Preservation taken seriously

(at least … trying to take it that way)

**PRE**servation **FORMA**ts for culture information/e-archives

Is an: EU-funded project (running since January 2014 until December 2017)

   … with very different partners:

   Partners comming from …

   - Center of Expertise
   - Cultural Heritage
   - Software-Evaluation

## Center of Expertise …

- PACKED EXPERTISECENTRUM DIGITAAL ERFGOED VZW, Belgium
- HOGSKOLAN I SKOVDE (University of Skovde), Sweden

## Cultural Heritage …

- RIKSARKIVET, Sweden
- STICHTING NEDERLANDS INSTITUUT VOOR BEELD EN GELUID, Netherlands
- KONINKLIJK INSTITUUT VOOR HET KUNSTPATRIMONIUM, Belgium
- GREEK FILM CENTRE AE, Greece
- LOCAL GOVERNMENT MANAGEMENT AGENCY, Ireland
- STIFTUNG PREUSSISCHER KULTURBESITZ, Germany
- AYUNTAMIENTO DE GIRONA, Spain
- EESTI VABARIIGI KULTUURMINISTEERIUM, Estonia
- KUNGLIGA BIBLIOTEKET, Sweden

## Software-Evaluation und – tests …

- UNIVERSITA DEGLI STUDI DI PADOVA, Italy
- FRAUNHOFER (Ilmenau), Germany

## Coordination / media-partner…

- PROMOTER, Italy

**PRE**servation **FORMA**ts for culture information/e-archives

What is it about ?

**PREFORMA**

1. Some information has to be preserved for a very long time

2. This information (often) is stored in files

3. These files use file-formats

4. These file-formats are (mostly) standardised

5. Usage of files (data, information, …) is done with programms that are build on the standards of file-formats

6. Files that should be readable in many years ahead have to follow the file-format standards

(If not: We cannot be sure that future programms can read/interpret the information contained in the files correctly)

Not everything that calles itself a TIF-file is a TIF-file … and
Not everything that is a TIF-file, is a TIF6.0-file … and
Not everything that is a TIF6.0-file follows the TIF6.0-Baseline-Standard

Not everything that comes along as PDF is a PDF … and
Not everything that is a PDF is a PDF/A … and
Not everything that is a PDF/A is a PDF/A-1b

…

…

…

… What's in a name ? – A closer look is needed …

To be able to „look more closely" we need tools !

**Development of good tools for file-format-validation – that is the objective of PREFORMA**

The development of the tools is coordinated and controlled by Preforma-project.

The software development is done by companies/consortia that have been choosen by Preforma (and that get nearly all of the finances of the project)

**PREFORMA**

**The development of good tools for file-format-validation – objective of PREFORMA**

1. PREFORMA-consortium had defined which formats should be validated

   Decision was for PDF, TIF, MKV/FFV1

2. PREFORMA-consortium had defined requirements for the tools

   Tools have to be easy-to-use, scalable, multilingual, … and (above all) OPEN SOURCE

3. PREFORMA-consortium organised a tender and choose companies/consortia

   The developers of the tools are …

… but …

**PREFORMA**

**The development of good tools for file-format-validation – objective of PREFORMA**

… of course, there was an evaluation of existing file-format validation tools:

- There are very few good open-source validators

- It has been shown that the quality of the results of validation differs very much! While one validator says a given file is valid others decides that the same file is not valid

- It is often very hard to integrate existing validators into already established technical workflows for digital preservation

… now the developers:

**Tools for file-format validation …**

**PDF** :: Is developed by a consortium of Open Preservation Foundation (OPF) and PDF Association – supported by Digital Preservation Coalition



http://verapdf.org/home/

**PREFORMA**

**Tools for file-format validation ...**

**TIF** :: Is developed by easyinnova (Barcelona) and Digital Humanities Lab der Uni Basel



http://www.dpfmanager.org/

**Tools for file-format validation ...**

**Matroska/FFV1** :: Is developed by mediaarea (developers of *mediainfo*) and supported by the developers of Matroska and of FFmpeg



https://mediaarea.net/MediaConch/

PREFORMA

All consortia agreed to create OPEN-SOURCE software (GPLv3+)

