research Data	Données de la recherche				
Software	Logiciel				
Spreadsheet	Feuille de calcul				
Video	Vidéo				
Web	Web				
Not content type specific	Hors classement				

Base de  
connaissance

Outils

Stratégie de  
conservation

Coopération

# 1<sup>ère</sup> version publiable

Acteurs (personnes physiques ou morales) ▾	[Pays] ▾	Lieu d'exercice actuel (depuis le) ▾	Content type ▾	Précision du sujet ▾	Contact ▾	Références/sources ▾	Public/privé
Centre Informatique National de l'Enseignement Supérieur	France ▾	Centre Informatique National de l'Enseignement Supérieur	Image	TIFF/GeoTIFF, PNG, JPEG, JPEG2000, 3D (.ply, .dae)	<a href="mailto:svp@cines.fr">svp@cines.fr</a>	<a href="https://www.cines.fr/archivage/">https://www.cines.fr/archivage/</a> <a href="https://alfresco.cines.fr/alfresco/faces/jsp/browse/browse.jsp">https://alfresco.cines.fr/alfresco/faces/jsp/browse/browse.jsp</a>	Personne morale secteur public
Institut national de l'audiovisuel	France	Institut national de l'audiovisuel	Video	Tout type de format	<a href="mailto:expertises@ina.fr">expertises@ina.fr</a>	<a href="https://www.ina.fr/">https://www.ina.fr/</a>	Personne morale

Base de connaissance

Outils

Stratégie de conservation

Coopération

# 1<sup>ère</sup> version publiable

Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Document	PDF/A
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Software	Design (.agdb .anf .axd)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Web	.html .htm .mht
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Software	programmation (.cpp), codes sources (.f)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Software	développement (.inc .list), implémentation (.m)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Document	feuilles de style (.css)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Document	mise en page (.fm)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Image	2D/3D (.igs)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Software	Unix (.unix) Archive Unix (.tar)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Document	LaTeX (.tex)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Research Data	HDF
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Geospatial	FITS
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Image	CAO : STEP (.step .stp .p21)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Image	CAO : .cate
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Image	
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Building information specific	building information modeling (BIM) : formats ICF
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Not content type specific	Imagerie médicale : format DICOM (.dcm)
Expertise recherchée	<input type="checkbox"/>	Expertise recherchée	Research Data	données brutes scientifiques (applications LIMS)

**NOUS AVONS BESOIN DE VOUS !**

Base de  
connaissance

Outils

Stratégie de  
conservation

Coopération

# Et après ?

Acteurs (personnes  
physiques ou morales) ▾

[Pays] ▾

Lieu d'exercice actuel  
(depuis le) ▾

Content type ▾

Précision du sujet ▾

Contact ▾

Références/sources ▾

Public/privé

## Les acteurs des autres pays

Base de  
connaissance

Outils

Stratégie  
de conser-  
vation

Coopération

# Connaissance des formats et critères d'obsolescence et de pérennisation

Composé des Archives nationales, du ministère de l'Europe et des Affaires étrangères, du Ministère des Armées, de la Bibliothèque nationale de France et du département de la Moselle.

- Contextes et enjeux très différents parmi les participants.
- Plusieurs constats : difficultés de définition, difficultés à formuler des recommandations.



## Questions

- Quels sont les critères de pérennisation adoptés par d'autres institutions ?
- Quelle synthèse peut-on en proposer ?

Projet de publication sur HAL.

# Critères de pérennisation

## Réponse

**Définir une politique formats : les neuf critères essentiels.**

CELLULE NATIONALE DE VEILLE SUR LES FORMATS

Sous-groupe « connaissance des formats existants et émergents et définition de critères d'obsolescence et de pérennisation »

Avril 2020

CCSD HAL Episciences.org Sciencesconf.org Support fr Humbert Marion

Accueil Dépôt Consultation Recherche Documentation Mon espace

Mon espace / Mes dépôts

Documents en attente de vérification 1

Identifiant	Référence	Date de dépôt
hal-02983527, v1	Bertrand Caron, Martine Sin Blima-Barru, Émeline Levasseur, Erwann Ramondenc, Isabelle Josse, et al.. Définir une politique formats : les neuf critères essentiels.. 2020. (hal-02983527)	2020-10-30



# Connaissance des formats conservés

## Réponse

## Questions

- Quels formats sont conservés au sein de nos établissements ?

