

## Empfehlungen zur Verwendung von Dateiformaten beim Einreichen von Forschungsdaten in HeFDI-Repositoryn

Forschende sind im Zuge der [Guten Wissenschaftlichen Praxis \(GWP\)](#) dazu verpflichtet, ihre Forschungsdaten für mindestens 10 Jahre aufzubewahren und damit einhergehend in der Verantwortung, sich Gedanken über die Verwendungsfähigkeit und Lesbarkeit der Daten in diesem Zeitraum und ggf. darüber hinaus zu machen. Dies umfasst eine Reihe von Problematiken, wie z.B. die stetig steigenden Datenmengen, die Entwicklungszyklen von Hard- und Software, die Alterung von Datenträgern, aber auch die Frage nach geeigneten Dateiformaten.

Forschungsdaten liegen, abhängig von Erhebung und Methode, in verschiedenen Dateiformaten vor, die unterschiedlich gut für eine langfristige Nutzung geeignet sind. Für die Sicherung von Daten im wissenschaftlichen Bereich gemäß GWP sollte vor allem auf die Kompatibilität mit verschiedenen Programmen und die verlustfreie Konvertierung in alternative Formate geachtet werden. Weitere gute Kriterien sind Transparenz in Form von Mensch- und Maschinenlesbarkeit sowie die anzunehmende Langzeitstabilität (genormter Standard - bekanntes und verbreitetes Format) eines Dateiformats.

Diese Handreichung empfiehlt Dateiformate zur Ablage in einem Repository und nimmt dabei den mittelfristigen Zeitraum von bis zu 10 Jahren in den Blick. Eine Langzeitarchivierung, die den Erhalt der Nutzbarkeit der Daten über 10 Jahre hinaus zum Ziel hat, erfordert eine aufwändige Vorbereitung und aktive Pflege der Daten. Die Auswahl eines geeigneten Dateiformats ist eine gute Vorbereitung, reicht aber allein nicht aus, um eine Eignung für die Langzeitarchivierung zu gewährleisten; ebenfalls deckt die Ablage in einem Forschungsdatenrepository nicht die Langzeitarchivierung der Daten ab.

Die folgende Tabelle gibt eine aktuelle Einschätzung der Fachstelle Forschungsdatenmanagement und Datenerhalt an der ETH Zürich (ETH-Bibliothek) zur Eignung von häufig verwendeten Datenformaten. Diese Tabelle basiert auf Erfahrung sowie einer ausführlichen Auswertung von Empfehlungen und Richtlinien internationaler Einrichtungen mit Archivierungsauftrag. In der vierten Spalte stehen Empfehlungen zur Konvertierung in geeignetere Formate. Falls eine Konvertierung zwar möglich ist, aber mit einem geringeren Funktionsumfang oder mit Informationsverlust einhergeht, wird empfohlen, die Daten in beiden Formaten abzulegen. Wenn es nicht möglich ist, ein empfohlenes Dateiformat zu verwenden, sind die Daten voraussichtlich 10 Jahre später nicht mehr nutzbar.

| Dateiart                         | Empfohlen  | Bedingt geeignet  | Nicht geeignet  |
|----------------------------------|--|---|---|
| <b>Text</b>                      | <ul style="list-style-type: none"> <li>• PDF/A (*.pdf, bevorzugte Subtypen 2b und 2u)</li> <li>• Unformatierter Text (*.txt oder Quellcode, usw.) kodiert als ASCII, UTF-8 oder UTF-16 mit Byte Order Mark (BOM)</li> <li>• XML (inklusive XSD/XSL/XHTML, etc.; Schema &amp; Buchstabenkodierung inklusive)</li> </ul> | <ul style="list-style-type: none"> <li>• (*.pdf) mit eingebetteten Fonts</li> <li>• Unformatierter Text (.txt, .asc, .c, .h, .cpp, .m, .py, .r usw.) (ISO 8859-1 kodiert)</li> <li>• Rich Text Format (*.rtf)</li> <li>• HTML und XML (ohne externe Inhalte)</li> <li>• Word *.docx</li> <li>• PowerPoint *.pptx</li> <li>• LaTeX und TeX (inkl. lizenzfreie Softwarepakete mit Spezialfonts und resultierendes PDF)</li> <li>• OpenDocument Formate (.odm, .odt, .odg, .odc, *.odf)</li> </ul> | <ul style="list-style-type: none"> <li>• Word *.doc</li> <li>• PowerPoint *.ppt</li> </ul>  |
| <b>Spreadsheets und Tabellen</b> | <ul style="list-style-type: none"> <li>• Komma- oder Tab-begrenzte Text Files (*.csv)</li> </ul>   | <ul style="list-style-type: none"> <li>• Excel *.xlsx</li> <li>• OpenDocument Formate (.odm, .odt, .odg, .odc, *.odf)</li> </ul>  | <ul style="list-style-type: none"> <li>• Excel .xls, .xlsb<br/><b>Konvertierung:</b> in .xlsx</li> </ul>  |
| <b>Rohdaten und Workspace</b>    |  | <ul style="list-style-type: none"> <li>• Unformatierter Text (ASCII-kodiert)</li> <li>• S-Plus (*.sdd)</li> <li>• Matlab (*.mat) ab v7.3 MAT-Datei</li> </ul>   | <ul style="list-style-type: none"> <li>• Matlab-Dateien *.mat (binär) <b>Konvertierung:</b> HDF5-Format.</li> <li>• R-Dateien *.RData<br/><b>Konvertierung:</b>HDF5-</li> </ul> |

