

ANF Huma-Num

Qualité des données, comprendre les enjeux des formats de fichiers

Michel Jacobson IR* Huma-Num

8-décembre-2023

La notion de format de fichier

Point terminologique

La norme ISO-14721:2012 OAIS (Open Archival Information System) distingue :

- **l'information**: connaissance que l'on peut échanger ;
- **la donnée**: représentation formalisée de la connaissance.



Pour manipuler une information, on lui donne une forme (la donnée)

Cette forme suit des règles qui permettent son codage et son décodage.

Un **format de fichier** constitue un ensemble de règles qui décrivent comment coder les informations sous forme de données et comment décoder celles-ci pour retrouver les informations. Ces règles sont généralement consignées dans des « spécifications » et sont mises en œuvre par des outils logiciels (éditeurs, convertisseurs...) pour permettre de manipuler l'information.

Point terminologique

- Les données DOIVENT être accompagnées de **métadonnées** afin de comprendre et réutiliser les informations qu'elles contiennent.

L'identification du format est une de ces métadonnées

- Mais les métadonnées sont aussi des informations, qu'il faut coder, formater et décrire...

- Les formats sont très liés au marché et plus globalement à l'ensemble de l'environnement numérique. Ils sont donc soumis à l'obsolescence plus ou moins rapidement.

Le choix des formats cibles

- Quels sont les enjeux/risques ?
- Comment faire le choix dans le cadre du projet ?
- Comment atteindre les formats choisis ?
- Comment vérifier que les données sont bien formatées ?

Quels sont les risques ?

Dépendance

ex.

Un format fermé lisible que par un logiciel propriétaire d'un éditeur dont on ne connaît pas la feuille de route

Non maîtrise de la qualité des données

ex.

Impossibilité de valider par un outil autre que celui de l'éditeur. Tout repose sur la confiance...

Perte de l'information

ex.

Un format propriétaire dont l'éditeur ne maintient plus les outils d'édition/lecture et qui ne fonctionne plus avec la nouvelle version de MS-Windows

Comment évaluer les risques ?

1) Des critères basés sur les spécifications

- Explicites ou non
- Secrètes / Publiques / Standardisées / Normalisées
- Avec ou sans entrave légale (brevet, copyright...)
- Qualité des spécifications
 - Complexité
 - Auditabilité

2) Des critères liés au marché

- Existence d'un outil d'édition et d'alternatives
- Existence d'outils de conversion
- Existence d'outils de validation
- Mais aussi évaluation des risques sur les logiciels (nombre d'utilisateurs, dynamisme des communautés d'utilisateurs, dépendances logicielles, fréquence des mises-à-jour...)

« Guide méthodologique pour le choix des formats numériques pérennes dans un contexte adaptés aux données orales et audiovisuelles » https://francearchives.fr/fr/circulaire/DGP_SIAF_2010_010

Spécifications

Absence de spécifications ou spécifications implicites

ex. SQL définit un langage d'interrogation pour base de données mais pas de format de représentation

Spécifications non publiées

ex. les anciens formats de Microsoft (doc, ppt...)

Spécifications standardisées/normalisées

ex.  OpenXML (pptx, docx...)

 XML

 CSV

Complexité 1

Taille des spécifications

CSV 8 pages

XML 35 pages

OpenXML + de 6000 pages

Complexité 2

Dépendances

ex.

La recommandation « **Extensible Markup Language (XML)** »

- est liée à d'autres recommandations « Namespaces in XML », « XML Inclusions (XInclude) », « Associating Stylesheets with XML », « Associating Schemas with XML »
- basé sur un format texte mais ne spécifie pas ce qu'est un format texte
- ne spécifie pas l'encodage des caractères mais juste comment l'explicitier, + sa valeur par défaut (UTF8), + les encodages que tous les outils doivent permettre (UTF-8 and UTF-16), mais les outils peuvent en proposer d'autres
- ne spécifie pas les langues mais juste la norme à utiliser pour les exprimer

ex.

Le format **OpenDocumentFormat** utilise une enveloppe ZIP contenant une arborescence de fichiers essentiellement formatés en XML

Complexité 3

Variabilité

ex. Le format **WAV** est un format conteneur qui permet de stocker des données sonores de différents encodages (PCM, ADPCM, μ -law...)

On peut avoir un fichier WAV avec un codex mp3. Ce qui n'est pas la même chose qu'un format mp3 dont le codec est forcément mp3

ex. Un format **MP4** est un format conteneur qui peut contenir des pistes audio, des pistes vidéo, des sous-titre, des métadonnées, etc.

On peut avoir un fichier MP4 qui ne contient que de l'audio ou qui ne contient aucune piste audio ni vidéo

Silence

ex. Un format **CSV** ne dit rien de l'encodage des caractères qu'il utilise

Connaître l'organisation du conteneur ne veut pas dire maîtriser ses contenus

Initiatives d'évaluation des risques

Preservation Risk matrix



NARA (Archives du gouvernement fédéral des États-Unis)

- ✓ https://github.com/usnationalarchives/digital-preservation/tree/master/Digital_Preservation_Risk_Matrix
- ✓ Grille d'évaluation pour déterminer le niveau de risque que représente l'utilisation d'un format
- ✓ Synthèse du risque en 3 niveaux : haut, moyen, bas
- ✓ Une cinquantaine de critères
 - ✓ Le format est-il propriétaire?
 - ✓ Les spécifications sont-elles publiées?
 - ✓ Existe-t-il des outils de validation du respect des spécifications par un fichier?
 - ✓ Les spécifications ont-elles été validées par un organisme de normalisation?
 - ✓ Le format est-il utilisé dans le "federal government"?
 - ✓ Le format est-il utilisé à l'extérieur du "federal government"?
 - ✓ Le format est-il activement maintenu et mis à jour par une organisation, une personne individuelle ou une communauté?
 - ✓ Existe-t-il plusieurs outils de visualisation?
 - ✓ [...]
 - ✓ Le format permet-il l'encapsulation d'information comme le watermarking?

