

# Format Registry Ontology

Rev. 2003/Mar/10 – SLA

The format registry maintains persistent, unambiguous bindings between public *identifiers* for digital representation formats and *representation information* for those formats. A format is a representation of an *information model*, a formal expression of exchangeable knowledge. Thus, the format registry is a typing mechanism for encoded digital content.

This specification defines two fundamental classes of formats:

1. *Content stream* formats
2. *Physical media* formats

A content stream format is a fixed byte-serialized encoding of an information model, independent of the physical medium underlying its manifestation. (The registry makes no assumption regarding the size of a byte, i.e., a byte is not necessarily an octet.) A content stream is not required to have homogeneous composition, but can be defined as an aggregation of other content streams.

A physical media format is a fixed encoding of a content stream in a tangible form on a physical storage structure.

Format representation information enables the mapping of typed formats into more meaningful concepts by capturing the significant *syntactic* and *semantic* properties of formats, with particular relevance towards the operational needs of digital repositories, including, but not limited to: object format identification, characterization, ingest validation, interchange, migration, emulation, and other archival preservation activities.

Within the registry, formats are classified functionally according to an *ontology* that hierarchically decomposes the format space based upon the fundamental nature of the data units representable by the formats. The top-level organization of the hierarchy is similar to the MIME media type classification.

- 0      **Information model** [ *expression of exchangeable knowledge* ]
- 1      **Content stream** [ *byte-serialized encoding of information model* ]
- 1.1      **Logical** [ *data representing truth and falsity* ]
  - XDR Boolean (RFC 1832)
- 1.2      **Numeric** [ *data representing mathematical ordinality or cardinality* ]
- 1.2.1      **Scalar**
- 1.2.1.1      **Integer**
  - XDR integer (RFC 1832)
- 1.2.1.1.1      **Unsigned integer**
  - XDR unsigned integer (RFC 1832)
- 1.2.1.2      **Real**
- 1.2.1.2.1      **Floating point**
  - IEEE 754
- 1.2.1.3      **Complex**
- 1.3      **Text** [ *character data directly interpretable by humans* ]
  - EBCDIC
  - ISO/IEC 646:1991 (ASCII) → UTF-8
  - ISO/IEC 8859-1:1999 (Latin 1)
  - Mac OS Roman
  - UTF-8 (Unicode, ISO/IEC 10646-1:2000)
  - Windows code page 1252
- 1.3.1      **Structured text** [ *text with structural constraints* ]
  - CSV
  - Tab delimited
- 1.3.1.1      **Mark-up language**
  - HTML
  - ISO 8879:1986 (SGML)
  - LaTeX
  - RTF

1.3.1.1.1	- XML <b>Definition</b> - SGML DTD - XML DTD - XSD → XML
1.3.1.1.2	<b>Transformation</b> - XSLT → XML
1.3.1.1.3	<b>Presentation</b> - CSS - XSL-FO → XML
1.3.1.2	<b>Programming language</b>
1.3.1.2.1	<b>Functional</b> - Lisp
1.3.1.2.2	<b>Declarative</b>
1.3.1.2.2.1	<b>Interpreted</b> - AppleScript - Perl - sh
1.3.1.2.2.2	<b>Compiled</b> - C++ - C# - Fortran - Java
1.3.1.3	<b>Message</b>
1.3.1.3.1	<b>Mail</b> - MIME (RFC 2045)
1.3.1.3.2	<b>News</b> - USENET (RFC 1036)
1.4	<b>Image</b> [ <i>visual data requiring rendering technology for human interpretability</i> ]
1.4.1	<b>Still</b>
1.4.1.1	<b>Font</b> [ <i>character glyph data</i> ]
1.4.1.1.1	<b>Outline</b> - Adobe Type 1 - OpenType - TrueType
1.4.1.1.2	<b>Raster</b>
1.4.1.2	<b>Graphic</b>
1.4.1.2.1	<b>Vector</b>
1.4.1.2.1.1	<b>2D</b> - SVG → XML
1.4.1.2.1.2	<b>3D</b> - VRML
1.4.1.2.2	<b>Raster</b> [ <i>rectilinear array of picture elements</i> ] - GIF - JFIF - PCD - TIFF
1.4.1.3	<b>Page description</b> - PDF - PostScript - Quark Xpress
1.4.2	<b>Motion</b> - AVI - MPEG - QuickTime
1.5	<b>Audio</b> [ <i>aural data requiring playback technology for human interpretability</i> ] - AIFF - MP3 - Real - WAV

