



**SLUB**

Wir führen Wissen.

# Langzeitarchivfähige Dateiformate

SLUB Dresden

Version 2.1, 2024-07-26

# Inhaltsverzeichnis

Überblick .....	1
Einleitung .....	1
Risikoanalyse .....	2
Vorbetrachtung .....	3
Kriterien zur Beurteilung .....	3
Risikobetrachtung ausgewählter Dateiformate des SLUBArchivs .....	8
Comma Separated Values (CSV) .....	9
Expert Witness Format (EWF) .....	10
Matroska/FFV1 Multimediacontainer .....	11
PDF/A & PDF/UA (Portable Documents Format) .....	12
Rows of String Values (RSV) .....	13
TIF (Tagged Image File Format) .....	14
XML basierte Dateiformate .....	15
ZIP basierte Dateiformate .....	16



- Neustrukturierung
- kleinere Darstellungsprobleme gefixt
- Aufnahme Diagramme Risikobewertung

## Überblick

Dieses Dokument ist ein Teil der Übernahmespezifikation für das SLUBArchiv. Zur Übernahmespezifikation gehören die folgenden Dokumente <sup>[1]</sup>:

- In der **Übernahmevereinbarung** zwischen SLUB und Dienstnehmer sind die Daten, Ansprechpartner und organisatorischen Randbedingungen beschrieben. Dies schliesst die zu verwendenden Handreichungen für Dateiformate bzw. Objektgruppen mit ein.
- In der **Langzeitarchivfähige Dateiformate** sind die Formate aufgeführt, die die SLUB als potenziell archivfähig bewertet und für die die Funktionalitäten der Formaterkennung, Formatvalidierung und Metadatenextraktion in einem ausreichenden Maß durch das SLUBArchiv gewährleistet werden könnten. Die genaue Festlegung erfolgt spezifisch für jeden Workflow und jede Objektgruppe auf Basis der ermittelten signifikanten Eigenschaften.
- Die **SIP Spezifikation für automatischen Ingest SLUBArchiv** beschreibt den Aufbau der Ablieferungspakete (englisch: Submission Information Package, SIP) mit denen der Dienstnehmer die zu archivierenden Dokumente für das SLUBArchiv bereitstellt.
- Die **DIP Spezifikation für automatischen Access SLUBArchiv** beschreibt den Aufbau eines Auslieferungspaketes (englisch: Dissemination Information Package, DIP), welches für die automatische Weiterverarbeitung zielgruppengerechter Ausspielungen (Access) von im SLUBArchiv archivierter digitaler Datenobjekte (IE) geeignet ist
- Die **Workflow Spezifikation für automatisierte Interaktionen mit dem SLUBArchiv** beschreibt den Prozess der Übergabe zu archivierender Dokumente in das SLUBArchiv (Ingest / AIP Update), das Fehlerprotokoll und den Zugriff auf die archivierten digitaler Objekte (Access).
- Das Dokument **Spezifikation Rechteauszeichnung SLUBArchiv** beschreibt, wie rechtliche Informationen zu einem Datenobjekt kodiert und abgelegt werden müssen.
- Das Dokument **Webservice SLUBArchiv** beschreibt Funktionen, die Dienstnehmer nutzen können, um Informationen über ihre Daten im SLUBArchiv über einen Webservice abzufragen.
- Vom SLUBArchiv verwendete Begriffe sind im **Glossar SLUBArchiv** definiert.

## Einleitung

Digitale Langzeitarchivierung stellt besondere Anforderungen an die verwendeten Dateiformate. Digitale Master müssen in offen spezifizierten Dateiformaten vorliegen, die eine leichte Konvertierbarkeit erlauben, um dem Formatwandel über die Jahre Rechnung zu tragen. Gleichzeitig müssen auch über mehrere Konvertierungen die signifikanten Eigenschaften erhalten bleiben, die eine zukünftige Nutzbarkeit erlauben.

Deswegen ist es wichtig, vor der permanenten Archivierung eine Menge an Dateiformaten

festzulegen, die für prinzipiell langzeitarchivfähig gehalten werden und deren Einspeisung in das Langzeitarchiv erlaubt werden darf.

Langzeitarchivfähige Dateiformate müssen folgenden Kriterien genügen:

- offen standardisiert (notfalls offen spezifiziert); kein proprietäres Format
- weit verbreitet
- geringe Komplexität
- ohne Zugriffsschutzmechanismen wie Kopierschutz, Verschlüsselung, DRM
- selbstdokumentierend
- robust
- keine Abhängigkeiten zu anderen Dateiformaten
- lizenzfrei
- validierbar

Die offene Standardisierung von Dateiformaten erlaubt es im Fall der Formatobsoleszenz, Lese- und Konvertierungsprogramme neu zu erstellen. Weitverbreitete Dateiformate bringen es mit sich, dass das Formatwissen über diese stärker verbreitet und für die Nachwelt eher dokumentiert ist. Hinzu kommt, dass die Verfügbarkeit von geeigneter Software mit hoher Wahrscheinlichkeit länger gegeben ist.

Die geringe Komplexität eines Dateiformates geht oft mit einer besseren Verständlichkeit und einer weniger fehleranfälligen Implementierung einher.

Unter selbstdokumentierenden Dateiformaten versteht man jene, die über eine bestimmte Signatur (magic byte, meist am Dateianfang) eine einfache Identifikation des Dateiformates erlauben. Dateiformate sind dann robust, wenn sich einzelne Bitfehler auf den Inhalt der Datei nur gering oder gar nicht auswirken. Aus diesem Grund wird oft von komprimierten Datenformaten abgeraten. Wenn, wie im Falle des Videodatenformates Matroska/FFV1, verlustfreie Datenkompression mit dem gezielten Hinzufügen von Redundanz, z. B. durch CRC, kombiniert wird, kann die Forderung nach Robustheit erfüllt werden.

Sinnvoll ist es, wenn Dateiformate vollständig spezifiziert sind und nicht auf andere Dokumente angewiesen sind. Dies erhöht die Chance Lese- und Konvertierungsprogramme nur anhand der vorhandenen Spezifikation zu erstellen. Lizenzbehaftete, proprietäre Spezifikationen erschweren den Erhalt, da eine Sicherung der Spezifikationsdokumente notwendig ist und Lizenzgeber über die Zeit nicht mehr verfügbar sein können. Dateiformate sollten aus dem selben Grund auch keine Zugriffsschutzmechanismen enthalten, da der Zugriff auf den Inhalt über die Zeit verloren gehen kann (oder gar nicht möglich ist).

Sinnvoll ist ebenso, dass für die Dateiformate in der digitalen Langzeitarchivierung geeignete Validatoren oder zumindest Referenzimplementierungen zur Verfügung stehen, die eine Beurteilung ermöglichen, inwieweit Dateien von Programmen interpretiert werden können.

Proprietäre Datenformate halten diesen Anforderungen nicht stand.

