

Submission Guidelines

Leitfaden für Datenlieferungen an das Langzeitarchiv EWIG des Zuse-Instituts Berlin (ZIB)

Stand: Dezember 2019

Das Dokument erläutert Optionen für die technisch-organisatorischen Details einer Übernahme von Daten in das digitale Langzeitarchiv des ZIB. Die Submission Guidelines dienen als Leitfaden für Datengeber. Sie orientieren sich an den Standards ISO 14721:2012 (Open Archival Information System – OAIS) und DIN 31645 (Leitfaden zur Informationsübernahme in digitale Langzeitarchive).

Erfahrungen aus Datenübernahmen in das Langzeitarchiv EWIG fließen in neue Versionen des Dokuments ein.

Folgende Grundbedingungen gelten für Datenübernahmen:

- Zusammengehörige Daten im Sinne einer Intellektuellen Einheit (IE) werden vom Datengeber als strukturierte Datenpakete (digitale Objekte) geliefert. Diese Objekte werden im Langzeitarchiv als Einheiten verwaltet. Für den Transfer können mehrere (vollständige) Einheiten zusammengefasst werden, wenn sie wie im Anhang beschrieben voneinander abgegrenzt sind.
- Die Repräsentationen der Informationsobjekte werden soweit möglich nur in archivfähigen Formaten geliefert. Das Langzeitarchiv kann im Vorfeld der Datenübernahme Empfehlungen für archivfähige Formate geben.
- Eine Transformation der Dateien durch Konvertierung (Migration bzw. Normalisierung) ist nicht im Zuge der Datenübernahme vorgesehen – sie erfolgt gegebenenfalls beim Datengeber.
- Metadaten werden soweit möglich in validen Metadatenformaten geliefert.
- Es werden ausschließlich Metadaten in das Langzeitarchiv übernommen, denen mindestens eine digitale Repräsentation eindeutig zugeordnet werden kann. Umgekehrt werden keine digitalen Repräsentationen übernommen, zu denen keine sie beschreibenden Metadaten vorliegen.
- Alle Datei- und Pfadnamen sollen nach dem gleichen Schema aufgebaut sein und keine Leer- und Sonderzeichen enthalten. Ferner wird eine einheitliche Ordnerstruktur angestrebt.
- Es werden bei einer Lieferung ausschließlich Transferpakete in das Langzeitarchiv übernommen. Eventuell zusätzlich übermittelte Dateien werden ignoriert.
- Die signifikanten Eigenschaften der Daten sollen gemäß den Erwartungen (zukünftiger) Nutzer der Daten berücksichtigt sein.

1. Notwendige Metadaten der Datenlieferung

Metadaten, die eine gesamte Lieferung beschreiben, werden in einem **Submission Manifest** zusammengestellt. Dieses Submission Manifest dient vorrangig der Identifizierung der Lieferung und der Ansprechpartner beim Datengeber im Zuge des Lieferprozesses.

Gemäß der nachfolgenden Aufstellung wird es je nach Lieferformat entweder innerhalb von METS (submission-manifest.xml) oder als separate Textdatei (submission-manifest.txt) im YAML-Format (<https://yaml.org/>) bereitgestellt (vgl. Abschnitt 2 und Anhang). Das Submission Manifest ersetzt nicht die individuellen Beschreibungen der Objekte in den Metadaten des Transferpakets.

Die Textdatei muss UTF8-codiert sein. Alle Felder sind Pflichtfelder, es sei denn sie sind als *optional* gekennzeichnet. Sie wird in enger Abstimmung zwischen den Partnern mit folgenden Angaben erstellt:

SubmissionManifestVersion: 2.0

SubmittingOrganization: Name der abliefernden Institution (Datengeber).

OrganizationIdentifier: z.B. ISIL-Nummer der abliefernden Institution. Dient der eindeutigen Zuordnung einer im Archiv bekannten Organisation.

ContractNumber: Vertragsnummer (Verweis auf zugrundeliegende vertragliche Regelung zw. Datengeber und Langzeitarchiv). Dient der administrativen Zuordnung der Daten etwa zu Laufzeiten und Speichermengen.

