



**DATA**  
**TERRA**

# Ecole Thématique DATA SDUE

**Guide de Survie dans la jungle des données  
en Sciences de l'Univers et de l'Environnement (SDUE) :  
Comment gérer les données pour les valoriser?**

## **Session « Choix et utilisation de formats de fichiers standards, ouverts et libres »**

Olivier Rouchon



# Les généralités et le contexte

Donnée, information, connaissance, compétence

*Crédits section : L.Duplouy(BnF) – Formation « pérennisation de l'information numérique » Aristote/PIN 2022*

# Que cherchons-nous à conserver ?

010101001000011010101010100011110101000011110101010101111  
1110110100000101011111110101010111010101010101010101011  
110101010000111010101010111111010110101010111111101000010  
010101010101010101111011010101010111111000010100101010101  
110000111110111101010011101010101011010101010111101010101  
101010111100101010010000110101010101000111101010000111101  
010101011111110110100000101011111110101010111010101010101  
010101010111101010100001110101010101111110101101010101111  
111010000100101010101010101011110110101010101111110000101  
001010101011100001111101111010100111010101010110101010101  
111010101011010101111000101010101110000111110111101010011

# Qu'est-ce que l'on veut pérenniser ?

- Suivant que l'on veut conserver/transmettre/pérenniser:
  - L'apparence visuelle d'un document dans un but de preuve et/ou de mémoire
  - Le contenu sans contrainte sur l'apparence visuelle du document qui sera restitué dans le futur
  - Des fonctionnalités complexes existant dans une base de données
- Les choix résultants
  - Choix de standards (formats de données, métadonnées)
  - Choix techniques associés

seront différents

# Quelle est la problématique ?

- But = Pérennisation des informations numériques
  - C'est l'information qui doit être pérennisée
- L'information (dans ce contexte) = donnée compréhensible par un être humain
  - Vidéos, images, sons, textes.
  - Attention importance forme/contenu pour le texte et les documents composites (texte + image)
- Quelle est la différence entre une donnée, une information, une connaissance ?

# Les notions

22

- Est une donnée

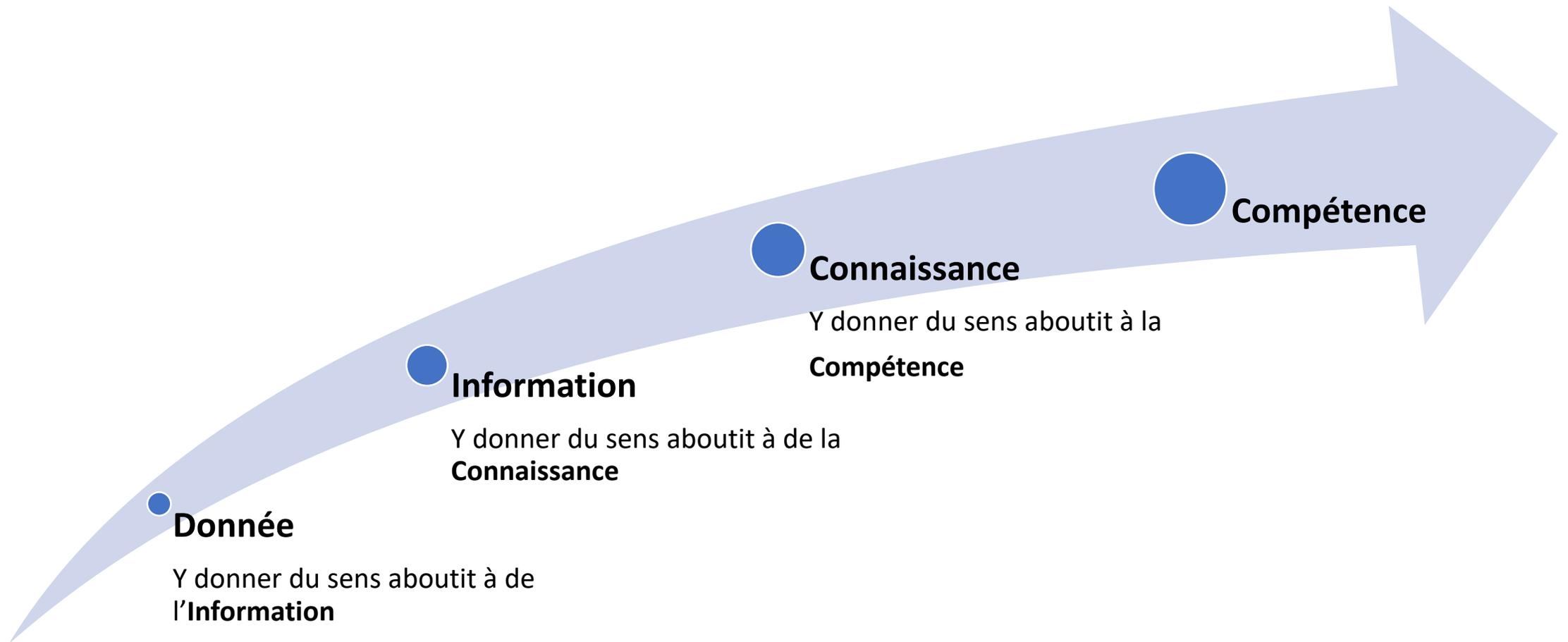
22 degrés Celsius

- Est une information

22 degrés Celsius est une température confortable

- Est une connaissance

# Donnée, information, connaissance



# Information vs données

- *Information*

- tout type de connaissance pouvant être échangée
- indépendante des formes (physique ou numérique) utilisées pour la représenter

- *Données*

- représentation de l'information interprétable de manière formelle permettant la communication, l'interprétation ou le traitement.

