

---

## Resources

Survey questionnaire is accessible at the following URL (as of 2012/04/01): [https://docs.google.com/document/d/1pP\\_qjPhlx-YjctcRvWsur2B2rBuV8DdoLCItTeoIBtg/edit](https://docs.google.com/document/d/1pP_qjPhlx-YjctcRvWsur2B2rBuV8DdoLCItTeoIBtg/edit)

Survey data (translated from the original Portuguese to English) is accessible at the following URL (as of 2012/04/01): <https://docs.google.com/spreadsheets/ccc?key=0Ai3bWoF9slIddGtKMVpEN293MVjxeVWd3WHFoUkppaFE>

Lucas do Rio Verde city government website: <http://www.lucasdoriverde.mt.gov.br> (accessed 2012-04-01).

## RECONSIDERING THE CHECKSUM FOR AUDIOVISUAL PRESERVATION: Detecting digital change in audiovisual data with decoders and checksums

Dave Rice

### Firstly, what are checksums for?

A checksum is small data value computed from a given amount of data, such as a file or bit-stream, for the purpose of facilitating the future ability to detect changes in that given data. The generation and verification of checksums for digital archival holdings is a central principle of digital preservation and enable archivists to trust that data held within an archive is the same data that was received by the archive. Although checksum wrangling is typically a behind-the-scenes process within digital storage systems and repositories, these values are worth a closer look. The checksum value is generally expressed in hexadecimal representation (aka base 16) comprised of the numbers 0 through 9 and the letters A through F. Several checksum algorithms, such as CRC32, MD5, and SHA-1, have been introduced offering varying degrees of processing efficiency, security, and collision resistance<sup>19</sup>. As an example, the CRC32 checksum value for a text file that contains the ASCII characters “checksum” would be de6dfd9a and the MD5 value for the same file would be 226190d94b21d1b0c7b1a42d855e419d. If the text file changed, whether through manipulation, bit rot, or data corruption, then further evaluations of the file would produce a different checksum value. The mismatch of a newly calculated checksum and a stored checksum produced earlier is an alert that data under care has been changed.

Within digital preservation environments, the generation and verification of checksums against digital files can aid the confirmation or denial of digital authenticity. A checksum mismatch is an alert that a file under care has changed from a prior state; potentially triggering retrieval of backups, review of hardware, or migration of content. Generally, if a given checksum algorithm is applied to a file, then as long as the same checksum can be regenerated from the file then the data is verified, else a mismatched checksum reveals a digital change. Further details such as the whereabouts, extent, or significance of the change in data are not revealed by the checksum mismatch but only that the data examined now is not the same as the data examined before.

A checksums alone is not very self-descriptive or actionable. As checksums are evaluated from digital files, the PREMIS data dictionary requires that the checksum value be stored alongside the name of the algorithm that was utilized (such as MD5) and optionally with the name of the original producer of the checksum<sup>20</sup>. Metadata standards such as reVTMD<sup>21</sup> also recommend that the date and time that the checksum was evaluated is store with the checksum as well. This information provides further context to enable future use of the checksum.

For digital preservation activities, the application of checksums is an essential recommendation; however the checksum of a file offers only a yes or no response to the question “Has this file changed since it produced an earlier documented checksum?” Within a digital audiovisual collection if the files has changed, follow-up questions could include:

- “Does the change affect the presentation of audiovisual data?”
- “Where is the change?”
- “How extensive is the change?” or
- “Do these two different files present the same audiovisual data, i.e. is one file a loss-less representation of the other?”

19 While this paper will not go into detail regarding collision resilience further information is available at [http://en.wikipedia.org/wiki/MD5#Collision\\_vulnerabilities](http://en.wikipedia.org/wiki/MD5#Collision_vulnerabilities)

20 See the PREMIS Data Dictionary for Preservation Metadata at <http://www.loc.gov/standards/premis/v2/premis-2-1.pdf>.

21 See the reVTMD metadata schema at <http://www.archives.gov/preservation/products/reVTMD.xsd>.

---

Additionally are the checksums of digital preservation best practices as effective when dealing with the large file sizes of audiovisual data? Should checksum strategies be extended for the preservation of video? Furthermore, what other roles have checksums played in audiovisual data? This paper will examine the relationship between audiovisual data and the objectives of checksums, will review tactics to allow the effectiveness and function of checksums to scale for audiovisual content, and will review the history of checksums in use within digital video to enable error detection and error concealment.