# Risikoanalyse

## Vorbetrachtung

Basierend auf der Risikodefinition "Risiko ist Eintrittswahrscheinlichkeit mal Schadenshöhe" ist eine Risiko-Bewertung der im Archiv befindlichen Dateiformate vorzunehmen. Arbeitsgrundlage ist der Entwurf der Digital Preservation Risk Matrix der US National Archives (NARA) [https://github.com/usnationalarchives/digital-preservation/tree/master/Digital\\_Preservation\\_Risk\\_Matrix](https://github.com/usnationalarchives/digital-preservation/tree/master/Digital_Preservation_Risk_Matrix).

Das Risiko besteht für das Archiv notwendige Formatmigrationen dann, wenn bei genutzten Dateiformaten Obsoleszenz droht. Die Schadenshöhe ist dabei an den Aufwand der Migration, bzw. Totalverlust der Dateien bei nicht stattfindender Migration gekoppelt. Die Eintrittswahrscheinlichkeit wird durch mehrere Parameter, wie Grad der Verbreitung, Komplexität des Dateiformates, Quelloffene Implementierung, offene Spezifikationen und Abkündigung von Softwareprodukten etc. bestimmt.

## Kriterien zur Beurteilung

### Bewertung und Gewichtung der Kriterien

Die Bewertung der Kriterien erfolgt zwecks einfacherem Handling nach einem Punktesystem. Die Punkte werden pro Kriterium addiert.

Die Bewertung erfolgt aus Sicht des vom SLUBArchiv gewählten Dateiformat-Profiles. Bei Zweifeln oder bei zusammengesetzten Dateiformaten, wie z. B. Matroska-FFV1, wird bei Widersprüchen in den Antworten die Punktezahl *Sonstige* bzw. die pessimistische Variante gewählt.



Die pessimistische Bewertung eines Kriteriums ist dem Umstand geschuldet, dass es sinniger ist ein Risiko zu über- statt zu unterschätzen.



Es können gleiche Fragen in unterschiedlichen Kriterien vorkommen. Der Hintergrund ist der, dass die zugrundeliegende Eigenschaft Auswirkung auf beide Kriterien haben kann.

Um die Nachvollziehbarkeit der Beurteilung zu gewährleisten, sollte eine kurze Begründung hinterlegt werden.



Die Bewertung eines Dateiformates anhand der vorliegenden Kriterien dient zur Bestandsaufnahme der möglichen Risiken, die durch die Verwendung von bestimmten Dateiformaten im Archiv bestehen.

Sie ersetzt nicht eine kritische Reflektion über die Ergebnisse der Bewertung und ist nicht hinreichend für die Auswahl eines archivtauglichen Dateiformates.

Der Fall, dass ein oder mehrere Kriterien mit 0 bewertet werden, ist daher nicht zwingend ein Ausschlusskriterium für die Ungeeignetheit eines Dateiformates. Genausowenig gilt der umgedrehte Fall.

Die Bewertung ist nur ein Indikator, der eine objektivierte Auseinandersetzung mit den möglichen Risiken ermöglicht und damit Basis für ein Risikomanagement.

## Offenlegungsgrad (offen)

Dieses Kriterium ermittelt, inwieweit ein Dateiformat offen dokumentiert ist. Je niedriger die Punktzahl, desto höher das Risiko, dass ein Dateiformat nicht mehr migrierbar ist, weil keine Konvertierungsprogramme entwickelt werden können.

- Ist das Dateiformat offen spezifiziert, d.h. das Recht die Spezifikation zu lesen, anzuwenden und basierend darauf Werkzeuge zu entwickeln ist nicht eingeschränkt?? (Ja=2 / Nein=0 / Sonstige=0)
- Gibt es aktuell verfügbare und lauffähige Validierungswerkzeuge? (Ja=2 / Nein=0 / Sonstige=1)
- Sind diese Validierungswerkzeuge Opensource? (Ja=1 / Nein=0 / Sonstige=0)
- Ist das Dateiformat durch eine internationale Standardisierungsorganisation geprüft und als Standard veröffentlicht? (Ja=2 / Nein=0 / Sonstige=1)
- Ist die Spezifikation kostenfrei zugänglich? (Ja=1 / Nein=0 / Sonstige=0)
- Ist die Spezifikation komplett verfügbar? (Ja=2 / Nein=0 / Sonstige=1)

## Adoptionsgrad (verbreitet)

Dieses Kriterium definiert inwieweit ein Dateiformat verbreitet ist. Je niedriger die Punktzahl, desto höher das Risiko, dass ein Dateiformat nicht mehr migrierbar ist, weil niemand das Dateiformat kennt und es nicht einmal mehr anzeigbar ist.

- Wird das Dateiformat durch Behörden und öffentliche Einrichtungen häufig verwendet? (Ja=1 / Nein=0 / Sonstige=0)
- Wird das Dateiformat auch außerhalb von Behörden und öffentlichen Einrichtungen häufig verwendet? (Ja=2 / Nein=0 / Sonstige=0)
- Wird das Dateiformat aktiv betreut? (durch eine Community=3 / einzelne Person=1 / einzelne öffentliche Einrichtung=2 / gar nicht=0)
- Sind für das Dateiformat verschiedene Anzeigeprogramme verfügbar? (Ja=2 / Nein=0 / Sonstige=0)
- Hat die Archiv- und Bibliothekscommunity dieses Dateiformat als Archivfähiges Datenformat allgemein anerkannt? (Ja=2 / Nein=0 / Sonstige=0)

## Transparenzgrad (transparent)

Dieses Kriterium definiert die einfache Durchführbarkeit von Analysen. Je niedriger die Punktezahl, desto höher ist das Risiko Dateien nur mit unverhältnismäßig hohem Aufwand analysieren zu können.

- Ist das Dateiformat menschenlesbar und mittels Texteditor anzeig- und bearbeitbar? (Ja=1 / Nein=0 / Sonstige=0)
- Ist die verfügbare Spezifikation des Dateiformates detailliert genug, so dass man das

Dateiformat mittels Hexeditor in Grundzügen analysieren kann? (Ja=2 / Nein=0 / Sonstige=0)

- Verwendet das Dateiformat Standardzeichen oder andere bekannte Kodierungsmethoden, z. B. UTF8 für Zeichen bzw. Zahlen nach IEEE? (Ja=1 / Nein=0 / Sonstige=0)
- Ist Software, die das Dateiformat erzeugen kann für wenig Geld oder gar kostenlos erhältlich? (Ja=2 / Nein= 0 / Sonstige=1)
- Ist Software, die das Dateiformat erzeugen kann für aktuelle Computerplattformen verfügbar? (Ja=1 / Nein=0 / Sonstige=0)
- Ist das Dateiformat komprimiert bzw. benutzt es Datenkompression? (Ja=0 & Nein=2 / Sonstige=1)
- Erlaubt das Dateiformat nutzerdefinierte Datenkompression, die zu Veränderungen der essentiellen Charakteristiken, z. B. Kompressionsartefakte, der Datei führen? (Ja=0 / Nein=2 / Sonstige=0)

### **Grad der Selbstdokumentation (selbstdokumentiert)**

Dieses Kriterium definiert, inwieweit das Dateiformat selbstdokumentiert ist. Je niedriger die Punktzahl, desto höher das Risiko, dass man ohne externe Metadaten eine Datei nicht zuordnen und analysieren kann.