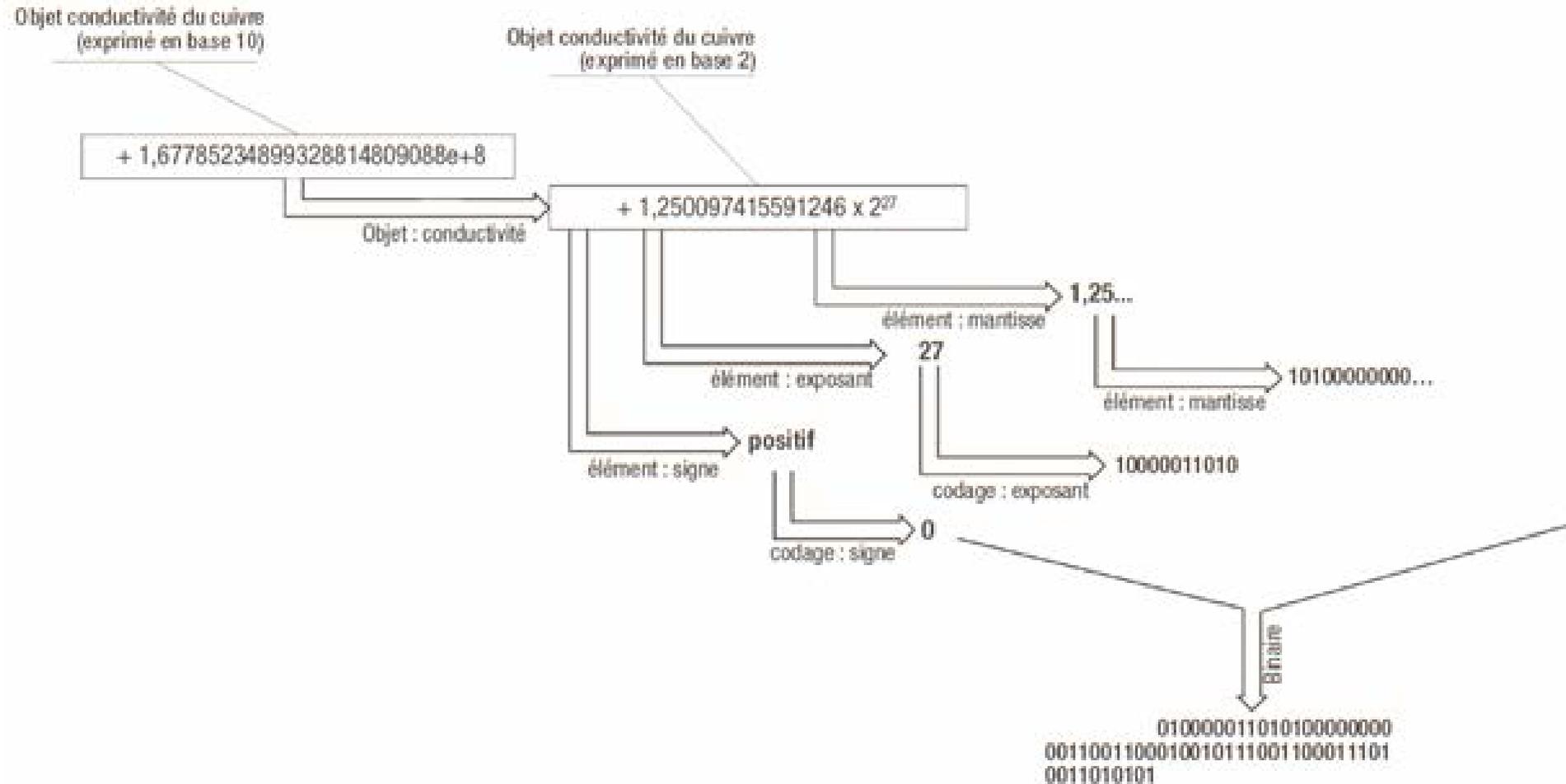


# En résumé

*Texte, son, vidéo* = *Information*

*Format de représentation* = *Support d'information*

# Exemple





# La définition d'une politique de formats

# Les critères de sélection d'un format de fichier

- Elaboré par la cellule nationale de veille sur les formats
- S'appuie sur les travaux de :
  - Library of Congress,
  - British Library,
  - Harvard Library,
  - National Archives and Records Administration (NARA)
  - Bibliothèque nationale de France (BnF).
- 9 critères retenus :
  - Adoption
  - Existence et disponibilité de la documentation
  - Transparence
  - Contenu additionnel embarqué
  - Dépendances externes
  - Incidence des brevets
  - Dispositifs techniques de protection
  - Résilience
  - Compacité



# Adoption

- Utilisation du format pour produire, conserver, diffuser ou échanger l'information.
- Plus l'adoption est large
  - Plus les outils de création, de lecture, de transformation, d'émulation, etc. sont nombreux
  - Plus le format est maintenu et mis à jour.
- Pour les services patrimoniaux, l'utilisation d'un format recommandé pour la conservation par d'autres institutions est un point important de l'adoption.
- L'âge du format et sa diffusion dans plusieurs pays sont notamment des facteurs permettant de mesurer le taux d'adoption.