Contact: Name der administrativ verantwortlichen Person in der abliefernden Institution, die als Kontakt zu grundsätzlichen Fragen der Datenlieferungen zuständig/befugt ist. Die Person ist befugt, Entscheidungen hinsichtlich Datenmanagement-Prozessen (Rücklieferung, Veränderung o.ä.) zu treffen. Der Kontakt erhält zudem das Übernahmeprotokoll nach Abschluss des Ingest-Prozesses (Format: Nachname, Vorname). Dient der Kommunikation mit dem Archiv, somit hier keine kurzzeitig befristeten Personen benennen.

ContactRole: Funktion von *Contact* in der abliefernden Institution (z.B. Abteilungsleitung Digitalisierung, Leitung Digitale Dienste, Sammlungsleitung o.ä.). Dient der Recherche nach Ansprechpartnern in der Institution nach Ausscheiden des Contact.

ContactEmail: E-Mail-Adresse von *Contact*. Dient als Adresse zur Kommunikation mit Datengeber.

TransferCurator: Name der Person, die die Daten bereitstellt (wenn abweichend von *Contact*). Wird bei konkreten Rückfragen bzgl. des aktuellen Transfers kontaktiert. Erhält Kopie des Übernahmeprotokolls. (Format: Nachname, Vorname).

TransferCuratorEmail: E-Mail-Adresse von *TransferCurator*.

SubmissionName: Eindeutiger, vom Datengeber vergebener Identifikator, der die Lieferung kennzeichnet. Bei aggregierenden Einrichtungen, die Datenlieferungen mehrerer Bereiche einer Institution abgeben, empfiehlt sich ein Zuordnungspräfix, wenn durch Signaturen/andere Nummer keine Zuordnung möglich ist. (Beispiel: GrSammlung_L_x42-2020 für eine Graphische Sammlung oder UA_S2716-2 bei einem Universitätsarchiv als Teil einer Universitätsbibliothek). Dient der eindeutigen Identifizierung einer Lieferung aus Sicht der einliefernden Einrichtung. Erlaubte Zeichen: A-Za-z0-9_()#-

SubmissionDescription: Selbstgewählte inhaltliche Beschreibung der Lieferung. Sollte in jedem Fall einen (beschreibenden) Titel der Daten beinhalten. Zum Beispiel einen Projekttitel: „Digitalisierungsprojekt Korrespondenzen 1995, Teil 1 bis 5. Siehe auch Artikel zum Projekt von Mustermann (1995): XYZ...“.

RightsHolder: Der Inhaber der nötigen Rechte zur Veröffentlichung oder Lizenzierung. Dies beinhaltet auch die Angabe zu Inhabern von Leistungsschutzrechten. Bei mehreren Rechteinhabern sollte angegeben sein, wer welche Rechte hält (z.B. bei Film: Regisseur, Tonmeister, Komponist, Drehbuchautor usw.). (Format: Nachname_1, Vorname_1 (Rolle_1) usw., mit Semikolon als Trennzeichen).

Falls keine Rechte (mehr) bestehen, erfolgt der Eintrag wie folgt: N/A

Rights: Eine Auszeichnung des Rechtstatus durch einen URI von rightsstatements.org oder loc.gov (Kennzeichnung von Gemeinfreiheit entweder durch <http://id.loc.gov/vocabulary/preservation/copyrightStatus/pub> oder <http://rightsstatements.org/vocab/NKC/1.0/>)

RightsDescription (optional): Dient der Beschreibung der Nutzungsrechte bei komplizierteren Sachverhalten (nicht maschinell auswertbar; siehe Beispiele unten).

License: Nutzungsbedingungen für die potentielle Nachnutzung der Objekte, zu denen in den gelieferten Metadaten *keine* Lizenzangaben vorhanden sind. Die hier angegebene Lizenz dient als „Fallback-Information“ für das Langzeitarchiv falls keine sonstigen Lizenzangaben vorhanden sind. Lizenzangaben müssen als URI erfolgen, damit sie maschinell auswertbar sind. Beispiel: Creative Commons Lizenzen. Gemeinfreiheit wird durch <https://creativecommons.org/publicdomain/mark/1.0/> gekennzeichnet. Kann keine URI angegeben werden, erfolgt der Eintrag wie folgt: N/A