|                              |   |  |   |
|------------------------------|---|--|---|
|                              |   | <ul style="list-style-type: none"> <li>• Network Common Data Format oder NetCDF (.nc, .cdf)</li> <li>• Hierarchical Data Format (HDF5) (.h5, .hdf5, *.he5)</li> </ul>  | Format (mit dem Paket rhdf)   |
| <b>Rastergrafik (Bitmap)</b> | <ul style="list-style-type: none"> <li>• TIFF (*.tif, unkomprimiert, möglichst TIFF 6.0+)</li> <li>• Portable Network Graphics (*.png, unkomprimiert)</li> <li>• JPEG2000 (*.jp2, verlustfreie Komprimierung)</li> <li>• Digital-Negative-Format (*.dng)</li> </ul> | <ul style="list-style-type: none"> <li>• TIFF (*.tif, komprimiert)</li> <li>• GIF (*.gif)</li> <li>• BMP (*.bmp)</li> <li>• JPEG/JFIF (*.jpg)</li> <li>• JPEG2000 (*.jp2, verlustbehaftete Komprimierung)</li> </ul> |   |
| <b>Vektorgrafik</b>          | <ul style="list-style-type: none"> <li>• SVG ohne JavaScript binding (*.svg)</li> </ul>   |  | <ul style="list-style-type: none"> <li>• Grafik InDesign (.indd), Illustrator (.ait)</li> <li>• Encapsulated Postscript (*.eps)</li> <li>• Photoshop (*.psd)</li> </ul> |
| <b>CAD</b>                   | <ul style="list-style-type: none"> <li>• AutoCAD Drawing (*.dwg)</li> <li>• Drawing Interchange Format, AutoCAD (*.dxf)</li> <li>• Extensible 3D, X3D (.x3d, .x3dv, *.x3db)</li> </ul>  |  |   |
| <b>Ton, Audio</b>            | <ul style="list-style-type: none"> <li>• WAV (*.wav) (unkomprimiert, pulse-code moduliert)</li> </ul>   | <ul style="list-style-type: none"> <li>• Advanced Audio Coding (*.mp4)</li> <li>• MP3 (*.mp3)</li> </ul>   |   |
| <b>Video<sup>1</sup></b>     | <ul style="list-style-type: none"> <li>• FFV1 Codec (ab Version 3) in Matroska Container (*.mkv)</li> </ul>   | <ul style="list-style-type: none"> <li>• MPEG-2 (.mpg, .mpeg)</li> </ul>   | <ul style="list-style-type: none"> <li>• Windows Media Video (*.wmv)</li> <li>• QuickTime Movie (*.mov)</li> </ul>  |

- MP4, heisst auch MPEG-4 Part 14 (\*.mp4)
- Audio Video Interleave (\*.avi)
- Motion JPEG 2000 (.mj2, .mjp2)

<sup>1</sup>Neben dem Dateiformat (bzw. Containerformat) spielen auch der verwendete Codec und die Kompressionsart eine wichtige Rolle.

Quelle: ETH Zürich, ETH-Bibliothek [Archivtaugliche Dateiformate](#), vereinfacht und kommentiert durch [HeFDI](#), [CC-BY 4.0](#)

## Weitere Informationen

Für Dateiformate, die entweder nicht in den Empfehlungen auftauchen oder als nicht geeignet bezeichnet werden, ist zunächst zu prüfen, ob als Alternative ein Format aus der Empfehlungsliste genutzt werden kann.

Um den Umgang mit Daten, insbesondere Daten in nicht empfohlenen Dateiformaten, zu erleichtern und einen möglichst langen Erhalt der Nutzbarkeit zu ermöglichen, empfiehlt es sich, eine sogenannte README-Datei zusammen mit den Daten abzulegen. In dieser einfachen Textdatei wird der Kontext der Erstellung der Daten beschrieben, v. a. mit welcher Software (inkl. Version) die Daten erstellt wurden, sowie Informationen zu bestimmten Einstellungen von Messinstrumenten, Codierung, und alle weiteren Hinweise, die helfen können, später Rückschlüsse zu ziehen, wie die Daten genutzt werden können.

Es erhöht die Chancen auf langfristige Nutzbarkeit, wenn eingebettete Objekte (wie z. B. Abbildungen, Tabellen, etc.) zusätzlich als separate Datei abgelegt werden.

Bei der Konvertierung empfiehlt es sich, die Qualität des Ergebnisses sorgfältig visuell zu überprüfen, beispielsweise bei Texten insbesondere die Formeln, Sonderzeichen, Umlaute, speziellen Schriftarten.

## Weiterführende Links

[Anleitung zum Konvertieren von .docx zu .pdf](#)

[Vortrag: Der PDF/A-Standard und seine verschiedenen Versionen](#)