

PreForma MediaInfo

Project Acronym: PREFORMA

Grant Agreement number: 619568

Project Title: PREservation FORMAts for culture information/e-archives

Prepared by:

- MediaArea SARL
- Jérôme Martinez
- Dave Rice
- Tessa Fallon
- Ashley Blewer
- Erik Piil
- Guillaume Roques

Prepared for:

Date: December 31, 2014

Licensed under: Creative Commons CC-BY v4.0

Summary: This reports serves as a snapshot of MediaArea's research, planning, and design work to develop a conformance checker, tentatively entitled PreForma MediaInfo.

- [Project Introduction](#)
- [Introduction of Featured Formats](#)
 - [Matroska](#)
 - [FFV1](#)
 - [Linear PCM](#)
- [Development of a Conformance Checker](#)
 - [Implementation Checker](#)
 - [Policy Checker](#)
 - [Reporter](#)
 - [Fixer](#)
 - [Shell](#)
 - [Optimization for Large File Size](#)
 - [Focus on Fixity](#)
 - [Reference and Test Files](#)
- [Intended Behavior by Use Case](#)
 - [Overview](#)
 - [Conformance Checking at Digitization Time](#)
 - [Conformance Checking at Migration Time](#)
- [Open Source Ecosystem](#)
 - [Cross Platform Support](#)
 - [Online Resources](#)
 - [Advance Improvement of Standard Specification](#)
- [Architectural Layers](#)

- Transport layer
 - Preforma MediaInfo: File on disk or direct memory mapping
 - Plugin integration proof of concept: libURL
- Container/Wrapper implementation checker
 - Preforma MediaInfo: Matroska checker
 - Plugin integration proof of concept: mkvalidator
- Container/Wrapper Demultiplexing
 - Preforma MediaInfo
 - Plugin integration proof of concept: FFmpeg
- Stream/Essence implementation checker
 - Preforma MediaInfo:
 - Plugin integration proof of concept: jpylyzer
 - Plugin integration proof of concept: DV Analyzer
 - Optional
- Stream/Essence decoder
 - Preforma MediaInfo
 - Plugin integration proof of concept: FFmpeg
 - Plugin integration proof of concept: OpenJPEG
- Container/Wrapper vs Stream/Essence Coherency Check
 - Preforma MediaInfo
- Baseband Analyzer
 - Preforma MediaInfo
 - Plugin integration proof of concept: QCTools
- Hypervisor (Controller)
- Human interface
- Style Guide
- Source Code Guide
 - Portability
 - Modularity
 - Deployment
 - API's
- Open Source Practices
 - Development
 - Open Source Platforms
- Contribution Guide
 - File Naming Conventions
 - Rules for Qt/C++ code:
 - Rules for contributing code
 - Rules for contributing feedback
 - Linking
- License
- Checker Design: Conformance and Coherency
- Conformance Check Registry

- Conformance Check Registry Requirements
- Elements
- Matroska Conformance Checks (Draft)
 - Extension Test
 - Extension Test MKV
 - Extension Test MKA
 - Extension Test MKS
 - Extension Test MK3D
 - EBML Element Start
 - EBML vint efficiency
 - Element ID Registered
 - Element Size 0x7F Reservation
 - Element Size Byte Length Limit
 - Element Size Unknown
 - Level 0 Segment
 - Only One EBML Header recommended
 - File Size Consistency
 - EBMLVersion Presence
 - EBMLReadVersion Presence
 - EBMLMaxIDLength Presence
 - EBMLMaxSizeLength Presence
 - DocType Presence
 - DocTypeVersion Presence
 - DocTypeReadVersion Presence
 - EBML Version Coherency
 - EBMLMaxIDLength Limits
 - EBMLMaxSizeLength Limit
 - EBMLMaxSizeLength Matches
 - DocType
 - DocTypeVersion Coherency
 - DocTypeVersion Limits
 - Top Elements Coded on 4 Octets
 - CRC Order
 - CRC-32 Size Coherency
 - CRC Validation
 - CRC Not Pointlessly Used
 - CRC-Presence
 - Single Segment Composition
 - Seek-Presence
 - SeekID-Presence
 - SeekID-Type
 - SeekPosition-Presence
 - Segment-Info-Presence
 - SegmentUID-Range
 - SegmentUID-Size
 - SegmentUID-Type
 - SegmentFilename-Type
 - PrevUID-Size
 - PrevUID-Type

- PrevFilename-Type
 - NextUID-Size
 - NextUID-Type
 - NextFilename-Type
 - SegmentFamily-Size
 - SegmentFamily-Type
 - TimecodeScale-Presence
 - Duration-Range
 - Duration-Type
 - DateUTC-Type
 - Title-Type
 - Tag Total Parts
 - Tag Part Number
 - Tag Part Offset
 - Tag Title
 - Tag Subtitle
 - Tag URL
 - Tag Sort_with
 - Tag Email
 - Tag Address
 - Tag Fax
 - Tag Phone
 - Tag Initial_Key
 - Tag Law_Rating
 - TAG ICRA
 - Tag DATE_RELEASED
 - Tag DATE_RECORDED
 - Tag DATE_ENCODED
 - Tag DATE_TAGGED
 - Tag DATE_DIGITIZED
 - Tag DATE_WRITTEN
 - Tag DATE_PURCHASED
 - Tag Play_Counter
 - Tag FPS
 - Tag BPM
 - Tag Measure
 - Tag Tuning
 - Tag Replay Gain (Gain)
 - Tag Replay Gain (Peak)
 - Tag (Identifiers) ISRC
 - Tag (Identifiers) MCDI
 - Tag (Identifiers) ISBN
 - Tag (Identifiers) Barcode
 - Tag (Identifiers) Catalog number
 - Tag (Identifiers) Label code
 - Tag (Identifiers) LCCN
 - Tag (Commercial) Purchase Item
 - Tag (Commercial) Purchase Price
 - Tag (Commercial) Purchase Currency
- FFV1 Conformance Checks (Draft)

- Missing header
- version
- version 2
- micro_version 2
- coder_type
- state_transition_delta
- colorspace_type
- bits_per_raw_sample
- (More header tests to add)
- crc_parity
- LPCM Conformance Checks (Draft)
 - formatType
 - bitsPerSample
 - bytesPerSecond
 - blockAlignment
 - channelCount
 - nChannels
 - sampleRate
- Container/Stream Coherency Checks (Draft)
 - Aspect Ratio Match
 - Width Match
 - Height Match

Project Introduction

The PreForma challenge illuminates and responds to a significant and real obstacle that faces the preservation community today. This report encompasses a snapshot of MediaArea’s design plans to create a toolset (tentatively entitled “PreForma MediaInfo”) as the conformance checker, policy checker, reporter, and fixer of a select list of formats.

As preservation workflows have incorporated digital technology, significant amounts of careful research have gone into the selection of file format recommendations, lists of codec specifications, and development of best practices; however, despite the existence of such recommendations, there remains a lack of assessment tools to verify and validate the implementation of such recommendations. A few validation tools (such as mkvalidator) are produced alongside the development of their associated standards; however, most file format specifications are not officially tied to any validation tool and are documented through human-readable narrative without equivalent computer-actionable code. Where a metadata standard may be described in both a data dictionary and a computer usable XML Schema, file formats standards often lack a computer-useable verification method. The PreForma project recognizes this discrepancy and the resulting long-term impacts on archival communities and seeks to fill in the gaps necessary to provide memory institutions with levels of control to verify, validate, assess and repair digital collections.

MediaArea’s approach to this challenge centers on Free Software, modular design, and interoperability and will rely strongly on MediaInfo (an open source MediaArea product) to meet this challenge. MediaInfo is often advised as the first tool to use when a media file is not playable, allowing the user to identify characteristics that would help find an appropriate playback or transcoding tools. MediaInfo’s open licensing and agility in technical metadata reporting have encouraged its integration into several archival repository systems and OAIS workflows to assist archival with technical control of collections.

MediaArea sees community involvement as a key factor of evaluating the success of the project. To encourage this, during the prototype phase MediaArea will perform the development work for command line utilities,

graphical user interfaces, and documentation in publicly accessible repositories at github.com. We will set up an online set of project resources such as public access to a corpus of test media, an IRC channel, and a responsive public issue tracker.

Introduction of Featured Formats

During the development phases MediaArea will focus on one container format, Matroska, and two streams, LPCM and FFV1. The design work of MediaArea will address formats and codecs through a modular architecture so that other formats or codecs may easily be added alongside or after development.

Matroska, FFV1, and LPCM describe very unique concepts of information including: - a container format, Matroska - an audio stream, LPCM - a video stream, FFV1

These three information concepts will inform distinct user interface and reporting design in order to process these concepts through differing strategies. For instance reporting on the status of 100,000 of video frames within a video recording may be done more efficiently in a different interface as one designed to communicate the hierarchical structure of a file format.

Additionally other formats currently being addressed by PreForma fit within these three conceptual categories; for instance, PDF and TIFF are formats (containers) and JPEG2000 is a video stream. These three concepts affect the design of an overall application shell as conformance information for each category can have its own optimized user interface.

Matroska

Matroska is a open-licensed audiovisual container format with extensive and flexible features and an active user community. The format is supported by a set of core utilities for manipulating and assessing Matroska files, such as `mkvtoolnix` and `mkvvalidator`.

Matroska is based on EBML, Extensible Binary Meta Language. An EBML file is comprised of one of many defined “Elements”. Each element is comprised of an identifier, a value that notes the size of the element’s data payload, and the data payload itself. The data payload can either be stored data or more nested elements. The top level elements of Matroska focus on track definition, chapters, attachment management, metadata tags, and encapsulation of audiovisual data. The Matroska element is analogous to QuickTime’s atom and AVI’s chunk.

Matroska integrates a flexible and semantically comprehensive hierarchical metadata structure as well as digital preservation features such as the ability to provide CRC checksums internally per selected elements.

FFV1

FFV1 is an efficient lossless video stream that is designed in a manner responsive to the requirements of digital preservation. Version 3 of this lossless codec is highly self-descriptive and stores its own information regarding field dominance, aspect ratio, and colorspace so that it is not reliant on a container format to store this information (other streams that rely heavily on its container for technical description often face interoperability challenges).

FFV1 version 3 mandates storage of CRCs in frame headers to allow verification of the encoded data and stores error status messages. FFV1 version 3 is also a very flexible codec allowing adjustments to the encoding process based on different priorities such as size efficiency, data resilience, or encoding speed.

Linear PCM

Linear Pulse Code Modulation (LPCM) is a ubiquitous and simple audio stream. PCM audio streams may be comprised of signed or unsigned samples, arranged in little-endian or big-endian arrangements, in any number

of audio channels. PCM is very flexible but is not self-descriptive. A raw PCM file can not properly be decoded without knowing the sample encoding, channel count, endianness, bit depth, and sample rate. PCM is typically dependent on its container format (such as WAV) to store sufficient metadata for its decoding.

Because PCM streams contain only audio samples without any codec structure or metadata within the stream, by itself any data could be considered valid PCM and decoded as audio. Determining the conformity or technical health of PCM data requires the context of information provided by its container format.

Development of a Conformance Checker

The design of the conformance checker is intended to allow interoperability between the conformance checkers with PreForma's other suppliers so that users may integrate multiple conformance checkers within a single 'shell'. The conformance checker is comprised of four components:

- implementation checker
- policy checker
- reporter
- metadata fixer

The conformance checker developed within the PreForma project must document and associate conformance rules with data types (such as formats, streams, frames, etc) and authorities (such as specifications, community practices, or the local rules of a memory institution). MediaArea recommends that communication between the conformance checker and the shell be performed through an API designed via collaboration of the PreForma suppliers.

Implementation Checker

Each implementation checker (for Matroska, FFV1, and LPCM) should assess compliance and/or deviation between files and a series of conformance checks which are written by dissecting rules and logic from each format's underlying specifications, including rules that may be deduced or inferred from a close reading of the specification or related authoritative documentation. MediaArea has drafted registries of conformance rules within the PreForma design phase and plans to collaborate with each format's specification communities to refine them. See the [Conformance Check Registry](#).

For streams such as FFV1 some implementation checks may be performed frame-by-frame to discover frame-specific issues such as CRC mismatches, invalid frame headers, or incomplete frames. Frame-by-frame implementation assessments will naturally be time consuming as nearly the entire file must be read. In order to accommodate user's various time priorities the implementation checker will use options to perform implementation checks on the first few frames of a stream, a percentage of the frames, or all of the frames.