AccessRights: Angaben zu Nutzungs- bzw. Zugriffsrechten, die für alle in der Lieferung enthaltenen Objekte gelten. Dient dem Archiv zur Entscheidung, ob Datensätze perspektivisch durch das Archiv öffentlich publiziert werden dürfen. Enthält ausschließlich einen der folgenden drei Begriffe:

- *institution*
Die Datensätze werden vom Archiv nicht veröffentlicht. Nur der Datengeber darf auf die archivierten Objekte zugreifen.
- *public*
Die archivierten Objekte dürfen potentiell (nach Rücksprache) durch das Archiv zugänglich gemacht werden. Es gelten die in den deskriptiven Metadaten der Transferpakete beschriebenen Lizenzen. Wenn dort nichts verzeichnet ist, gilt die im submission-manifest unter *License* aufgeführte Lizenz.
- *embargoUntil YYYY-MM-DD*
Mischform aus den beiden oben genannten. Bis zum angegebenen Datum darf nur der Datengeber auf die archivierten Objekte zugreifen, danach gelten im Fall einer Veröffentlichung durch das Langzeitarchiv die in den deskriptiven Metadaten der transferierten Pakete beschriebenen Lizenzen. Wenn dort nichts verzeichnet ist, gilt die im submission-manifest unter *License* aufgeführte Lizenz.

DataSourceSystem: Software, aus der Daten exportiert wurden (mit Versionsangabe). Dient dem Archiv zur Auswahl der weiteren Verarbeitung und Datennutzern als Kontext.

MetadataFile: Name und Pfad zur Metadaten-datei im Transferpaket. Wenn es mehrere Metadaten-dateien mit unterschiedlichen Dateinamen pro Lieferung gibt, werden entsprechend Sternchen (*) als Platzhalter verwendet. Sind diese weiterhin in unterschiedlichen Verzeichnissen enthalten, werden auch für diese Verzeichnisnamen Platzhalter verwendet. Die Verzeichnistiefe wird hiermit ebenfalls bestimmt und ist pro Lieferung eindeutig. Beispiele: meta.xml oder **/meta.xml* (für gleichlautende Metadaten-dateien in unterschiedlichen Ordnern) bzw. ***/*.xml* (für Dateien mit verschiedenen Dateinamen in unterschiedlichen Ordnern). Im zweiten Fall darf nur eine Datei mit dieser Endung im Verzeichnis enthalten sein. Dient zur Verarbeitung und Indizierung von Metadaten der Datenlieferung.

MetadataFileFormat: URI des Metadaten-Namespaces (zum Verständnis von Struktur/Format der Metadaten wie dcterms, datacite, lido, ead, mods). Beispiele: <http://www.lido-schema.org> oder <http://www.loc.gov/mods/v3>. Sind diese in einen METS-Container eingebettet, so ist hier <http://www.loc.gov/METS/> anzugeben.

CallbackParams (optional): Für die Übermittlung von Parametern während einer automatisierten Datenübertragung an EWIG. Diese können genutzt werden, um während des Ingestprozesses automatisch Rückmeldungen an das abliefernde System geben zu können. Die genaue Ausgestaltung erfolgt in Absprache mit EWIG.

1.1 Beispiele und Erläuterungen

```
SubmissionManifestVersion: 2.0
SubmittingOrganization: Küchenbibliothek Berlin
OrganizationIdentifier: DE-0815
ContractNumber: ZIB-XXXX
Contact: Bonnhofer, Ingo
ContactRole: Leiter Digitale Dienste
ContactEmail: bingo@example.com
TransferCurator: Inionski, Manfred
```

TransferCuratorEmail: minion@example.com
SubmissionName: ABT1_WERKZEUGE_MESSER_A-F

SubmissionName „hierarchie abbildend“: Der Transfer beinhaltet Daten aus der hauseigenen Hierarchie Abteilung1/Werkzeuge/Messer/A-F

SubmissionName: Projekt-FOOD-2019-S1001

SubmissionName „intellektuell verständlich“: Der Transfer beinhaltet Daten aus dem Projekt „Food“ des Jahres 2019 mit Signatur S1001.

SubmissionDescription: 12 Mecky Messer Kochbücher von 1927. Kooperationsprojekt mit Rezeptbuchsammlung Brecht.