- Enthält das Format eine Selbstbeschreibung durch deskriptive Metadaten? (Ja=2 / Nein=0 / Sonstige= 0)
- Enthält das Format eine Selbstbeschreibung durch technische Metadaten? (Ja=2 / Nein=0 / Sonstige=0)
- Enthält das Format eine Selbstbeschreibung durch administrative Metadaten? (Ja=1 / Nein=0 / Sonstige=0)
- Liegen diese Metadaten in einem anerkannten, internationalen Metadatenstandard vor? (Ja=2 / Nein=0 / Sonstige=0)
- Sind die Formatmetadaten robust genug für ein akkurate Dateianalyse? (Ja=2 / Nein=0 / Sonstige=1)
- Erlaubt das Format den Nachweis der Authentizität durch eingebettete elektronische Signaturen ohne die Nutzbarkeit einzuschränken? (Ja=2 / Nein=0 / Sonstige=1)

### **Grad der Abhängigkeit von externer Hard- und Software (technologieunabhängig)**

Dieses Kriterium definiert die Unabhängigkeit von spezifischer Soft- und Hardware und ist ein Indiz für proprietäre Formate. Je niedriger die Punktezahl, desto höher das Risiko, dass ein Erhalt der Nutzbarkeit durch Emulation und Formatmigration wegen fehlender Abspielumgebung nicht mehr durchführbar sind.

- Benötigt das Dateiformat eine ganz spezielle Hardwareumgebung (z. B. Grafikkarten, Soundkarten, Speicheranforderungen, Blu-rayplayer, Spielekonsole), um damit zu interagieren, es nutzen und migrieren zu können? (Ja=0 / Nein=3 / Sonstige=1)
- Benötigt das Dateiformat proprietäre Software, um es anschauen, nutzen und migrieren zu

können? (Ja=0 / Nein=3 / Sonstige=1)

- Benötigt das Dateiformat spezielle Plugins oder Scripte, um es anschauen, nutzen und migrieren zu können? (Ja=0 / Nein=2 / Sonstige=0)
- Benötigt das Dateiformat ein bestimmtes Betriebssystem, um es anschauen, nutzen und migrieren zu können? (Ja=0 / Nein=1 / Sonstige=0)

### **Grad der Lizenz- und Patentabhängigkeit (rechtefrei)**

Dieses Kriterium definiert das Risiko der Behinderungen von Eigenentwicklungen für die Formatmigration. Je niedriger die Punktezahl, desto höher das Risiko, dass ein Dateiformat nicht mehr migrierbar ist, weil keine Konvertierungsprogramme entwickelt werden können.

- Ist das Format Gegenstand von Lizenz- oder Patentansprüchen, die die Entwicklung von Open-Source-Tools zum Öffnen und Verwalten der Dateien behindern könnten? (Ja=0 / Nein=3 / Sonstige=1)
- Sind mit dem Format Gebühren verbunden, die sich aus den Lizenz- oder Patentansprüchen ergeben? (Ja=0 / Nein=2 / Sonstige=1)
- Unterliegt die Formatspezifikation Open-Source-Lizenzbedingungen? (Ja=3 / Nein=0 / Sonstige=1)

### **Grad der technischen Schutzmaßnahmen (drmfrei)**

Dieses Kriterium definiert die Gefährdung der Nutzbarkeit durch nicht mehr entschlüsselbare oder verwendbare Inhalte. Je niedriger die Punktezahl, desto höher das Risiko dass eine Datei nicht mehr nutzbar ist.

- Verfügt das Format über die Möglichkeit, die gesamte oder einen Teil der resultierenden Datei zu verschlüsseln? (Ja=0 / Nein=2 / Sonstige=0)
- Erfordert das Format die Verwendung von Verschlüsselung? (Ja=0 / Nein=2 / Sonstige=1)
- Können technische Schutzmaßnahmen (z. B. Digital Rights Management) angewendet werden? (Ja=0 / Nein=2 / Sonstige=0)
- Erlaubt das Format eingebettete Informationen zur Durchsetzbarkeit von Rechtsansprüchen, die die zukünftige Nutzbarkeit gefährden könnten, wie z. B. Wasserzeichen, zeitlich befristete Lizenzschlüssel? (Ja=0 / Nein=2 / Sonstige=0)

### **Grad der Alterung (aktuell)**

Dieses Kriterium ist sowohl ein Indiz für mögliche Formatobsoleszenz, als auch für die Stabilität von Dateiformaten. Je niedriger die Punktezahl, desto höher das Risiko, dass eine Datei nicht gelesen werden kann, weil die Dateiformatspezifikation veraltet oder zu frisch ist.

- Wann wurde die Formatspezifikation erstmals erstellt? (>30Jahre=0 / >20Jahre=1 / >10 Jahre=2 | <10 Jahre =1)
- Wann wurde die Formatspezifikation letztmalig aktualisiert? (>30 Jahre=0 / >20 Jahre=1 / >10Jahre=2 | <10 Jahre =3)

- Ist die Formatspezifikation unter einem üblichen persistenten Identifier bzw. persistenter URL auffindbar? (Ja=2 / Nein=0 / Sonstiges=0)

## Komplexitätsgrad (einfach)

Dieses Kriterium definiert, wie umfangreich Formatwissen aufgebaut werden muss und wieviel Aufwand für die Fehleranalyse und Entwicklung von Migrationslösungen aufgebracht werden muss. Je niedriger die Punktezahl, desto höher das Risiko, dass Formatwissen nur mit unverhältnismäßig hohem Aufwand aufgebaut und Dateien nicht hinreichend analysiert werden können.

- Verweist die Spezifikation explizit auf andere Dateiformat-Standards? (Ja=0 / Nein=2 / Sonstige=1)
- Erlaubt die Spezifikation Varianten des Dateiformats? (Ja=0 / Nein=3 / Sonstige=1)
- Besteht die Spezifikation aus mehreren Teilspezifikationen, die für sich veröffentlicht wurden und nicht alle aufeinander verweisen? (Ja=0 / Nein=2 / Sonstige=1)
- Verwendet das Datenformat für Daten und Metadaten ein eingeschränktes, kontrolliertes Vokabular (Ja=1 / Nein=0 / Sonstiges=1)
- Ist die Spezifikation umfangreich? (<100 Seiten=3 / < 1000 Seiten=2 / < 5000 Seiten=1 / >5000 Seiten=0)
- Erlaubt die Spezifikation mehrere Kodierungsvarianten von Daten? (Ja=0 / Nein=2 / Sonstige=1)
- Ist das Datenformat multiplex fähig, sprich: Mit anderen Datenformaten im Bitstrom mischbar, z. B. weil es Datenpaket orientiert aufgebaut ist (wie z. B. MP3)? (Ja=0 / Nein=2 / Sonstige=0)

## Grad der Robustheit (robust)

Dieses Kriterium definiert wie gut ein Dateiformat mit fehlerhaften Übertragungen umgehen kann. Je niedriger die Punktezahl, desto höher das Risiko, dass Inhalte von beschädigten Dateien nicht wieder hergestellt werden können.