# Existence et disponibilité de la documentation

- Spécifications disponibles pour le format :
  - Complète
  - Précise
  - Compréhensible
- Pouvant provenir d'une source institutionnelle mais pas forcément
- Existence d'outils accessibles pour identifier et valider le format.
  - Le facteur peut par exemple être corroboré par l'existence d'outils de validation élaborés par différentes sources, preuve que les spécifications sont accessibles et complètes.
  - Les formats ouverts présentent généralement une mise à disposition de la documentation plus complète que les formats propriétaires mais ce n'est pas forcément systématique.

# Transparence

- Accessibilité du format à l'aide d'outils non spécialisés
  - Possibilité pour un humain de lire et de comprendre le code du format grâce à des outils de type éditeur de texte.
  - Outils qui servent à produire le format doivent par ailleurs être utilisables sur des machines standard.
- Suppose également que le format soit non compilé et encodé avec des caractères standardisés.
- La compression, si elle est inévitable, doit utiliser des algorithmes largement répandus ou précisément documentés par l'utilisateur.

# Contenu additionnel embarqué

- Capacité d'un format à être autoporteur
  - Comprendre au sein de son code les flux de données nécessaires à sa compréhension et son exploitation.
- Ces flux sont de plusieurs types :
  - Représentation (permet de transformer des données en information exploitable),
  - Identification,
  - Contexte (lien entre une information et son environnement),
  - Intégrité (traçabilité des modifications),
  - Provenance (historique de l'information)
- Inverse d'un format qui obligerait l'utilisateur à stocker ces flux séparément de l'information primaire.

# Dépendances externes

- Nécessité d'utiliser
  - Des logiciels (notamment propriétaires)
  - Du matériel informatique spécifiques (joystick, microphone, etc.) pour exploiter le format.
- Ce facteur concerne particulièrement les contenus dynamiques, comme les jeux vidéo ou certaines œuvres d'art numérique

# Incidence des brevets

- Présence de brevets ou de toute contrainte juridique sur un format
  - Peut entraîner des coûts de conservation importants et limiter le développement d'outils d'analyse et d'exploitation du format.
- Présence d'un brevet (ou d'une réglementation similaire) doit être associée à l'analyse de l'adoption du format
  - Bon indicateur pour savoir si le brevet empêche la large diffusion du format et la création d'outils d'exploitation.

# Dispositifs techniques de protection

- Manipulation de l'information par les services de conservation à des fins de
  - Diffusion,
  - Transformation,
  - Réparation, etc.
- Empêchée par
  - Le chiffrement,
  - La protection par mot de passe,
  - Eventuellement la présence d'informations de copyright sous forme de watermark,
- Attention d'autres dispositifs comme la signature électronique, n'empêchent pas la manipulation du contenu, mais sont conçus pour assurer son authenticité.

# Résilience / Compacité

- Résilience

- Capacité du format à permettre la lisibilité de son contenu, même en cas d'altération partielle de celui-ci.

- Compacité

- Capacité du format à contenir un grand nombre de données dans un volume réduit.
  - Si la compacité est liée à la compression, s'assurer de la transparence et de la réversibilité de la compression.

- Format maîtrisé = identifié et vérifiable
  - Format publié ; ex. WAVE, SVG ;
  - Format largement utilisé ; ex. XML, MPEG4 ;
  - Format normalisé si possible ; ex. PDF (ISO 32000-1:2008), PNG (ISO 15948:2004).
- Respect des spécifications du format
- Existence d'outils libres permettant une identification, validation et caractérisation.

Type	Format
Texte	TXT, XML, HTML, PDF, ODT
Image	GIF, JPEG, TIFF, PNG, SVG, JPEG200
Son	WAV, AIFF, AAC, OGG (VORBIS)
Vidéo	MPEG4, OGG (THEORA), MKV
Données brutes	CSV, XML, ODS, HDF5, NetCDF, FITS

# Quelques exemples

PDF, NetCDF, FITS

# PDF - Portable Data Format

- Langage de description de page présenté par la société Adobe Systems en 1992
- Norme ISO depuis 2008
- Préservation des
  - polices de caractères, images, objets graphiques,
  - mise en forme

de tout document source, quelles que soient le logiciel de lecture

- Format utilisé dans de nombreux logiciels
  - Exportation dans les suites bureautiques grand public,
  - Manipulations par des programmes spécialisés
  - Génération de documents officiels
- Format binaire



# PDF - Portable Data Format

- Spécifications : <https://www.iso.org/fr/standard/51502.html> (v1.7, 2008)
- PUID : fmt/14  
(<https://www.nationalarchives.gov.uk/pronom/fmt/14>), fmt/276
- Type MIME : application/pdf, application/x-pdf, application/x-bzpdf, application/x-gzpdf
- Signature : 25 50 44 46