More or less every month new „releases" are published. The lastest versions of the three validators can be downloaded from the Preforma-Open Source Portal



http://www.preforma-project.eu/open-source-portal.html

Testing the tools in their respective status of development was done perpetually. Everyone was invited to test the software and report errors or whishes …

An intense phase of tests coordinated and done by Preforma-partners started after the publication of „Release Candidates" in december 2016. The test were done with big and small, valid and corrupted, real world files and synthetic test files, …

By december 2017 the development of the tools should be finished.

Each of the three tools works with APIs that are geared to each other. This way it is possible to create a „Meta-Tool" incorporating all validators developed (and others to come)

PREFORMA

It's quite an effort to develop such tools. Some examples ...

PDF/A might contain images, annotations and signatures → Have to be validated too
PDF/A might contain font-definitions, scripts, forms etc. → Have to be validated too
PDF/A might appear as PDF/A-1a, PDF/A-1b, PDF/A-2a, PDF/A-2b, PDF/A-2u, PDF/A-3 →
The respective specifications have to be taken into account

TIFF might be based on different color-space-definitions → Has to be validated
TIFF might appear as TIFF-EP, LibTIFF, BigTIFF, TIFF-IT, GeoTIFF, ... → Has to be validated
TIFF has a large amount of Tags, TIFF-Tags might be missing, contain wrong
information, contain right information in a wrong way, might be placed at a wrong
place → Each Tag has to be validated

It's quite an effort to develop such tools. Some examples …

(Footnote from http://www.digitalpreservation.gov/formats/content/tiff_tags.shtml)
TIFF image classes are described in the 1992 TIFF 6.0 [specification](#) and may be summarized as follows:
• Class B. Baseline bilevel.
• Class G. Baseline grayscale.
• Class P. Baseline palette-color.
• Class R. Baseline RGB.
• Class Y. Extension YCbCr.

The TIFF/IT specification (ISO 12639, 2004) defines the following image categories:
• CT. Color continuous-tone picture.
• LW. Color line art.
• HC. High-resolution continuous-tone.
• MP. Monochrome continuous-tone picture.
• BP. Binary picture.
• BL. Binary line art.
• SD. Screened data image.
• FP. Final page.

**PREFORMA**

It's quite an effort to develop such tools. Some examples …

Matroska/FFV1 has the problem that these format-codec-combination is just on the way to become widely used and to become a standard

Matroska actually is in the process of formal standardisation with the IETF (The Internet Engineering Task Force)

(One can validate the compliance to a standard only if a standard is well documented and widely used …)

**Important**: Validating if a standard is followed ... that cannot be the only thing

Standards (if existing) are – as shown – in a way flexible, they might be interpreted very strict of (in parts) more freely

To enable cultural heritage institutions to use the tools – the institutions have to be able to influence the validation:

Examples:
- Some cultural heritage institutions might define PDF/A-3 as the format of choice for preservation of text (allowing container-elements in the PDF), another institution decides their format of choice for text is PDF/A-1b (no container-elements allowed)

- One museum thinks it very important that in their TIFF files for each time-entry also the time zone is stored (TIFF/EP), another museum considers this as not so important and wants to check against the baseline standard only

Validating if a standard is followed … that cannot be the only thing

Rules …
- Cultural heritage institutions have to be able to check against their own policies (interpretations of the standards). This implies that the tools must offer the policies als option (or must be able to store them as options)

- It has to be made easy for cultural heritage institutions to define their „rules" and implement them in the „tools"

Fixer …
- In some cases missing or wrongly used tags can be reconstructed automatically from values stored in other tags … this way (sometimes) the compliance with a standard (and policy) might be created automatically. The preforma tools have a basic „metadata fixer" component