MediaArea has drafted a [registry](#) of metadata elements to be used in described an implementation conformance check, which provides unique identifier, the scope, and underlying rationale and authority for the check. Code created to perform implementation checks will be internally documented with references to conformance check's unique identifiers, so that MediaArea may create resources for each conformance check that relate the identity of the check, its underlying authority, sample files, and associated code.

Policy Checker

For each format addressed through a conformance checker MediaArea will create a vocabulary list of technical metadata and contextual descriptions. Additionally MediaArea will define a list of operators to enable various comparators between the actual technical metadata and the user-provided expected metadata. Such defined language will allow users to make policy check expressions such as:

- FFV1.gop_size MUST_EQUAL “1”
- FFV1.slice_crc MUST_BE_ENABLED
- FFV1.version GREATER_THAN_OR_EQUAL “3”
- MKV.tag.BARCODE MUST_START_WITH “ABC”
- MKV.tag.DATE_DIGITIZED IS_BEFORE “2014-01-01”
- MKV.tag.ISBN MATCHES_REGEX “(?=[-0-9xX]{13})(? : [0 - 9] + [-])3[0 - 9] * [xX0 - 9]”

MediaArea proposes that PreForma suppliers collaborate to define a common expression for sets of policy checks via an XML Schema, associated data dictionary, and vocabulary of comparative operator. The collaboration would include agreement and definition on the operators (“Greater Than”, “Starts With”, etc) of the policy checks and attempts to normalize technical metadata between common formats where they have overlapping concepts. Each policy checker would produce a vocabulary of technical metadata specific to its format for policies to be checked against as well as inclusion within an API so that the Shell can access the possible operators of any enabled conformance checker.

MediaArea will provide sample sets of policy checks based on interviews with memory institutions and community practice.

Reporter

MediaArea proposes that PreForma suppliers collaborate to create a common XML Schema to define the expression of PreForma reporting (referred to here as “PreFormaXML”). The schema should define methods to express technical metadata and relates checks to formats/streams (including components of formats and streams such as frames or attachments). The XML Schema should encompass not only information about the files being assessed but also the conditions and context of the particular use (which shell was used, with what policy sets, at what verbosity, etc). The XML Schema should be supported by a data dictionary that is also collaboratively written by PreForma suppliers. MediaArea anticipates that the implementations and features performed upon the basis of a common XML Schema may vary from supplier to supplier or per conformance checker, but that adherence to a common schema is essential to interoperability and consistent user experience amongst conformance checkers.

The PreFormaXML schema should accommodate the expression of results from multiple conformance checkers upon a single file. For instance a Matroska file that contains a JPEG2000 stream, an FFV1 stream, and an LPCM stream should be able to express one XML element about the file with sub-elements about each conformity check to reduce redundancy.

MediaArea plans to include these features commonly within MKV, FFV1, and LPCM reporters:

- Export of a standardized PreFormaXML
- Export PreFormaXML with gzip compression (to reduce the impact of large and highly verbose XML files)
- Export of the same data within a semantically equivalent JSON format
- Other functions based on PreFormaXML (such as generation of PDF formats or summarization of multiple collections of PreFormaXML) will happen within the “Shell” component

Fixer

MediaArea will produce a fixer that allows for editing the file. Enabling this function will be performed with a substantial amount of caution as in some cases a user could use it to change a file considered a preservation master. The fixer will support assessing a file first to determine the risk of editing a structurally unhealthy file and provide suitable levels of warning to the user.

The metadata fixer shall support both direct editing on the input file (with warning) or producing a new output file as a copy that the metadata change as requested.

The metadata fixer will support comprehensive logging of the change and offer options to log the performance of the edit itself with the file if it has a means to accommodate it (such as Matroska).

In addition to metadata manipulation the fixer will accommodate structural fixes to improve the structural health of the file, such as repairing Matroska Element order if ordered incorrectly, or validating or adding Matroska CRC Elements at selected levels, or fixing EBML structures of truncated Matroska files.

Substantial care should be exercised to ensure that the Conformance Checker properly associates risk, user warnings, and assessments with each fix allowed. In order to allow a fix the software must properly understand and classify what may be fixed and be aware of how the result may be an improvement. Adjustments directly to a preservation file must be handled programmatically with great caution with diligent levels of information provided to the user.

An example of a fix that could be enabled in the RIFF format could be verifying that any odd-byte length chunk is properly followed by a null-filled byte. Odd-byte length chunks that do not adhere to this rule cause substantial interoperability issues with data in any chunk following the odd-byte length one (this is particularly found in 24 bit mono WAV files). If the odd-byte length chunk is not followed by a null-filled padding byte, then most typically the next chunk starts where the padding byte is and the padding byte may be inserted so that other following chunks increase their offset by one byte. This scenario can be verified by testing chunk id and size declaration of all following bytes so that the software may know beforehand if the fix (inserting the null-filled padding byte) will succeed in correcting the RIFF chunk structure's adherence to its specification.

Fixes for Matroska files could include fixing metadata tags that don't include a SimpleTag element or re-clustering frames if a cluster does not start on a keyframe.

Because many files focused on with FFV1 and Matroska conformance checkers will be quite large, MediaArea plans to provide options to either rewrite the original file with the check or edit the file in place so that the file is only changed according to the fix that is request. With the latter option is the user is 'fixing' the metadata in a 50 gigabyte Matroska file only the last few megabytes of the Matroska tagging element may be rewritten without a requirement to rewrite the non tagging elements at the beginning of the file (MediaArea deployed a similar feature within BWF MetaEdit).

Shell

The Shell shall coordinate the actions of the implementation checker, policy checker, reporter and fixer. As PreForma seeks that the Shell developed by each supplier supports each supplier's conformance checker(s), MediaArea encourages all suppliers to work collaboratively to negotiate API documentation to support not only our own interoperability but to support third-party development of additional conformance checkers to utilize the produced shells.

The development of the shell will strive to facilitate an intuitive and informed use by memory institutions at both expert and non-expert levels. The shell will include substantial internal documentation that mimics the online resources that we will provide so that the shell and conformance checker function well offline.

MediaArea will implement a scheduling service within the shell so that large tasks may be performed overnight or according to a defined schedule. MediaArea will enable the Shell to load queues of files from lists of filepaths or URLs. Because of the size of data involved in audiovisual checkers MediaArea will give priority to designing the shell and conformance checker to perform multi-threaded and optimized processing.

Implementation Checker (Shell) The shell produced will support all functions and requirements of the implementation checker as described as an independent utility and also support:

- Allow the user to open one or many files at a time.
- Allow the user to queue simultaneous or consecutive file analysis.

- Allow the user to select how comprehensive or verbose an implementation check may be (for instance, samples frames or all frames of video).
- Enable the user to select sections of conformance checks or sets of conformance checks that they may wish to ignore.
- Enable the user to associate certain actions or warnings with the occurrence of particular checks.
- Provide feedback and status information live during the file analysis.
- (For Matroska) Present a user interface that displays the hierarchical EBML structure of the file with the corresponding policy outcome for each policy check.
- (For FFV1) Present a user interface that displays frame metadata in a table and enable the user to filter the presentation of frame metadata according to warnings or coherency events (for example, discontinuous aspect ratio).

Policy Checker (Shell) The shell produced will support all functions and requirements of the policy checker as described as an independent utility and also support:

- Allow PreForma-support technical metadata vocabularies to be imported or synchronized against an online registry.
- Provide an interface for the user to import, create, edit, or export a valid set of policy checks.
- Implement selected set of policy checks on all open files or selected files.
- Present the outcome of policy checks in a manner that allows comparison and sorting of the policy status of many files.
- Allows particular sets of policy to be associated with particular sets of files, based on file identification or naming patterns.
- (For Matroska) Present a user interface that displays the hierarchical EBML structure of the file with the corresponding policy outcome for each policy check.

Reporter (Shell) The shell produced will support all functions and requirements of the reporter as described as an independent utility and also support:

- Export of the PreFormaXML data at user-selected verbosity levels in a PDF format, which data visualizations supplied where helpful.
- Ability to read a collection of PreForma XMLs and provide a comprehensive summary and technical statistics of a collection to allow for prioritization, comparison, and planning.

Fixer (Shell) The shell produced will support all functions and requirements of the reporter as described as an independent utility and also support:

- Allow for single file or batch editing of file format metadata.
- Allow for selected metadata to be copied from one file to another or from one file to all other open files.
- Allow for file format metadata to be exported and imported to CSV or XML to enable metadata manipulation in other programs to then be imported back into the Shell and applied to the associated files.
- (For Matroska) Present a user interface that displays the hierarchical EBML structure of the file and allows the user to create, edit, or remove (with warning) any EBML element and display the associated policy or implementation check that corresponds with such actions.

Interfaces The selected formats (MKV, FFV1, and LPCM) represent substantially distinct concepts: container, video, and audio. The optimization of a conformance checker should utilize distinct interfaces to address the conformance issues of these formats, but allow the resulting information to be summarized together.

File View Import Export Help							
	File offset	Frame number	Tool	Type	Layer	Issue	Details
FFV1.mkv	0x00000123		Conformance	Error	Matroska	DocTypeVersion too high	DocTypeReadVersion=4, must be lower than DocTypeVersion=2
FFV1.mkv	0x00000123		Policy	Error	Matroska	DocTypeVersion	DocTypeReadVersion=4, Policy is DocTypeReadVersion<=3
FFV1.mkv	0x00004567	0	Conformance	Error	FFV1	CRC not valid	CRC is present but is not valid, frame corrupted
FFV1.mkv	0x00004567	0	Policy	Error	FFV1	micro_version	version=3, micro_version=3, this micro-version is pre-standard
FFV1.mkv			Conformance	Warning	Matroska	CRC-32 missing	Segment at offset 0x00000123 is missing mandatory CRC-32
FFV1.mkv			Conformance	Error	x-check	Width Coherency	Matroska width=704, FFV1 width=720
FFV1.mkv			Policy	Error	FFV1	Chroma subsampling	is 4:2:2, Policy is 4:4:4
FFV1.mkv			Policy	Error	FFV1	bit depth	is 8 bit, Policy is >=10
FFV1 - File 2.mkv			Policy	Error	FFV1	CRC missing	is not present, Policy is presence mandatory

Figure 1: frame scrolling mockup

```

Testing FFV1.mkv...
Conformance checking errors:
Matroska      ! DocTypeVersion too high      ! DocTypeReadVersion=4, must be lower than DocTypeVersion=2
FFV1          ! CRC not valid                    ! CRC is present but is not valid, frame corrupted
x-check       ! Width Coherency                  ! Matroska width=704, FFV1 width=720

Conformance checking warnings:
Matroska      ! CRC-32 missing                    ! Segment at offset 0x00000123 is missing mandatory CRC-32
FFV1          ! micro_version                    ! version=3, micro_version=3, this micro-version is pre-standard

Policy checking errors:
Matroska      ! DocTypeVersion                    ! DocTypeReadVersion=4, Policy is DocTypeReadVersion<=3
FFV1          ! Chroma subsampling                ! is 4:2:2, Policy is 4:4:4
FFV1          ! bit depth                        ! is 8 bit, Policy is >=10
FFV1          ! CRC missing                      ! is not presence, Policy is presence mandatory

Total: 3 conformance errors, 2 conformance warnings, 4 policy errors

```

Figure 2: CLI conformance output sample

Assessment of file conformance.

Assessment of file conformance can be displayed via a command line interface.

An interface for assessing conformance of FFV1 video could enable review of the decoded FFV1 frames (via a plugin) in association with conformance data so that inconsistencies or conformity issues may be reviewed in association of the presentation issues it may cause.