### **Little Checksums for Large Files**

Often in digital preservation and repository design, the whole of each digital file is documented by a checksum. One file gets one checksum. A digital collection of office documents, emails, and images may contain hundreds of thousands of files in the space of a terabyte. With the large file sizes associated with archival audiovisual data a similar terabyte may contain only tens of files. The effectiveness and function of checksums do not scale for collections with larger file sizes. In a collection of documents a checksum mismatch may identify an issue within the space of a megabyte whereas in a collection of audiovisual material the checksum mismatch may potentially make suspect dozens of gigabytes occupied by a single file. As a real world analogy a checksum mismatch of an electronic document could compare to a fire alarm triggering in an office building. The fire department is alerted to the address and floor at risk. A checksum mismatch on an archival audiovisual file covers a vast area of data – like triggering a fire alarm that alerts firefighters only to the general neighborhood.

White papers and best practice documents for the digital preservation community recommend to regularly checksum our digital collections and identify mismatches in order to monitor digital integrity, but advice on how to respond to the mismatch is hard to find. There may be two copies of the digital material, one with a checksum mismatch and one without. A stable one could replace the changed one, but a backup may not always be available. An audiovisual file with digital corruption could still look and perform as well as the original or the change may affect the container or key codec attributes that cause the process of demultiplexing or decoding the file to fail completely.

A changed bit will affect audiovisual data very differently depending on the codecs involved. A tweaked bit within uncompressed audio and video would most likely be imperceptible to notice, changing at most one pixel or one audio sample. Audiovisual data with compression, whether lossless or lossy, could be affected dramatically; depending on the type of codec a changed bit could alter anywhere from a section of a frame to several frames of video<sup>22</sup>.

### **Checksums Within Audiovisual Data**

Outside of digital preservation environments, checksum functions and features have been integrated into many audiovisual codecs and containers to suit different purposes. There are checksum applications existing within some types of audiovisual data that do enable a more specific response than whole file checksums can allow. Some of these applications deal with error detection and error concealment and can identify which frames have changed in order to minimize the visual effect of the error during presentation. Other applications of checksums enable the confirmation of whether a losslessly transcoded file represents the original encoded data authentically or confirm whether the contained audiovisual data is maintained accurately even when other embedded metadata tags within the same file have been added or altered.

---

22 Resilience of various audiovisual codecs and containers to digital corruption can be tested with an input fuzzer. See, for example the zzuf application (use cautiously): <http://caca.zoy.org/wiki/zzuf>

---

## DV Tape and DIF Block Parity Data

DV tapes, such as DVCam, miniDV, and DVCPro, contain encoded digital streams of audiovisual data. Given the small size of the digital tape, damage to the data held on the tape is commonplace and the DV codec and DV hardware are designed to anticipate handling damaged data through video error concealment and controlled audio dropouts. When a DV tape is recorded, audiovisual data is divided into DIF blocks and written onto the tape digitally. The data of each block is analyzed through an algorithm to produce a very small checksum (aka parity data) that is also written with the DIF block. With this parity data, a DV deck can read the DV tape and verify if the data is read properly or identify audiovisual data that mismatches with its corresponding parity data. Rather than play damaged or unexpected data, the DV hardware will either copy from the corresponding DIF block of the last valid frame or drop out the audio. In these cases the digital damage to the tape is corrected but is concealed to minimize the visual or aural impact and make the damage less noticeable. Storing parity data during recording to tape and verifying parity data during reading from tape, allows DV hardware to incorporate restoration principles in playback so that inaccurate audiovisual data may be concealed with an approximate of the data lost. If the DV data is transferred from the tape to a file without undergoing any transcoding the resulting file-based DV stream will contain data that notes the location of frames and DIF blocks that contained checksum mismatches<sup>23</sup>. Since DV on a tape contains one checksum per video block of a frame, this data can enable the preservationist to know how extensive any digital damage is and where it is located. Although the checksum, in this case, can not correct the data the use and assessment of these mismatches can inform the archivist's response so that the tape may be transferred under adjusted circumstances in hope to produce a more digitally accurate data transfer. In this case the embedded parity data within a DV stream, evaluated at the time of the recording, enables a more precise response to the digital damage than with a whole file checksum.