Validating if a standard is followed … that cannot be the only thing

Reports …
- It is very important that the tools create easily understandable reports and analyses. Even non-IT-people should be able to understand what and where the problems are

- The reports should be principally available in the language of the user

- The reports have to be available in machine-readable language too, to be passed to other programs that are eventually able to do more corrections

**PREFORMA**

Validating if a standard is followed … that cannot be the only thing

Integration …
- The tools have to be available as single (offline) version, they also have to allow shared use in a LAN or via web

- The tools have to make the integration into existing workflows for digital preservation easy

Scalability…
- The tools must be able to check very small and very big files and also be able to validate small or very big groups of files (e.g. folders with 10000 images)

**PREFORMA**

Status of development…

- The development is nearly finished. Tests revealed that at the end of the preforma project the tools are all working well.

- Everyone is free to continue the development, either by creating validators for other formats or by enhancing or improving the now existing validators.

- The software is created to enable multilinguality but the translations are not done yet (were not part of the project)

**… take it, use it, improve it, share it …**

**!**

Finally an example from the online-validator for TIF-files. The tested version of this form of the Tif-tool does not have the possibilities for defining / setting own „rules" (explained above) …

# DPF Manager

## CONFORMANCE CHECKER

### File ⓘ

Browse …

### Configuration ⓘ

○ Baseline HTML.dpf
○ Baseline JSON.dpf
○ Baseline PDF.dpf
○ Baseline XML.dpf
○ Custom config...

**Check files**

PREFORMA

The (limited) Online-validator (as of  2017-09-21)

A tiff-file …

The report …

🏷️ **IFD Tags**

☐ **Expert mode**  ☐ **Default values**

⛓️ **File structure**

| | Tag Id | Tag Name | Value |
|---|---|---|---|
| ↔️ | 256 | ImageWidth | 120 |
| ↕️ | 257 | ImageLength | 240 |
| ⠿ | 258 | BitsPerSample | 16 |
| ↙️ | 259 | Compression | None |
| 🌢 | 262 | PhotometricInterpretation | Bilevel |
| 🎬 | 274 | Orientation | TopLeft |
| ⠿ | 277 | SamplesPerPixel | 1 |
| 🏷️ | 284 | PlanarConfiguration | Chunky |
| 📄 | 296 | ResolutionUnit | 2 |

**Elements**

📄 **IFD0 - Main image**

📄 **Metadata analysis**

| | Description |
|---|---|
| ✅ | No metadata incoherencies found |

---

**An image IFD must have a X Resolution value**

An Image File Directory(IFD) that contains and image data must have a X Resolution value

*TIFF Baseline 6: Section 3: Bilevel Images. Page 21 TIFF Baseline 6: Section 4: Grayscale Images. Page 22 TIFF Baseline 6: Section 5: Palette-color Images. Page 23 TIFF Baseline 6: Section 6: RGB Full Color Images. Page 23*

---

| | | | |
|---|---|---|---|
| ❌ | IFDI-0004 | IFD1 | Image IFD must have tag X Resolution |
| ❌ | IFDI-0005 | IFD1 | Image IFD must have tag Y Resolution |
| ℹ️ | TAG-281-0005 | IFD1 | MaxSampleValue Tag is not defined. Then 2**(BitsPerSample) - 1 value is assumed |
| ℹ️ | TAG-280-0005 | IFD1 | MinSampleValue Tag is not defined. Then 0 value is assumed |
| ℹ️ | TAG-254-0009 | IFD1 | NewSubfileType Tag is not defined. Then a full resolution, single image, no transparency is assumed |
| ℹ️ | TAG-274-0006 | IFD1 | Orientation Tag is not defined. Then 0th row represents the visual top of the image, and the 0th column represents the visual left-hand side |

Explanation might be simpler!

# ... take it, use it, improve it, share it ...



Final conference: Tallinn, 11-12. October 2017

Thank you very much !

Stefan Rohde-Enslin | Institut für Museumsforschung (SMB-PK) | s.rohde-enslin@smb.spk-berlin.de