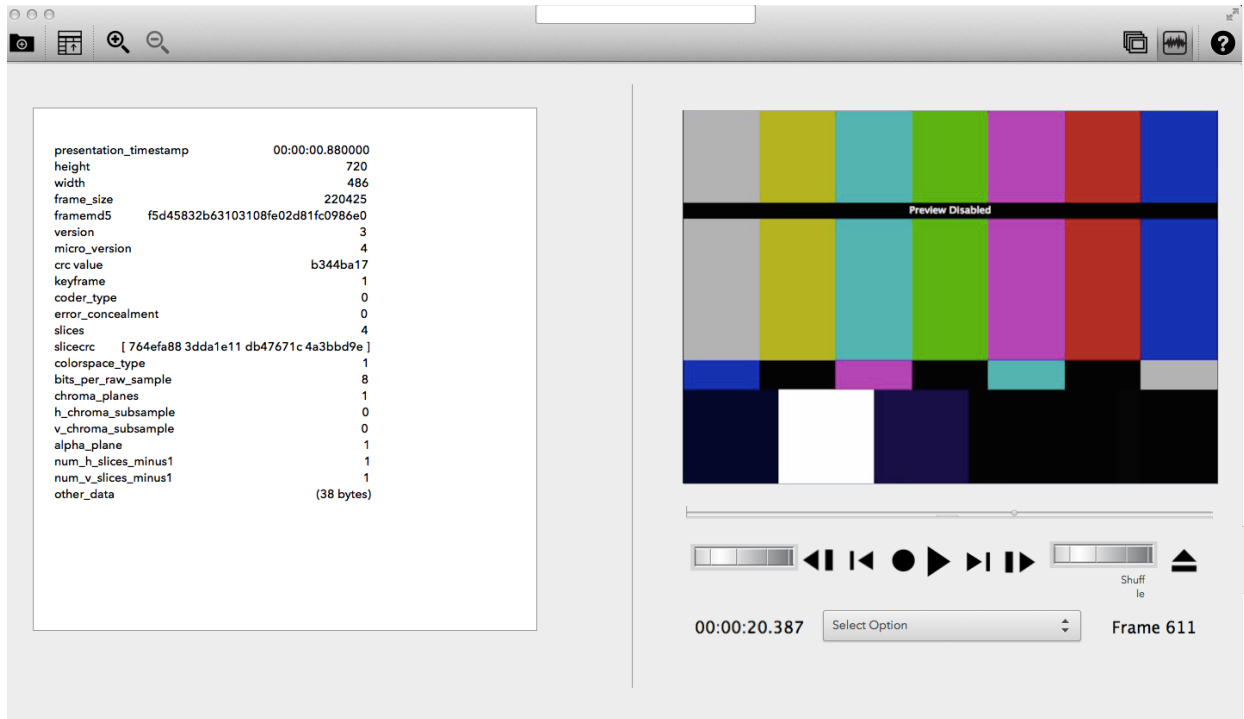


Figure 3: frame view mockup

MediaArea proposes an interface to present conformity issues for audio and video streams (FFV1 and LPCM) on a timeline, so that conformance events, such as error concealment or crc validation issues may be reviewed effectively according to presentation, parent Matroska block element, or video frame.

An interface for assessing conformance of FFV1 video could enable review of the decoded FFV1 frames (via a plugin) in association with conformance data so that inconsistencies or conformity issues may be reviewed in association with the presentation issues it may cause.

For deep analysis, a distinct interface that allows for its hierarchical structure to be reviewed and navigated. The presentation should allow for MKV and FFV1 elements to be expanded, condensed, or filtered according to element id or associated conformity issues.

A summary of the file properties can also be displayed via a command line interface for quick reference or export.

Optimization for Large File Size

Design of a conformance checker and shell should be considerate of the large file sizes associated with video. For instance, an hour-long PAL FFV1 file (which contains 90,000 frames per hour) should provide efficient access if cases where one FFV1 frame contains a CRC validation error.

A video conformance checker should be well optimized and multi-threaded to allow for multiple simultaneous processes on video files. Additionally the conformance checker should allow a file to be reviewed even as it is being processed by the conformance checker and allow assessment of files even as they are being written.

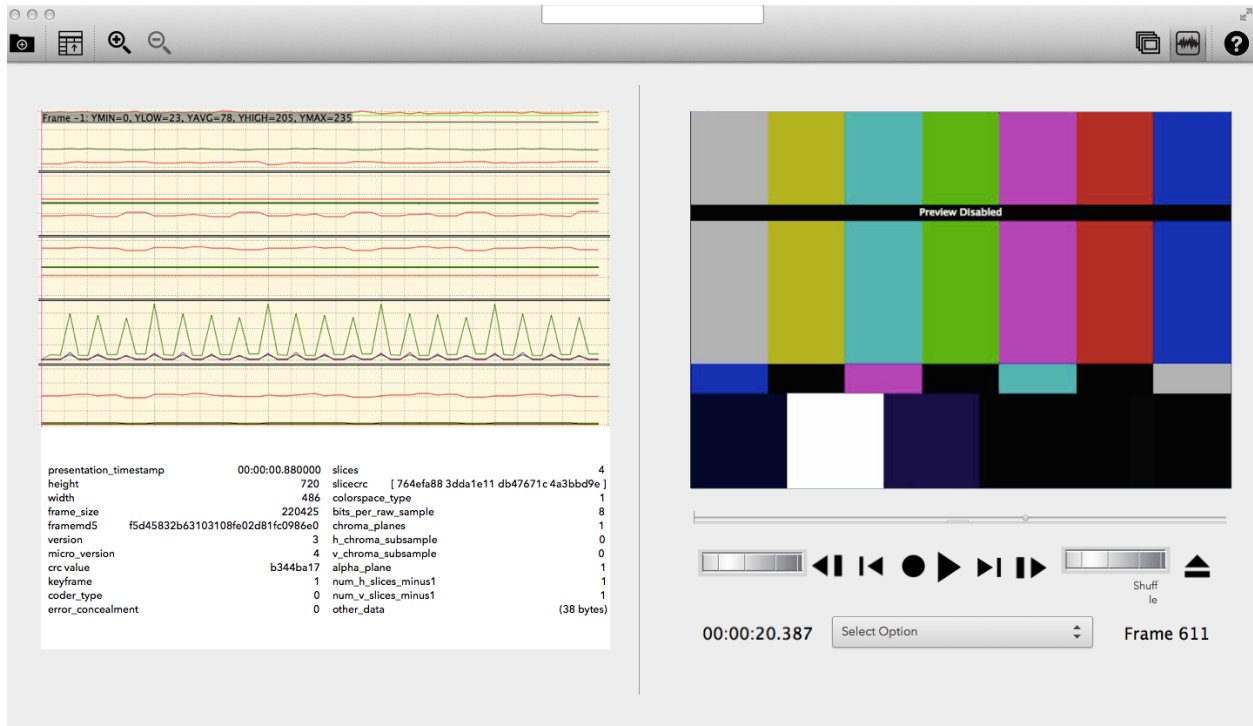


Figure 4: frame scrolling mockup

Focus on Fixity

Both FFMpeg and Matroska provide fixity features that serve the objectives of digital preservation by allowing data to be independently validated without the requirement of managing an external checksumming process. FFMpeg version 3 mandates CRCs on each frame. Matroska documents methods to embed checksums in Matroska elements to allow for future validation of any content.

Although the Matroska specification states that “All level 1 elements should include a CRC-32” this is not the practice of most Matroska multiplexers. As part of the Fixer aspect of this project, MediaArea proposes to develop a conformance checker that allows users to add CRC-32 to selected elements.

The advantages of embedded fixity in preservation media files are significant. The use of traditional external checksums does not scale fairly for audiovisual files, because since the file sizes are larger than non-audiovisual files there are less checksums per byte, which creates challenges in addressing corruption. By utilizing many checksums to protect smaller amounts of data within a preservation file format, the impact of any corruption may be associated to a much smaller digital area than the entire file (as the case with most external checksum workflows).

Reference and Test Files

MediaArea anticipates creating a large library of reference files that highlight many of the issues documented in our conformance check registry. Such files will be released under PreForma’s selected open licenses. Additionally we will curate references to relevant files in other reference libraries:

- Selections from <http://samples.ffmpeg.org>
- PDF/A files buggy files: <http://www.pdfa.org/2011/08/isartor-test-suite/>
- JPEG 2000 files: <https://github.com/openplanets/jpylyzer-test-files>