## MPEG and CRC Checks

MPEG streams are designed to decode gracefully even when interference or digital corruption hinders the transmission process. Formats such as MP3 or MPEG-2 transport streams may contain CRC (cyclic redundancy check) checksums that verify the data of the previously transmitted packet of audiovisual data. When an MPEG decoder, such as a software application or digital television, receives a section of data that does not validate to the attached CRC then the decoder can discard the invalid data and await the next valid section. The CRC checksums may occur frequently within an MPEG stream and mismatches between embedded CRC and a newly interpreted one make it feasible to note the extensiveness or location of digital change or corruption. MP3val<sup>24</sup> is one such validator for mpeg audio streams that make use of CRC checksums. The MP3val manual notes of the

“MPEG is a streamable format, that is, it is optimized for quick and easy recovery from errors. MP3 decoders are very tolerant to inconsistencies in the input file. Most players even don't report to user about stream errors. So, as a rule, the user doesn't know whether his files are valid or broken. But using broken files can eventually lead to problems during playback on certain software/hardware.”<sup>25</sup>

## Lossless FLAC and MD5 Signatures

FLAC, Free Lossless Audio Codec, uses another approach for embedded checksums. FLAC is a lossless format that compresses raw audio data to a smaller size while maintaining the ability to decode to the exact same audio samples that it encoded. Thus a 300 megabyte BWF

---

23 See the DV Analyzer project by AudioVisual Preservation Solutions which can evaluate this data: <http://www.avpreserve.com/dvanalyzer/>

24 <http://mp3val.sourceforge.net/>

25 <http://mp3val.sourceforge.net/docs/manual.html> accessed on April 1, 2012

---

(Broadcast Wave Format) audio file may be losslessly compressed to a FLAC audio file at about 100 megabytes. As the lossless nature of the codec implies, the FLAC file should then decode to the exact same audio samples as the original BWF. The header of a FLAC file contains an MD5 checksum or signature that represents the original audio data that is encoded. Since FLAC is a lossless format this MD5 also represents the audio data that should be decoded from the FLAC file.

As an example, suppose a Broadcast Wave Format file called “original.wav” is processed with Ffmpeg<sup>26</sup>, an open source set of tools to process audiovisual data. The raw audio data of the file can be decoded and sent to a checksum utility. With ffmpeg installed on a computer a command could be entered within a terminal window to decode the audio samples of the BWF file and evaluate the samples with the MD5 checksum algorithm, like this<sup>27</sup>:

```
ffmpeg -i original.wav -f md5 decoded_wav.md5
```

The resulting file “decoded\_wav.md5” would contain the following text:  
MD5=6d0b18032fd23dfbf400b092eb9642f7

Using the compress utility called flac, also available as a command line application, the Broadcast Wave Format file “original.wav” may be converted to a lossless FLAC file like this<sup>28</sup>:

```
flac original.wav --keep-foreign-metadata -V -o lossless.flac
```

Within the header of the resulting FLAC file, “lossless.flac”, is contained the MD5 signature of the original audio data: 6d0b18032fd23dfbf400b092eb9642f7 (the same checksum as produced when applying a checksum to the decoded audio data of the original file).

Although both files, “original.wav” and “lossless.flac”, decode to identical audio data the two files are very distinct structurally. The whole-file checksums of the two files would be different and the FLAC file, “lossless.flac”, would be a third of the size of “original.wav”; however both files should both decode to identical audio presentations. FLAC’s embedded checksum offers a functional advantage of the FLAC file over the BWF file in that the specification’s incorporation of the embedded MD5 checksums mean that the audio data may be verified without additional external checksum workflows. The encoding of audio data and generation of a checksum occurred in the same action and are packaged together in a single file.

The wiki at etree.org offers further explanation on the utility of using embedded checksums to document what the lossless encoded data should decode to rather than to document what the data actually is.

*A FLAC Fingerprint is generated only for the audio data portion of the file. (Therefore, changing the filename or the tags or FlacMetadata does not change the fingerprint calculation.) In contrast, an .md5 is generated against the whole file, including header portions. ... Under FLAC, you are allowed to change the compression ratio and add/remove meta data to .flac files without changing the actual audio. The audio may be identical, but the extra data will completely change the .md5 checksum. Checking these*

---

26 <http://ffmpeg.org>

27 Documentation and downloads for Ffmpeg are available at <http://www.ffmpeg.org>. This command means the “original.wav” file is used as the input (-i) and to be processed into the “md5” format (-f) then the output is directed to the file called “ decoded\_wav.md5”. With Ffmpeg’s default use of the md5 format the first audio and video streams of the input (in this case just a single PCM audio stream) are uncompressed to raw video and signed 16 bit audio to be processed by the md5 muxer.