RightsHolder: N/A

Rights: <http://id.loc.gov/vocabulary/preservation/copyrightStatus/pub>

RightsDescription: Die Digitalisate sind gemeinfrei.

License: <https://creativecommons.org/publicdomain/mark/1.0/>

AccessRights: public

Die Daten können ggf. nach Rücksprache mit dem Datengeber vom Archiv zugänglich gemacht werden. Sollten in den Daten selbst keine identifizierbaren, maschinenverarbeitbaren oder unvollständige Rechteinformationen enthalten sein, werden sie in dem o.g. Fall als gemeinfrei ausgezeichnet.

RightsHolder: Messer, Mecky (Autor)

Rights: <http://rightsstatements.org/vocab/InC/1.0/>

RightsDescription: Alle Rechte vorbehalten.

License: N/A

AccessRights: institution

Die Daten werden grundsätzlich nicht vom Archiv zugänglich gemacht.

RightsHolder: Messer, Mecky (Autor)

Rights: <http://rightsstatements.org/vocab/InC/1.0/>

RightsDescription: Das Digitalisat ist gemeinfrei nutzbar.

License: <https://creativecommons.org/publicdomain/zero/1.0/>

AccessRights: public

Die Daten können ggf. nach Rücksprache mit dem Datengeber vom Archiv zugänglich gemacht werden. Sollten in den Daten selbst keine identifizierbaren, maschinenverarbeitbaren oder unvollständige Rechteinformationen enthalten sein, werden sie in dem o.g. Fall als gemeinfrei nutzbar ausgezeichnet. Hierbei ist zu beachten, dass alle CreativeCommons-Lizenzen (außer die Kennzeichnung PublicDomainMark) die Nennung mindestens eines Rechteinhabers benötigen, als denjenigen, der das Nutzungsrecht eingeräumt hat (auch bei PublicDomainDedication mit CC0). Diese Differenzierung ist wichtig, weil viele Metadatenstandards Rechteinformationen nur eingeschränkt abbilden können.

DataSourceSystem: Kitodo Archive Plugin 1.0.0

MetadataFile: metadata.xml

MetadataFileFormat: <http://www.loc.gov/METS/>

2. Beschreibung des Transferpaketformats

Das Langzeitarchiv nimmt ausschließlich Daten an, die mit Integritätsinformationen versehen sind. Diese werden idealerweise im Zuge der Erstellung eines Transferpaketes automatisch erzeugt. Transferpakete sollen bevorzugt in Form eines mit **METS** beschriebenen Verzeichnisses übermittelt werden. Alternativ können Transferpakete in Form eines **bag**¹ oder als **komprimierter Container** (z.B. ZIP) übermittelt werden.

In Bezug auf die Verzeichnisstruktur hat der Datengeber verschiedene Möglichkeiten (s. Anhang). Das Langzeitarchiv unterstützt den Datengeber bei der Erstellung eines Transferpaketes. Grundsätzlich empfohlen wird die Übermittlung der strukturierenden Metadaten im METS-Format (Metadata Encoding and Transmission Standard).

Die Zuordnung zu Intellektuellen Einheiten erfolgt hierbei im METS und es kann in diesem Fall auf andere Strukturierungen, wie in den anderen Beispielen gezeigt, vollständig verzichtet werden.

¹ vgl. BagIt der Library of Congress, <https://tools.ietf.org/html/draft-kunze-bagit-17>

Ein Transferpaket beinhaltet in der obersten Verzeichnisebene das Submission Manifest mit den Liefermetadaten. Das Paket wird identifiziert durch den in dieser Datei enthaltenen *SubmissionName*. Dieser muss eindeutig sein in Bezug auf alle Lieferungen des Datengebers. Ein Transferpaket soll ausschließlich die zur Interpretation der Informationsobjekte notwendigen Daten und keine Referenzen zu externen Daten enthalten. Hierbei gelten Daten als extern, wenn sie nur für den Datengeber verfügbar sind. Referenzen auf persistent veröffentlichte Daten dürfen im Datenpaket enthalten sein.

Ein einzelnes Transferpaket besteht aus maximal 1,8 TB (Terabyte)² Daten. Der Datengeber soll in einem Transferpaket eine möglichst große Datenmenge zusammenfassen, um die Zahl der Transferpakete gering zu halten.