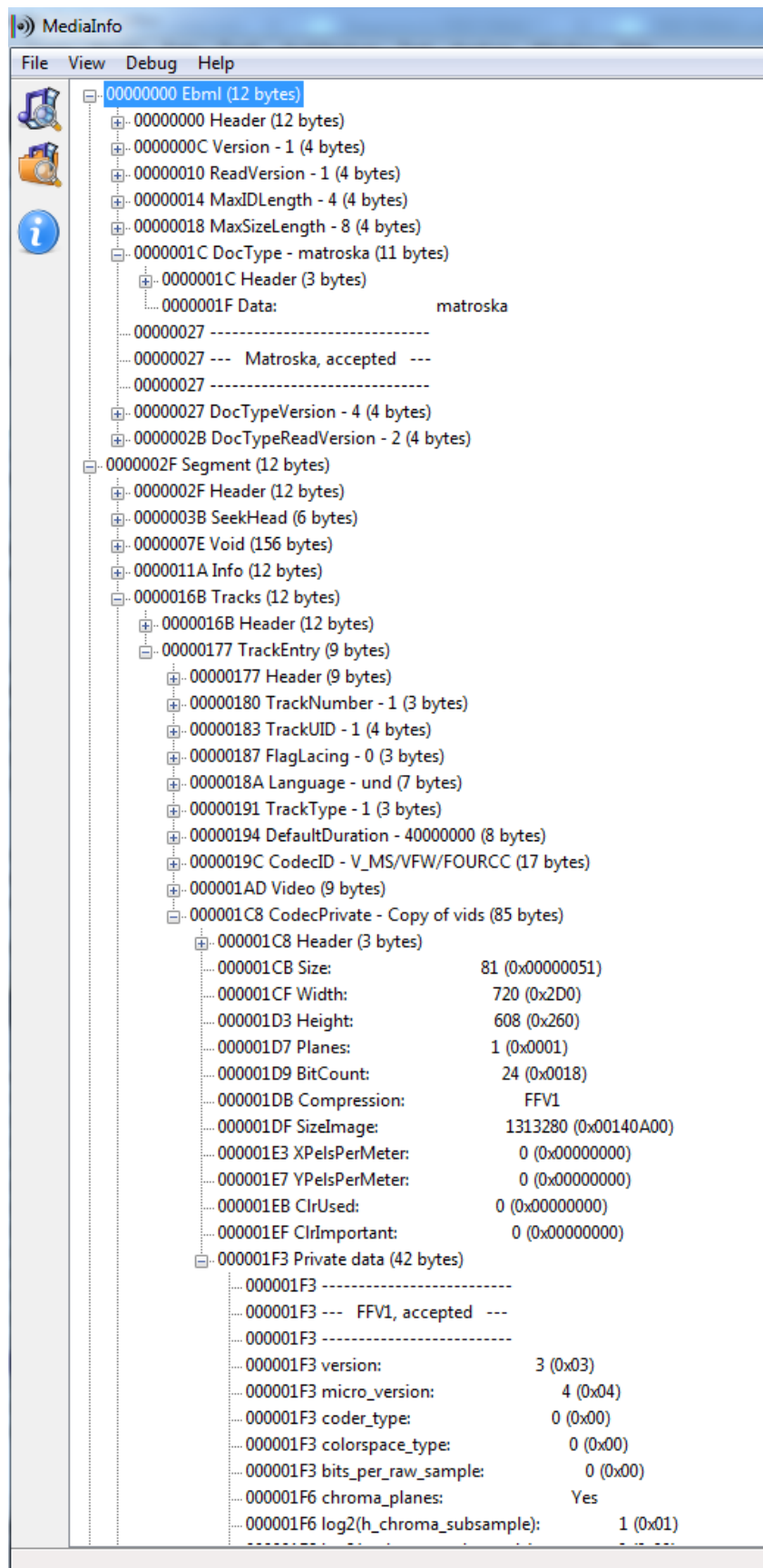


Figure 5: MediaInfo Windows

```

General
Unique ID                : 205266298566315392049931599426139420988 (0x9A6CD9E193016EB230C85E0E7E23793C)
Complete name            : /Users/ashley/Downloads/pres_metadata_sample_20110608.mkv
Format                   : Matroska
Format version           : Version 2
File size                : 69.2 MiB
Duration                 : 6s 430ms
Overall bit rate        : 90.1 Mbps
Encoded by               : {Name of the technician or organization responsible for the encoding and file creation process}
Encoded date             : UTC 2011-05-03 02:12:03
Mastered date            : UTC {Date and time of the file creation process. Enter in ISO 8601 format.}
Writing application      : mkclean 0.8.2 r from libebml v1.2.0 + libmatroska v1.1.0 + mkvmerge v4.7.0 ('I Got The...') built on Jun  6 2011 16:40:20
Writing library          : Lavf53.2.0
Original source form     : {Physical format of the source tape. Vocabulary?}
OriginalSourceForm/BarCode : {Barcode or other identifier from the tape.}
OriginalSourceForm/BarCode/source : {name of barcode authority or creator; example "XYZ Archive Barcode"}
OriginalSourceForm/LABEL  : {The record label or imprint on the source media object.}
OriginalSourceForm/DATE_RECORDED : {The time that the recording began. This is akin to the TDRC tag in ID3.}
OriginalSourceForm/condition : {Comments on the condition of the source, any preparation for playback of the source tape (if relevant).}
OriginalSourceForm/initial_source_timecode : 01:02:15;10
OriginalSourceForm/initial_source_timecode : DropFrame=Yes / 24HourMax=No / IsVisual=No
EncodedBy/Url             : {URL for the technician or organization listed in ENCODED_BY.}
EncodedBy/processing_actions : {Description of any actions performed during processing, such as trimming silence.}
EncodedBy/capture_software : Blackmagic Media Express
EncodedBy/capture_software/version : 2.0.3
EncodedBy/capture_operating_system : Ubuntu
EncodedBy/capture_operating_system/version : 10.10
EncodedBy/capture_device   : Decklink Studio SDI
EncodedBy/capture_device/manufacture : Blackmagic-Design
EncodedBy/capture_device/serial_no : xyz6789
EncodedBy/capture_device/settings : {notes on adjustments or settings}
EncodedBy/playback_device  : SV0-5800
EncodedBy/playback_device/manufacture : Sony
EncodedBy/playback_device/serial_no : abc1234
EncodedBy/playback_device/settings : {notes on adjustments or settings}
EncodedBy/playback_device/playback_signal : {Protocols used to transfer audiovisual data between the playback deck and the capture device. example: Component}
Attachment                : Yes / Yes

Video
ID                         : 1
Format                     : FFV1
Codec ID                   : V_MS/VFW/FOURCC / FFV1
Duration                   : 6s 440ms
Width                      : 720 pixels
Height                     : 486 pixels
Display aspect ratio       : 4:3
Frame rate mode            : Constant
Frame rate                 : 29.970 fps
Standard                   : NTSC
Compression mode           : Lossless
Language                   : English
Default                    : Yes
Forced                     : No
Encoded date               : UTC 2011-05-03 02:12:03

```

Figure 6: CLI output sample

- Matroska buggy files: Homemade + request to Matroska mailing list
- FFV1 buggy files: Homemade + request to FFmpeg mailing list
- LPCM files: Homemade

Intended Behavior by Use Case

Overview

The following use cases are presented to describe intended behaviors of the conformance checker:

Conformance Checking in an Open Archival Information System (OAIS) MediaArea recognizes the recommended practices described in the Open Archival Information System (OAIS) reference model, intended to provide long term preservation of digital information (CCSDS 650.0.-M-2). The stated conformance checker should be developed to address all relevant areas of activity within the OAIS model, as well as minimize and/or mitigate any incompatibility with additional OAIS-related standards.

Conformance checking plays a major role in OAIS-type services through the Information Packages created by Producers, maintained by Management in the Archive, and later, retrieved and put to use by Consumers and/or Designated Community. To take an example of such a service model, the Ingest Functional Entity of OAIS, or Ingest, performs quality assurance on incoming Submission Information Packages (SIPs). The conformance checker, as both an implementation and policy checker, would serve as the primary tool to verify SIPs and other submitted Preservation Description Information (PDI), with rules and specifications defined by the Archive. As Archival Information Packages (AIPs) are generated for the Archive, a conformance checker would map all Transformations through the collection of associated Representation, Content, and PDI information. And upon a dissemination request, reports created by the conformance checker would be used as descriptive information needed for the processing of objects for the Dissemination Information Package (DIP).

The conformance checker in effect resembles a complex Representation Net in the OAIS Archive, providing information needed to adequately understand the proper playback of the associated data object.

Conformance Checking at Digitization Time

Verification of Digitization Policy Archival digitization workflows are generally highly defined and consistent so that various analog source objects are associated with particularly digitization requirements. Generally digitization specifications are selected in order to reduce alterations to the significant characteristics of the analog source material. Example of such digitization scenarios may be:

- A PAL Betacam SX tape is digitized to a Matroska/FFV1 file at PAL specifications with YUV 4:2:2 8 bit video and 4 channels of 24 bit LPCM audio
- A NTSC U-Matic tape is digitized to a Matroska/FFV1 file at NTSC specifications with YUV 4:2:2 10 bit video and 2 channels of 24 bit LPCM audio
- A 1/4" audio reel is digitized to a 2 channel 96000 Hz, 24 bit audio LPCM file
- A CD-R is ripped to a 44100 Hz, 16 bit, 2 channel LPCM file
- A DAT tape is ripped to either a 32000, 44100, or 48000 Hz 16 bit file

Such digitization requirements may be expressed into a policy checker set through the shell or policy checker to verify that the results of digitization are consistent with the archive's specifications. This includes both sets of technical metadata and specification as well as anticipated embedded descriptive, preservation, or administrative metadata.

Verification of Lossless Digitization Until recently audiovisual digitization required a fairly inflexible set of hardware requirements and extremely limited possibilities for an open source approach to video digitization. Due to the bandwidth and processing requirements for the digitization of standard definition video required the installation of PCI cards and often the use of hardware encoders that were designed to encode video as fast as the video was being received to codecs like MPEG2 or JPEG2000. With modern connectivity options such as USB 3 and Thunderbolt it is easier to add video digitization capabilities to modern computers and more archive are performing this internally. Additionally modern computer processors can now transcode video losslessly in software from a video input without the need to rely on proprietary hardware-based encoders. Open source solutions such as DVA Profession, bmdcapture, and FFmpeg along with the open provision of video digitization software development kits, such as the Blackmagic SDK are facilitating new open development projects for archival video digitization.

As vendors and memory institutions are increasing considering and implementing digitization workflows that encode video directly to lossless codecs without the use of an intermediate file-based uncompressed audiovisual data, it is increasingly crucial to assess this lossless file soon after creation to detect any flaws within the digitization process.

For those digitizing video through processes that incorporate libav or FFmpeg such as bmdcapture or FFmpeg's decklink integration, a separate framemd5 may be written alongside the encoded ffv1 data. The resulting ffv1 data may then be verified against the framemd5 to verify that the correct bits were written to disk.

An inspiration for the use of framemd5 reports within a digitization workflow is inspired by the verify option with the flac utility available at <http://flac.sourceforge.net/>. The '-V' or --verify command is used to decode the encoded stream in parallel to the encoding process to double-check the losslessness of the transcoding. With this method any discrepancy between what data is read and transcoded versus what data is written to disk could be identified in a subsequent verification process. The use of framemd5 data within a digitization workflow enables verification in cases where an option similar to flac's --verify argument isn't available.

Assessment of Vendor/Producer Deliverables For archives that clarify specifications for audiovisual digitization projects, the conformance checker should facilitate a workflow for the archivist to express those specifications and verify received material against them. In addition to testing for the presence and order of required metadata tags the conformance checker should also be able to verify that they adhere to particular patterns as expressed through regular expressions.

The conformance checker should be able to verify that files were transferred completely and that the delivered material does not contain any partial files from an incomplete or aborted transfer.

The implementation and policy checker's reporting on deliverables will enable the user to provide specific feedback to the vendor or producer to create files with greater compliance or coherency.

Conformance Checking at Migration Time

Fixity Verification Migration of large amounts of data introduces risk for digital corruption and/or sector loss. Ongoing data migration is essential for digital preservation but can require a time consuming verification process. Both Matroska and FFV1 contain features for internal fixity so that a file copied from point A to point B can be assessed at point B alone to verify the data integrity of the frames. MediaArea recommends using Matroska's CRC features for use in digital preservation to allow for fixity verification to be more stable and achievable with the file alone without necessarily depending on external databases or records of checksums.

Obsolescence Monitoring Migration is typically an ideal time to perform obsolescence monitoring and preparing actions to limit complications in obsolescence status. Just as memory institutions must maintain the technology that their physical collections are dependent upon, this is equally true for digital collections. As this maintenance becomes more complex, costly, or unlikely archives will typically reformat material (with

as little compromise to the content and characteristics of the source as possible) to a format that has more sustainable characteristics.

To counteract arising obsolescence challenges it is critical to have access to thorough sets of technical metadata in order to associate certain codecs, formats, or technologies with sustainability risks or to identify what one format should be superseded by another in a particular digital preservation. For instance an institution that utilized FFV1 version 0 as a lossless preservation codec may wish to identify such files to reformat them to FFV1 version 3 (now that it is non-experimental) in order to take advantage of version 3's additional advantages. In our research one archive found that some digitized material received from a vendor was missing technical metadata about field dominance and had to identify exactly which materials were affected to order to rectify the issue.

Open Source Ecosystem

MediaArea has long been an open source natives and has an open source business model based on sponsored support (bug correction and feature requests). For instance, MediaInfo is officially provided by multiple open source distributions.

- Debian: <https://packages.debian.org/wheezy/mediainfo>
- Ubuntu: <http://packages.ubuntu.com/utopic/mediainfo>
- RedHat / Fedora: <https://apps.fedoraproject.org/packages/mediainfo>
- OpenSuse: <http://packman.links2linux.org/package/mediainfo>
- Arch Linux: <https://www.archlinux.org/packages/?q=mediainfo>
- FreeBSD: <http://www.freshports.org/multimedia/mediainfo/>
- Homebrew (open source software package management system for Mac): <http://brewformulas.org/MediaInfo>

Cross Platform Support

MediaArea excels in open source development for cross-platform support and chooses development frameworks and tools that enable cross-platform support to be maintained. Several applications developed by MediaArea such as QCTools, MediaInfo, and DVAnalyzer are available under nearly all major operating systems. To achieve this we will program in C++ and use the Qt application framework.

For an impression of MediaArea's focus on cross platform usability please see our download page: <http://mediaarea.net/en/MediaInfo/Download>.

Online Resources

MediaArea will utilize GitHub as a social and development center for PreForma MediaInfo development and uses GitHub's issue tracker and wiki features alongside development.

For communication MediaArea will establish public mailing lists and an IRC channel for foster support and involvement from memory institutions.

MediaArea will solicit, create, and accept test files and reference files that highlight various features of the conformance checker and illustrate likely preservation issues that may occur within the selected formats.

Advance Improvement of Standard Specification

FFV1 Specification Efforts to create an FFV1 specification began in April 2012, continuing through the August 2013 release of FFV1 version 3. Currently the specification remains in development at <http://github.com/ffmpeg/ffv1>. Ideally a specification should fully inform the development of a decoder

or parser without the need to reference existing implementations (such as the ffv1 implementations within ffmpeg and libav); however MediaArea's initial research and prototyping efforts with FFV1 found the current specification insufficient to create a decoder. As a result MediaArea utilized ffmpeg's FFV1 implementation to fully interpret the specification. Several threads on the ffmpeg-devel and libav-devel listserv reference discussions about the development of the FFV1 specification and consideration of efforts to standardize the specification through a standards organization, such as IETF (Internet Engineering Task Force) [1](#).

In consideration of FFV1's utilization within preservation contexts, the standardization of the codec through an open standards organization would better establish FFV1 as a trustworthy, stable, and documented option. At the moment FFV1 can be seen at a tipping point in its use within preservation context. It's speed, accessibility, and digital preservation features make it an increasingly attractive option for lossless video encoding that can be found in more and more large scale projects; the standardization of FFV1 through an open standards organization would be of broad interest to digital preservation communities and facilitate greater accessibility of lossless encoding options that are both efficient and standardized.

MediaArea proposes working closely with the lead authors of the FFV1 specification in order to update the current FFV1 specification to increase its self-reliance and clarity. Development of the FFV1 specification early within the PreForma project will generate substantial feedback to the authors of the specification which could then be offered through the specification's github page via pull requests or the issue tracker. MediaArea proposes at a later stage of development that the Preforma project serve as a catalyst to organize, facilitate, and sponsor the IETF standardization process for FFV1.

Considering the 2 year timeline of the PreForma project and usual pace of IETF standardization projects, we propose at least submitting FFV1 as an Independent Submission to IETF that could provide workable timeline, encourage a detailed review process, and assign a formal RFC number to the specification.