Vers un recensement des formats conservés : échanges, réflexions, documents de travail

CELLULE NATIONALE DE VEILLE SUR LES FORMATS

Sous-groupe « connaissance des formats existants et émergents et définition de critères d'obsolescence et de pérennisation »

Février 2020

DOMAINES (prendre la catégorisation PRONOM)	Intitulé du format	Extension	Type MIME	PUID	Wikidata	Caractéristiques	Outil particulier pour les manipuler au sein de l'institution	outil pour générer un fichier (en fonction des contextes institutionnels) <i>facultatif</i>	Nombre d'objets par institution	Volume en % par institution	Contexte de production (éléments de compréhension propre à chaque institution) <i>facultatif</i>	Commentaire sur l'entrée <i>facultatif</i>
	OpenDocument Text (ODF)	.odt				Ouvert						
	Portable Document Format (Adobe Systems)	.pdf				Propriétaire ouvert						



## Questions

- Quelles sont les caractéristiques d'un format précis ?

## Connaissance technique des formats

### Réponse



Travaux en cours : TIFF, Mp3, formats bureautiques, messageries.

# Composition de la cellule et programme de travail 2019-2020



## Sous-groupe « Outils et Corpus »

Composé de la BnF, du CINES, du TGIR HumaNum, de Mintika, du programme VITAM

- Recensement des corpus de fichiers librement réutilisables
- Utilisation de Wikidata pour recenser les outils

Base de connaissance

Outils

Stratégie de conservation

Coopération

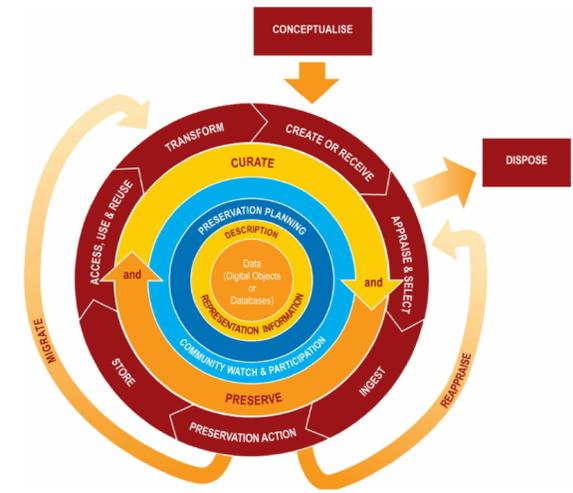
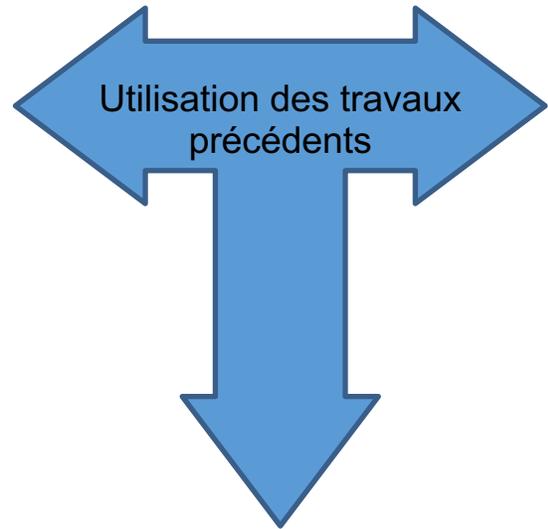
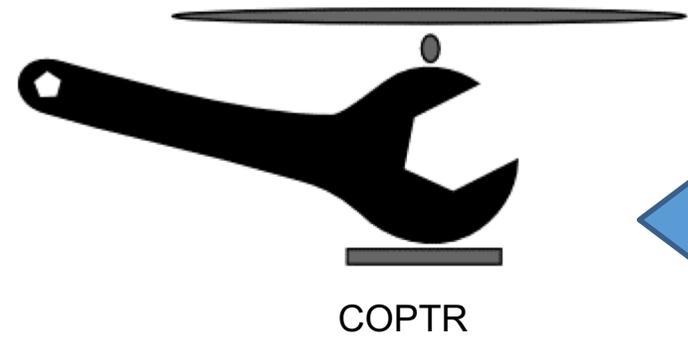
- Première version
- 12 corpus génériques repérés