3. Beschreibung der Transferoptionen

Zum Transfer der Daten in das Langzeitarchiv stehen zwei Optionen zur Verfügung:

- Internet: Der Transfer des vom Datengeber erstellten Transferpaketes erfolgt mittels sftp (secure file transfer protocol) oder scp (secure copy) in einen Transferbereich des Langzeitarchivs oder er wird vom Langzeitarchiv per scp aus einem vom Datengeber verwalteten Speicherbereich heruntergeladen.
- Externe Festplatten: Die Transferpakete werden auf externen Festplatten geliefert und die Daten von MitarbeiterInnen des Langzeitarchivs manuell in den Transferbereich übertragen. Integritätsinformationen werden vor der Übertragung durch den Datengeber erzeugt.

4. Beschreibung des Transferprozesses

Im Transferbereich des Langzeitarchivs werden die übermittelten Daten anhand der beim Datengeber erzeugten Prüfsummen auf Übertragungsfehler untersucht. Im Fehlerfall werden alle Dateien verworfen und eine erneute Übertragung initiiert. Bei fehlerfreier Übertragung und Validität von Format und Verzeichnisstruktur der Transferpakete werden die Daten in das Langzeitarchiv übertragen und anschließend die Transferkopie gelöscht.

Transferpakete werden nur als Ganzes übernommen. Aus den Transferpaketen werden im Langzeitarchiv nach definierten Verfahren Submission Information Packages (SIP), Archival Information Packages (AIP) und (intern verwertete) Dissemination Information Packages (DIP) erzeugt. Nachdem die AIPs in den Speicherbereich des Langzeitarchivs übertragen wurden, ist der Transferprozess technisch abgeschlossen.

Das Langzeitarchiv erstellt abschließend ein Übernahmeprotokoll, das dem Datengeber bereitgestellt wird. Der Transferprozess ist nach Bereitstellung des Übernahmeprotokolls auch organisatorisch abgeschlossen. Die Verantwortung über die Bewahrung der archivierten Objekte liegt nun beim Langzeitarchiv.

² 1 TB (Terabyte) = 10¹² Byte

ANHANG: Mögliche Verzeichnisstrukturen für Transferpakete

Bei der Erstellung eines Transferpakets hat der Datengeber die Möglichkeit, Daten unterschiedlich strukturiert zusammenzustellen. Im Folgenden werden die Alternativen für die Verzeichnisstruktur beschrieben. Sie stellen einen Kompromiss dar zwischen größtmöglicher Flexibilität für den Datengeber und Minimierung des Aufwands für das Langzeitarchiv.

In einem Transfer kann eine beliebige Anzahl von Objekten mit Metadaten enthalten sein, sofern das zuvor genannte maximale Datenvolumen eines Transferpaketes nicht überschritten wird. Die Objekte können aus einer oder mehreren Dateien (digitalen Repräsentationen, z.B. tiff-Dateien) bestehen und bilden zusammen mit den zugehörigen Metadaten jeweils eine Intellektuelle Einheit (IE).

Eine IE besteht zum Beispiel aus den Digitalisaten eines Buches mit zugehörigen beschreibenden und strukturierenden Metadaten. Das Konzept der Intellektuellen Einheit dient dazu, dass zusammengehörige Inhalte auf unbestimmte Zeit verständlich beschrieben bleiben und identifizierbar, verwaltbar sowie nutzbar bleiben.

Jede IE für sich kann beim Transfer nahezu beliebig bezeichnet werden, jedoch darf diese Bezeichnung keine Leer- und Sonderzeichen enthalten. Innerhalb eines Transfers werden in der Regel mehrere IE enthalten sein, jede IE erhält dann einen eigenen Bezeichner, der keinem festen Schema folgen muss.

Nachfolgend sind einzelne *Dateien* nach diesem Prinzip bezeichnet: *****.*****
Die ersten drei Sterne symbolisieren dabei den Dateinamen, die letzten drei die Dateierdung. Die Dateinamen können dabei ebenfalls beliebig gestaltet sein (jedoch ohne Leer- und Sonderzeichen), die Dateierdung muss dem jeweiligen Dateiformat entsprechen, z.B. ab45234.tif, 20348.tif, m29384.xml. Wenn in einer Variante ein bestimmtes Dateiformat verlangt wird, ist die Dateierdung entsprechend dargestellt, z.B. *****.xml**. Schrägstriche (/) kennzeichnen jeweils Unterverzeichnisse. Datei- und Ordernamen sollten generell nicht zu lang sein.