# PDF - Portable Data Format

```

%PDF-1.5
%µµµµ
1 0 obj
<</Type/Catalog/Pages 2 0 R/Lang(fr-FR) /StructTreeRoot 106 0 R/MarkInfo<</Marked true>>>
endobj
2 0 obj
<</Type/Pages/Count 11/Kids[ 3 0 R 28 0 R 48 0 R 68 0 R 83 0 R 90 0 R 92 0 R 97 0 R 99 0 R 101 0 R 103 0 R ] >>
endobj
3 0 obj
<</Type/Page/Parent 2 0 R/Resources<</Font<</F1 5 0 R/F2 12 0 R/F3 14 0 R/F4 16 0 R/F5 21 0 R/F6 23 0 R>>/XObject<</Image7 7 0 R/Image9 9 0 R/Image11 11 0 R>>/ProcSet[/PDF/Text/ImageB/ImageC/ImageI] >>/MediaBox[ 0 0 594.96 842.04]
/Contents 4 0 R/Group<</Type/Group/S/Transparency/CS/DeviceRGB>>/Tabs/S/StructParents 0>>
endobj
4 0 obj
<</Filter/FlateDecode/Length 4764>>
stream
xœÖ][!~7äyPRUiu
£æYix±i|Çil OdYôôÄ-ôHè óoóó7æi!YW^aoaä TM »UEò<[!ÔÄ†Óíçôæ<{vñÁé'pü²y"} x¼?ôw□» üý-(q¼pr»[Yn÷» «wOèèÖ-□. /"ç/_ $íož>'xEBÔCE'7ÿÿ>|lÿ'DÑ4£<aE?oiž>É'/oãOOÿ|X$È¿%7ÿyóÉTYié"-ú$àY$X£□
Á$Kóá")úäyþ!2$R²Á"wá"Šce!<xh>,HÁÚ@64"Ô$ô-üýÉÖ!lr□óáôÉ$œ,,$$)Mh0ÉS""%Dw9O6ðÄè»ð—J^i_cpnLQYüäØ!sZ°èFSP#`3¶Nws™™œª )5ý5o_É¿~rFÜjÁýíró2£@IA°ÓR`T4ñiè!±L%+èiÜ4Éqz3K3à4)è>Y'ó²"èpÖÿÛsQn `·
ç`i"§T(²4Kq+4{+.bg'½ú',¶»ÖÝž/è@7[ä!¥GuP!)J@°É²èðÚvpU7$ô¼ % ,Oùb'b·è¥XlÝäÄ),çM-%,Ð^o.ìõ"Ä"™-ÁÉ(/ èæÄ,fÔpSa"<¼0>,V³ÉÁGö_Í
f"ÿ¶FÔ™Ö` ,%J@½|2Éú¼
™@J«2`b¼¿½â«7KBW¿æ°ws_—1D½\□æ-µµp-™L`_joÐÄx`^¿]¶x¿NÁU$«ÍH*1Ý]S`6,a²d°OS,ô9Kú"Q...¶Q¥?ahâè"4i$|'6èÈuKCE²~•ùLéIa¿[èøâ-;ñm²"·O0ðYhN... (Ô=à8□ÝžÓÁÔ†g"NSì.™ª&¡³□
ç...6"ç'•,|~U"€ñQÁ"oRAQa[±?Ú?Át:&?.W!;Lú»óíánÉðóaaTbó^L_f)2",¿°èñkÐ( @ò)ð ò,#ú·†òù·èžÚ"!£pòé6TQä.b'X--Híßµ □¶T5€P³W³CEÁ]fÖ,; ...¹s!™¥Úa•
¶¶¶¶èiò,°DyÚ³±d"-|Á ·•<Ôè@Z@ÚMj"ò,æ4•$"þ$¶]wks`-...V,†_ç #XUELxp Ž»yp:|—ú!Ý@ñ'LÿyX@òÁpWP>Ú<¥úéV¿ßúié E'¥9ç¶|_ ^@äb-Oø†Z'tvHh:'¶A"ÓL" àø:
hÓ□
ÔóÄo[q"²" vU·b:l:±y~@ø)□□L"Ð"?,*P·óáW1LZE@Q0;f·GαÔ_ô 44)ÁÑÓá¶%s%GÔó=áyZðv          p%xxÁ,3i(ÉCEV      Á~p5C□          ÛÉryªÁðíH          C¥. ©tÉt$"hVÔk†çj™K]lf
"°" »gYÆÉ%†_Á¿pÇ.W9]æpá—«á™—&ðN+yé"½!a
pâ2$,`òÉ$<e%□_°ÜÖ_ÚNZ`ç8{ØœÐÓ]IÁ¹!¥<KÁ·Úí6e«ai¾mñxÚ`Ð`Éyx
&éq¼)5â"iáçúloèx1¼\,+è£•
xð`x#çÉy0óN/Æ?Ú; &Áè±?;ÚÇwkwKi>äüð'ÉU1_ú
É—@JbjÚÓá¶H-Á[...ó1ÖNC`A'Z".ÖRÓ4SÉé."
Úp@Æ@v °dª=¥°niÉj'0½I1@VKgáÄ+;@à>Š          %w·0áNÚ$tw±PJ`±tešÉti÷:(QÐf□@?k;ž'Ä.Æ(ý¾øè¿-E+ú43bRpCe{>-SicTšcøVEGHwæZúCdniÈÜxá¿;,->É6•x¶žÿÖúL·[TN_7ž[Z·hüYpè-Ä8Ú@xvB_¡AØ¼aóI5!M%00
òèYú+v[±ðóúRèn¿1KJhÿ!9rî5Y%`x¾-i3-ÄÄÉJž¼xÿyðfí~`¼°ñórvZ"j·?I-z—zNÑ¿ %øÿYfH$á-Ôtç$.../ÓNÁp ²éJ5
ÖSÑ4á
2YD-s·3Ó]oóRªx•«R¾"ÁH[#ÉB,s%h-nÿ áR¥ŠDèOÚá)¥VÇöÿø`_B:è=ÖH£ÐpysÒ¼_üòD`ñÄ·p±°QáU^¿zò²zuXí-Ááÿ[ZVk(ú;°G Ž $™™kgM —=P-zèD
ÚÉáŠ
v`Í<atí1-9Zø]x±R`|·µñí_Á
•fÚ=,uCEzT»_zÓðqexðá,íX ?FA`ýí,ÁLÉ °æœL.É]Úf(©æ1r¾3"Z°ipèñúwöð

```

# NetCDF - Network Common Data Form

- Format de fichiers pour données scientifiques stockées sous la forme de tableaux de nombres multi-dimensionnels
- Pluridisciplinaire, très utilisé en météorologie, océanographie
- Format binaire
- Structure de données « auto-documentée »
  - En-tête
    - Décrivant la disposition des données dans le reste du fichier,
    - Contenant une liste arbitraire de métadonnées présentées sous la forme d'attribut de type nom/valeur.
  - Tableaux de données
- Indépendant de l'architecture matérielle



# NetCDF - Network Common Data Form

- Spécifications : [https://docs.unidata.ucar.edu/netcdf-c/current/file\\_format\\_specifications.html](https://docs.unidata.ucar.edu/netcdf-c/current/file_format_specifications.html) (v4.0, 2008)
- PUID : fmt/282 (<https://www.nationalarchives.gov.uk/pronom/fmt/282>), fmt-283
- Type MIME : application/netcdf, application/x-netcdf
- Signature : 43 44 46 01
- Bibliothèques fournies par l'UCAR : <https://www.unidata.ucar.edu/software/netcdf/>
- File format validation : <https://www.seanoe.org/data/00344/45538/>



# FITS - Flexible Image Transport System

- Format de fichiers pour la communauté Astronomie

- Conçu spécifiquement pour des données scientifiques
- Inclut beaucoup de dispositifs pour décrire l'information photométrique et spatiale de calibrage, ainsi que les métadonnées de l'origine de l'image.



- Structure multiple

- Successions de HDU (« Header » + « Data Unit »)
- Les métadonnées de l'image ou du tableau multidimensionnel sont stockées dans un en-tête lisible par un humain, au format ASCII
- 3 modèles de données
  - Image,
  - Table ASCII
  - Table binaire



# FITS - Flexible Image Transport System

SIMPLE = T / file does conform to FITS standard BITPIX = -32 / number of bits per data pixel NAXIS = 3 / number of data axes NAXIS1 = 200 / length of data axis 1 NAXIS2 = 200 / length of data axis 2 NAXIS3 = 4 / length of data axis 3 EXTEND = T / FITS dataset may contain extensions COMMENT FITS (Flexible Image Transport System) format is defined in 'AstronomyCOMMENT and Astrophysics', volume 376, page 359; bibcode: 2001A&A...376...359H BSCALE = 1.0E0 / REAL = TAPE\*BSCALE + BZERO BZERO = 0.0E0 / OPSIZE = 2112 / PSIZE of original image ORIGIN = 'STScI-STSDAS' / Fitsio version 21-Feb-1996  
FITSDATE= '2004-01-09' / Date FITS file was created FILENAME= 'u5780205r\_cvt.c0h' / Original filename ALLG-MAX= 3.777701E3 / Data max in all groups ALLG-MIN= -7.319537E1 / Data min in all groups ODATTYPE= 'FLOATING' / Original datatype: Single precision real SDASMGNU= 4 / Number of groups in original image CRVAL1 = 182.6311886308 CRVAL2 = 39.39633673411 CRPIX1 = 420.  
CRPIX2 = 424.5 CD1\_1 = -1.067040E-6 CD1\_2 = -1.259580E-5 CD2\_1 = -1.260160E-5 CD2\_2 = 1.066550E-6  
DATAMIN = -7.319537E1 / DATA MIN DATAMAX = 3.777701E3 / DATA MAX MIR\_REVR= T ORIENTAT= -85.16 FILLCNT = 0  
ERRCNT = 0 FPKTTIME= 51229.798574 LPKTTIME= 51229.798742 CTYPE1 = 'RA---TAN' CTYPE2 = 'DEC--TAN'  
DETECTOR= 1 DEZERO = 316.6452 BIASSEVEN= 316.6715 BIASODD = 316.6189 GOODMIN = -5.064006  
GOODMAX = 2552.17 DATAMEAN= 0.4182382 GPIXELS = 632387 SOFTERRS= 0 CALIBDEF= 1466 STATICD  
= 0 ATODSAT = 16 DATALOST= 0 BADPIXEL= 0 OVERLAP = 0 PHOTMODE=  
'WFPC2,1,A2D7,LRF#4877.0,CAL' PHOTFLAM= 3.447460E-16 PHOTZPT = -21.1 PHOTPLAM= 4884.258 PHOTBW = 20.20996  
MEDIAN = -0.175651 MEDSHADO= -0.121681 HISTWIDE= 1.033711 SKEWNESS= -1.983727 MEANC10 = 0.12958  
MEANC25 = 0.3129676 MEANC50 = 0.4577668 MEANC100= 0.3916293 MEANC200= 0.3115222 MEANC300= 0.3295493  
BACKGRND= -0.3676353 ORIGIN = 'NOAO-IRAF FITS Image Kernel December 2001' / FITS file originator DATE = '2004-01-09T03:26:36' IRAF-TLM= '03:26:36 (09/01/2004)' FILETYPE= 'SCI' /  
type of data found in data file TELESCOPE= 'HST' / telescope used to acquire data INSTRUME= 'WFPC2' / identifier for instrument used to acquire data EQUINOX = 2000.0 / equinox of celestial coord.  
system / WFPC-II DATA DESCRIPTOR KEYWORDS ROOTNAME= 'u5780205r' / rootname of the observation set PROCTIME= 5.301314019676E+04 /  
Pipeline processing time (MJD) OPUS\_VER= 'OPUS 14.5a' / OPUS software system version number CAL\_VER = ' / CALWP2 code version / SCIENCE INSTRUMENT CONFIGURATION  
MODE = 'FULL' / instr. mode: FULL (full res.), AREA (area int.), SERIALS = 'OFF' / serial clocks: ON, OFF / IMAGE TYPE CHARACTERISTICS  
IMAGETYP= 'EXT' / DARK/BIAS/IFLAT/UFLAT/VFLAT/KSPOT/EXT/ECAL CDBSFILE= 'NO' / GENERIC/BIAS/DARK/PREF/FLAT/MASK/ATOD/NO PKTFMT = 96 / packet format code / FILTER  
CONFIGURATION FILTNAM1= 'FR533P15' / first filter name FILTNAM2= ' / second filter name FILTER1 = 69 / first filter number (0-48) FILTER2 =  
0 / second filter number (0-48) FILTROT = 15.0 / partial filter rotation angle (degrees) LRFWAVE = 4877.000000 / linear ramp filter wavelength / INSTRUMENT STATUS USED IN DATA  
PROCESSING UCH1CJTM= -88.2569 / TEC cold junction #1 temperature (Celsius) UCH2CJTM= -88.6697 / TEC cold junction #2 temperature (Celsius) UCH3CJTM= -88.3028 / TEC cold junction #3  
temperature (Celsius) UCH4CJTM= -88.7671 / TEC cold junction #4 temperature (Celsius) UBAY3TMP= 13.2302 / bay 3 A1 temperature (deg C) KSPOTS = 'OFF' / Status of Kelsall spot lamps: ON, OFF SHUTTER = 'A' / Shutter in place at  
beginning of the exposure ATODGAIN= 7.0 / Analog to Digital Gain (Electrons/DN) / RSDP CONTROL KEYWORDS MASKCORR= 'COMPLETE' /  
Do mask correction: PERFORM, OMIT, COMPLETE ATODCORR= 'COMPLETE' / Do A-to-D correction: PERFORM, OMIT, COMPLETE BLEVCORR= 'COMPLETE' / Do bias level correction BIASCORR= 'COMPLETE' / Do bias correction: PERFORM, OMIT,  
COMPLETE DARKCORR= 'COMPLETE' / Do dark correction: PERFORM, OMIT, COMPLETE FLATCORR= 'SKIPPED' / Do flat field correction SHADCORR= 'OMIT' / Do shaded shutter correction DOSATMAP= 'OMIT' / Output saturated pixel  
map DOPHOTOM= 'COMPLETE' / Fill photometry keywords DOHISTOS= 'OMIT' / Make histograms: PERFORM, OMIT, COMPLETE OUTDTYPE= 'REAL' / Output image datatype: REAL, LONG, SHORT  
LINENUM = '02.030' / PEP proposal line number SEQLINE = ' / PEP line number of defined sequence SEQNAME = ' / PEP define/use sequence name HISTORY MASKFILE=uref\$8213081u.r0h MASKCORR=COMPLETED  
HISTORY PEDIGREE=INFLIGHT 01/01/1994 - 15/05/1995 HISTORY DESCRIP=STATIC MASK - INCLUDES CHARGE TRANSFER TRAPS HISTORY BIASFILE=uref\$9a1612mu.r2h BIASCORR=COMPLETED HISTORY PEDIGREE=INFLIGHT 29/08/98 - 21/08/99  
HISTORY DESCRIP=not significantly different from j6e16008u. HISTORY DARKFILE=uref\$2g1549cu.r3h DARKCORR=COMPLETED HISTORY PEDIGREE=INFLIGHT 16/02/1999 - 16/02/1999 HISTORY DESCRIP=Pipeline dark: 120 frame superdark  
with hotpixels from HISTORY 16/02/99 HISTORY FLATFILE=uref\$4i1559cu.r4h FLATCORR=SKIPPED HISTORY PEDIGREE=DUMMY 18/04/1995 HISTORY DESCRIP=All pixels set to value of 1. Not flat-fielded.  
HISTORY PC1: bias jump level ~0.100 DN. HISTORY The following throughput tables were used: HISTORY crotacomp\$shst\_ota\_007\_syn.fits, crwfpc2comp\$wfpc2\_optics\_006\_syn.fits, HISTORY crwfpc2comp\$wfpc2\_lrf\_004\_syn.fits[wave#],  
HISTORY crwfpc2comp\$wfpc2\_dqepc1\_005\_syn.fits, HISTORY crwfpc2comp\$wfpc2\_a2d7pc1\_004\_syn.fits, HISTORY crwfpc2comp\$wfpc2\_flatpc1\_003\_syn.fits, HISTORY The following throughput tables were used:  
HISTORY crotacomp\$shst\_ota\_007\_syn.fits, crwfpc2comp\$wfpc2\_optics\_006\_syn.fits, HISTORY crwfpc2comp\$wfpc2\_lrf\_004\_syn.fits[wave#], HISTORY crwfpc2comp\$wfpc2\_dqewfc2\_005\_syn.fits, HISTORY  
crwfpc2comp\$wfpc2\_a2d7wf2\_004\_syn.fits, HISTORY crwfpc2comp\$wfpc2\_flatwf2\_003\_syn.fits CTYPE3 = 'GROUP\_NUMBER' / Extra dimension axis name CD3\_3 = 1 / CD3\_1 = 0 /  
CD1\_3 = 0 / CD2\_3 = 0 / CD3\_2 = 0 / END  
?Pcö½º;º% ??Jú€?[2?Yüð½Ä³ % #! ?gýÄ?RE%ÄÉ¿¿...U?]=¼B¿'w'¿¿¿'ò§%5¼v¼""ò%V||1"º½!ST¿÷xv?è%¼\$§Ä?UD%¼"ò%9 Ä¿-#t¿\$! ?O«A è¿?K\*L¼i08?ft@ CE!¿ãÑ¿ã":¿¿d¿fN@?6|p¼cF?%OóX¿¿S¿¿¿...  
Ö¿?;O-?\_¿?a+3¿B-?!Ü=|)(%|úú%eÉb?0¿¿¿¿{?@?Ä?+R×¿B%M"2?b¼k¿J)¿%P¿=Eh?W>¿-¿?¿%S"¿?G>?@¿Z¿"éi>cl¿Ám'¿e½".Ü¿¿e¿TX¿¼%,J¿,f2?¿z?R>¿¿B¿

# L'assurance qualité sur les formats

Les contrôles, les outils

*Crédits section : L.Duplouy(BnF) – Formation « pérennisation de l'information numérique » Aristote/PIN 2022*

# Les types de contrôle possible

- Identification
  - Déterminer le type de format présumé de l'objet numérique par signature
- Validation
  - Déterminer la conformité du format par rapport à sa spécification (standard ou norme) ou par rapport à un profil d'application du format (contexte spécifique/local-use)
- Caractérisation
  - Déterminer les propriétés (nécessite des outils) de l'objet numérique
- Évaluation
  - Déterminer le niveau de risque du format (peut varier selon propriétés)

# Les outils

- Commande « file » sous \*nix
  - Basé sur la signature
  - Paramétrable
  - Ne permet que l'identification

- DROID

- <http://www.nationalarchives.gov.uk/aboutapps/pronom/droid.htm>
- identification automatique des formats
- utilise les signatures internes et externes référencées dans la base Pronom (PUID, ex :fmt/40)
- beaucoup de formats sont justes signalés et pas décrits
- les spécifications sont externes (= problème pour la préservation, pas d'évaluation possible)



# Les outils

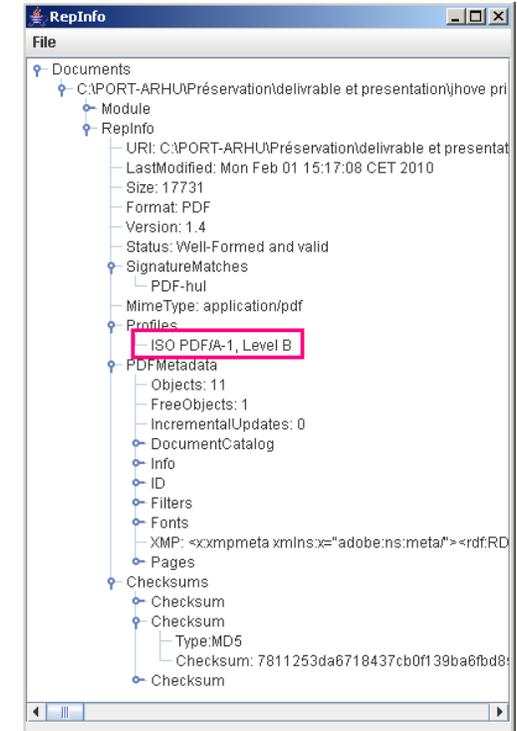
- FIDO (Format Identification for Digital Objects)
  - porté par openplanets foundation (financé par le projet Eu SCAPE )
  - <http://www.openplanetsfoundation.org/software/fido/>
  - Outil d'identification de format
  - Fournit des identifiants pronom
  - Va plus loin que DROID : analyse complète des fichiers
    - Pas uniquement l'entête de fichier
    - Détecte les formats imbriqués



# Les outils

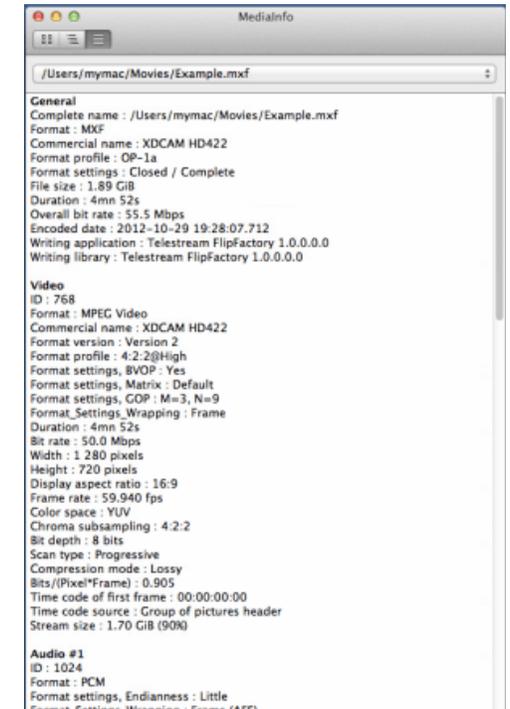
- JHOVE

- Initialement porté par l'université d'Harvard
- <http://jhove.sourceforge.net/>
- Outil de caractérisation
- Fournit une description complète des formats supportés
- Très populaire notamment pour l'analyse de fichier image (TIFF, JPEG, ...) remplacé par JHOVE2
  - <https://bitbucket.org/jhove2/main/wiki/Home>
  - Porté par la Library of Congress
- Amélioration du framework pour faciliter l'intégration de nouveaux modules



# Les outils

- mediainfo
  - Projet open-source, créé par Jérôme Martinez
  - <http://mediaarea.net/>
  - Outil d'analyse des formats audio-vidéo
    - Conteneur
      - MPEG-4, QuickTime, Matroska, AVI, MPEG-PS (y compris les DVD non protégés), MPEG-TS (y compris les Blu-ray non protégés), MXF, GXF, LXF, WMV, FLV, Real...
    - Tags
      - Id3v1, Id3v2, Vorbis comments, APE tags...
    - Vidéo
      - MPEG-1/2 Video, H.263, MPEG-4 Visual (DivX, XviD compris), H.264/AVC, Dirac...
    - Audio
      - MPEG Audio (MP3 compris), AC3, DTS, AAC, Dolby E, AES3, FLAC...
    - Sous-titres
      - CEA-608, CEA-708, DTVCC, SCTE-20, SCTE-128, ATSC/53, CDP, DVB Subtitle, Teletext, SRT, SSA, ASS, SAMI...



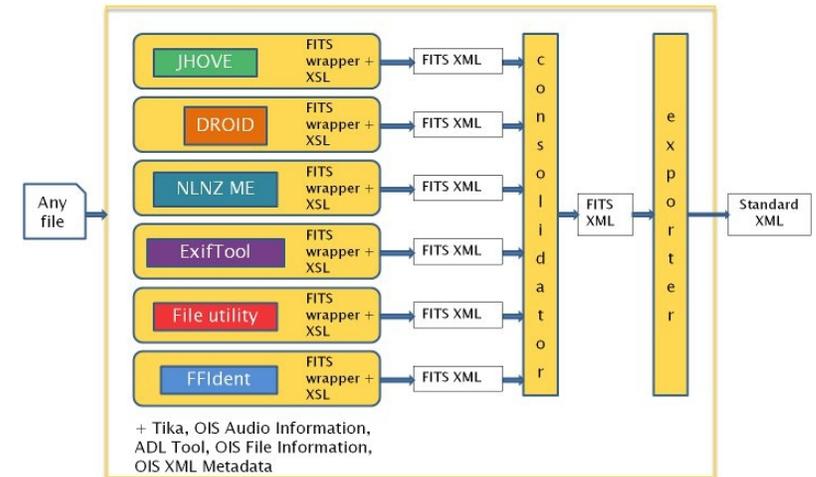
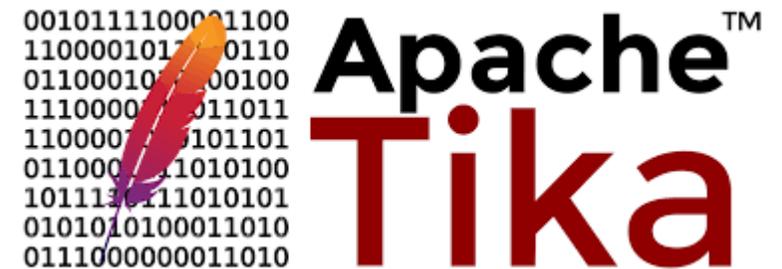
# Les outils

- Tika

- Porté par Apache software foundation
- <http://tika.apache.org>
- Utilise la technique des parsers
- Utilise le mime/type pour l'identification
- Permet la détection et l'extraction de métadonnées de différents formats de fichier (texte, image, audio, vidéo)
  - Permet l'extraction de métadonnées des formats bureautique ODF et OOXML

- File Information Tool Set (FITS)

- Projet porté par l'université d'Harvard
- <http://projects.iq.harvard.edu/fits>
- Fédérateur de différents outils
- Propose un format XML de sortie unifié



# L'exemple FACILE au CINES

Format d'Archivage du Cines par anaLyse et Expertise

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



## FACILE - Service de validation de formats

Vérifier l'éligibilité de vos documents à un archivage sur la plateforme PAC du CINES.

Validation

Correction PDF

Tutoriels

Web Service

+ Choisissez des fichiers - Taille max 2,5 Go

→ Valider

⊗ Annuler

[Cliquez ici pour demander l'aide d'un expert du 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.

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.	✗
JPEG RAW	Joint Photographic Experts Group - Raw JPEG Stream	[fmt/41]	[image/jpeg]	Format de représentation compressée d'une image fixe.	✗
JPEG2000	JPEG 2000	[fmt/151, x-fmt/392]	[image/jp2]	Extension du format JPEG.	✓
JPEG 1.00	Joint Photographic Experts Group	[fmt/42]	[image/jpeg]	Format de représentation compressée d'une image fixe.	✓
JPEG 1.01	Joint Photographic Experts Group	[fmt/42]	[image/jpeg]	Format de représentation compressée d'une image fixe.	✓

# FACILE

- Liste les formats de fichiers éligibles à l'archivage au CINES
- Permet
  - L'identification d'un format
  - La validation (par rapport au format identifié)
  - La correction éventuelle d'erreurs (format PDF uniquement)
  - L'assistance des experts du CINES
  - La formation aux bonnes pratiques (via des tutoriels)

# Discussion



# Vos formats

- Données numériques :
  - CSV (17) / Excel (10)
  - BDD relationnelles (5)
  - netCDF, XML (4)
  - geoTIFF, HDF, JSON ... (2)
  - ames, miniseed, FITS, et autres formats divers et variés ... (1)
- Données multimédias :
  - JPG, TIFF (4)
  - wav, MP3, MP4 (4)