Preservation Risk matrix

	A	B	C	D	E	F	G	H	I	J
1	NARA Guidance : Preferred	NARA Guidance: Acceptable	Numeric Risk Rating	Risk Level	NARA Format ID	Format Name	File Extension(s)	Category/Plan(s)	Is the format proprietary?	Does the format have a published open specification?
11			-26,00	High Risk	NF00182	Executable file	exe	Software and Code	-1	-1
12			-25,00	High Risk	NF00702	BlackBerry Binary Executable	cod	Software and Code	-1	-1
17			-25,00	High Risk	NF00473	Nikon RAW (NRW)	nrw	Digital Still Image	-1	-1
18			-25,00	High Risk	NF00474	Olympus RAW	orf	Digital Still Image	-1	-1
19			-25,00	High Risk	NF00356	PCPaint Image	pic; clp	Digital Still Image	-1	-1
20			-25,00	High Risk	NF00475	Pentax RAW	pef; ptx	Digital Still Image	-1	-1
21			-25,00	High Risk	NF00476	Sigma RAW	x3f	Digital Still Image	-1	-1
22			-25,00	High Risk	NF00477	Sony RAW	srf	Digital Still Image	-1	-1
634	X		45,00	Low Risk	NF00650	Broadcast Wave (BWF) unspecified version	wav	Digital Audio	-1	2
635	X		45,00	Low Risk	NF00561	eXtensible Markup Language 1.1	xml	Web Records; Software	-1	2
636			45,00	Low Risk	NF00646	Portable Document Format/Archiving (PDF/A-4)	pdf	Presentation and Publication	2	2
637	X		45,00	Low Risk	NF00406	SIARD 1.0	siard	Databases	-1	2
647			47,00	Low Risk	NF00618	SIARD 2.0	siard	Databases	-1	2
648	X	X	48,00	Low Risk	NF00602	Portable Document Format/Archiving (PDF/A-2a) accessible	pdf	Presentation and Publication	2	2
649	X	X	48,00	Low Risk	NF00634	Portable Document Format/Archiving (PDF/A-2b) basic	pdf	Presentation and Publication	2	2
650			48,00	Low Risk	NF00642	Portable Document Format/Archiving (PDF/A-2u) unicode	pdf	Presentation and Publication	2	2
651			49,00	Low Risk	NF00619	SIARD 2.1.1	siard	Databases	-1	2
652										

Liste des formats archivables au CINES

Liste des formats validables

⚠ Attention : le validateur de formats permet de valider certains formats qui ne sont pas pris en charge par la plateforme d'archivage du CINES.

- <https://facile.cines.fr/>

Format	Nom	PRONOM PIUD	Type MIME	Commentaire	Archivable dans PAC
AAC AAC	Advanced Audio Codings	[fmt/199]		Format Mpeg-4 contenant uniquement un flux audio au format AAC.	✓
AIFF PCM	Audio Interchange File Format	[fmt/414]	[audio/x-aif, audio/x-aiff]	Format audio contenant uniquement un flux PCM.	✓
APNG	Animated Portable Network Graphics	[fmt/935]	[image/vnd.mozilla.apng, image/apng]	L'APNG est une extension du format PNG permettant de réaliser des animations graphiques.	✗
DAE UTF-8 1.4.1	Collada		[application/xml]	Format permettant de stocker des données géométriques sous forme de scènes (plusieurs objets combinés dans le même référentiel), et d'y ajouter des informations supplémentaires pour décrire la scène et les objets (matériaux, environnement lumineux, animations, ...) ou pour ajouter des notions sémantiques (relations entre les objets, découpage d'un objet en plusieurs éléments fonctionnels, etc...).	✗
FLAC FLAC 1.2.1	Free Lossless Audio Codec	[fmt/279]	[audio/ogg, audio/x-flac]	Format audio compressé sans perte.	✓
GIF 87a	Graphics Interchange Format	[fmt/3]	[image/gif]	Format image pouvant contenir également des animations.	✓
GIF 89a	Graphics Interchange Format	[fmt/4]	[image/gif]	Format image pouvant contenir également des animations.	✓
GeoTIFF	Geographic Tagged Image File Format	[fmt/155]	[image/tiff]	Format dérivé du TIFF contenant des informations de géoréférencement et de géolocalisation.	✓
HDF5 1.0	Hierarchical Data Format	[fmt/286]		Format de données à caractère scientifique.	✗
HDF5 2.0	Hierarchical Data Format	[fmt/287]		Format de données à caractère scientifique.	✗

Gestion des fichiers

Des formats riches en FAIR



Pour qu'un fichier de données soit lisible et réutilisable

- 1) Son format doit être identifié explicitement et précisément
- 2) Son format doit présenter le moins possible de risque (multicritères)
- 3) Sa forme doit être vérifiée et valide (conforme aux spécifications des formats et codages utilisés)
- 4) La donnée dans ce format doit être dans un circuit de maintenance qui surveille l'obsolescence

Quelques bonnes pratiques de gestion

nommage des fichiers

rangement des fichiers

éviter/repérer de doublons

Objectifs de gestion

Le nommage des fichiers et l'organisation des dossiers sont deux opérations importantes pour

- faciliter l'accès aux données (rôle des extensions, du plan de classement)
- faciliter le tri des documents
- faciliter la compréhension de l'ensemble des documents du projet
- éviter la perte et la duplication erronée de fichiers

Convention de nommage

Exemple 1 : Doranum* :

<https://doranum.fr/stockage-archivage/comment-nommer-fichiers/>

- Donner un nom bref et explicite
- Ne pas mettre d'espace ni de caractères spéciaux
- Indiquer les dates au bon format
- Placer l'élément important en premier
- Indiquer les versions des documents

Convention de nommage

Exemple 2 : Alain Rivet, Marie-Laure Bachèlerie, Auriane Denis-Meyere et Delphine Tisserand « Traçabilité des activités de recherche et gestion des connaissances : Guide pratique de mise en place », 2018

https://qualite-en-recherche.cnrs.fr/wp-content/uploads/2021/08/guide_tracabilite_activites_recherche_gestion_connaissances.pdf

Les intitulés des fichiers :

- doivent être succincts et précis ;
- doivent être uniques ;
- peuvent être caractérisés a minima par une date, un sujet et un type de document ;
- ne doivent pas excéder 31 caractères (extension comprise).