28 The flac utility and documentation are available at <http://flac.sourceforge.net/>. This command converts the input file, “original.wav”, into the output file (-o) “lossless.flac”. The “-V” or “--verify” command is used to decode the encoded stream parallel to the encoding process to double-check the losslessness of the transcoding. If FLAC is encoded and decoded with the “--keep-foreign-metadata” then the FLAC may also maintain non-audio data from the source file, such as a bext chunk or other types of embedded metadata.

.md5s against the new .flac files will report failure, even though there is nothing actually wrong with the new fileset. That can cause major confusion.<sup>29</sup>

## Codec Independent Implementations of Frame Based Checksums

FFmpeg incorporates formats called framecrc<sup>30</sup> and framemd5<sup>31</sup>. These formats were integrated in order to facilitate testing functions such as verifying that an adjusted decoder maintains intended results or that an FFmpeg decoder decodes a proprietary stream to the same data as a decoder by another vendor. Although originally intended as a testing format, the framemd5 format helps enable some of the same goals performed by checksums in the cases of DV, MPEG, and FLAC by providing a checksum of each decoded frame.

By producing checksums on a more granular level, such as per frame, it is more feasible to assess the extent or location of digital change in the event of a checksum mismatch. In recalling the analogy of a whole-file checksum for audiovisual data as a neighborhood file alarm, the additional generation of more granular checksums enables tracking of digital change or alarms with greater precision.

By decoding a file and processing the decoded data to generate a framemd5 document, each decoded audio and video frame is documented according to its timestamp, digital size, and MD5 checksum. For the first three frames of video, the framemd5 output could be<sup>32</sup>:

```
#tb 0: 1001/24000
0,    0,    0,    1,  518400, 5bc19af1a75adb8bda9d79390981a0ea
0,    1,    1,    1,  518400, bb485b0d6fd001358aa7dbe76031ff4d
0,    2,    2,    1,  518400, 30dc414cd4487dd58b0d16a5ddafba35
```

In this output the columns refer to the stream number, counting from zero, (column 1), the decoding and presentation timestamps (column 2 and 3), the samples duration (column 4), the size of the data checksummed in bytes<sup>33</sup> (column 5), and the MD5 checksum for that data.

Storing a framemd5 file along with each audiovisual file can not replace the function of a traditional whole-file checksum. That is, it is useful to have both the fire alarm for all the buildings and floors within a neighborhood, as well as the fire alarm for the whole neighborhood. It is still possible for a file to be changed in a way that would result in a mismatch for a future whole-file checksum analysis, but not create any difference between a stored framemd5 output and a newly created framemd5 output. This could occur when embedded metadata is edited but the stored audiovisual data remains the same.

For audiovisual data, storing both a whole-file checksum and a framemd5 output enables greater awareness of digital change in managed files, a more strategic and aware response to change, and the ability to verify lossless transcoding. If an audiovisual file is found to have a mismatch between a newly generated whole-file checksum and one generated previously, indicative of digital change, then comparison between a stored framemd5 document and a newly generated one could facilitate in pinpointing the digital change as it affects audiovisual presentation if at all.

29 <http://wiki. etree.org/index.php?page=FlacFingerprint> accessed on April 1, 2012

30 <http://ffmpeg.org/ffmpeg.html#framecrc>

31 <http://ffmpeg.org/ffmpeg.html#framemd5>

32 The output framemd5 is produced by using this FFmpeg command: "ffmpeg -i a\_movie\_file.mov -f framemd5 a\_movie\_file.framemd5". The (-i) refers to the input file, (-f) refers to the output frame which in this case is 'framemd5' and then 'a\_movie\_file.framemd5' is the name of the output file.

33 518400 is the number of bytes used for a frame of uncompressed YUV 4:2:0 video at 720x480.