- Ist das Dateiformat menschenlesbar und mittels Texteditor anzeig- und bearbeitbar? (Ja=2 / Nein=0 / Sonstige=0)
- Ist das Dateiformat XML basiert? (Ja=2 / Nein=0 / Sonstige=0)
- Verwendet das Dateiformat einen eingeschränkten Bytewerte-Bereich für die Daten? (Ja=1 / Nein=0 / Sonstige=0)
- Verwendet das Dateiformat in Voreinstellung Kompression? (Ja=0 / Nein=2 / Sonstige=0)
- Verwendet das Dateiformat als zentralen Einstiegspunkt den Dateianfang das Dateiende oder ist es Datenblock-orientiert? (Dateianfang=2 / Datenblock=3 / Dateiende=0 / Unbekannt=1)
- Sind die Headerinformationen Prüfsummen gesichert? (Ja=2 / Nein=0 / Sonstige=0)
- Sind alle Daten außerhalb des Headers Prüfsummen gesichert? (Ja, alle=2 / Ja, nur Metadaten=1 / Ja, nur Nutzdaten=1 / Nein=0 / Sonstige=0)
- Sind 1-Bitfehler erkennbar? (Ja=2 / Nein=0 / Sonstige=1)
- Sind 1-Bitfehler korrigierbar (z. B. bei CRC32)? (Ja=2 / Nein=0 / Sonstige=1)

- Sind n-Bitfehler erkennbar? (Ja=2 / Nein=0 / Sonstige=1)

## Risikobetrachtung ausgewählter Dateiformate des SLUBArchivs

Dieses Dokument listet Dateiformate auf, die nach **aktuellem** Stand der Bearbeitung durch die SLUB grundsätzlich für die dauerhafte Aufnahme in das Langzeitarchiv der SLUB in Frage kommen.



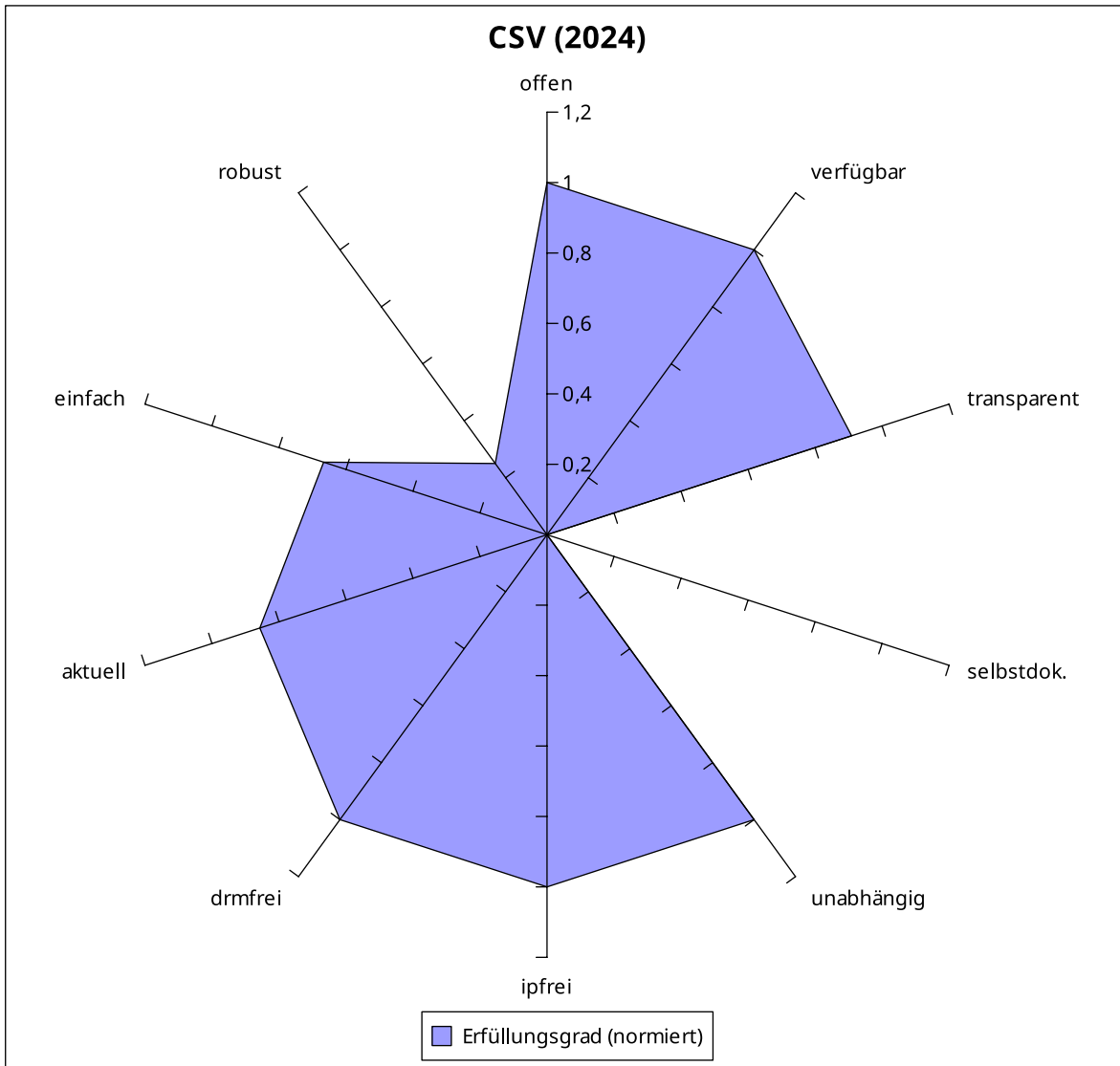
Da einige Kriterien Abhängigkeiten zum Zeitpunkt der Erstellung aufweisen, sind in Klammern im Titel der jeweiligen Diagramme die Bewertungszeitpunkte angegeben.