Éviter :

- les accents (é, è ë, ê), cédille (ç), caractères spéciaux (, ; . : ! ? * »%...@ &) ;
- les mots vides : le, la, les, un, une, des, et, ou... ;
- les dénominations vagues, par exemple « divers », « autres », « à classer » ;
- les espaces que l'on peut remplacer par des « _ » (tiret 8).

Convention de nommage

Quelques monstres rencontrés

« vidéo de l'entretien de Monsieur Pierre Martin en 2002 à son domicile du 43 rue Paul Doumer.mp4 »

- Des accents : ne vont pas passer indemne à travers tous les systèmes
- Des espaces : sans précaution certains logiciels vont considérer qu'il y a autant de fichiers que de mots
- Mélange minuscules majuscules : risque de confusion
- Long : risque de troncation
- Contient des métadonnées : c'est un rôle pris en charge par d'autres outils que le simple nommage des fichiers
- Contient des informations à caractère personnel

Mais aussi :

✓ « 马孝珍采访音频.doc » → « é© ¬å?ç??é??è®¿é?³éç?.doc »

✓ « toto.bin »

✓ Des caractères retour chariot en fin de nom

✓ Des caractères BOM « byte order mark » en milieu de mot. Le BOM est un caractère invisible 'espace insécable sans chasse' U+FEFF ou U+FFFE qui indique l'ordre des octets pour composer les caractères multi-octets

Quelques outils de traitement de fichiers de données

- ✓ Outils d'identification de format
- ✓ Outils de validation du formatage
- ✓ Outils de conversion de format

Identification des format des fichiers

Comment identifier un format ?

- Quels indices pour « deviner » le format d'un fichier ?
 - Son extension (.jpg, .wav, .pdf...)
 - Un outil de lecture
 - Les métadonnées embarquées
 - Des outils spécialisés : DROID, FIDO, file,...
- Où l'indiquer ?
 - Dans le DMP ; dans les métadonnées
- Comment l'exprimer ?
 - En utilisant des identifiants de référentiels connus (type-mime, PRONOM)
- Pourquoi ?
 - Guider la bonne lecture des données

Un exemple d'outil d'identification



« Digital Record Object Identification » DROID

- Logiciel d'identification de format de fichiers (Java, New BSD License)
- Lié au registre de formats de fichier PRONOM
- Éditeur : Les Archives nationales du Royaume-Uni

Un vrac numérique



1archive-yann-bi
siou_02-42097.jp
g



3-IMG202111131
21429.jpg



06-07_Accueil_Ire
sco.pdf



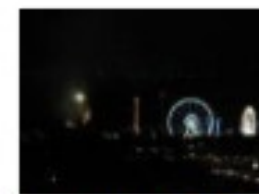
7savigny_50_01-0
cadcdESSIN
JEUNE.jpg



14_01_zoom.jpg



14_02_zoom.jpg



20100501_FoireD
uTrone_003.JPG



20100501_FoireD
uTrone_018.JPG



20100501_FoireD
uTrone_021.png



20210920_104035.
jpg



20210920_104035.
pdf



20210920_104107.
jpg



20210920_104107.
pdf



9327631989_9d4e
e03ed1_o.jpg



1636737300376.jp
g



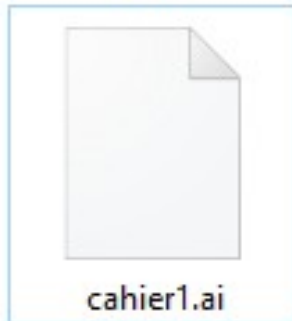
1636737320330.jp
g



AAC_Table-ronde
_APRE.pdf



Adieu aux
Fantastiques.jpg



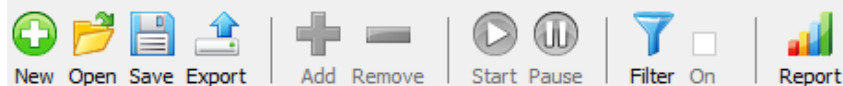
cahier1.ai

Rapport DROID

DROID v6.5



File Edit Run Filter Report Tools Help



Untitled-2 x Untitled-3 x Untitled-4 x Untitled-5 x

Resource	Extension	Format	Version	Mime type	△ PUID
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\cahier.1.ai	ai ⚠	Acrobat PDF 1.4 - Portable Document Format	1.4	application/pdf	fmt/18
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\20210920_104035.pdf	pdf	Acrobat PDF 1.5 - Portable Document Format	1.5	application/pdf	fmt/19
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\20210920_104107.pdf	pdf	Acrobat PDF 1.5 - Portable Document Format	1.5	application/pdf	fmt/19
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\06-07_Accueil_Iresco.pdf	pdf	Acrobat PDF 1.5 - Portable Document Format	1.5	application/pdf	fmt/19
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\AAC_Table-ronde_APRE.pdf	pdf	Acrobat PDF 1.6 - Portable Document Format	1.6	application/pdf	fmt/20
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\20210920_104035.jpg	jpg	JPEG File Interchange Format	1.01	image/jpeg	fmt/43
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\20210920_104107.jpg	jpg	JPEG File Interchange Format	1.01	image/jpeg	fmt/43
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\9327631989_9d4ee03ed1_o.jpg	jpg	JPEG File Interchange Format	1.01	image/jpeg	fmt/43
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\1636737300376.jpg	jpg	JPEG File Interchange Format	1.01	image/jpeg	fmt/43
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\1636737320330.jpg	jpg	JPEG File Interchange Format	1.01	image/jpeg	fmt/43
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\Adieu aux Fantastiques.jpg	jpg	JPEG File Interchange Format	1.01	image/jpeg	fmt/43
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\1archive-yann-bisiou_02-42097.jpg	jpg	JPEG File Interchange Format	1.01	image/jpeg	fmt/43
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\3-IMG20211113121429.jpg	jpg	JPEG File Interchange Format	1.01	image/jpeg	fmt/43
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\7savigny_50_01-0cadcdESSIN JEUNE.jpg	jpg	JPEG File Interchange Format	1.01	image/jpeg	fmt/43
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\14_01_zoom.jpg	jpg	JPEG File Interchange Format	1.02	image/jpeg	fmt/44
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\14_02_zoom.jpg	jpg	JPEG File Interchange Format	1.02	image/jpeg	fmt/44
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\20100501_FoireDuTrone_021.png	png ⚠	Exchangeable Image File Format (Compressed)	2.21	image/jpeg	fmt/645
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\20100501_FoireDuTrone_003.JPG	jpg	Exchangeable Image File Format (Compressed)	2.21	image/jpeg	fmt/645
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\data\20100501_FoireDuTrone_018.JPG	jpg	Exchangeable Image File Format (Compressed)	2.21	image/jpeg	fmt/645