1.5.1	<b>Music</b>
	- MIDI
1.6	<b>Application</b> [ <i>arbitrary data requiring technological mediation for human interpretability</i> ]
1.6.1	<b>Communication</b>
1.6.2	<b>Database</b>
1.6.2.1	<b>Hierarchical</b>
1.6.2.2	<b>Relational</b>
1.6.2.2.1	<b>Schema</b>
	- SQL DDL
1.6.2.2.2	<b>Query</b>
	- SQL DML
1.6.2.2.3	<b>Data</b>
	- MySQL
	- Oracle
	- Postgres
1.6.3	<b>Executable</b>
	- ELF
	- EXE
	- Java byte code
1.6.4	<b>Presentation</b>
	- PowerPoint
1.6.5	<b>Spreadsheet</b>
	- Excel
1.6.5.1	<b>Macros</b>
1.6.6	<b>Word processing</b>
	- Word
	- WordPerfect
1.7	<b>Transformation</b> [ <i>composable encodings applied against primary formats</i> ]
1.7.1	<b>Compression</b> [ <i>size minimizing encodings</i> ]
1.7.1.1	<b>Lossless</b>
	- CCITT T.4
1.7.1.1.2	<b>Lempel-Ziv</b>
	- compress
	- gzip
	- LZ77
	- LZW
1.7.1.2	<b>Lossy</b>
	- JPEG
1.7.2	<b>Container</b> [ <i>aggregations of atomic information model units</i> ]
	- jar
	- PKZIP
	- Stuffit
	- tar
1.7.3	<b>Transfer</b> [ <i>data encodings for transmission over a communication channel</i> ]
1.7.3.1	<b>7-bit safe</b>
	- Base64
	- BinHex
	- uuencode
2	<b>Physical media</b> [ <i>tangible encoding of content stream on physical storage structure</i> ]
2.1	<b>Magnetic</b>
2.1.1	<b>Disk</b>
2.1.2	<b>Tape</b>
2.1.2.1	<b>Reel</b>
2.1.2.1.1	<b>9 track</b>
	- ANSI X3.54-1986
2.1.2.2	<b>Cartridge</b>
2.1.2.2.1	<b>3480 class</b>
	- ANSI X3.180-1990
2.1.2.2.2	<b>DLT</b>
	- ISO/IEC 15307:1997
	- ISO/IEC 16382:2000

2.2	Optical
2.2.1	Disk
2.2.1.1	CD-ROM
	- ISO 9660:1988
2.2.1.2	DVD
2.2.2	Film
2.3	Paper
2.3.1	Card
	- Holerith
2.3.2	Tape

1.6.3	Domain [ <i>domain-specific data</i> ]
1.6.3.1	Agriculture
1.6.3.2	Commerce
1.6.3.2.1	EDI
	- ASC X12
	- UN/CEFACT
1.6.3.2.2	Finance
1.6.3.3	Education
1.6.3.4	Engineering
	- ISO 10303 (STEP)
1.6.3.4.1	CAD
	- DXF
	- IGES
1.6.3.5	Humanities
1.6.3.6	Law
1.6.3.7	Medicine
1.6.3.7.1	Pharmaceutical
1.6.3.8	Science
	- HDF
	- netCDF
1.6.3.8.1	Information
1.6.3.8.1.1	Archive
1.6.3.8.1.2	Computer
1.6.3.8.1.3	Library
1.6.3.8.1.3.1	Catalog
	- MARC
1.6.3.8.2	Physical
1.6.3.8.2.1	Astronomy
	- FITS
1.6.3.8.2.2	Biology
	- PDB
1.6.3.8.2.2.1	Genetics
	- GenBank
1.6.3.8.2.3.2	Zoology
1.6.3.8.2.3	Chemistry
1.6.3.8.2.4	Earth Science
1.6.3.8.2.4.1	Geography
1.6.3.8.2.4.1.1	GIS
	- GeoTIFF → TIFF
	- SDTS
1.6.3.8.2.4.2	Geology
1.6.3.8.2.4.3	Meteorology
1.6.3.8.2.4.4	Oceanography
1.6.3.8.2.5	Mathematics
	- Mathematica
1.6.3.8.2.5.1	Statistics
	- SAS
	- SPSS

<b>1.6.3.8.2.6</b>	<b>Physics</b>
<b>1.6.3.8.3</b>	<b>Social</b>
<b>1.6.3.8.3.1</b>	<b>Anthropology</b>
<b>1.6.3.8.3.2</b>	<b>Economics</b>
<b>1.6.3.8.3.3</b>	<b>Political Science</b>
<b>1.6.3.8.3.4</b>	<b>Psychology</b>
<b>1.6.3.8.3.5</b>	<b>Sociology</b>