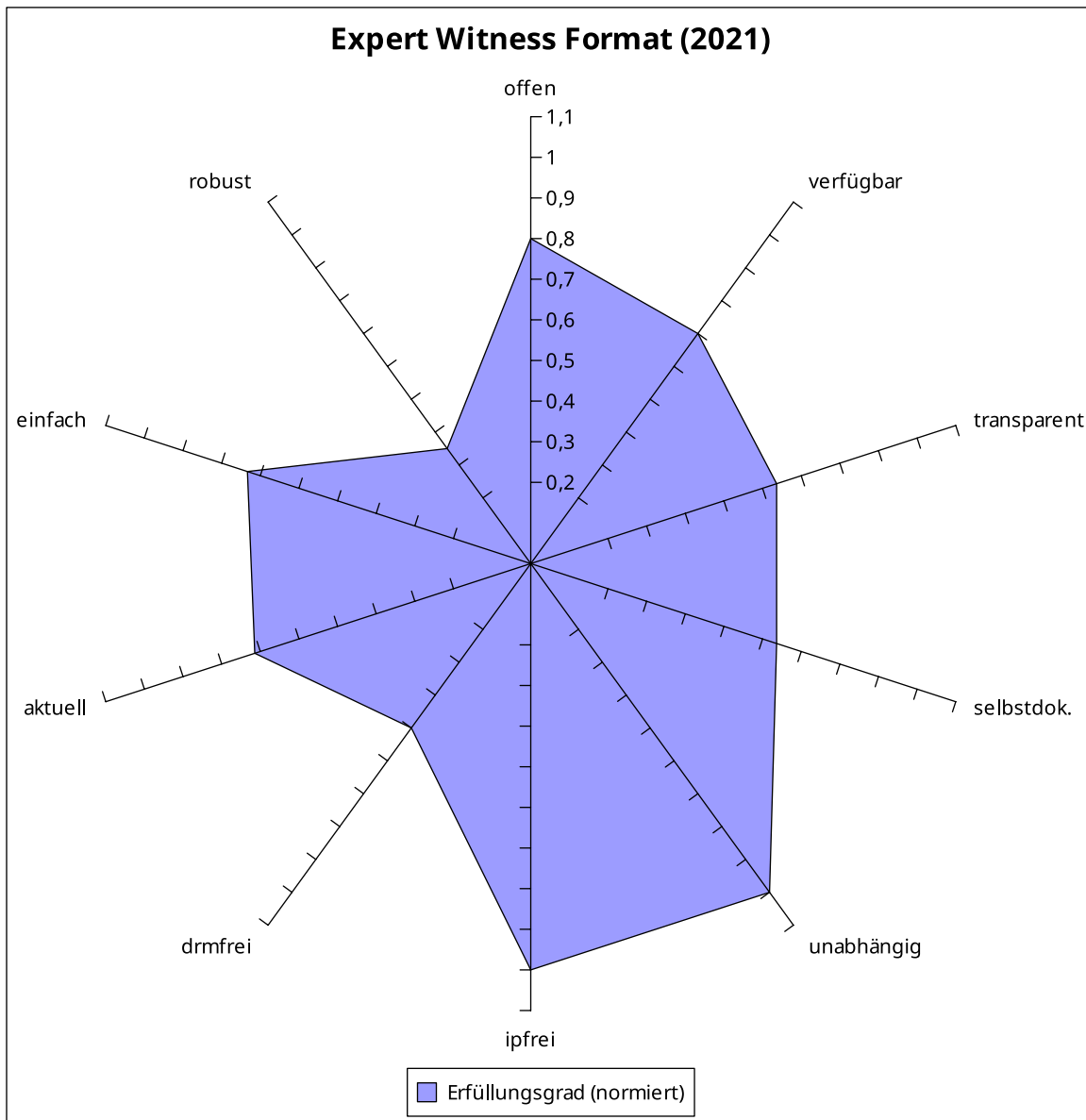
Die **exakte Festlegung**, welches Dateiformat im SLUBArchiv verwendet wird, erfolgt durch das SLUBArchiv **spezifisch** für jeden Workflow und jede Objektgruppe auf Basis der zu diesem Zweck ermittelten und dokumentierten **signifikanten Eigenschaften**.

Für bestimmte Dateiformate hat die SLUB Handreichungen herausgegeben, die die Kriterien für die Aufnahme in das SLUBArchiv genauer spezifizieren.

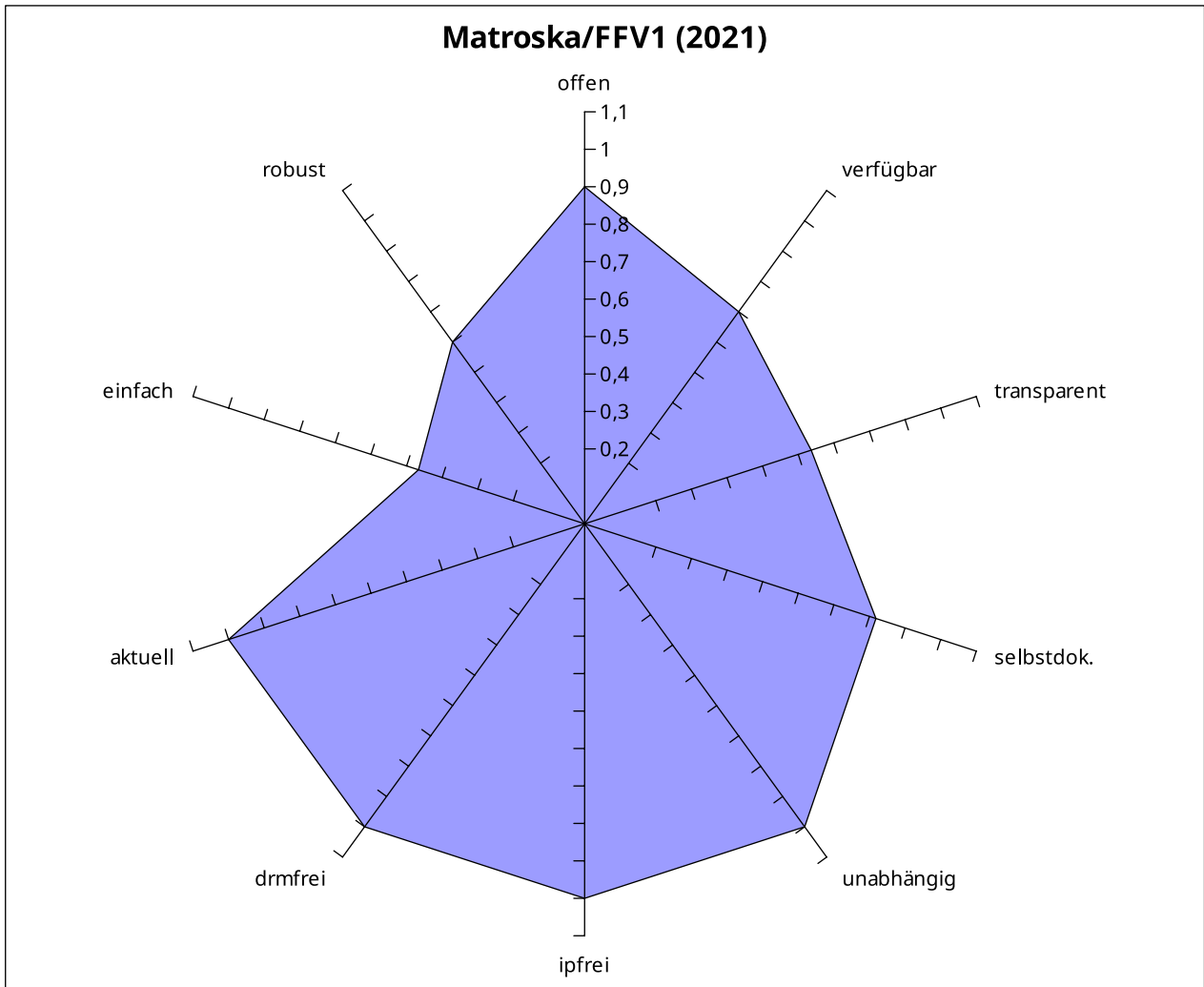
# Comma Separated Values (CSV)