Lecture du rapport DROID

- 5 extensions : (pdf, jpg, JPG, ai, png) dont 2 avec warnings
 - le fichier avec l'extension png est en fait un format jpeg
 - l'extension ai (Adobe Illustrator) est en fait un format pdf
- 2 types mime : application/pdf, image/jpeg
- 6 formats :
 - 3 formats PDF : « PDF version 1.4 »(fmt/18) , « PDF version1.5 » (fmt/19) et « PDF version 1.6 » (fmt/20)
 - 3 formats JPEG « JPEG File Interchange Format 1.01 » (fmt/), « JPEG File Interchange Format 1.02 » (fmt/) et « Exchangeable Image File Format (Compressed) 2.21 »

Retour d'expérience sur Nakala

Situation d'octobre 2023 sur un peu plus d'un million de fichiers

- 77 types mime
 - 68 % image/jpeg
 - 81 % image/*
 - 10 % application/pdf
- 186 formats différents au registre PRONOM

Retour d'expérience sur Nakala

- 18 formats de type application/pdf

	A	B	C	D
1	fmt/276	Acrobat PDF 1.7 - Portable Document Format (1.7)	44953	36,16 %
2	fmt/20	Acrobat PDF 1.6 - Portable Document Format (1.6)	40023	32,20 %
3	fmt/95	Acrobat PDF/A - Portable Document Format (1a)	12997	10,46 %
4	fmt/18	Acrobat PDF 1.4 - Portable Document Format (1.4)	10294	8,28 %
5	fmt/17	Acrobat PDF 1.3 - Portable Document Format (1.3)	9370	7,54 %
6	fmt/19	Acrobat PDF 1.5 - Portable Document Format (1.5)	5449	4,38 %
7	fmt/477	Acrobat PDF/A - Portable Document Format (2b)	564	0,45 %
8	fmt/354	Acrobat PDF/A - Portable Document Format (1b)	389	0,31 %
9	fmt/16	Acrobat PDF 1.2 - Portable Document Format (1.2)	143	0,12 %
10	fmt/478	Acrobat PDF/A - Portable Document Format (2u)	41	0,03 %
11	fmt/157	Acrobat PDF/X - Portable Document Format - Exchange 1a:2001	31	0,02 %
12	fmt/488	Acrobat PDF/X - Portable Document Format - Exchange PDF/X-4	15	0,01 %
13	fmt/476	Acrobat PDF/A - Portable Document Format (2a)	13	0,01 %
14	fmt/479	Acrobat PDF/A - Portable Document Format (3a)	9	0,01 %
15	fmt/480	Acrobat PDF/A - Portable Document Format (3b)	5	0,00 %
16	fmt/15	Acrobat PDF 1.1 - Portable Document Format (1.1)	2	0,00 %
17	fmt/158	Acrobat PDF/X - Portable Document Format - Exchange 3:2002	2	0,00 %
18	fmt/14	Acrobat PDF 1.0 - Portable Document Format (1.0)	1	0,00 %

Les registres de formats

Les types MIME

En 1996, le standard de l'IETF, RFC-2046 « Multipurpose Internet Mail Extensions » définit des types de média dans le cadre des échanges mail (SMTP)

L'organisme d'enregistrement de ces types de média est IANA « Internet Assigned Numbers Authority »

- ✓ Les identifiants sont structurés en 2 parties
 - ✓ Types ou grandes familles : application, audio, example, font, image, message, model, multipart, text et video
 - ✓ Sous-types : par ex. audio/mpeg ; audio/aac ; image/png ; text/csv....
- ✓ Ils contribuent par exemple dans Nakala à orienter le choix d'une visionneuse adaptée

Le registre types MIME



Registries included below

- [application](#)
- [audio](#)
- [font](#)
- [example](#)
- [image](#)
- [message](#)
- [model](#)
- [multipart](#)
- [text](#)
- [video](#)

application

Available Formats



CSV

Name	Template	Reference
1d-interleaved-parityfec	application/1d-interleaved-parityfec	[RFC6015]
3gpdash-qoe-report+xml	application/3gpdash-qoe-report+xml	[3GPP][Ozgur_Oyman]
3gppHal+json	application/3gppHal+json	[3GPP][Ulrich_Wiehe]
3gppHalForms+json	application/3gppHalForms+json	[3GPP][Ulrich_Wiehe]
3gpp-ims+xml	application/3gpp-ims+xml	[3GPP][John_M_Meredith]
A2L	application/A2L	[ASAM][Thomas_Thomsen]
ace+cbor	application/ace+cbor	[RFC9200]
ace+json	application/ace+json	[RFC9431]

Le registre PRONOM

<http://www.nationalarchives.gov.uk/pronom/>

- ✓ Éditeur : Les Archives nationales du Royaume-Uni (TNA)
- ✓ Identifie les formats avec un PUID « PRONOM unique identifier
- ✓ Décrit les formats de fichiers avec de nombreuses propriétés
- ✓ Tout le monde peut proposer l'ajout de nouveaux formats par un formulaire. Par ex. suite aux travaux effectués avec le CINES sur la validation de la TEI, ajout de nouveaux formats fmt/1474, fmt/1475, fmt/1476 et fmt/1477
- ✓ Un portail web pour accéder à la documentation structurée sur les formats
- ✓ Un fichier de « signatures » de formats à utiliser dans le cadre de DROID pour l'identification
- ✓ Liens directs par l'outil DROID au registre en ligne de PRONOM par le PUID

Le format « fmt/1 »

Details for: Broadcast WAVE 0 Generic

 Save as... XML | CSV  Print

Go to: [Summary](#) | [Documentation](#) > | [Signatures](#) > | [Compression](#) > | [Character encoding](#) > | [Rights](#) > | [Reference files](#)
> [Properties](#) >