Matroska Specification Both the Matroska specification and its underlying specification for EBML are at mature and stable stage with thorough documentation and existing validators, but several efforts of the PreForma project can serve as contributions to this specifications. The underlying EBML specification [2](#) has already been drafted into RFC format but is has not yet been submitted to IETF as an Independent Submission or otherwise. MediaArea recommends that PreForma play a similar catalyst role for further standardization with Matroska as well, helping enable the refinement of the current RFC draft and coordinating an IETF process.

Matroska has a detailed metadata specification at <http://www.matroska.org/technical/specs/tagging/index.html>. Each tag has an official name and description while provides rules and recommendations for use. Many of these tags could be associated with validation rules, such as expressed by regular expression to assure that the content of the tag conforms to expectations. For instance tags such as URL, EMAIL, or ISBN have specific allowable patterns for what may be contained. As part of build a conformance tool for Matroska, MediaArea will generate conformance tests for individual tags and these tests may be contributed back to the Matroska specification in a list of regex values, an XML schematron file, or other acceptable contribution method.

Other Suggested Improvements or Contributions to Standard Specifications

- Register an official mime type via IETF for Matroska.
- Register dedicated FFV1 codecid with Matroska (current use is via fourcc).
- Proposal of a tagging extension to Matroska based on the requirements of the digital preservation community.
- Feedback for features and functions of FFV1 version 4, which is currently under development.
- Creation of metadata translators to convert common descriptive metadata formats within memory institution. For instance convert EBUCore into the XML representation of the Matroska tagging specification so that such metadata may be easily imported and exported between EBUCore and Matroska. # Architectural Layers

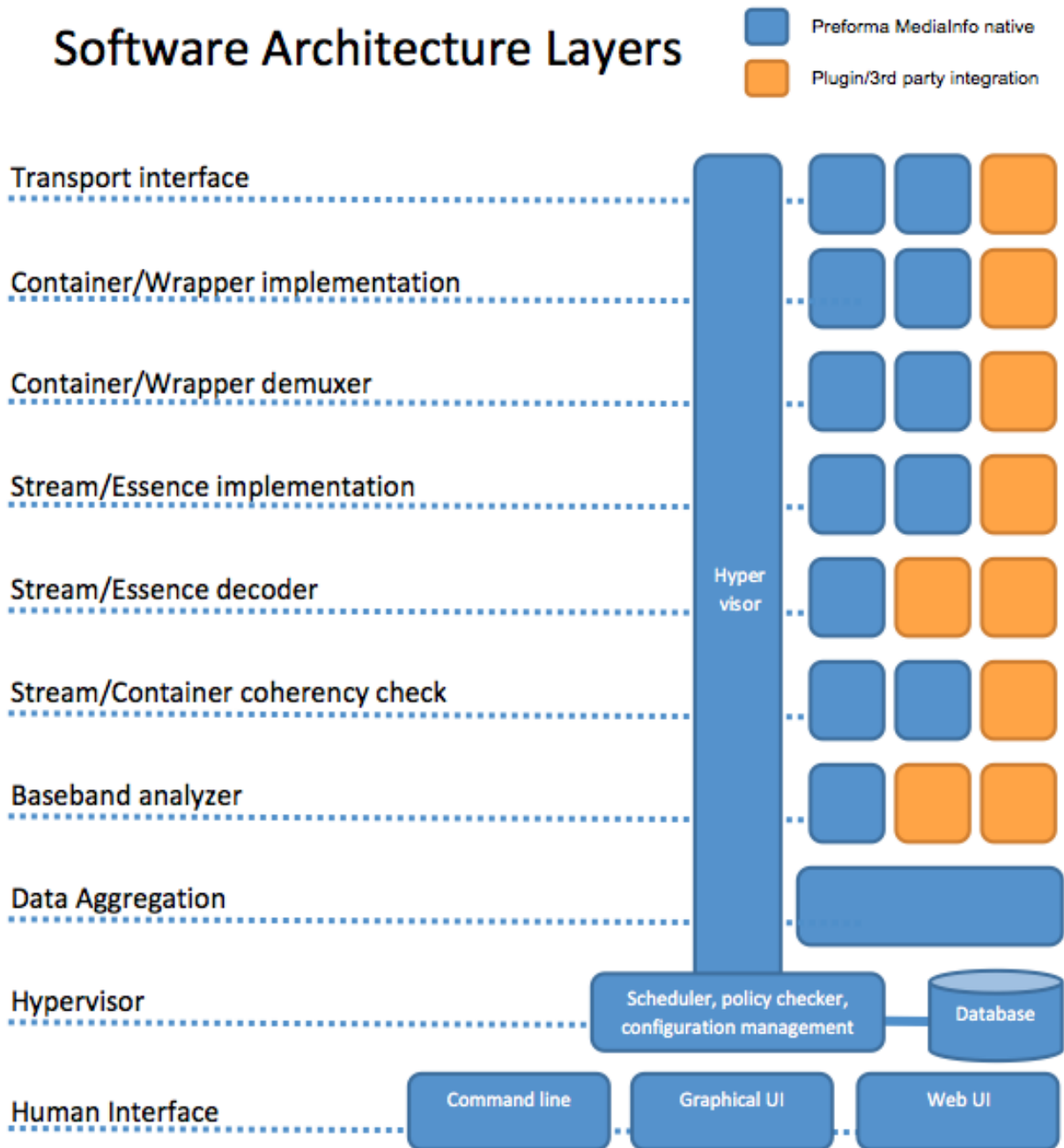


Figure 7: software architecture layers

The design of the conformance checker portion of the PreForma MediaInfo application will be comprised of several layers which will communicate via a Hypervisor. The layers shall include:

- Transport Layer
- Container Implementation (Conformance Check)
- Container Demuxer
- Stream Implementation (Conformance Check)
- Stream Decoding (through plugin)
- Stream/Container Coherency Check
- Baseband Analysis (through plugin)
- Data Aggregation
- Playback and Playback Analysis (through plugin)
- Human Interface

Note that the PreForma tender does not require decoding or subsequent baseband analysis or playback; however, from our experience in conformance checker design with DV Analyzer and QCTools and through discovery interviews, we've found that users are quick to require some form of playback in order to facilitate decision-making, response, and strategies for fixing. For instance if the conformance checker warns that the Matroska container and FFV1 codec note contradictory aspect ratios or a single FFV1 frame registers a CRC mismatch it is intuitive that the user would need to decode the video to determine which aspect ratio is correct or to assess the impact of the CRC mismatch. These layers can be supporting by designing a conformance checker and shell that is prepared to utilize FFmpeg as an optional plugin to enable additional analysis features and playback. Our overall proposal is not dependent on supporting an FFmpeg plugin but we believe that preparing a conformance checker that could support FFmpeg as an optional plugin could create a more intuitive, comprehensive, and informed user experience.

We propose incorporating several compatible utilities into PreForma MediaInfo to extend functionality and add immediate convenience for users. Each component is built as a plugin and can be replaced by a third party tool.

Transport layer

Preforma MediaInfo: File on disk or direct memory mapping

Preforma MediaInfo natively offers a file API for each operating system to enable direct file access, including files that are still in the process or being written. The inclusion of MediaInfo also offers features for direct memory mapping which will be useful for third-party development or plugins.

Plugin integration proof of concept: libcurl

libcurl is licensed under an MIT license that is compatible with both GPLv3+ and MPLv2+. curl offers extensive support for transferring data through many protocols. By incorporating curl into PreForma MediaInfo the tool will be able to assess files that may be accessible online by providing a URL (or list of URLs) in place of a filepath.

Since we will be generating a library of reference and sample files that will include large audiovisual files, users will be able to assess reference files without necessarily needing to download them.

Used as a proof of concept of plugin integration: HTTP/HTTPS/FTP/FTPS support via MediaInfo open source GPLv3+/MPL2+ and libcurl (MIT license, compatible with GPLv3+/MPL2+)

Container/Wrapper implementation checker

Preforma MediaInfo: Matroska checker

Plugin integration proof of concept: mkvalidator

mkvalidator is a basic and no more maintained Matroska checker (BSD license, compatible with GPLv3+/MPL2+) which will be used mostly for demonstration of the plugin integration.

Container/Wrapper Demultiplexing

Preforma MediaInfo

PreForma MediaInfo will utilize MediaInfo's existing demuxing libraries which will allow for PreForma's selected video codecs, FFV1 and JPEG2000, to be assessed from within many formats found within archives although these container formats themselves aren't the focus of the current PreForma project. Through discovery interviews with archives and vendors we have found FFV1's archival implementations to use a variety of container formats such as AVI and QuickTime as well as Matroska. In order to allow developed tools to support FFV1 even if not contained within Matroska, PreForma MediaInfo will support the following formats for demuxing (though not necessarily for conformity (yet)):

- MXF (commonly found within memory institutions)
- MOV/MP4 (often found containing FFV1, JPEG2000, and LPCM)
- DV (video stream format which uses LPCM)
- AVI (used with FFV1 by DV Profession, NOA, Austria Mediathek)
- WAV (a common container for LPCM)

By supporting the demultiplexing of these formats through MediaInfo, the developed tools will be applicable to a wide variety of files that contain PreForma's selected codecs: FFV1, JPEG2000, and LPCM. This demultiplexing support can be available through MediaInfo's existing libraries in a manner that is compatible with PreForma's licensing requirements.

Plugin integration proof of concept: FFmpeg

FFmpeg is one of the most ubiquitous, comprehensive, and open tools for demultiplexing and decoding audiovisual data; however, although FFmpeg's GPLv2+ license is compatible with PreForma's selected GPLv3+ license, it is not compatible with PreForma's other selected license, MPLv2+. As the PreForma conformance project evolves to support additional formats and codecs through plugins the use of FFmpeg's features are expected to becoming more and more appealing. For instance the integration of FFmpeg can provide integration of very comprehensive decoding and demultiplexing support beyond what can be easily provided with MediaInfo's demuxing libraries. FFmpeg's libavfilter library also provides access to waveform monitoring, vectorscope, audio meters, and other essential audiovisual inspection tools.

Although PreForma MediaInfo won't incorporate FFmpeg in order to comply with the MPLv2+ licensing requirement, we would like to design plugin support for FFmpeg. In this way a memory institution using PreForma MediaInfo could separately download FFmpeg and link the two together to enable additional tools such as:

- Video Waveform Monitor
- Vectorscope
- Ability to inspect luminance and chroma planes separately

- Audio Meters

We anticipate that the implementation of FFmpeg plugin support will substantially simplify the development of other plugins for broader codec and format support so that an entire decoder or demuxer does not need to be written from scratch in order to extend support.

Stream/Essence implementation checker

Preforma MediaInfo:

- FFV1
- PCM (including D-10 Audio, AES3)

Plugin integration proof of concept: jpylyzer

For JPEG 2000 (GPLv3+ license, compatible with GPLv3+ but not with MPL2+)

Plugin integration proof of concept: DV Analyzer

For DV (BSD license, compatible with GPLv3+ and MPL2+)

Optional

- MPEG-1/2 Video (including IMX, AS-07, D-10 Video, FIMS...)
- H.264/AVC (including AS-07)
- Dirac
- AC-3 (including AS-07)
- MPEG 1/2 Audio
- AAC
- Any other essence format on sponsor request (we have skills in DV, VC-1, VC-3, MPEG-4 Visual, H.263, H.265/HEVC, FLAC, Musepack, Wavepack, , BMP, DPX, EXR, JPEG, PNG, SubRip, WebVTT, N19/STL, TTML...)

Stream/Essence decoder

Preforma MediaInfo

- PCM (including D-10 Audio, AES3)

Plugin integration proof of concept: FFmpeg

FFmpeg decoder (GPLv2+ license, compatible with GPLv3+ but not with MPL2+)

Plugin integration proof of concept: OpenJPEG

OpenJPEG decoder (BSD license, compatible with GPLv3+/MPL2+)

Container/Wrapper vs Stream/Essence Coherency Check

PreForma MediaInfo

PreForma MediaInfo will support the coherency check between all supported formats (see Container/Wrapper implementation checker and Stream/Essence implementation checker parts)

Baseband Analyzer

PreForma MediaInfo

- None (only creation of the API)

Plugin integration proof of concept: QCTools

QCTools graphs (report on and graph data documenting video signal loss, flag errors in digitization, identify which errors and artifacts are in original format and which resulted from the digital transfer based on all the data collected in the past.)