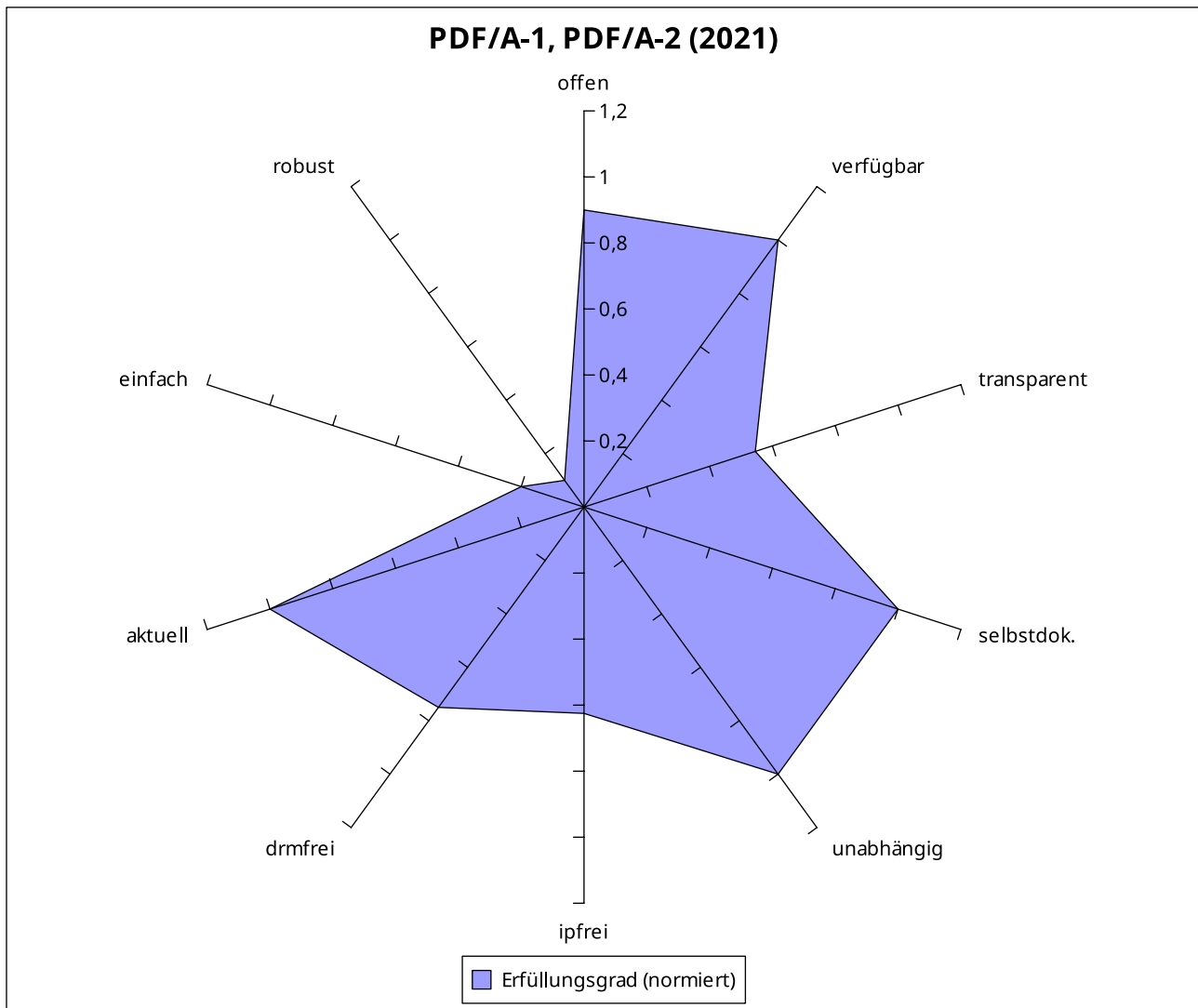
# Expert Witness Format (EWF)