Summary

Name	Broadcast WAVE
Version	0 Generic
Other names	BWAVE (0), BWF (0)
Identifiers	MIME: audio/x-wav PUID: fmt/1
Family	WAVE
Classification	Audio
Disclosure	Full
Description	<p>Broadcast WAVE is a chunk-based audio format developed by the European Broadcasting Union based on the Microsoft WAVE format, which is in turn based on the generic Resource Interchange File Format (RIFF) specification developed by Microsoft and IBM. Structurally, a BWAVE file consists of a number of chunks, each comprising a four character code chunk identifier, the chunk identifier, and the chunk data. It comprises a RIFF header with a WAVE data type identifier, followed by a Broadcast Audio Extension chunk, containing meta-information, a Format chunk, which describes the audio data, and a Data chunk, containing the audio data itself. BWAVE files which contain compressed audio data must also include a Fact chunk, containing file-dependent information. A BWAVE identified as generic by DROID likely uses an encoding other than PRONOM, or perhaps a structural difference, and users are encouraged to let the format maintainers know about this error.</p>

Orientation	Binary
Byte order	Little-endian (Intel)
Related file formats	Has lower priority than Broadcast WAVE (0 PCM Encoding) Has lower priority than Broadcast WAVE (0 MPEG Encoding) Has lower priority than Broadcast WAVE (0 WAVEFORMATEXTENSIBLE Encoding) Has priority over Waveform Audio Has priority over Waveform Audio (PCMWAVEFORMAT) Has priority over Waveform Audio (WAVEFORMATEXTENSIBLE) Is previous version of Broadcast WAVE (1 Generic)
Technical Environment	
Released	01 Jan 1997
Supported until	01 Jul 2001
Format Risk	
Developed by	None.
Supported by	None.
Source	 Digital Preservation Department / The National Archives
Source date	11 Mar 2005
Source description	
Last updated	19 Jul 2013
Note	

Validation du bon formatage des fichiers

Comment tester la validité d'un fichier ?

- La lecture avec son éditeur ne suffit pas
 - Peu fiable : les éditeurs sont souvent très tolérants aux erreurs afin de ne pas bloquer la lecture
 - Problème de « juge et partie »
- Il existe des outils dont c'est la finalité
 - Des outils spécialisés sur un format précis : Par ex. jpylyzer (Open Planets Foundation) pour le format JPEG2000
 - Des outils généralistes : Jhove (JSTOR and the Harvard Library)
 - Des intégrateurs d'outils : facile (CINES), FITS (Harvard Library)

Pourquoi valider ?

- ✓ Un fichier mal formé n'a pas toutes les garanties d'être correctement lu par tous les outils
- ✓ Il est hasardeux d'effectuer des conversions sur des fichiers mal formés : pas de garantie du résultat
- ✓ Les résultats aux tests de validation permet à des services d'archives de guider leur décision de prise en charge ou de niveau de service.
 - ✓ Prise en charge de la conservation de l'information
 - ✓ Prise en charge de l'intégrité de la donnée

Valider les données

Test du fichier « EF circulaire NMPP.jpg » avec Jhove

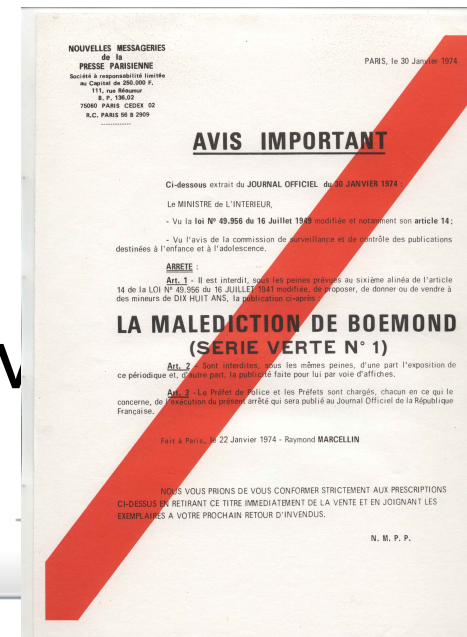


Format : JPEG 1.02
Statut : not well-formed
Message : Unexpected end of file à l'offset 364954 sur un fichier de 370857 octets

ReplInfo

File

- Documents
 - C:\jacobson\Formations\2020s\2021\2021-12-09_ANF_ISOREJeuDeData_Selection_TP-PMEF circulaire NMPP.jpg
 - Module
 - JPEG-hul
 - Release: 1.5.2
 - Date: 5 nov. 2019
 - ReplInfo
 - URI: C:\jacobson\Formations\2020s\2021\2021-12-09_ANF_ISOREJeuDeData_Selection_TP-PMEF circulaire NMPP.jpg
 - LastModified: Tue Nov 30 15:15:12 CET 2021
 - Size: 370857
 - Format: JPEG
 - Version: 1.02
 - Status: Not well-formed
 - SignatureMatches
 - JPEG-hul
 - Messages
 - ErrorMessage: Unexpected end of file
 - ID: JPEG-HUL-2
 - Offset: 364954
 - MimeType: image/jpeg



