



Digital **Preservation** Coalition

Tool demo: making sense of a collection





Digital **Preservation** Coalition

Scale = automation = software tools

- Digital typically means working at scale
- Automation becomes a necessity
- Software tools
- Techie stuff
- But...
- Hands on experience





Digital **P**reservation Coalition

1010101010111010001010101000	SOI
10010100101110100101010101001	APPO JFIF
00100101010100100000101010101	1.2
01111101000101010100010010100	APP13 IPTC
10111010010101010100100100101	APP2 ICC
0101001000001010101010111101	DQT
000101010001001010010111010	SOFO 200x392
0000101010101011010101011101	DRI
000101010001001010010111010	DHT
010101010010010010101010010	SOS
00001010101011110100010101	ECS0
010001001001011101001010101	RST0
010010010010101001000001010	ECS1
1010101111010001010101000100	RST1
1010010111010010101010100100...	ECS2...





Making sense of a collection

- Understand the data, then assess risks, plan, take action to preserve
- Characterisation:
 - How many files?
 - How big are the files?
 - What file formats?
 - Is the data dynamic or interactive?
 - Does it contain personal information?
 - Is it encrypted?



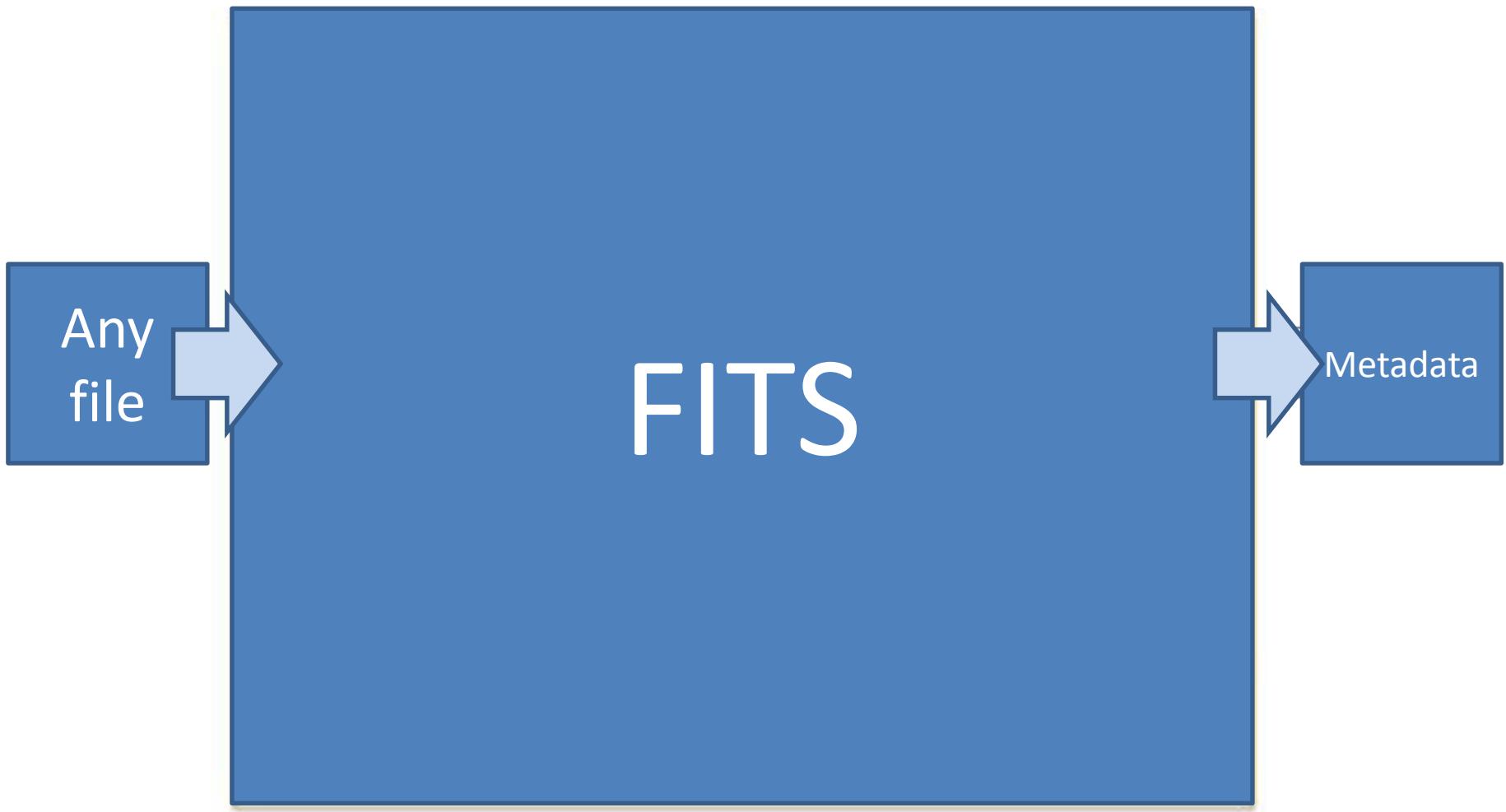
File formats: what flavour is your data?

- What are they likely to contain?
- What are the likely preservation risks associated with that file format?
- What software do I need to render, interpret or make sense of that file format?
- Caveat: Fuzzy concept, doesn't give all the answers
 - [Abusing file formats](#) by Ange Albertini (chapter 6)



Digital **Preservation** Coalition

FITS – File Information Tool Set





- FITS: Content profiling tool
- Works with FITS to visualize FITS output metadata
- Setup can be a little more taxing



The background of the image is a dark, atmospheric landscape. It features a sky filled with heavy, swirling clouds in shades of dark blue, purple, and black. In the lower right corner, there is a dark, silhouetted shape that looks like the side of a mountain or a large hill. The overall mood is mysterious and foreboding.

TRUST NO ONE

Assume nothing, validate everything



Stuff happens

- Whenever a digital collection is moved, processed, curated or altered in any way.... things can go wrong!
 - Network dropouts at critical times
 - Disks get full, subsequent data copied there is lost
 - Software bugs lead to unexpected results
 - Human error leads to all sorts of issues
- Stuff happens a lot more at scale!