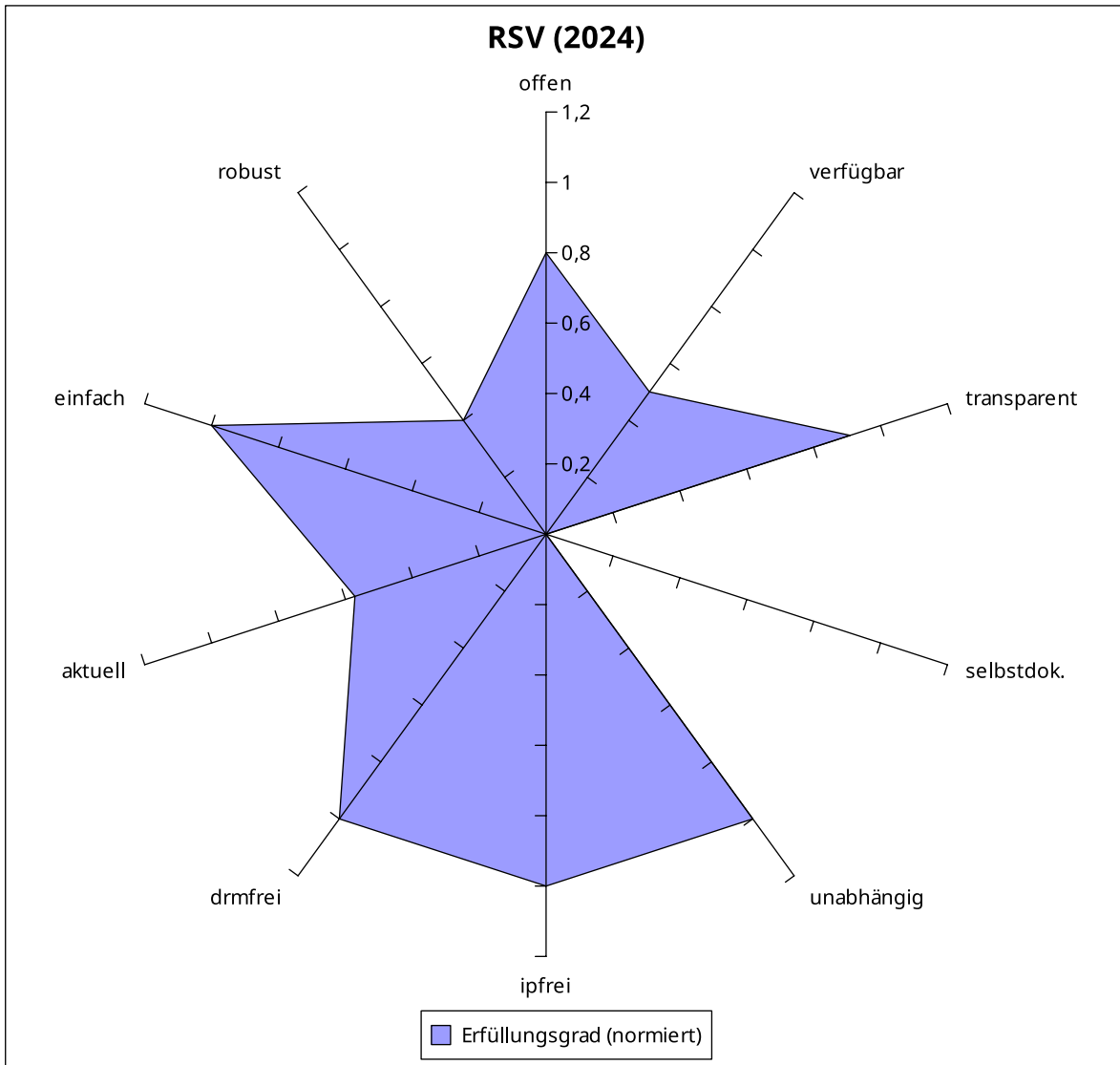
# Matroska/FFV1 Multimediacontainer



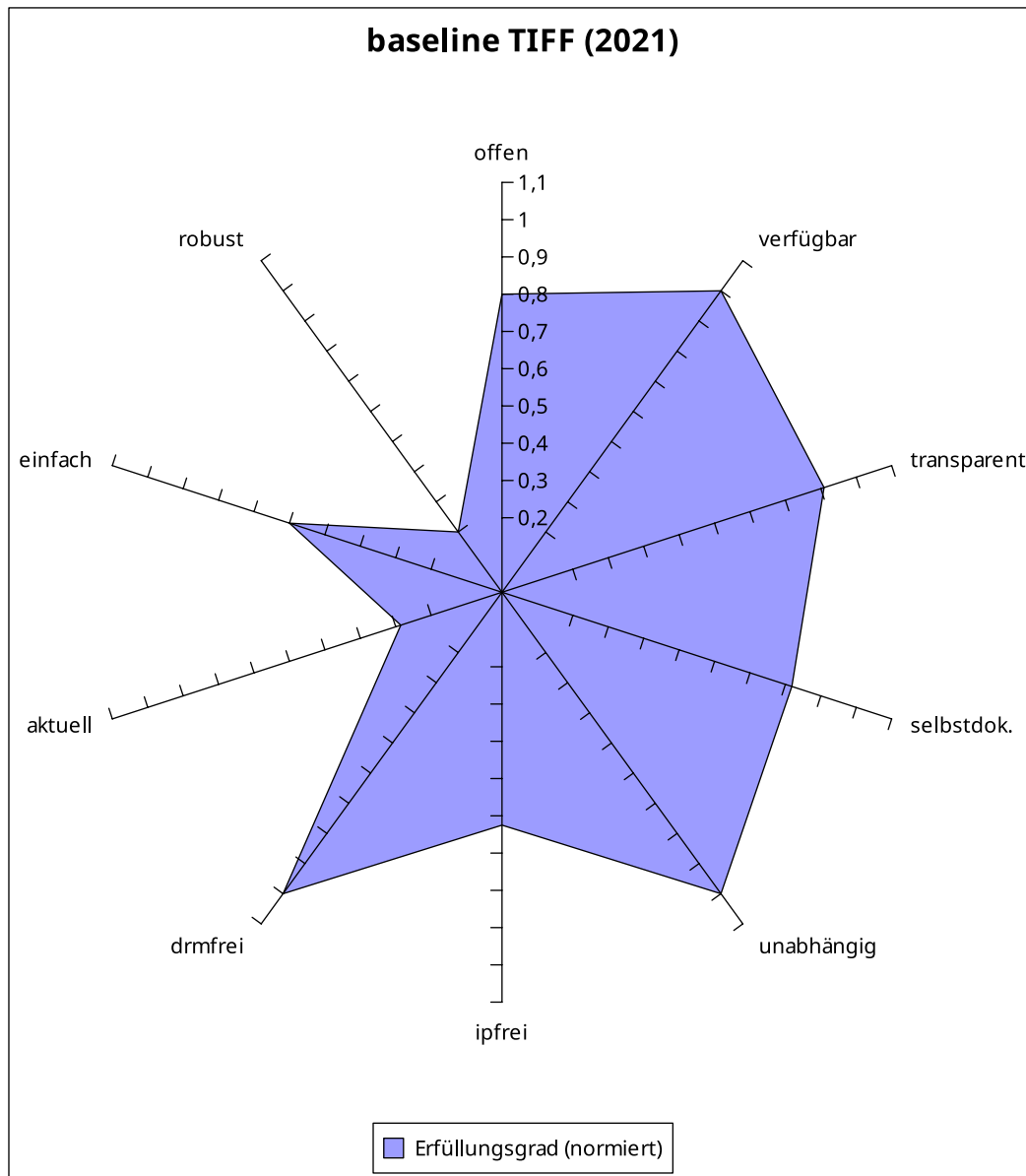
# PDF/A & PDF/UA (Portable Documents Format)