Hypervisor (Controller)

The Hypervisor serves as communication between all plugins within and outside of the PreForma MediaInfo system and between all layers. The Hypervisor layer supports the following requirements:

- Scheduling
- Statistics
- Reporting
- User management
- Policies management

Human interface

PreForma MediaInfo will provide three different options for a human interface for maximum flexibility. These three interfaces are:

- Command line interface
A command line interface will be functional on nearly all kinds of operating systems, including those with very little graphical interface support.
- GUI (based on Qt)
The GUI, being based on Qt, has the strength of being versatile between operating systems and does not require additional development time to provide support for multiple platforms.
- Web UI (server/client)
The web interface will provide access to conformance checks without software installation.

Style Guide

Source Code Guide

Portability

Source code **MUST** be built for portability between technical deployment platforms.

Modularity

Source code **MUST** be built in a modular fashion for improved maintainability.

Deployment

The Conformance Checker **MUST** allow for deployment in these five infrastructures/environments: The PREFORMA website, an evaluation framework, a stand-alone system, a network-based system, and legacy systems.

To sufficiently demonstrate the scope and functionality of the Conformance Checker, it, along with associated documentation and guidelines, must be made available at the PREFORMA project website. The PREFORMA website will be considered as the deliverable for the PREFORMA project.

In order to gather sufficient structured feedback on the conformance checking process, the Conformance Checker will require deployment within the DIRECT infrastructure for test and evaluation of the tool in the PCP procedure.

The Conformance Checker must have the capability to be packaged and run as an executable on a PC running any standard operating system (at least for: MS Windows 7, Mac OSX, common Linux distributions such as Ubuntu, Fedora, Debian, and Suse). This ensures the conformance checker can be used in small-scale institutions without centralized IT infrastructure.

The Conformance Checker must allow for deployment in network-based solutions (dedicated server, cloud solutions) for digital repositories.

The Conformance Checker must have the capability of being plugged into legacy systems via written API integration.

API's

The Conformance Checker **MUST** interface with other software systems via API's.

Open Source Practices

Development

Development of software in open source projects in PREFORMA **MUST** utilise effective open source work practices. Effective open source work practices include:

- use of nightly builds Nightly builds are automated neutral builds that reflect the current, most up-to-date status of a piece of developed software. Access made to nightly builds allow for groups of developers to work collaboratively and always assess the most current state of the software, with consideration to potential bugs or other hazards that could occur during the development process. Programmers are able to confidently determine if they “broke the build” (made the software inoperable) with their code and prevent or correct changes quickly, as needed.

- use of software configuration management systems (e.g. Git) Using a software configuration management system allow for version and revision control, an essential component to developers working collaboratively. A version control system allows multiple people to work on same or similar sections of the source code base at the same time with awareness and prevention of overlapping or conflicting work. Git will be used as the software configuration management system for this project.
- use of an open platform for open development (e.g. Github) An open platform on which to develop software facilitates the open development of that software. Public visibility and ability to contribute to the software by anyone allows for heartier, more reliable software. Feedback is more easily sought and more readily provided with the use of an open platform. Github will be used as the open platform for open development of this project.

Open Source Platforms

All development of software and all development of digital assets (related to developed open source software) in PREFORMA MUST be conducted and provided in open source projects at open development platforms.

Contribution Guide

File Naming Conventions

Files related to documentation should be named in CamelCase. Sample data should be added in snake_case with a sufficiently descriptive title.

Commit messages should concisely summarize the contribution. Commits should be cohesive and only include changes to relevant files (e.g. do not fix a typo in the Style Guide, change scope parameters, and fix a bug all in the same commit).

Rules for Qt/C++ code:

4 spaces are used for indentation. Tabs are never used.

For more guidelines, refer to the Qt Coding Style guide: http://qt-project.org/wiki/Qt_Coding_Style For even more guidelines, Google guide on C++: <http://google-styleguide.googlecode.com/svn/trunk/cppguide.html>

Rules for contributing code

Contributions of code or additions to documentation must be written with Qt and must be made in the form of a branch submitted as a pull request.

1. Fork this repository (<https://github.com/MediaArea/PreFormaMediaInfo/fork>)
2. Create your feature branch (`git checkout -b my-new-feature`)
3. Commit your changes (`git commit -am 'Added some feature'`)
4. Push to the branch (`git push origin my-new-feature`)
5. Create a new Pull Request with a more verbose description of the proposed changes

Rules for contributing feedback

Feedback of all kind is encouraged and can either be made through [opening an issue](#) or by contacting the team directly at info@mediaarea.net

Linking

In order to facilitate self-description, intuitive discovery, and use of resulting code and documentation it is highly encouraged to utilize linking through documentation, tickets, commit messages, and within the code. For instance the registry itemizes individual conformance checks should link to code blocks and/or commits as software is developed that is associated to that conformance check. In this manner it should be feasible to easily review both human-readable descriptions of conformity checks and associated programmatic implementations.

License

The software and digital assets delivered by tenderer are made available under the following IPR conditions:

All software developed during the PREFORMA project MUST be provided under the two specific open source licenses: “GPLv3 or later” and “MPLv2 or later”.

All source code for all software developed during the PREFORMA project MUST always be identical between the two specific open source licenses (“GPLv3 or later” and “MPLv2 or later”).

All digital assets developed during the PREFORMA project MUST be provided under the open access license: Creative Commons CC-BY v4.0 and in open file formats, i.e. an open standard as defined in the European Interoperability Framework for Pan-European eGovernment Service (version 1.0 2004)

Checker Design: Conformance and Coherency

Conformance checks for both container formats (such as Matroska) and streams (such as LPCM and FFV1) shall be defined, registered, and associated with the code that performs the check. PreForma MediaInfo will perform and report on a growing list of tests per format. Many of these tests will be derived directly from the specifications or standard documents of a given file format, but other tests will derive from expected patterns and structural incoherency. Some checks focus on coherency between a stream and a container (such as if the container and stream utilize contradictory aspect ratios).

These checks shall have logical cause and effect or conditional relationships and shall be documented by the citation of external standards documentation or by the project’s own research and development. MediaArea plans to provide guidance for user communities to develop and explain their own ruleset in shareable form. An XML schema for conformance definition is provided. MediaArea’s development of conformance and policy checkers will involve several categories of tests. In addition to supporting conformity checks and logical interpretation of selected file formats, there is user desire for checks performed based on internal or institutional policy that are not necessarily embedded in the file format technical specifications. A PreForma MediaArea ‘shell’ shall be able to load multiple sets of conformity/coherency rulesets so that users may select which rulesets they choose to adhere to as well as create their own.

Conformance and coherency rulesets specifically targeting specification compliance of FFV1, Matroska, and LPCM are currently under development.

The Conformance Check Registry defines the basis of how a conformance check may be expressed. MediaArea proposes developing online resources that define checks and relate them to sample files, associated code, rationale, potential responses and community discussion. The Conformance Check Registry provides a basis on how information on the implementation checks themselves can be communicated to users.

Conformance Check Registry

This documentation explains the elements, structure, and intent of the Preforma Mediainfo Conformance Check Registry.

A conformance check is a particular test applied to a file format, stream, or section of a format or stream in order to quantify the adherence of such data to associated sets of rules and practices. The registry refers specifically to checks performed by the implementation checker. Rules performed by the policy checker are defined elsewhere.

MediaArea plans to maintain an identity and open documentation for each Conformance Check in both a public online space and within the internal help documentation of a PreForma MediaInfo Shell. As the implementation checker assesses given files against a series of checks and the reporter presents the findings to the user, MediaArea plans for the shell to facilitate the user to discover more underlying information, advice, or responses to conformance checks that appear as problematic from the implementation checker.

Conformance Check Registry Requirements

- Each conformance check must be identified by a unique identifier.
- Documentation for Conformance Checks must offer hierarchical relationships between related checks.
- Conformance Checks must be documented according to their CCID (Conformance Check Identifier) and Version number.
- Any revised Conformance Check must maintain a changelog as well as records of all past versions of the conformity check.
- The following keywords to indicate requirement levels when used in a conformance check description MUST be used according to their definitions provided by (RFC2119)[<http://tools.ietf.org/html/rfc2119>]: “MUST”, “MUST NOT”, “REQUIRED”, “SHALL”, “SHALL NOT”, “SHOULD”, “SHOULD NOT”, “RECOMMENDED”, “MAY”, and “OPTIONAL”.

Elements

A draft list of elements to be used in the definition of an implementation conformance check are provided.

Name A human-readable name for the conformance check.

CCID Conformance Check Identifier. An alphanumeric identifier used to reference or identify the conformance check to order to relate code, documentation, and reports that reference the check.

Version The version number of the reference Conformance Check. This value should be expressed as a standard GNU version number in major.minor.revision format. A value of “0” may be used to indicate draft status.

Authority The authority associates each conformance check with a standards organization, community, or logic from which the conformance check is derived. Examples: EBU, Microsoft. Within the use of a Conformance Checker, the user may enable or disable check from certain authorities; for instance to perform checks against specifications of Standards Organization A and Community Practice B, but not Standards Organization C or Community Practice D.

Target Format The name of the file format, codec, or bit stream that is to be test.

Target Format Version Identify the version of range of versions of the target format which are eligible for the conformance test. A numeric range should be used or the word “all” if the rule applies to all known versions.

Target Format Part The name of a chunk, atom, element, bitstream, or other smaller component of the target format that the conformance check relates to.

Citation A reference of the specific document, specification, or reference from which the conformance check is derived. Typically the citation will be a publication or expression of the ‘Authority’.

Quote A quote from the authority that demonstrates the logic or reasons for the check.

Rule Clarity Expresses whether the check is based of an explicit statement of an underlying specification or based on deductive logic or inference from a reading of the specification.

Definition A clear description of conformance check.

Regex Parameters A human-readable description of what the regular expression must check for.

Regular Expression Code that checks for conformance to the suggested parameters. `## Matroska Conformance Checks (Draft) ### Extension Test`

Descriptor	Value
CCID	MKV-EXT
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	File name
Citation	http://www.matroska.org/node/2/revisions/153/view

Rule Clarity: Inferred

Quote:

Definition: The file extension SHOULD be one of the following (MKV, MKA, MKS, MK3D, WEBM)

Extension Test MKV

Descriptor	Value
CCID	MKV-EXT-MKV
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	File name
Citation	http://www.matroska.org/node/2/revisions/153/view

Rule Clarity: Inferred

Quote:

Definition: If the file extension is MKV, the file SHOULD contain at least one video track.

Extension Test MKA

Descriptor	Value
CCID	MKV-EXT-MKA
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	File name
Citation	http://www.matroska.org/node/7/revisions/214/view

Rule Clarity: Inferred

Quote:

Definition: If the file extension is MKA, the file SHOULD contain at least one audio track and no other type of track, i.e. “audio-only”.

Extension Test MKS

Descriptor	Value
CCID	MKV-EXT-MKS
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	File name
Citation	http://www.matroska.org/node/2/revisions/153/view

Rule Clarity: Inferred

Quote:

Definition: If the file extension is MKS, the file SHOULD contain at least one subtitle track and no other type of track, i.e. “subtitle-only”.

Extension Test MK3D

Descriptor	Value
CCID	MKV-EXT-MK3D
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	?
Target Format Part	File name, StereoMode element
Citation	http://www.matroska.org/node/2/revisions/153/view

Rule Clarity: Inferred

Quote:

Definition: If the file extension is MKV3D the file SHOULD contain at least one video track AND SHOULD contain at least one StereoMode element.

EBML Element Start

Descriptor	Value
CCID	MKV-EBML-ELEM-START
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	specdata.xml

Rule Clarity:

Quote: “Set the EBML characteristics of the data to follow. Each EBML document has to start with this.”

Definition: An Matroska file MUST start with an EBML element id, ie. 0x1A45DFA3.

EBML vint efficiency

Descriptor	Value
CCID	EBML-VINT-EFF
Version	0
Authority	EBML Specification
Target Format	EBML
Target Format Version	all
Target Format Part	EBML Structure
Citation	http://matroska.org/technical/specs/rfc/index.html

Rule Clarity:

Quote: Section 2.2 “IDs are always encoded in their shortest form, e.g. 1 is always encoded as 0x81 and never as 0x4001.”

Definition: The bits following the Element ID’s Length Descriptor are not more than $(8 - \{\text{bit-length-of-length-descriptor}\})$ successive 0 bits, i.e. vint is expressed as efficiently as feasible.