# Recensement des corpus de fichiers

	A	B	C	D	E	F	G	H
1	Recensement de corpus de fichiers librement utilisables							
2	Code	Nom	Organisation	URL	Type	Volumétrie	Licence	Repérage des fichiers valides ou non
3								
4	CF_01	Format corpus	OPF	<a href="https://github.com/openpreserve/format-corpus">https://github.com/openpreserve/format-corpus</a>	Tous		CC0	?
5	CF_02	veraPDF-Corpus	OPF, DualLabs	<a href="https://github.com/veraPDF/veraPDF-corpus">https://github.com/veraPDF/veraPDF-corpus</a>	PDF		CC-BY 4.0	Oui
6	CF_03	Isartor	PDF Association	<a href="https://www.pdfa.org/resource/isartor-test-suite/">https://www.pdfa.org/resource/isartor-test-suite/</a>	PDF/A		Freely downloadable and usable without restriction	Oui
7	CF_04	Bavaria	PDFLib	<a href="https://github.com/bfosupport/pdfa-testsuite">https://github.com/bfosupport/pdfa-testsuite</a>	PDF/A		Creative Commons Public License	
8	CF_05	Jpylyzer test Suite	OPF	<a href="https://github.com/openpreserve/jpylyzer-test-files">https://github.com/openpreserve/jpylyzer-test-files</a>	JP2000		CC-BY	Oui
9	CF_06	Google image test suite		<a href="https://code.google.com/archive/p/imagetestsuite/">https://code.google.com/archive/p/imagetestsuite/</a>	Images : TIFF, PNG, GIF, JPEG	143 Mo	CC-BY 3.0	



# Registre des outils de préservation



Le modèle du DCC sur le cycle de vie de la préservation numérique



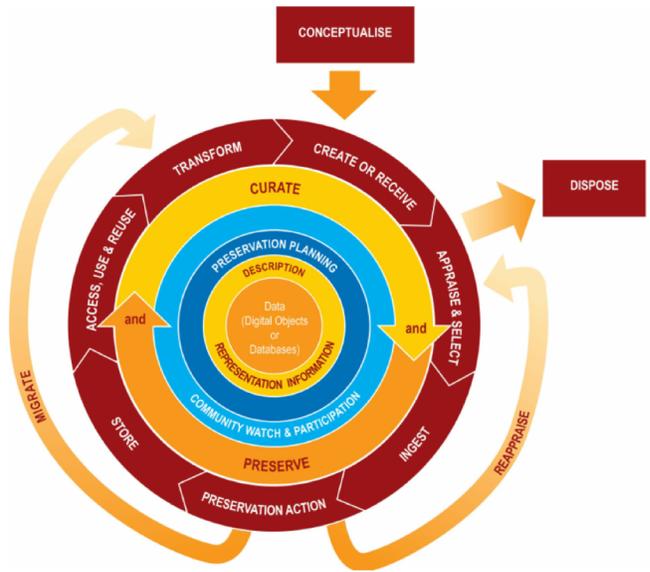
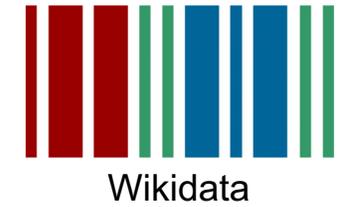
Wikidata