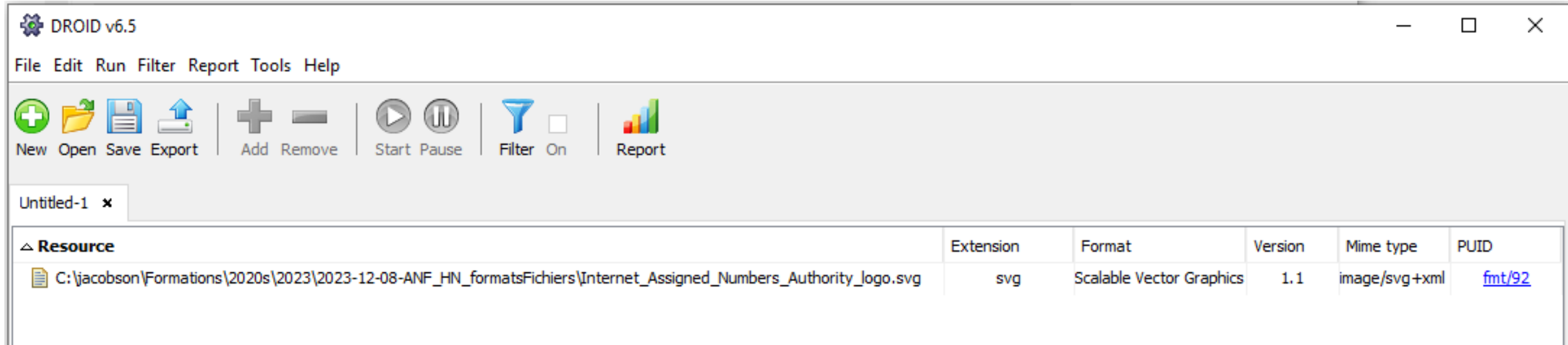
Test sur un fichier iana.svg

La lecture du fichier avec LibreOffice



Test sur un fichier iana.svg

Identification du fichier avec DROID




The screenshot shows the DROID v6.5 application window. The title bar reads "DROID v6.5". The menu bar includes "File", "Edit", "Run", "Filter", "Report", "Tools", and "Help". The toolbar contains icons for "New", "Open", "Save", "Export", "Add", "Remove", "Start", "Pause", "Filter", "On", and "Report". The main area displays a table with the following data:

Resource	Extension	Format	Version	Mime type	PUID
C:\jacobson\Formations\2020s\2023\2023-12-08-ANF_HN_formatsFichiers\Internet_Assigned_Numbers_Authority_logo.svg	svg	Scalable Vector Graphics	1.1	image/svg+xml	fmt/92

Test sur un fichier iana.svg

Validation du fichier avec Markup Validation Service

**Markup Validation Service**
Check the markup (HTML, XHTML, ...) of Web documents

Error found while checking this document as SVG 1.1!

Result:	1 Error	
File:	<input type="text" value="Parcourir..."/> Aucun fichier sélectionné. <small>Use the file selection box above if you wish to re-validate the uploaded file Internet_Assigned_Numbers_Authority_Logo.svg</small>	
Encoding:	utf-8	<input type="button" value="(detect automatically)"/>
Doctype:	SVG 1.1	<input type="button" value="(detect automatically)"/>
Root Element:	svg	
Root Namespace:	http://www.w3.org/2000/svg	

Options
 Show Source Show Outline List Messages Sequentially Group Error Messages by Type
 Validate error pages Verbose Output Clean up Markup with HTML-Tidy
[Help](#) on the options is available.

Validation Output: 1 Error

 **Line 89, Column 9: element "ligne" undefined**

```
<ligne/ >
```

You have used the element named above in your document, but the document type you are using does not define an element of that name. This error is often caused by:

- incorrect use of the "Strict" document type with a document that uses frames (e.g. you must use the "Frameset" document type to get the "<frameset>" element),
- by using vendor proprietary extensions such as "<spacer>" or "<marquee>" (this is usually fixed by using CSS to achieve the desired effect instead),
- by using upper-case tags in XHTML (in XHTML attributes and elements must be all lower-case).

Stratégies mises en œuvre dans des services

Le service d'archives du CINES

CINES (Centre informatique national de l'enseignement supérieur)

- ✓ Un guide méthodologique pour choisir un format
- ✓ Une liste de formats acceptés
- ✓ Un outil de test : Facile



La validation, c'est facile

Exemples de résultats

- Formats non pris en charge
 - jpeg-raw, md et mww, pptx
- 1 fichier pdf non valide
 - Message :
« Nombre d'objet ou de flux d'objet invalide »
- 1 format obsolete
 - PDF v. 1.3

Détails	Fichier	Format identifié	Bien formé	Valide	Archivable dans PAC	Commentaire
🔍	nakala-faq.md		✘	✘	✘	
🔍	Poster_ISIDORE_DH2018.pdf	PDF 1.3	✔	✔	✘	
🔍	Human-numGBverticale.png	PNG 1.0	✔	✔	✔	
🔍	equipe-HN-2021.csv	TXT UTF-8	✔	✔	✔	
🔍	VIGNETTE_NAKALA_et_IIIF.png	PNG 1.2	✔	✔	✔	
🔍	Poster_ISIDORE_DH2018.png	PNG 1.0	✔	✔	✔	
🔍	logo-grand-nakala-rvb.png	PNG 1.2	✔	✔	✔	
🔍	logo_recontre_HN-2026.png	PNG 1.2	✔	✔	✔	
🔍	Human-numGBverticale.jpg	JPEG RAW	✔	✔	✘	
🔍	guide-module-isidore-suggestions.pdf	PDF 1.6	✔	✔	✔	
🔍	ISIDORE_widget_2-0.png	PNG 1.2	✔	✔	✔	
🔍	Human-numGBverticale.ai	PDF 1.5	✔	✔	✔	
🔍	ANF2015-MOREL-C.pdf	PDF 1.5	✘	✘	✘	Corriger automatiquement votre fichier Consulter nos tutoriels pour corriger votre fichier
🔍	lettre_info25-iINSHS_TribuneHumaNum.pdf	PDF 1.4	✔	✔	✔	
🔍	ANF2015-Francart.zip		✘	✘	✘	
🔍	ISIDORE Suggestions_low.wmv		✘	✘	✘	
🔍	HUMA-NUM_PPTX_prez.pptx		✘	✘	✘	
🔍	NAKALA_et_IIIF.mp4	MPEG-4 AVC/AAC LC	✔	✔	✔	

L'entrepôt Nakala



- ✓ Nakala accepte tous les fichiers (= sans contrôle)
 - ✓ Dans tous les cas, les fichiers déposés peuvent être téléchargés (au détail près des droits d'accès)
 - ✓ Certains formats peuvent être visualisés directement sur la page de présentation

Les visionneuses proposées

- ✓ Fichier image : OpenSeadragon (tif, jpg, jp2, png, pdf)
- ✓ Fichier CSV : DataTables
- ✓ Fichier audio ou vidéo : Plyr
- ✓ Fichier PDF : PDF.js
- ✓ Fichier Markdown
- ✓ Archives : .zip, .rar, .phar, .tar, .tgz, .gz, .bz2
- ✓ Fichier de code : XML, HTML, JSON, etc.