Element ID Registered

Descriptor	Value
CCID	MKV-KNOWN-ELEM
Version	0
Authority	
Target Format	EBML

Descriptor	Value
Target Format Version	all
Target Format Part	
Citation	

Rule Clarity: Inferred

Quote:

Definition: Ensure MKV Element ID is registered in specdata.xml (as of Dec. 13, 2014 this is 224 registered Element IDs)

Element Size 0x7F Reservation

Descriptor	Value
CCID	EBML-ELEM-SIZE-7F
Version	0
Authority	EBML Specification
Target Format	EBML
Target Format Version	all
Target Format Part	EBML Element Size
Citation	http://matroska.org/technical/specs/rfc/index.html

Rule Clarity: Warning, since it is possible (though unlikely) element size is unknown but then happens to be 127 bytes.

Quote: “Note that the shortest encoding form for 127 is 0x407f since 0x7f is reserved.”

Definition: If Element Size is set to 0x11111111 but element size is actually 127 bytes provide a warning.

Element Size Byte Length Limit

Descriptor	Value
CCID	EBML-ELEM-SIZE-CAP
Version	0
Authority	EBML Specification
Target Format	EBML
Target Format Version	all
Target Format Part	EBML Element Size
Citation	http://matroska.org/technical/specs/rfc/index.html

Rule Clarity:

Quote: Section 2.3: “The EBML element data size is encoded as a variable size integer with, by default, widths up to 8.”

Definition: The first eight bits of any Element Size may not start with 0b00000000.

Element Size Unknown

Descriptor	Value
CCID	EBML-ELEM-SIZE-UNK
Version	0
Authority	EBML Specification
Target Format	EBML
Target Format Version	all
Target Format Part	EBML Element Size
Citation	Dave

Rule Clarity: Warning

Quote: “Values with all data bits set to 1 means size unknown, which allows for dynamically generated EBML streams where the final size isn’t known beforehand.”

Definition: Warning on unknown element sizes, unoptimized MKV.

Level 0 Segment

Descriptor	Value
CCID	MKV-LEVEL-0
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	specdata.xml

Rule Clarity:

Quote: Inferred: EBML and Segment are the only level 0 elements, both are allowed to occur multiple times.

Definition: The EBML Header MUST be immediately followed by another EBML Header Element, 0x1A45DFA3, or a Segment Element, 0x18538067. {{Can global Elements exist at level 0?!}}

Only One EBML Header recommended

Descriptor	Value
CCID	MKV-1-EBML
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Matroska structure
Citation	is there a rule to prevent this?

Rule Clarity: Warning

Quote: Assumed: Two EBML Headers in one MKV file seems contradictory.

Definition: There SHOULD only occur one EBML level 0 element within an MKV file. (EBML Headers

could recur if an MKV file is an attachment of an MKV file).

File Size Consistency

Descriptor	Value
CCID	MKV-FILESIZE-MATCH
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Matroska structure
Citation	http://www.matroska.org/technical/specs/index.html#block_structure

Rule Clarity:

Quote: Inferred

Definition: The actual file size should be the sum of all level 0 Element Size declarations plus the sum of the byte sizes of level 0 Element IDs and Element Sizes.

EBMLVersion Presence

Descriptor	Value
CCID	MKV-EBML-VER
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity:

Quote:

Definition: Within any EBML Header exactly one EMBL Version element must be present.

EBMLReadVersion Presence

Descriptor	Value
CCID	MKV-EBML-RV
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity:

Quote:

Definition:

EBMLMaxIDLength Presence

Descriptor	Value
CCID	MKV-EBML-MAXIDL
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity:

Quote:

Definition:

EBMLMaxSizeLength Presence

Descriptor	Value
CCID	MKV-EBML-MAXSL
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity:

Quote:

Definition:

DocType Presence

Descriptor	Value
CCID	MKV-EBML-DOCT
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity:

Quote:

Definition:

DocTypeVersion Presence

Descriptor	Value
CCID	MKV-EBML-DOCTV
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity:

Quote:

Definition:

DocTypeReadVersion Presence

Descriptor	Value
CCID	MKV-EBML-DOCTRV
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity:

Quote:

Definition:

EBML Version Coherency

Descriptor	Value
CCID	MKV-VER-COH
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	http://www.matroska.org/technical/specs/index.html#block_structure

Rule Clarity: Inferred

Quote:

Definition: The value of EBMLVersion MUST be greater than or equal to the value of EBMLReadVersion.

EBMLMaxIDLength Limits

Descriptor	Value
CCID	MKV-MAXID-LIMIT
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	specdata.xml

Rule Clarity: Spec says “4 or less”, but since the EBML ID length itself is 4, the EBMLMaxIDLength has no other valid value.

Quote:

Definition: MUST equal 4

EBMLMaxSizeLength Limit

Descriptor	Value
CCID	MKV-MAXSL-LIMIT
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	specdata.xml

Rule Clarity: “The maximum length of the sizes you’ll find in this file (8 or less in Matroska).”

Quote:

Definition: Must be less than or equal to 8 and greater than or equal to 1.

EBMLMaxSizeLength Matches

Descriptor	Value
CCID	MKV-MAXSL-MATCH
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header

Descriptor	Value
Citation	specdata.xml

Rule Clarity:

Quote:

Definition: No Element Size Length exceeds the length noted in EBMLMaxSizeLength

DocType

Descriptor	Value
CCID	MKV-DOCT-KNOWN
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity:

Quote:

Definition: MUST equal either “matroska” or “webm”

DocTypeVersion Coherency

Descriptor	Value
CCID	MKV-DOCTV-COH
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity:

Quote:

Definition: The value of DocTypeVersion MUST be greater than or equal to the vale of DocTypeReadVersion.

DocTypeVersion Limits

Descriptor	Value
CCID	MKV-DOCTV-LIMIT
Version	0

Descriptor	Value
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	EBML Header
Citation	

Rule Clarity: Warning

Quote:

Definition: Values for DocTypeVersion and DocTypeReadVersion must be either 1, 2, 3, or 4.

Top Elements Coded on 4 Octets

Descriptor	Value
CCID	MKV-TOP-ELEM-4CODE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Matroska structure
Citation	http://www.matroska.org/technical/specs/index.html#block_structure

Rule Clarity: “All top-levels elements (Segment and direct sub-elements) are coded on 4 octets, i.e. class D elements.”

Quote:

Definition: Note: this seems to contradict EBML rule to use most efficient element size, but perhaps this is an intention deviation of MKV to achieve top elements starting on multiples of 4 octets. ?

CRC Order

Descriptor	Value
CCID	MKV-CRC-ORDER
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	CRC Element
Citation	http://www.matroska.org/technical/specs/index.html#block_structure

Rule Clarity: “The CRC element should be the first in it’s parent master for easier reading.”

Quote:

Definition: CRC Elements SHOULD be the first sub-Element of its parent Element.

CRC-32 Size Coherency

Descriptor	Value
CCID	MKV-CRC-COH
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	CRC Element
Citation	http://www.matroska.org/technical/specs/index.html#block_structure

Rule Clarity: Inferred: “The CRC in use is the IEEE CRC32 Little Endian”

Quote:

Definition: The Element Size of the CRC-32 Element MUST be 4 bytes (aka 32 bit).

CRC Validation

Descriptor	Value
CCID	MKV-CRC-VAL
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	CRC Element
Citation	

Rule Clarity:

Quote:

Definition: The crc hash of the CRC-32 element MUST validate the subsequent data of the parent Element, from the Element that follows the CRC-32 element to the end of the parent Element.

CRC Not Pointlessly Used

Descriptor	Value
CCID	MKV-CRC-REASON
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	CRC Element
Citation	author

Rule Clarity: Recommended

Quote:

Definition: A CRC-32 element should not be the only child Element of its parent Element (ie hashing no data).

CRC-Presence

Descriptor	Value
CCID	MKV-CRC-PRES
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	CRC Element
Citation	

Rule Clarity: “All level 1 elements should include a CRC-32.” but CRC-32 Element is NOT Mandatory. ?

Quote:

Definition: Warning when Level 1 elements have no CRC-32. Very common.

Single Segment Composition

Descriptor	Value
CCID	
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	
Citation	specdata.xml

Rule Clarity:

Quote: “Typically a Matroska file is composed of 1 segment.”

Definition: File MUST contain at least one segment.

Seek-Presence

Descriptor	Value
CCID	MKV-SEEK-PRES
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Meta Seek Element
Citation	specdata.xml

Rule Clarity:

Quote:

Definition: File MUST contain at least one Seek element.

SeekID-Presence

Descriptor	Value
CCID	MKV-SEEKID-PRES
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Meta Seek Element
Citation	specdata.xml

Rule Clarity:

Quote:

Definition: File MUST contain at least one SeekID element.

SeekID-Type

Descriptor	Value
CCID	MKV-SEEKID-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Meta Seek Element
Citation	specdata.xml

Rule Clarity:

Quote:

Definition: SeekID MUST be in binary format.

SeekPosition-Presence

Descriptor	Value
CCID	MKV-SEEKPOSITION-PRES
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Meta Seek Element
Citation	specdata.xml

Rule Clarity:

Quote:

Definition: File MUST contain at least one SeekPosition element.

Segment-Info-Presence

Descriptor	Value
CCID	MKV-SEGMENTINFO-PRES
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	specdata.xml

Rule Clarity:

Quote:

Definition: Segment information MUST contain at least one Info element.

SegmentUID-Range

Descriptor	Value
CCID	MKV-SEGMENTUID-RNG
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity: Range cannot be zero.

Quote:

Definition: SegmentUID MUST be greater than zero.

SegmentUID-Size

Descriptor	Value
CCID	MKV-SEGMENTUID-SIZE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, SegmentUID MUST be 128 bits (16 bytes) in size.

SegmentUID-Type

Descriptor	Value
CCID	MKV-SEGMENTUID-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, SegmentUID MUST be in binary format.

SegmentFilename-Type

Descriptor	Value
CCID	MKV-SEGMENTFILENAME-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, SegmentFilename MUST be in UTF-8 format.

PrevUID-Size

Descriptor	Value
CCID	MKV-PREVUID-SIZE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, PrevUID MUST be in binary format.

PrevUID-Type

Descriptor	Value
CCID	MKV-PREVUID-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, PrevUID MUST be 128 bits (16 bytes) in size.

PrevFilename-Type

Descriptor	Value
CCID	MKV-PREVFILENAME-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, PrevFilename MUST be in UTF-8 format.

NextUID-Size

Descriptor	Value
CCID	MKV-NEXTUID-SIZE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, NextUID MUST be in binary format.

NextUID-Type

Descriptor	Value
CCID	MKV-NEXTUID-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, NextUID MUST be 128 bits (16 bytes) in size.

NextFilename-Type

Descriptor	Value
CCID	MKV-NEXTFILENAME-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, NextFilename MUST be in UTF-8 format.

SegmentFamily-Size

Descriptor	Value
CCID	MKV-SEGMENTFAMILY-SIZE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, SegmentFamily MUST be in binary format.

SegmentFamily-Type

Descriptor	Value
CCID	MKV-SEGMENTFAMILY-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, SegmentFamily MUST be 128 bits (16 bytes) in size.

TimecodeScale-Presence

Descriptor	Value
CCID	MKV-TIMECODESCALE-PRES
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: File MUST contain at least one TimecodeScale element.

Duration-Range

Descriptor	Value
CCID	MKV-DURATION-RANG
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, duration range MUST be greater than 0

Duration-Type

Descriptor	Value
CCID	MKV-DURATION-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	http://www.matroska.org/technical/specs/index.html

Rule Clarity:

Quote:

Definition: If present, duration type MUST be float integer.

DateUTC-Type

Descriptor	Value
CCID	MKV-DATEUTC-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	specdata.xml

Rule Clarity: UTC standards inferred.

Quote:

Definition: If present, DateUTC MUST be in date format and follow UTC standards.

Title-Type

Descriptor	Value
CCID	MKV-TITLE-TYPE
Version	0
Authority	Matroska Specification
Target Format	Matroska
Target Format Version	all
Target Format Part	Segment Element
Citation	specdata.xml

Rule Clarity:

Quote:

Definition: If present, Title MUST be in UTF-8 format.