Wikipedias



# Etape 1 : créer les catégories du registre



Cycle de vie du DCC

Chaque étape du cycle de vie est modélisé par une entité WD

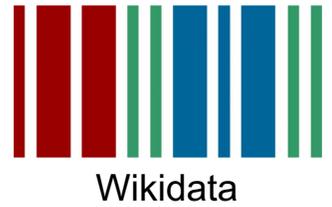
P31 (instance of - *nature de*) : Q714737 (category - *catégorie*)  
 P361 (part of - *partie de*) : COPTR  
 P1545 (series ordinal - *rang dans la série*) = 1

rang	qid	etape	stage
1	<a href="#">Q90878156</a>	Accéder, Utiliser et Réutiliser	Access, Use and Reuse
2	<a href="#">Q90887043</a>	Créer ou Réceptionner (Acquérir)	Create or Receive (Acquire)
3	<a href="#">Q89021580</a>	Agir transversalement	Cross-Lifecycle Functions
4	<a href="#">Q90889758</a>	Éliminer	Dispose
5	<a href="#">Q90885460</a>	Verser	Ingest
6	<a href="#">Q90890255</a>	Agir pour préserver	Preservation Action
7	<a href="#">Q90890401</a>	Planifier la préservation	Preservation Planning
8	<a href="#">Q90890485</a>	Stocker	Store

<https://w.wiki/N5r>



# Etape 2 : lier les processus aux catégories



Category: [Discussion](#) Read Edit View history  Go Search

Category: Create or Receive (Acquire)

**Subcategories**  
This category has the following 7 subcategories, out of 7 total.

**D**

- Data capture and Deposit
- Disk Imaging

**F**

- File Copy

**O**

- OCR

**W**

- Web Crawl
- Web Snapshot

**W cont.**

- Workflow and Lab Notebook Management

**Lifecycle stage definition:** Functions that support the DCC Lifecycle Stage defined as "Create data including administrative, descriptive, structural and technical metadata. Preservation metadata may also be added at the time of creation. Receive data, in accordance with documented collecting policies, from data creators, other archives, repositories or data centres, and if required assign appropriate metadata."

of - facette de) : ?etape

**NOUS AVONS BESOIN DE VOUS !**

	process	processus
741970	file copying	Copie de fichier
4929239	data collection	collecte de données
id:Q61466324	web crawling	exploration du Web
wd:Q167555	optical character recognition	reconnaissance optique de caractères

Lier les processus aux catégories

**Activités:**

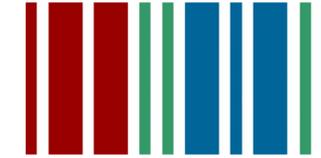
- Chercher une entité WD préexistante appropriée
- En créer une nouvelle, le cas échéant

Constat : Très souvent une entité représentant le résultat du processus existe mais pas le processus lui-même.  
Par exemple : "disk image file" ([Q592312](#)) au lieu de "disk imaging"

<https://w.wiki/N6M>



# Etape 3 : donner accès aux outils



Wikidata

P31 (instance of - *nature de*)/P279(subclass - *sous-classe*)\* Q7397 (software)

P366 (use - *usage*) : ?processus

<https://w.wiki/S4q>

stage	etape	process	processus	tool	toolLabel
Ingest	Verser	file format identification	identification du format de fichier	<a href="#">Q wd:Q24305922</a>	Digital Record Object IDentification
Ingest	Verser	file format identification	identification du format de fichier	<a href="#">Q wd:Q16927990</a>	JHOVE
Ingest	Verser	file format identification	identification du format de fichier	<a href="#">Q wd:Q2858088</a>	Apache Tika
Ingest	Verser	file format identification	identification du format de fichier	<a href="#">Q wd:Q181820</a>	file
Ingest	Verser	data integrity	intégrité	<a href="#">Q wd:Q6802825</a>	Md5deep
Ingest	Verser	data integrity	intégrité	<a href="#">Q wd:Q1932013</a>	Md5sum

**Les outils sont classifiés**  
**De nouveaux outils sont identifiés et découverts**

Base de  
connaissance

Outils

Stratégie de  
conservation

Coopé-  
ration

## Sous-groupe « Traduction »

Composé de la BnF, du  
CD Moselle, de l'ENC, de  
Mintika et du Ministère  
des Armées (SHD)

- Traduction de la grille d'évaluation de la NDSA
- Traduction du Rapid Assessment Model du DPC
- Traduction du Handbook sur la préservation numérique élaboré par le DPC

Base de connaissance

Outils

Stratégie de conservation

Coopération

# Document 1 : le NDSA Level of Preservation



Domaine fonctionnel	Niveaux			
	Niveau 1 (connaître vos contenus)	Niveau 2 (protéger vos contenus)	Niveau 3 (surveiller vos contenus)	Niveau 4 (pérenniser vos contenus)
<b>Stockage</b>	<p>Posséder deux copies complètes dans des lieux distincts.</p> <p>Documenter tous les supports de stockage où les contenus sont stockés.</p> <p>Utiliser des supports de stockage stables.</p>	<p>Posséder trois copies complètes avec au moins une copie à un emplacement géographique distinct.</p> <p>Documenter le stockage et les supports de stockage en indiquant les ressources et dépendances nécessaires à leur fonctionnement.</p>	<p>Posséder au moins une copie à un emplacement géographique présentant un type de menace différent de ceux des autres emplacements.</p> <p>Posséder au moins une copie sur un support de stockage différent.</p> <p>Surveiller l'obsolescence du stockage et des supports.</p>	<p>Posséder au moins trois copies dans des emplacements géographiques présentant des types de menaces différents.</p> <p>Augmenter la variété des supports de stockage pour éviter les points de défaillance uniques.</p> <p>Avoir un plan et mener des actions pour remédier à l'obsolescence des supports de stockage, des logiciels et du matériel informatique.</p>
<b>Intégrité</b>	<p>Vérifier l'information d'intégrité si celle-ci a été fournie avec les contenus.</p> <p>Générer une information d'intégrité si aucune information n'est disponible.</p> <p>Contrôler la présence de virus. Le cas échéant, mettre les contenus en quarantaine.</p>	<p>Vérifier l'information d'intégrité lors de la migration ou de la copie des contenus.</p> <p>Utiliser des bloqueurs d'écriture lors des travaux sur les supports originaux.</p> <p>Sauvegarder l'information d'intégrité et stocker la copie dans un emplacement distinct de celui des contenus.</p>	<p>Vérifier l'information d'intégrité à intervalles réguliers.</p> <p>Documenter les processus et les résultats des vérifications de l'information d'intégrité.</p> <p>Mener des audits d'intégrité à la demande.</p>	<p>Vérifier l'information d'intégrité à la suite d'événements ou d'activités spécifiques.</p> <p>Remplacer ou réparer les contenus corrompus le cas échéant.</p>
<b>Contrôle</b>	<p>Déterminer les agents humains et logiciels autorisés à lire, écrire, mettre à jour et supprimer les contenus.</p>	<p>Documenter les droits de lecture, d'écriture, de mise à jour et de suppression des agents humains et logiciels.</p>	<p>Identifier les agents humains et logiciels qui mènent des actions sur les contenus et journaliser ces actions.</p>	<p>Examiner périodiquement les journaux des opérations et des accès.</p>
<b>Métadonnées</b>	<p>Créer un inventaire des contenus. Y documenter les emplacements utilisés pour le stockage.</p> <p>Sauvegarder cet inventaire et en conserver au moins une copie à part des contenus eux-mêmes.</p>	<p>Stocker suffisamment de métadonnées pour connaître les contenus (possibilité de combiner les métadonnées administratives, techniques, descriptives, de préservation et structurelles).</p>	<p>Déterminer quel standard de métadonnées appliquer.</p> <p>Trouver et combler les lacunes dans les métadonnées pour se conformer à ces standards.</p>	<p>Archiver les actions de préservation associées au contenu et les occurrences de ces actions.</p> <p>Choisir et implémenter des standards de métadonnées.</p>
<b>Contenu</b>	<p>Documenter les formats de fichiers et toutes les autres propriétés essentielles (<i>significant properties</i>) des contenus, y compris les modalités et la date d'acquisition de cette documentation.</p>	<p>Vérifier les formats de fichiers et les autres propriétés essentielles (<i>significant properties</i>) des contenus.</p> <p>Développer des relations avec les créateurs de contenus pour encourager des choix de formats de fichiers durables.</p>	<p>Surveiller l'obsolescence et les évolutions des technologies dont dépendent les contenus.</p>	<p>Mener des opérations de migration, de normalisation, d'émulation, etc. pour s'assurer que les contenus restent accessibles.</p>