Stratégie dans Nakala

✓ Extension non reconnue vs. Extension reconnue et prise en charge

Fichiers

- aaa.truc
- bad.png
- canvas.png
- aaa.csv

Visualisation

Aperçu non disponible

Copier l'ID Copier l'url d'intégration Copier l'url de téléchargement

Fichiers

- aaa.truc
- bad.png
- canvas.png
- aaa.csv

Visualisation

Afficher 250 entrées Chercher

fichiers ↑↓	*dc:publisher ↑↓	*dcterms:created[dcterms:W3CDTF] ↑↓	dc:contributor[olac]
144259.wav	Laboratoire de langues et civilisations à tradition orale	1973	Ozanne-Rivierre, Fra

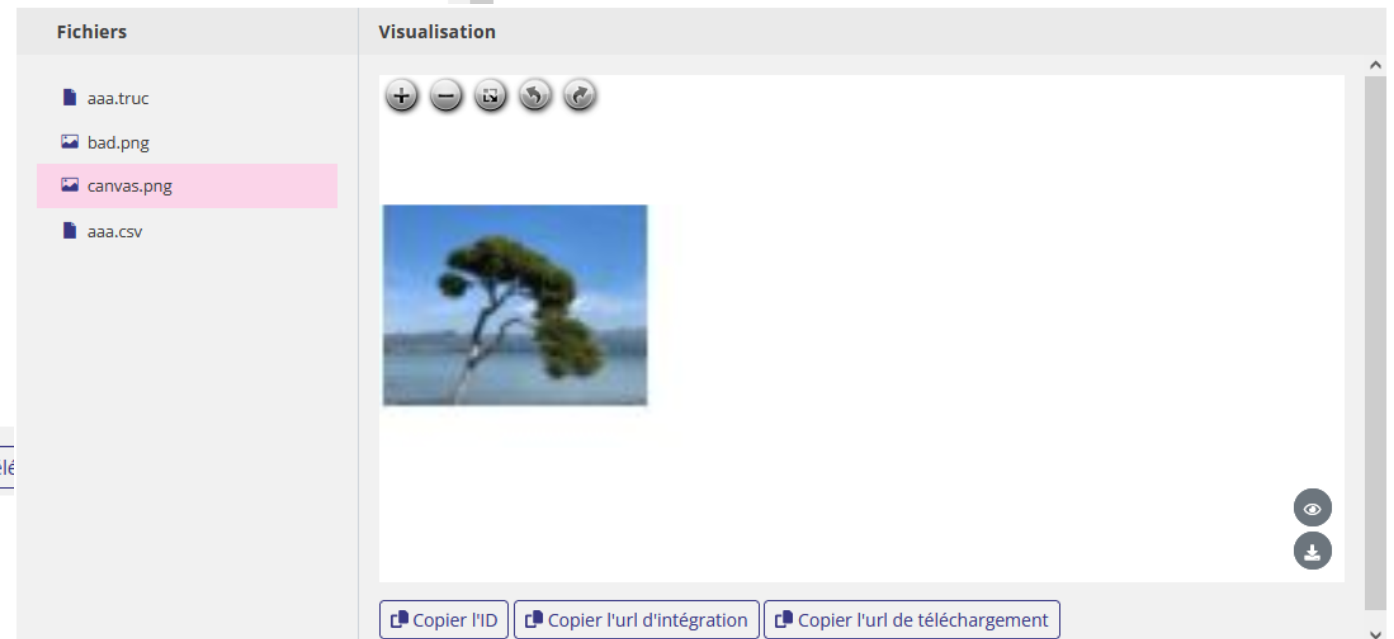
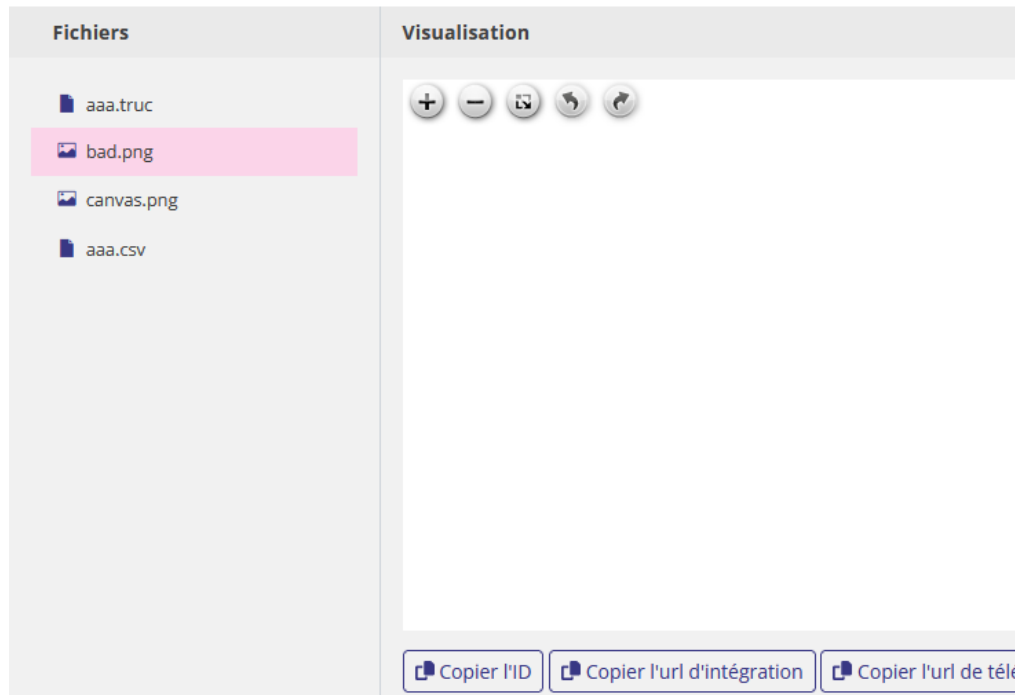
Copier l'ID Copier l'url d'intégration Copier l'url de téléchargement