Tag Total Parts

Descriptor	Value
CCID	MKV-TAG-TOTALPARTS
Version	0
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Part Number

Descriptor	Value
CCID	MKV-TAG-PARTNUMBER
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Part Offset

Descriptor	Value
CCID	MKV-TAG-PARTOFFSET
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Title

Descriptor	Value
CCID	MKV-TAG-TITLE
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Subtitle

Descriptor	Value
CCID	MKV-TAG-SUBTITLE
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag URL

Descriptor	Value
CCID	MKV-TAG-URL
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Sort_with

Descriptor	Value
CCID	MKV-TAG-SORT_WITH
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Email

Descriptor	Value
CCID	MKV-TAG-EMAIL
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Address

Descriptor	Value
CCID	MKV-TAG-ADDRESS
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Fax

Descriptor	Value
CCID	MKV-TAG-FAX
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Phone

Descriptor	Value
CCID	MKV-TAG-PHONE
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Initial_Key

Descriptor	Value
CCID	MKV-TAG-INITIAL_KEY
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Law_Rating

Descriptor	Value
CCID	MKV-TAG-LAW_RATING
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

TAG ICRA

Descriptor	Value
CCID	MKV-TAG-ICRA
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag DATE_RELEASED

Descriptor	Value
CCID	MKV-TAG-DATE_RELEASED
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag DATE_RECORDED

Descriptor	Value
CCID	MKV-TAG-DATE_RECORDED
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag DATE_ENCODED

Descriptor	Value
CCID	MKV-TAG-DATE_ENCODED
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag DATE_TAGGED

Descriptor	Value
CCID	MKV-TAG-DATE_TAGGED
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag DATE_DIGITIZED

Descriptor	Value
CCID	MKV-TAG-DATE_DIGITIZED
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag DATE_WRITTEN

Descriptor	Value
CCID	MKV-TAG-DATE_WRITTEN
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag DATE_PURCHASED

Descriptor	Value
CCID	MKV-TAG-DATE_PURCHASED
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Play_Counter

Descriptor	Value
CCID	MKV-TAG-PLAY_COUNTER
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag FPS

Descriptor	Value
CCID	MKV-TAG-FPS
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag BPM

Descriptor	Value
CCID	MKV-TAG-BPM
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Measure

Descriptor	Value
CCID	MKV-TAG-MEASURE
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Tuning

Descriptor	Value
CCID	MKV-TAG-TUNING
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Replay Gain (Gain)

Descriptor	Value
CCID	MKV-TAG-REPLAYGAIN_GAIN
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag Replay Gain (Peak)

Descriptor	Value
CCID	MKV-TAG-REPLAYGAIN_PEAK
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Identifiers) ISRC

Descriptor	Value
CCID	MKV-TAG-ISRC
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Identifiers) MCDI

Descriptor	Value
CCID	MKV-TAG-MCDI
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Identifiers) ISBN

Descriptor	Value
CCID	MKV-TAG-ISBN
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Identifiers) Barcode

Descriptor	Value
CCID	MKV-TAG-BARCODE
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Identifiers) Catalog number

Descriptor	Value
CCID	MKV-TAG-CATALOG_NUMBERA
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Identifiers) Label code

Descriptor	Value
CCID	MKV-TAG-LABEL_CODE
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Identifiers) LCCN

Descriptor	Value
CCID	MKV-TAG-LCCN
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Commercial) Purchase Item

Descriptor	Value
CCID	MKV-TAG-PURCHASE_ITEM
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Commercial) Purchase Price

Descriptor	Value
CCID	MKV-TAG-PURCHASE_PRICE
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

Tag (Commercial) Purchase Currency

Descriptor	Value
CCID	MKV-TAG-PURCHASE_CURRENCY
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

FFV1 Conformance Checks (Draft)

Missing header

Descriptor	Value
CCID	OUTOFBAND-HEADER-MISSING
Version	0
Authority	FFV1 Specification
Target Format	FFV1 >=2
Target Format Version	all

Descriptor	Value
Target Format Part	Header
Citation	http://www.ffmpeg.org/~michael/ffv1.html

Rule Clarity:

Quote: “Version 2 and later files use a global header”

Definition: If version is 2 or more, there should be a global header in the container private data

version

Descriptor	Value
CCID	FFV1-HEADER-version
Version	0
Authority	FFV1 Specification
Target Format	FFV1
Target Format Version	all
Target Format Part	Header
Citation	http://www.ffmpeg.org/~michael/ffv1.html

Rule Clarity: Warning

Quote: “version 0, 1 or 3”

Definition: Maximum known version is 3, analysis stops (note: doc sometimes indicates version 4)

version 2

Descriptor	Value
CCID	FFV1-HEADER-version2
Version	0
Authority	FFV1 Specification
Target Format	FFV1
Target Format Version	all
Target Format Part	Header
Citation	http://www.ffmpeg.org/~michael/ffv1.html

Rule Clarity: Warning

Quote: “Version 2 was never enabled in the encoder thus version 2 files should not exist”

Definition: Version 2 is forbidden, analysis stops

micro_version 2

Descriptor	Value
CCID	FFV1-HEADER-micro_version
Version	0

Descriptor	Value
Authority	FFV1 Specification
Target Format	FFV1
Target Format Version	all
Target Format Part	Header
Citation	http://www.ffmpeg.org/~michael/ffv1.html

Rule Clarity: Warning

Quote: “For version 3, micro_version is 4, micro versions prior to this represent pre standard”

Definition: Not supported version, high risk of decoding issue

coder_type

Descriptor	Value
CCID	FFV1-HEADER-coder_type
Version	0
Authority	FFV1 Specification
Target Format	FFV1
Target Format Version	all
Target Format Part	Header
Citation	http://www.ffmpeg.org/~michael/ffv1.html

Rule Clarity:

Quote: “0 (Golomb Rice), 1 (Range coder), 2 (Range coder with custom state transition table)”

Definition: coder_type >2 is not supported

state_transition_delta

Descriptor	Value
CCID	FFV1-HEADER-state_transition_delta
Version	0
Authority	FFV1 Specification
Target Format	FFV1
Target Format Version	all
Target Format Part	Header
Citation	http://www.ffmpeg.org/~michael/ffv1.html

Rule Clarity:

Quote:

Definition: (To be defined)

colospace_type

Descriptor	Value
CCID	FFV1-HEADER-colorspace_type
Version	0
Authority	FFV1 Specification
Target Format	FFV1
Target Format Version	all
Target Format Part	Header
Citation	http://www.ffmpeg.org/~michael/ffv1.html

Rule Clarity:

Quote: “0 (YCbCr), 1 (JPEG2000_RCT)”

Definition: colorspace_type >1 is not supported

bits_per_raw_sample

Descriptor	Value
CCID	FFV1-HEADER-bits_per_raw_sample
Version	0
Authority	FFV1 Specification
Target Format	FFV1
Target Format Version	all
Target Format Part	Header
Citation	http://www.ffmpeg.org/~michael/ffv1.html

Rule Clarity: Are other values valid?

Quote: “commonly 8, 9, 10 or 16”

Definition:

(More header tests to add)

Descriptor	Value
CCID	
Version	
Authority	
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition:

crc_parity

Descriptor	Value
CCID	FFV1-HEADER-crc_parity
Version	0
Authority	FFV1 Specification
Target Format	
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote: “32bit that are choosen so that the global header as a whole or slice as a whole has a crc”

Definition:

LPCM Conformance Checks (Draft)

formatType

Descriptor	Value
CCID	BWF-LPCM-FMT
Version	0
Authority	EBU BWAV v2 Specification
Target Format	BWF
Target Format Version	2
Target Format Part	FormatChunk ‘fmt’
Citation	EBU Tech 3285 v2, pg. 16

Rule Clarity: Inferred

Quote: “If the field of the is set to WAVE_FORMAT_PCM, then the waveform data consists of samples represented in pulse code modulation (PCM) format.”

Definition: WAVE_FORMAT_PCM = 0x0001

bitsPerSample

Descriptor	Value
CCID	BWF-LPCM-BPS
Version	0
Authority	EBU BWAV v2 Specification
Target Format	BWF
Target Format Version	2
Target Format Part	FormatChunk ‘fmt’
Citation	EBU Tech 3285 v2, pg. 17

Rule Clarity: Inferred

Quote: “The field specifies the number of bits of data used to represent each sample of each channel. If there are multiple channels, the sample size is the same for each channel.”

Definition: valid bits per sample 16, 20 or 24

bytesPerSecond

Descriptor	Value
CCID	BWF-LPCM-BYT
Version	0
Authority	EBU BWAV v2 Specification
Target Format	BWF
Target Format Version	2
Target Format Part	FormatChunk ‘fmt’
Citation	EBU Tech 3285 v2, pg. 17

Rule Clarity: Inferred

Quote: “For PCM data, the field of the ‘fmt’ chunk should be equal to the following formula rounded up to the next whole number: $(nChannels \times nSamplesPerSecond \times nBitsPerSample) / 8$ ”

Definition: MUST equal $(nChannels \times nSamplesPerSecond \times nBitsPerSample) / 8$ **only important for compressed formats

blockAlignment

Descriptor	Value
CCID	BWF-LPCM-BLK
Version	0
Authority	EBU BWAV v2 Specification
Target Format	BWF
Target Format Version	2
Target Format Part	FormatChunk ‘fmt’
Citation	EBU Tech 3285 v2, pg. 17

Rule Clarity: Inferred

Quote: “The field should be equal to the following formula, rounded to the next whole number: $(nChannels \times nBitsPerSample) / 8$ ”

Definition: C ontainer size (in bytes) of one set of samples. MUST equal $(nChannels \times nBitsPerSample) / 8$ EBU
**Note: The above formulae do not always give the correct answer. Strictly speaking, the number of bytes per sample $(nBitsPerSample / 8)$ should be rounded first. Then this integer should be multiplied by (which is always an integer) to give . This in turn should be multiplied by to give].

channelCount

Descriptor	Value
CCID	BWF-LPCM-CHN
Version	0
Authority	EBU BWAV v2 Specification
Target Format	BWF
Target Format Version	2
Target Format Part	FormatChunk ‘fmt’
Citation	EBU Tech 3285 v2, pg. 17

Rule Clarity: Inferred

Quote: 1 = mono, 2 = stereo, etc.

Definition: 1 = mono, 2 = stereo, etc.

nChannels

Descriptor	Value
CCID	BWF-LPCM-CHN
Version	0
Authority	EBU BWAV v2 Specification
Target Format	BWF
Target Format Version	2
Target Format Part	FormatChunk ‘fmt’
Citation	EBU Tech 3285 v2, pg. 17

Rule Clarity: Inferred

Quote: “Number of channels in the wave, 1 for mono, 2 for stereo”

Definition: 1 = mono, 2 = stereo, etc.

sampleRate

Descriptor	Value
CCID	BWF-LPCM-SRT
Version	0
Authority	EBU BWAV v2 Specification
Target Format	BWF
Target Format Version	2
Target Format Part	FormatChunk ‘fmt’
Citation	EBU Tech 3285 v2, pg. 17

Rule Clarity: Inferred

Quote: “Frequency of the sample rate of the wave file. This should be 48000 or 44100 etc. This rate is also used by the sample size entry in the fact chunk to determine the length in time of the data.”

Definition: 32000, 44100, 48000, etc.

Container/Stream Coherency Checks (Draft)

Aspect Ratio Match

Descriptor	Value
CCID	COHERENCY-DAR
Version	0
Authority	
Target Format	FFV1/MKV
Target Format Version	
Target Format Part	
Citation	

Rule Clarity:

Quote:

Definition: Display Aspect Ratio indicated in the container (e.g. Matroska) is not the Display Aspect Ratio indicated in the FFV1 stream

Width Match

Descriptor	Value
CCID	COHERENCY-WIDTH
Version	0
Authority	FFV1 and Container Specifications
Target Format	FFV1
Target Format Version	Container
Target Format Part	all
Citation	Header

Rule Clarity:

Quote:

Definition: Width indicated in the container (e.g. Matroska) is not the width indicated in the FFV1 stream

Height Match

Descriptor	Value
CCID	COHERENCY-HEIGHT
Version	0
Authority	FFV1 and Container Specifications
Target Format	FFV1
Target Format Version	Container
Target Format Part	all
Citation	Header

Rule Clarity:

Quote:

Definition: Height indicated in the container (e.g. Matroska) is not the height indicated in the FFV1 stream