Es ist möglich, im **optionalen Verzeichnis „submissionDocumentation“** *zusätzliche* inhaltliche Informationen zu den im jeweiligen Transferpaket enthaltenen Intellektuellen Einheiten abzulegen, die bei einer (zukünftigen) Nutzung helfen können, die Objekte zu interpretieren. Diese werden ausdrücklich *nicht* die beschreibenden Metadaten der IE ersetzen; sie ergänzen die IE lediglich um solche Informationen, die in die eigentlichen beschreibenden Metadaten nicht aufgenommen werden können. Diese Daten können auch unstrukturiert vorliegen. *Beispiele*: Bilder vom Experimentaufbau; Skizze vom Fotosetup; Emails mit Informationen, die auf die Herkunft der Daten/Objekte schließen lässt.

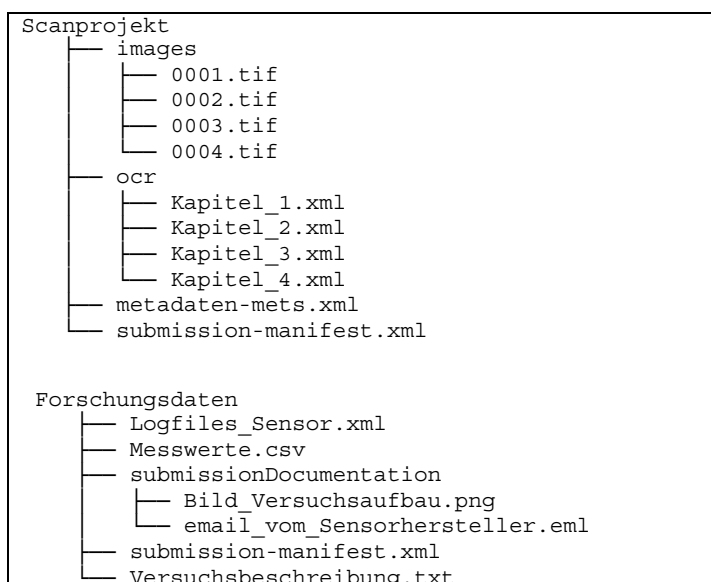
An den Inhalten des Verzeichnisses submissionDocumentation werden im Langzeitarchiv keine Erhaltungsmaßnahmen durchgeführt.

Der Datengeber kann sich bezüglich der Verzeichnisstruktur für eine der folgenden Varianten entscheiden, wobei das Langzeitarchiv ihn bei der Auswahl bei Bedarf beratend unterstützt:

1. Struktur im METS ordnet jede Datei einer Intellektuellen Einheit zu

Datencontainer (Bilder, Filme etc.) und Metadatencontainer werden entsprechend ihrer IE in einer vom Langzeitarchiv vorgegebenen METS-Struktur verzeichnet. Diese METS-Datei wird als submission-manifest.xml in der obersten Dateiebene mitgeliefert. Der Speicherort der jeweiligen Dateien ergibt sich relativ zu dieser METS-Struktur. Die Erzeugung von METS-Dateien erfolgt in enger Absprache zwischen den Partnern.

Beispiele:



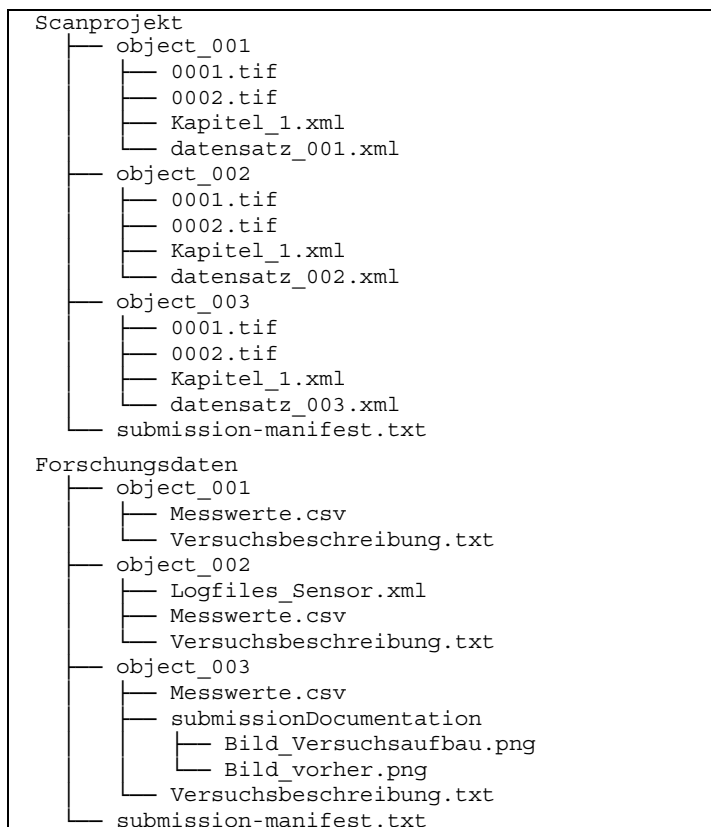
oder

2. Einzelne Unterverzeichnisse je IE

Die einzelnen Dateien einer IE werden in jeweils einem Unterverzeichnis zusammengefasst. Die Lieferung erfolgt als Container (ZIP) oder bag.

- /BezeichnerName/***.*** – Primärdaten (beliebige Anzahl an Dateien, mindestens eine)
- /BezeichnerName/***.xml -Datei mit Metadaten (insbesondere wenn auch die Primärdaten in XML abgelegt werden, muss diese Datei im Feld *MetadataFile* in der Datei *submission-manifest.txt* genau spezifiziert werden)
- /BezeichnerName/submissionDocumentation/ (optional)
- /submission-manifest.txt

Beispiele:

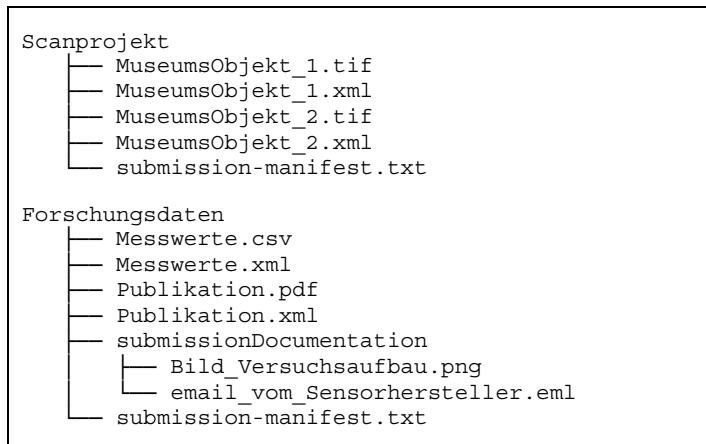


oder

3. Alle Dateien der verschiedenen IE in einem Verzeichnis

- Primärdateien (Bilder, Filme etc.) und Metadatendateien mit gleicher Benennung (Beispiel lv1102a.tif, dann lv1102a.xml usw.) in einem Hauptverzeichnis. Dies ist nur möglich bei Intellektuellen Einheiten, die aus einer Digitalisat- und einer Metadatendatei mit jeweils unterschiedlicher Dateierdung bestehen, also z.B. nicht für XML-Primärdaten, wenn die Metadaten ebenfalls in einem XML-Format vorliegen. Die Lieferung erfolgt als Container (ZIP) oder bag.
- /submissionDocumentation/ (optional) – als Unterverzeichnis des Hauptverzeichnisses, gilt nur für die gesamte Datenlieferung.
- /submission-manifest.txt – im Hauptverzeichnis

Beispiele:



oder

4. Einzelne ZIP-Container für einzelne IE

Die einzelnen Dateien einer IE werden in jeweils einer ZIP-Datei zusammengefasst. Die Lieferung erfolgt als Container (ZIP) oder bag.

- BezeichnerName.zip – Einzelne Zip-Container mit Daten und Metadaten als Einheiten
- /submission-manifest.txt

Beispiele:

