



# The Digital Preservation of Audiovisual Content – or – what's so special about audiovisual?

Richard Wright, PrestoPRIME

BBC Research & Development, UK

SUN Preservation and Archiving Special Interest Group

Malta 24-26 June 2009

## Stone, papyrus, film, hard drives

Medium	bits/cm <sup>2</sup>	life
Stone	10	10 000
Paper	10 <sup>4</sup>	1000
Film	10 <sup>7</sup>	100
Disc	10 <sup>10</sup>	10

⇒ Each change 1000 times more dense and hence 1000 times *cheaper*

but the media lasts 0.1 times as long

⇒ Interval between changes also 0.1 times as long

# The problems of AV content

- Traditional media are obsolete (except film, and that's "endangered")
- So we've been digitising
  - Probably 5 to 10 million hours of digitisation by broadcasters in last decade, in Europe
  - National audiovisual collections add millions more
  - Born digital: archiving 300 hrs/week = 30 TB (BBC)
    - Uncompressed Standard Definition; HD will be ???
- So: moving from media that lasts a few decades, to files – and to technology that lasts a few years
- So we'd like to use digital library technology
- If we could

# PrestoPRIME: problems of using digital library technology in broadcasting

- Because MXF is the primary professional broadcast wrapper format
  - In Europe, at LOC, in US Public Broadcasting and for digital cinema
- Because MXF isn't supported by any digital library tools: JHOVE, PRONOM, metadata extractors
- Because OAIS is little-known in broadcasting
- Because broadcasters use MAM, not digital archives much less digital repositories
- Because AIP etc has few broadcast exemplars
- Hence: a shotgun marriage of OAIS and MXF

# The National Archives

Search the archives  
[Advanced search](#)

You are here: [Home](#) > [Services for professionals](#) > [Preservation](#) > [PRONOM](#) > [Search by format](#) > Results



## The technical registry PRONOM

- Welcome
- About
- Add an entry
- Search
- Help
- Information resources

[? Help : report on file format](#)

### Search Results

- Simple search
- File format
- PRONOM Unique Identifier
- Software
- Vendor
- Lifecycles

You searched for: "mxf"

Save as... XML | CSV Print

No data matched the search criteria you entered.

#### Getting in touch

- Contact us
- Press office
- Visit us

#### Site help

- A-Z index
- Accessibility
- Site map

#### About us

- Jobs
- Terms of use
- Freedom of information

#### Websites

- Office of Public Sector Information
- Learning Curve
- Directgov

[? Help : report on file format](#)

## Search Results

Simple search File format PRONOM Unique Identifier Software Vendor Lifecycles

You searched for: "wav"

Save as... XML | CSV Print

page 1

PRONOM Unique ID	Format Name	Format Version	Extension	Format Risk
fmt/1	<a href="#">Broadcast WAVE</a>	0	wav	
fmt/2	<a href="#">Broadcast WAVE</a>	1	wav	
x-fmt/397	<a href="#">Exchangeable Image File Format (Audio)</a>	2.0	wav	
x-fmt/389	<a href="#">Exchangeable Image File Format (Audio)</a>	2.1	wav	
x-fmt/396	<a href="#">Exchangeable Image File Format (Audio)</a>	2.2	wav	
fmt/6	<a href="#">Waveform Audio</a>		wav	
fmt/141	<a href="#">Waveform Audio (PCMWAVEFORMAT)</a>		wav wave	
fmt/142	<a href="#">Waveform Audio (WAVEFORMATEX)</a>		wav wave	
fmt/143	<a href="#">Waveform Audio (WAVEFORMATEXTENSIBLE)</a>		wav wave	

Page: 1 of 1

page 1

Getting in touch

Site help

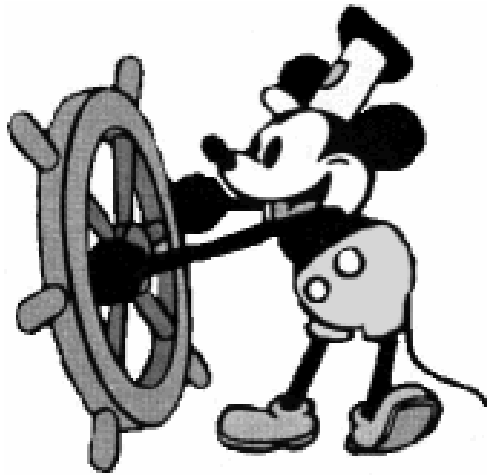
About us

Websites

# Policy: Compressed vs Uncompressed

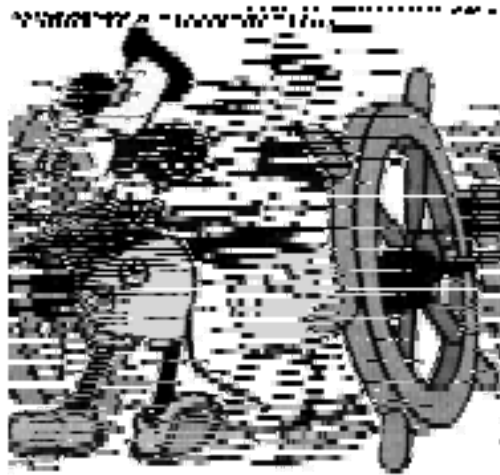
- Why not go from one high-quality compressed format to the next, forever?
  - Because there will not, in general, be a transcoder
  - So the material c1 has to first be decompressed to u1, then encoded to c2
  - At the next cycle, c2 is decoded to u2
  - $u1 \neq u2 (\neq u3 \neq u4 \dots)$
  - So the “real thing” drifts away
- And- an issue of file resilience (mitigation of loss):

# Mitigation of Loss



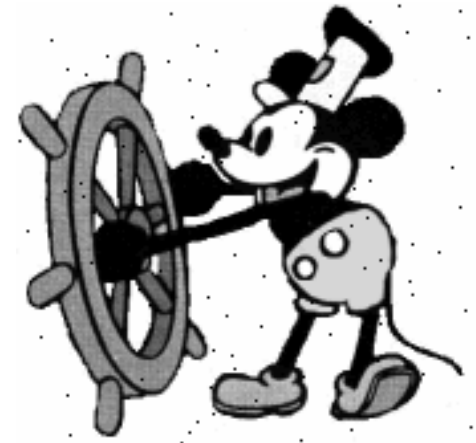
Not quite ready for retirement:  
*Steamboat Willie*, 1928  
Used without permission

BMP no errors 40k



Not quite ready for retirement:  
*Steamboat Willie*, 1928  
Used without permission

GIF 3 errors in 10k



Not quite ready for retirement:  
*Steamboat Willie*, 1928  
Used without permission

BMP 160 errors in  
40k

Heydegger, V. (2008). Analysing the impact of file formats on data integrity. *Proceedings of Archiving 2008*, Bern, Switzerland, June 24-27.



# Format Roadmap: low quality media

Ingest Format	Migration format	Notes
VHS tape	DVD	<b>Access</b> Perfectly adequate for VHS playback
VHS tape	MPEG-4 files	<b>Access</b> Adequate for quality. Minimum data rates (MPEG-4): 500k b/s. There are MANY potential access formats, and they come and go.
VHS tape	DV files	<b>Archive</b> (temporary) 25 M b/s, 12 GB/hr. Migrate to lossless for preservation.
'low end' digital files	Save as is, AND save as DV or lossless	<b>Archive</b> (temporary) Before format or DV format becomes obsolete, migrate to lossless for preservation.
DVD	DV files	<b>Archive</b> (temporary) 25 M b/s, 12 GB/hr. Migrate to lossless for preservation.

# Conclusions about compression

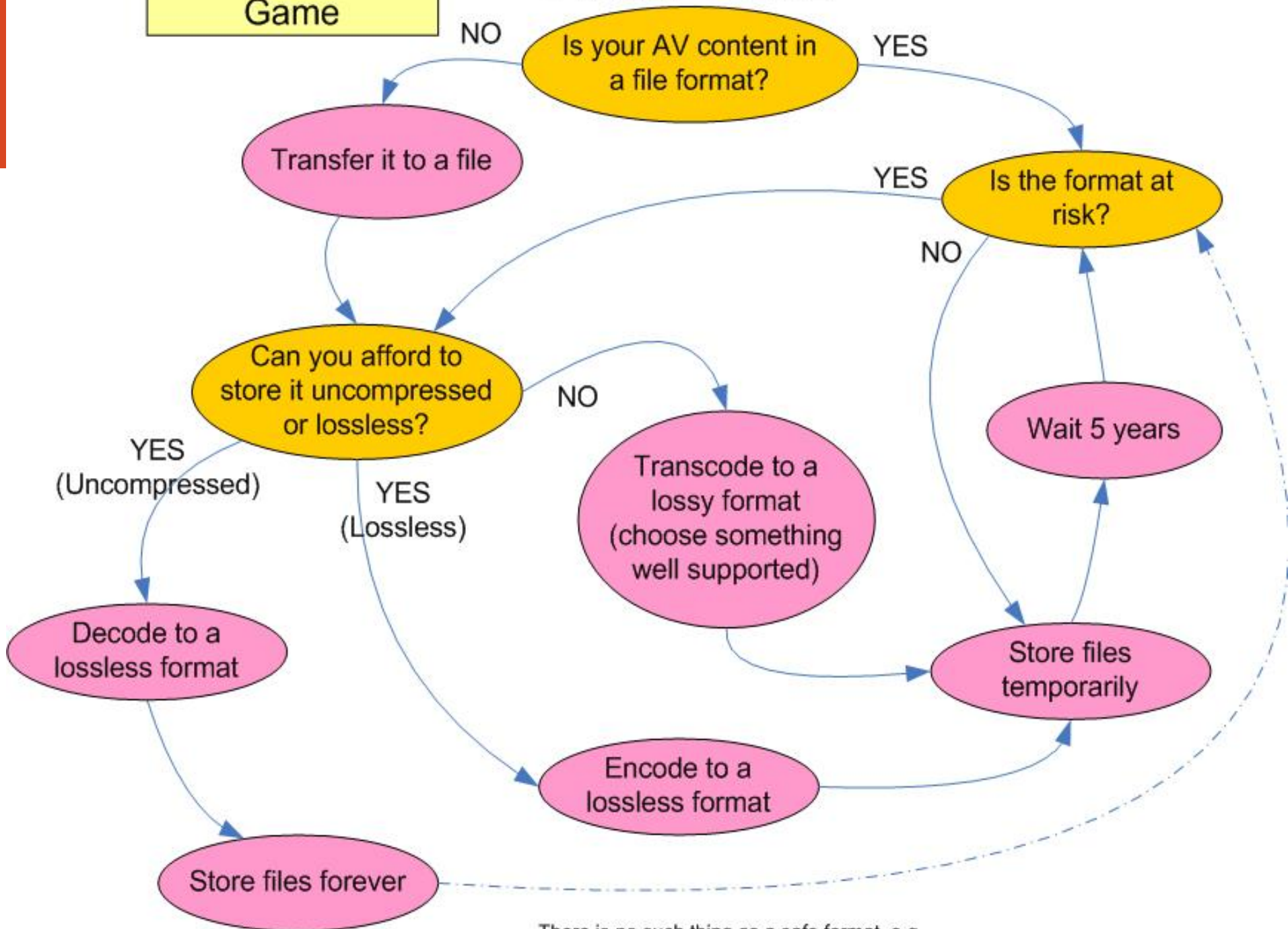
- **Should only be used once** in a ‘life history’ of audiovisual content
- Then have to move to uncompressed to avoid the ‘drift’ of never having a fixed reference
- But: compressed is less resilient
- So: only use interframe compression (as in DV and high-level MPEG profiles)
  - Which would allow individual frames to be lost without affecting the rest of the file
  - If only partial recovery of files were possible!

# Strategy: Migration, Emulation

- Migrate to uncompressed as soon as you can
- Check every format for risks, obsolescence, every N years
- Keep the original bitstream – if you can!
- Can use emulation to extend the usable life on an at-risk format
- Never move from lossy A to lossy B
  - Possible special pleading for genuine transcoding
- Implies possibility of keeping a poor format until it is obsolete – then migrating once to a not-so-poor format – then migrating to uncompressed
- Lossless compression (interframe compression only) is nearly as robust as uncompressed, for 1/3 the storage (but adds complexity and dependencies). Total risk unknown!

# The Preservation Game

## START HERE



There is no such thing as a safe format, e.g. the idea of 'files' will one day be obsolete

## FINISH HERE

# A Brief Reminder of Time

- AV content is time based (you knew that)
- And the files are very large (you knew that too)
- So: important for users to have a description of the internal structure of a file – so the user can navigate to ‘the interesting bit’ eg storyboard, “stripe-image”
- Requires: time-based metadata
  - Time-based presentation of the file’s contents
  - An ability to stream / download just ‘the interesting bit’
- Building communities: requires time-based *citation* and *annotation* of content (and sharing them)

Simple | Avançada

buscant en DIGITON

fixa  estrats  tot

ID  Data  +

Domini

Resultats (65 - DIGITON) Llistats | Recerca

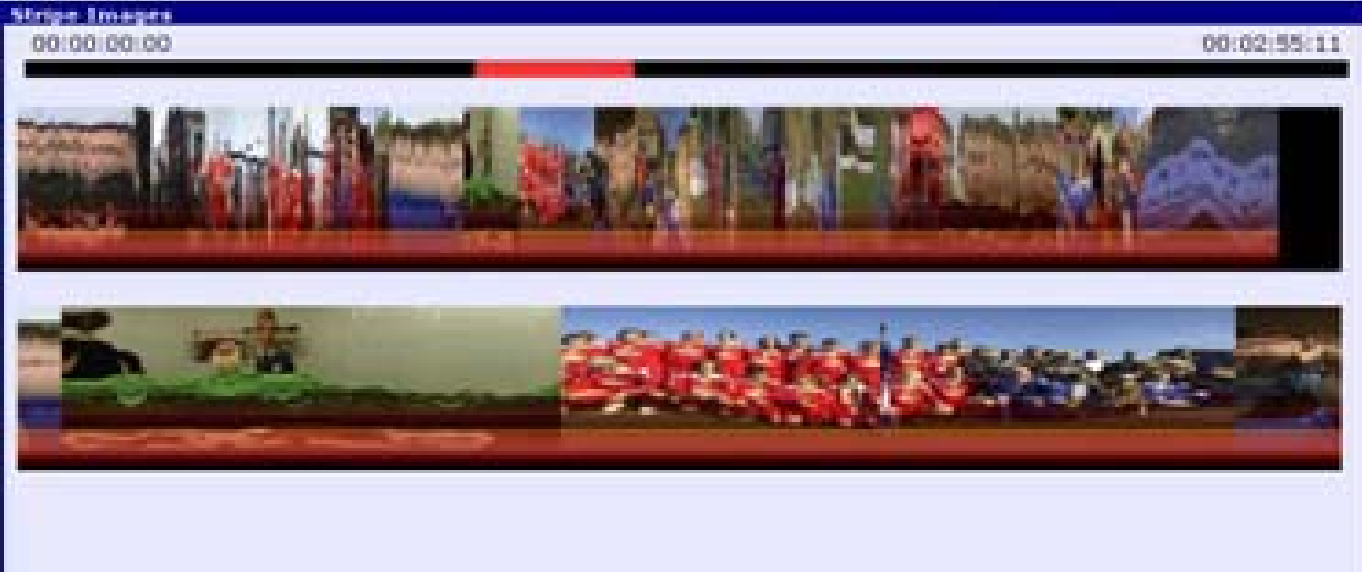
Id	Títol	Data
1319	FUTBOL: REAL MADRID - OSASUNA (2-0). RESUM DEL PARTIT	EDI 2007-12-16 23:23:08
1396	FUTBOL: VALÈNCIA - BARÇA (0-3). DECLARACIONS DE JOAN LAPORTA SOBRE LA LESIÓ DE LEO MESSI I SOBRE L'ESTAT DE FORMA DE RONALDINHO	EDI 2007-12-16 23:18:54
1312	FUTBOL: PARTIT DE LA MARATÓ DE TV3. FIATC - ESTRELLA DAMM (3-1). RESUM POSTPRODUÏT DEL PARTIT	EDI 2007-12-16 21:44:54

1312 (3/65)

Vídeo 1312

Inici: 00:00:00:00 | Final: 00:02:55:11 | 00:01:14:23

In: 00:00:00:00 | Out: 00:00:00:00



Ficha

ID	GENERAL
1312	
Durada	00:02:55:11
Títol	FUTBOL: PARTIT DE LA MARATÓ DE TV3. FIATC - ESTRELLA DAMM (3-1). RESUM POSTPRODUÏT DEL PARTIT
Data	2007-12-16 21:44:54
TEMA/NOM	SANT BOI DE LLOBREGAT * FUTBOL * MARATÓ TELEVISIVA * TVC * TV3
Resum	Partit disputat a Sant Boi entre els equips FIATC i Estrella Damm, en el que han participat personatges coneguts de diversos àmbits, amb motiu de la Marató de TV3 dedicada a les malalties cardiovasculars.
Text	quan truquem ve pleu durada 02'43" (ID=0) (ID=1) (ID=2) (ID=3) (ID=4) (ID=5)
	Així de simple és la tàctica que s'aplica al partit de futbol de la marató. Una cita clàssica que revivim el

# Access work in PrestoPRIME

- Expressing legacy metadata in RDF
- Combining legacy ontologies
  - And using the result to support search (invisibly)
- And time-based annotation tools
  - Trying to follow new W3C work
- And applying access technology to Europeana, the European Digital Library
  - (plus an effort to quantify, codify and automate rights management – for broadcast content)

# Other things I haven't said about PrestoPRIME:

- It's an EC Integrated Project
- Running from Jan 2009 for 42 months
- It's the only project in its 'cohort' about broadcast content
- It's about digital libraries, rights, access
- Ex Libris an important partner: real-world contact with digital library technology
- But also delivering open source implementations of tools developed by PrestoPRIME
- Launching a Networked Competence Centre to live on after the project, to support AV preservation



# Thank you

- PrestoPRIME
  - INA, RAI, BBC, B&G, ORF
  - Univ's of Amsterdam, Innsbruck, Liverpool, Southampton
  - Europeana, Joanneum Research
  - Ex Libris, Eurix, Doremi, Technicolor
  
- [www.prestoprime.eu](http://www.prestoprime.eu)
  - (but not much there yet!)

richard.wright@bbc.co.uk