---

## Whole-File Checksum plus framemd5

### Determining Extent and Relevance of Digital Change

Suppose a QuickTime file is documented by both a whole-file checksum and a framemd5 output and then a user inadvertently opens the QuickTime file, clicks File>Save, and then closes the file. In this process the QuickTime file will remain largely the same except that an embedded timestamp will update to the last saved date changing a couple bytes of the file. Upon the generation of a new whole-file checksum on the QuickTime file there will be a mismatch between the new checksum and the stored one, confirming the change. This checksum mismatch notes that an issue exists but offers little clue to what it is. If the QuickTime file is again used to generate a new framemd5 output, it will show that the decoded frames of the QuickTime file are identical to what was decoded before the change. With this information the nature of the change can be clarified; in this case digital change did occur but did not affect the audiovisual presentation produced by the file.

Within another scenario, a QuickTime file is migrated from one server over a wireless connection to another server and within the data transmission process a few bits are inaccurately copied. In this scenario too, the resulting whole-file's checksum would mismatch the checksum of the original file. Upon comparison of framemd5 documents between the original file and the copy, individual lines of the framemd5 documents would mismatch, revealing exactly which frames were affected by the changes and informing how extensive the change is.

The use of lossless codecs alone does not guarantee that the resulting encoded lossless audiovisual file could be used to reconstruct the original audiovisual data. A preservation-suitable lossless audiovisual encoding should decode to the same data that the original source would decode to, meaning that each pixel, frame, and timing decoded from the lossless version should be the same as the decoded original.

An original uncompressed digital audiovisual file called "uncompressed.mov" could produce this framemd5 output (the first four video frames are listed):

```
0, 0, 0, 1, 518400,6a6b8be7dbac428b86669992cc740d10
0, 1, 1, 1, 518400,6a6b8be7dbac428b86669992cc740d10
0, 2, 2, 1, 518400,2eb3a9a53f42d5c6b21e3b6b7e267414
0, 3, 3, 1, 518400,e67d03e253f3997ae2db46aa28d8c749
```

Let's imagine that an archive decides to transcode the uncompressed file to a lossless codec in order to reduce storage requirements. This `ffmpeg` command generates a lossless `ffv1` encoding from the original file and copies the audio data as-is.

```
ffmpeg -i uncompressed.mov -map 0 -c:v ffv1 -c:a copy lossless.mov
```

Generating a framemd5 report on the "lossless.mov" should produce the same output because both files, although utilizing different codecs, both decode to identical audiovisual presentations. If the two files do not decode to produce identical framemd5 documents then it is likely that the transcoding from the uncompressed codec to the lossless codec was not truly mathematically lossless<sup>34</sup>.

---

34 Within some transcoding environments it can be easy to inadvertently cause a change in colorspace, pixel format, chroma subsampling or other technical audiovisual specifications between the input and the output that would compromise the losslessness of the process although the resulting codec itself may be lossless.

---

## Under Development

At the time of this writing, the technical documentation<sup>35</sup> and a new version of FFV1, FFmpeg's lossless video codec, are under development. Version 1.3 of FFV1, currently marked as an experimental encoder, adds a mandatory CRC checksum to each frame header. This feature facilitates lossless encoding of video data in a manner that binds checksum protections into the encoding process. This feature allows digital lossless video collections to effectively assess holdings to identify stored data that mismatches the data originally encoded, enabling a lossless format to be self-checking.

## In Summary...

The preservation of audiovisual collections requires optimization of preservation systems and practices that can scale to and endure the large file sizes required. Once within the territory of audiovisual data, the traditional whole-file checksums can appear to be slow and vague. In looking as well to checksum applications found within compression technology we can identify tactics that enable precise responses to data corruption in more granular ways. While the traditional users of the embedded checksums of DV, MPEG, and FLAC have been decoding hardware and software, the functions and opportunities of these more granular checksum approaches serve needs within digital preservation environments as well. By utilizing FFmpeg's framemd5 format those advantages do not need to be restricted to a particular set of codecs. Utilizing frame-based or more granular checksumming in addition to whole-file checksums allows for more targeted discovery and could provide the archivist with an additional tool to analyze (and potentially repair) digital corruption within archival collections.

Thanks to:  
Michael Niedermayer, FFmpeg  
George Blood, George Blood, L.P.  
Jimi Jones, Library of Congress

This work is licensed under the Creative Commons Attribution-NoDerivs 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nd/3.0/> or send a letter to Creative Commons, 444 Castro Street, Suite 900, Mountain View, California, 94041, USA.

---

35 Technical documentation is being drafted here: <https://github.com/FFmpeg/FFV1>, accessed April 1, 2012