→ ne mettez pas des extensions exotiques, non-standards

Stratégie de Nakala

✓ Fichier bien formé vs. Fichier mal formé

Le fichier bad.png lu par LibreOffice



→ vérifiez vos données avant de les publier

Stratégie de Nakala

- ✓ Enfin ce n'est pas parce qu'une visionneuse de Nakala affiche correctement votre fichier, que celui-ci est valide
- ✓ Si votre fichier est mal-formé bien que bien-représenté
 - ✓ En cas d'obsolescence vous serez amené à convertir votre fichier pour maintenir sa lisibilité, mais cette conversion sera risquée car votre outil de conversion risque de ne pas savoir lire le fichier ou ne sera pas comment corriger l'anomalie
 - ✓ En cas d'archivage le fichier sera testé et donc rejeté pour cause de malformation

→ vérifiez vos données avant de les publier

Nakala (bonnes pratiques)

- ✓ Que faut-il déposer dans Nakala comme fichier ?
 - ✓ un format peu risqué (cf. matrice des risques de la NARA)
 - ✓ un format pris en charge par le CINES (cf. liste sur facile)
 - ✓ un fichier dont le formatage est valide (cf. facile)
 - ✓ un fichier bien nommé (extension standard)
 - ✓ le fichier le plus qualitatif que vous avez qui répond aux précédents critères (le protocole IIIF pour les images permet par exemple de générer des formats de consultation plus léger)
- ✓ Que faut-il décrire du format d'un fichier dans Nakala
 - ✓ Le système va extraire automatiquement des informations techniques de vos fichiers (l'extension, le type mime, le PUID (bientôt), la taille en octets. Il est donc inutile, risqué et inadapté de redéfinir ces informations dans les métadonnées DublinCore.

Conversion de format de fichiers

un format maintenable

- Pour lutter contre l'obsolescence, pour maîtriser les risques, pour faciliter des usages
 - Trouver un format équivalent pour exprimer toute l'information et permettre de conserver les usages peut s'avérer difficile. Cela se fait parfois avec une perte d'information qu'il faut évaluer
 - Trouver des outils de conversion (il en existe plein) et les qualifier (résultat fidèle et valide)
 - Par ex. ffmpeg pour l'audio-visuel (utilisable à travers Sharedocs)
 - Par ex. ImageMagic pour les formats image
 - ...

Conversion de format pour normalisation WMV → MP4/H264

The screenshot shows a file explorer interface with a green header bar. The left sidebar displays a folder tree under 'Mes fichiers', with 'hnTools_watchFolder' expanded to show 'Video' > 'ffmpeg' > 'toMP4_h264'. The main pane shows a table of files:

Nom...	Etiquette	Taille	Date de modifica...	Image pr...
ISID...	WMV	30.3 MB	17:43	


The file 'ISID...' is highlighted in green. The right pane shows a video player for 'ISIDORE Suggestions_low...' with a black video frame. Below the player, the following properties are listed:

- Type: Windows Media Video
- Taille: 30.3 MB
- Date Création: 17:43
- Evaluation: ☆☆☆☆☆

Audio properties are also shown:


- Durée: 8:12
- Codec: Windows Media Audio V8


Conversion de format pour normalisation WMV → MP4/H264

[HUMA-NUM] Video transcoding is finished  Boîte de réception x



sysadmin@huma-num.fr

 À moi ▾

 anglais ▾ > français ▾ [Traduire le message](#)

HN-Tools : job information

--- INPUT

File name : ISIDORE Suggestions_low.wmv
File path : mjacobson/Video/ffmpeg/toMP4_h264/IN
File size : 30.26mo
File extension : wmv
File type : video/x-ms-asf
File hash : 7f6af8e56647ea7079e21af04dc65dc337572ef254613fbf8005d14e
File date : 1631202208
Submit by : Jacobson Michel (michel.jacobson@gmail.com)

--- OUIPUT

Output file : mjacobson/Video/ffmpeg/toMP4_h264/OUT/ISIDORE Suggestions_low.mp4
Tool : Video
Engine : ffmpeg
Preset 1 : toMP4_h264
Execution time : 2mn 24.191s

Conversion de format pour normalisation WMV → MP4/H264

The screenshot shows a file manager interface with a green header bar. On the left is a sidebar with a tree view of folders: 'Mes fichiers', 'depot', 'test', 'Bibliothèques', 'hnTools_software', 'hnTools_watchFolder', 'Audio', 'OCR', 'PDF', 'Video', 'ffmpeg', 'toMP4_h264', 'IN', 'OUT', and 'toMP4_h264_0720p'. The 'OUT' folder is highlighted. The main pane shows the path 'hnTools_watchFolder > Video > ffmpeg > toMP4_h264 > ...'. A table lists a file 'ISID...' with a green checkmark, 'MP4' format, a lightning bolt icon, '16.4 MB' size, and '17:50' duration. The right pane shows a video player for 'ISIDORE Suggestions_low...' with a black video frame. Below the player, metadata is displayed: 'Type MPEG-4 Video', 'Taille 16.4 MB', 'Date Création 17:50', and 'Evaluation ☆☆☆☆☆'. An 'Add details..' link is at the bottom right of the metadata section.

Nom...	Etiquette	Taille	Date de modifica...	Image pr...
✓ ISID...	MP4 ⚡	16.4 MB	17:50	

Type MPEG-4 Video
Taille 16.4 MB
Date Création 17:50
Evaluation ☆☆☆☆☆
Add details..

LibreOffice

Options PDF

Général | Vue initiale | Interface utilisateur | Liens | Sécurité | Signatures numériques

Plage

Tout

Diapos :

Sélection

Afficher le PDF après export

Images

Compression sans perte

Compression JPEG Qualité :

Réduire la résolution des images

Filigrane

Signer avec un filigrane

Général

PDF hybride (fichier ODF incorporé)

Archive PDF/A, (ISO 19005)

PDF/A version :

PDF marqué (ajouter la structure du document)

Créer un formulaire PDF

Format d'envoi :

Autoriser les doublons de noms de champ

Structure

Exporter les repères de texte

Commentaires en tant qu'annotations PDF

Exporter les pages de notes

Exporter seulement les pages de notes

Exporter les pages masquées

Exporter les pages vides insérées automatiquement

Utiliser les XObjets de référence

Aide

Exporter

Annuler

Conclusion

Le formatage des fichiers numériques n'est pas qu'un enjeu d'interopérabilité

C'est aussi un enjeu de pérennité, dans la mesure où c'est sur lui que se repose le codage de l'information sous forme de données et donc leur lisibilité