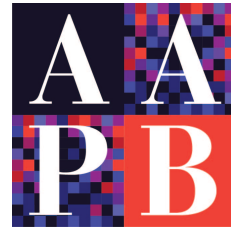# LTO Tapes at WGBH

and the American Archive of
Public Broadcasting

# The Project

- American Archive of Public Broadcasting (AAPB) – 40,000 hours of a/v content from around the country

- Partnership between WGBH and the Library of Congress to work with stations to deliver material, with Crawford Media Services as vendor

- At WGBH, 300 TB of digital a/v material (11,561 files) identified for inclusion in the American Archive

# The Data Transition

- WGBH backup data stored on LTO 4

- Administered through a SAMFS/QFS storage management system

- Metadata in Artesia Digital Asset Management system

Use of Artesia DAM being discontinued at WGBH – opportunity to pull media and data out of Artesia-SAMFS/QFS system & store it locally on the archive's dedicated LTO-6 tapes

# The Data Transition

- List of files characterized as media reported by Artesia

- SAMFS/QFS queried to generate report on which LTO-4 tapes those files were stored on

- Files organized into 'batches' for download to local storage based on reported storage location

- Files transferred remotely from LTO-4 via ssh and downloaded onto external hard drives

- Drives shipped to Crawford Media Services for transcode and upload to Archival Management System

# Initial Problems

- 57% of files in first large batch of video (2069 files) sent to Crawford proved to be incomplete or unreadable

- Several types of failure: 0-byte files, files that failed analysis by standard tools such as ffmpeg and mediainfo, files that could be analyzed but failed QC



QC Failure: File reads as normal when tested, but content is replaced with repeated glitching image or green screen a portion of the way through the video

# Deeper Analysis

- Most failed files transferred a certain amount of data successfully, then eventually cut out and replaced with nonsense 'filler' data



Media data for uncorrupted file

Media data at point of corruption

# Process Evolution

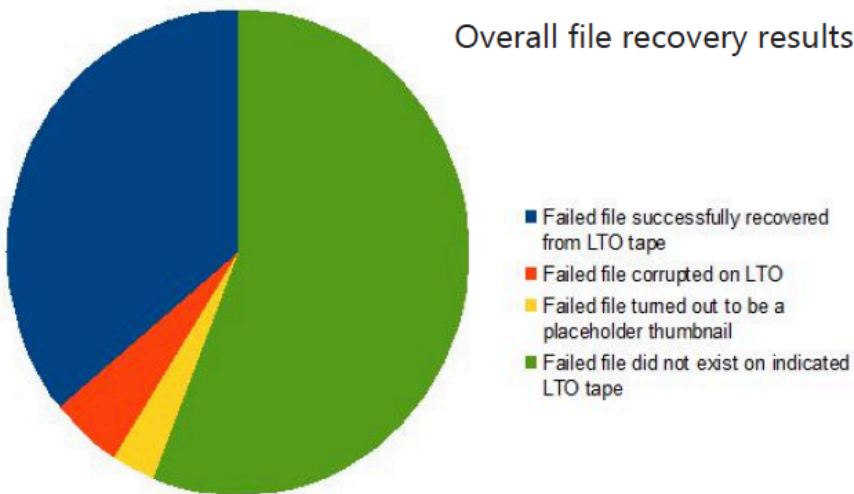- Each 'batch' of files was analyzed by ffmpeg immediately after download to determine which files failed

- Many files that initially failed could be successfully downloaded on a second or third try

- Other failed files (QC failure) were not detected until re-analysis after the drives returned to WGBH

- Final tally: 9957 out of 13492 were successfully re-ingested at WGBH -- 74% of files overall

# LTO Investigation

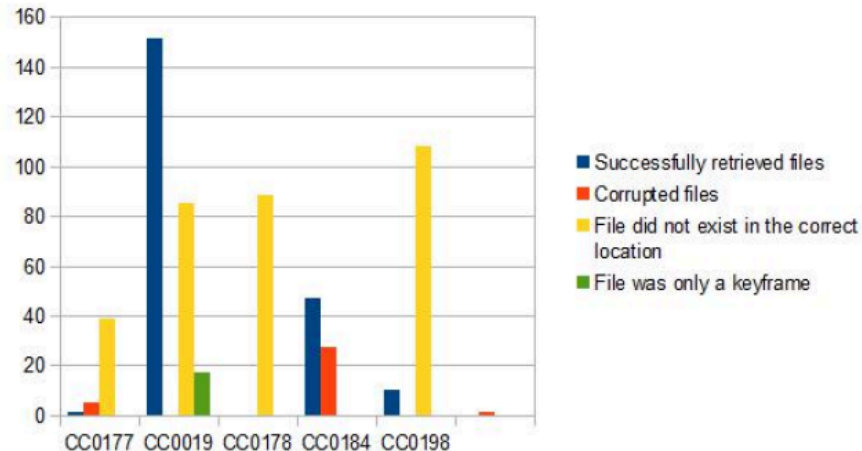- Ran checksums on all the files that had existing checksums recorded – 1005 failed files passed checksum analysis when checked in storage, 20 checksums could not be generated

- Generated analysis of failed files by LTO-4 storage tape and identified tapes with greatest number of failures, then requested then from IT

- Used MLA LTO-6 decks to do a direct data dump with dd of all the tar files written to the LTO-4 tapes

Overall file recovery results

- Failed file successfully recovered from LTO tape
- Failed file corrupted on LTO
- Failed file turned out to be a placeholder thumbnail
- Failed file did not exist on indicated LTO tape

Results by LTO tape

- Successfully retrieved files
- Corrupted files
- File did not exist in the correct location
- File was only a keyframe

# Inconclusive Conclusions

- Most files are still good on tape

- Some of the failures may be due to corruption on the LTO tapes themselves, but only a small percentage

- Processing problems may have been caused by inaccurate reporting from SAMFS/QFS

# Next Steps

- Re-ingesting successful files onto local LTO-6 through direct connection

- Tracking location in XML files entered into local Filemaker-based DAM system

- Re-attempting networked transfer of the other 26% and analyzing results

# Ideas or discussion welcome!

rebecca_fraimow@wgbh.org