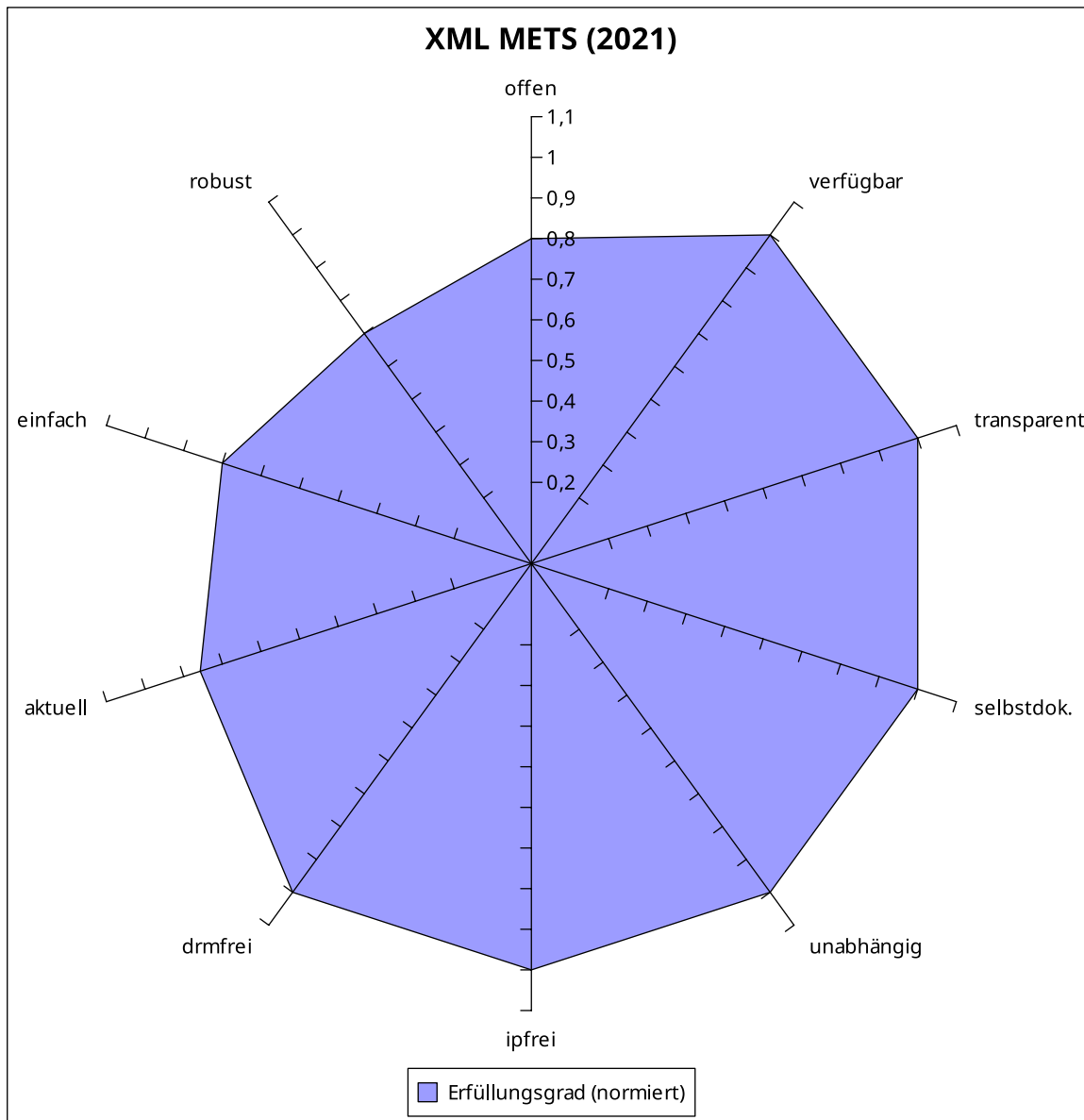
# Rows of String Values (RSV)



# TIF (Tagged Image File Format)

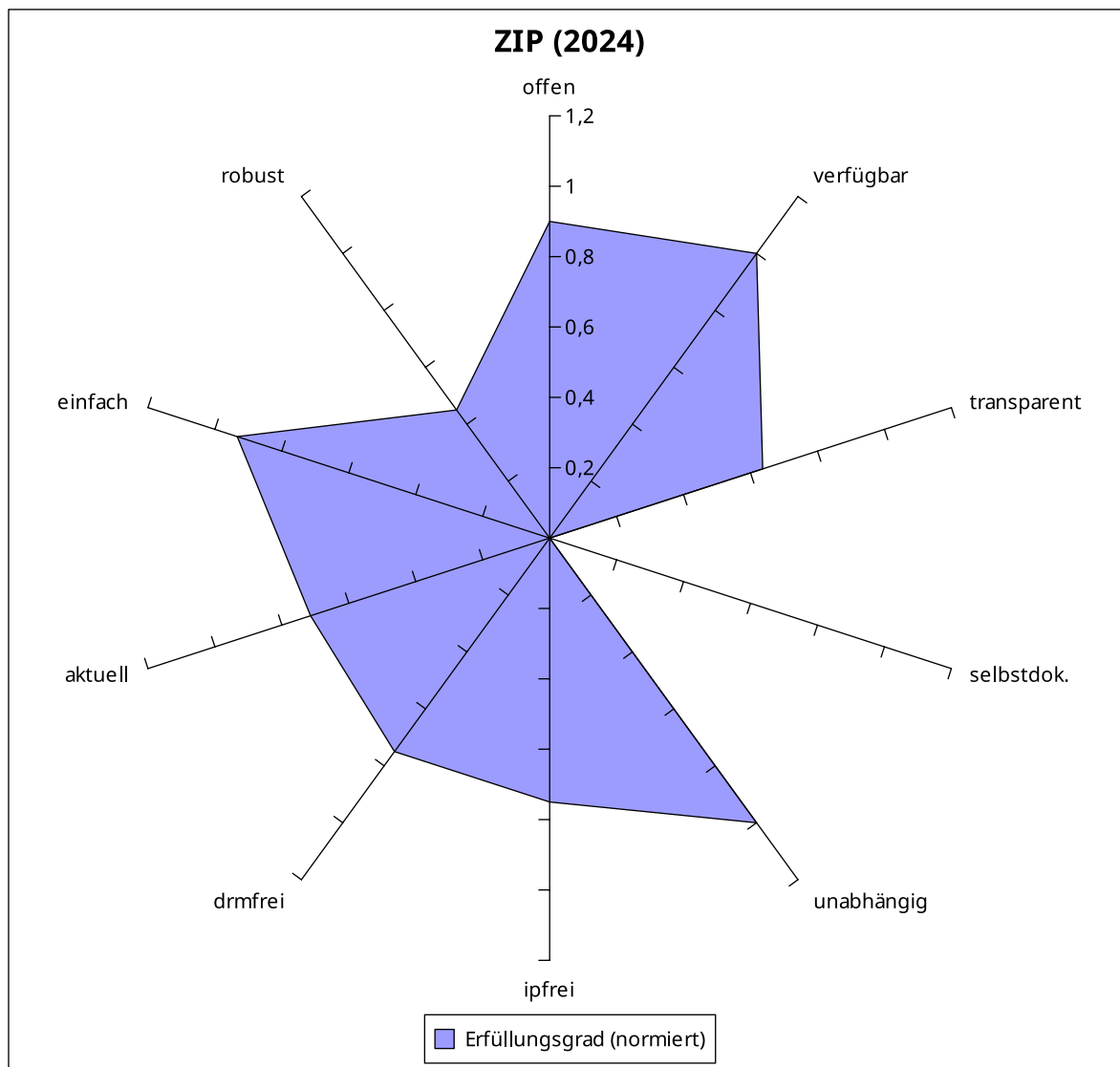


# XML basierte Dateiformate



Für Dateiformate, die auf XML basieren, kann es vereinzelt zu leichten Abweichungen bei den Bewertungskriterien kommen.

## ZIP basierte Dateiformate



Für Dateiformate, die auf ZIP basieren, kann es vereinzelt zu leichten Abweichungen bei den Bewertungskriterien kommen.

[1] Die genannten Dokumente sind auf der Webseite des SLUBArchivs unter <https://slubarchiv.slub-dresden.de/technische-standards-fuer-die-ablieferung-von-digitalen-dokumenten/> veröffentlicht und sind dort, technisch bedingt, spezifischer benannt.