**Grille d'évaluation des niveaux de préservation avec les domaines fonctionnels en ligne**

<https://hal-bnf.archives-ouvertes.fr/hal-02552208v1> et <https://hal-bnf.archives-ouvertes.fr/hal-02551807v1>

Base de  
connaissance

Outils

Stratégie de  
conservation

Coopé-  
ration

# Document 2 : le Rapid Assessment Model de DPC

Grille d'évaluation rapide de la Digital Preservation Coalition

## Annexe 1 – tableur DPC RAM

<b>Organisation</b>	
<b>Responsable de l'évaluation</b>	
<b>Date de l'évaluation</b>	
<b>Observations sur le périmètre de l'évaluation</b>	

Niveau stratégique				
	Niveau actuel	Justification	Niveau cible	Observations
<b>A - Viabilité de l'organisation :</b> Gouvernance, structure organisationnelle, dotation en personnel et en ressources des activités de préservation numérique				
<b>B - Politique et stratégie :</b> Politiques, stratégies et procédures qui régissent le fonctionnement et la gestion des archives numériques				



Base de  
connaissance

Outils

Stratégie de  
conservation

Coopé-  
ration

# Document 3 : le Handbook de DPC



Digital**Preservation**Coalition

Digital Preservation **Handbook**

## Explore the Handbook

[Home](#)

[Contents](#)

[Introduction](#)

[Digital preservation briefing](#)

[Getting started](#)

[Institutional strategies](#)

[Organisational activities](#)

[Technical solutions and tools](#)



En cours

## Digital Preservation Handbook

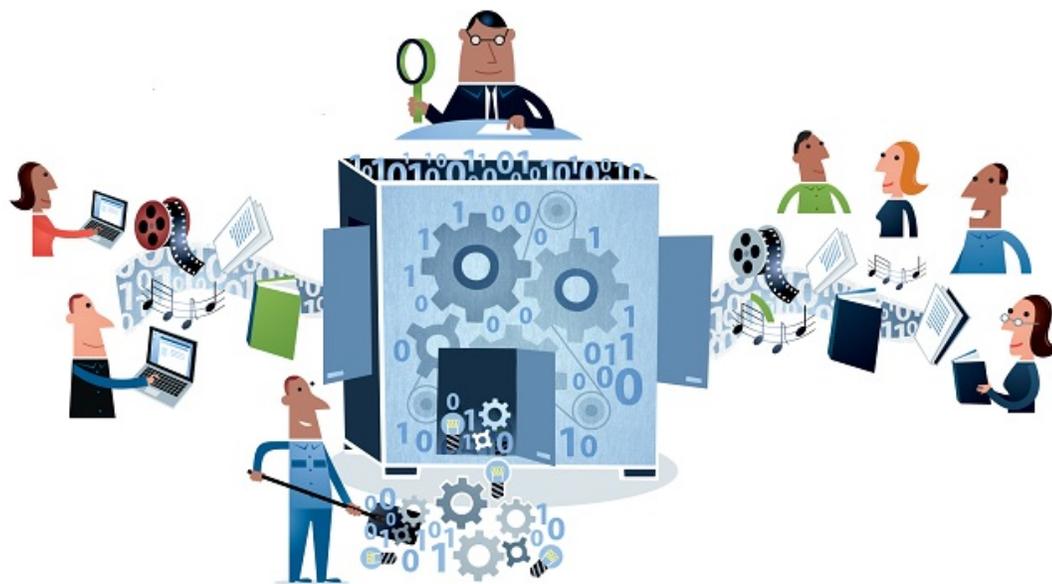


Illustration by Jørgen Stamp digitalbevaring.dk CC BY 2.5 Denmark

Welcome to the revised 2nd edition of the Digital Preservation Handbook. A key knowledge base for

**Merci de votre attention !**

**Contact :** [archivage.numerique.siaf@culture.gouv.fr](mailto:archivage.numerique.siaf@culture.gouv.fr)

LES FORMATS, C'EST LONG,  
C'EST TECHNIQUE, ÇA  
N'INTÉRESSE PERSONNE...  
MAIS ÇA PEUT SAUVER VOS  
DOCUMENTS ET DONNÉES !

