

# METADATEN FÜR OBJEKTE DES KULTURELLEN ERBES - QUALITÄTSANFORDERUNGEN. VORAUSSETZUNGEN FÜR DIE NACHNUTZBARKEIT UND DIE VERNETZUNG IN KULTURPORTALEN

Karolin Schmahl

Deutsche Fotothek, Sächsische Landesbibliothek – Staats- und Universitätsbibliothek  
Dresden, karolin.schmahl@slub-dresden.de

**KURZDARSTELLUNG:** Basierend auf den Erfahrungen der Datenaggregation und Datenverarbeitung für die Deutsche Digitale Bibliothek (DDB) und die Europeana sollen in Bezug auf die unterschiedliche Erfassungs- und Erschließungspraxis von Bildmedien Ansprüche an die Qualität von Metadaten erläutert werden. Mit dem Fokus auf der Qualitätssicherung werden außerdem Regelwerke, verbindliche Standards sowie Datenbanken und Tools thematisiert. Ausgehend von einer Begriffs- und Kriterienbestimmung zu Metadatenqualität und einer knappen Analyse des Status quo werden anhand spezifischer Use Cases im Portal der DDB (einfache Suche, Facettensuche, Relevanzbeurteilung der Ergebnisse) einige der wichtigsten Anforderungen – ohne Anspruch auf Vollständigkeit – zur optimalen Nachnutzbarkeit und Weitergabe von Metadaten zur Präsentation in Nachweisportalen vorgestellt und mögliche Lösungsansätze aufgezeigt. Ein Fokus liegt dabei auf der Optimierung von Erfassungsdaten im Sinne von Linked Data, welche die Kultur- und Wissenschaftseinrichtungen vor neue Herausforderungen hinsichtlich der Interoperabilität und Qualität ihrer Metadaten stellt. Denn nur über die Vernetzung der Daten mit anderen Objekt- und Normdatenbeständen können die Bestände unterschiedlicher Institutionen in Zukunft sowohl individuell zugänglich, als auch maschinell auswertbar und kontextuell erfahrbar sein.

## 1. EINFÜHRUNG

Die Deutsche Digitale Bibliothek (DDB) ist die Plattform für Kultur und Wissen in Deutschland, mit der die verteilten Bestände und Sammlungen des kulturellen Erbes virtuell zusammengeführt und über das Portal als gemeinsamen Zugangspunkt sichtbar gemacht werden. Als zentrales Zugangportal zu digitalen Objekten aus Kultur und Wissenschaft arbeitet die DDB spartenübergreifend (Archive, Bibliotheken, Museen, Mediatheken, Denkmalpflege, Wissenschaft) und interdisziplinär. Sie ist weiterhin der nationale Aggregator für die Europeana und versteht sich sowohl als kooperatives Netzwerk von Kultur- und Wissenschaftseinrichtungen (KWE) in Deutschland als auch als Plattform für Daten und Dienste. [1]

Kultur- und Wissenschaftsportale, die Informationsobjekte aus verschiedenen Sparten und in unterschiedlichen Formaten aggregieren sind in besonderer Weise auf die Verlässlichkeit der Daten angewiesen, wenn die Nutzersuche zu aussagekräftigen und verlässlichen Ergebnissen führen soll. Die Daten müssen zum einen so weit

vereinheitlicht sein, dass sie sich in die übergeordnete Struktur des Portals einfügen und zum anderen müssen sie hinreichende Spezifität und Unterscheidungskraft besitzen, um die Dokumente gezielt auffindbar zu machen.

Das Funktionieren der Deutschen Digitalen Bibliothek setzt also ein hohes Maß an Informations- und Datenqualität voraus. Das heißt nicht nur, dass ein Mindestmaß an beschreibenden inhaltlichen und administrativen Informationen mitgegeben werden muss, sondern auch, dass die Datenwerte und -inhalte korrekt sind. Beides ist nicht immer zuverlässig der Fall. Die DDB strebt eine Qualitätssicherung in allen Stufen des Aggregationsprozesses an, eine detaillierte Qualitätsprüfung der Metadaten auf Einzelebene ist jedoch nicht möglich. Aus diesem Grund ist eine Unterstützung seitens der datengebenden Institutionen unerlässlich in dem Sinne, dass diese möglichst tief erschlossene, standardisierte und korrekte Daten übermitteln.

Der Beitrag zeigt anhand von Beispielen die gravierendsten Probleme in Bezug auf eine effektive Suche auf, leitet aus diesen

Anforderungen an die Erfassung von Bildmedien aus Sicht eines spartenübergreifenden Kulturportals ab und spricht diesbezügliche Empfehlungen aus.

## **2. METADATENQUALITÄT – BEGRIFFSBESTIMMUNG**

Wie ist Metadatenqualität zu definieren und welche Kriterien können zur Bewertung herangezogen werden? Aus der Vielzahl der Definitionen sei hier beispielhaft die der Europeana Task Force on Metadata Quality erwähnt, welche Metadatenqualität wie folgt definiert: “Metadata quality is controlled by a set of processes which ensure that cultural heritage objects are described in such a way that they can be identified, discovered and seen in context by end-users, in a manner appropriate to the context in which the data provider created them. Metadata must include information on the potential re-use of cultural heritage objects.” Als Qualitätskriterien für die Metadaten wird bestimmt, dass diese das Ergebnis vertrauenswürdiger Prozesse, auffindbar, sinnvoll, lesbar, standardisiert, sichtbar und nachnutzbar sein sollen. [2]

Eine ähnliche Positionsbestimmung nehmen Thomas Bruce und Diane Hillmann vor, in dem sie fordern, dass Metadaten vollständig, genau, logisch konsistent und zusammenhängend, die Erwartungen erfüllend, zugänglich, aktuell und in ihrer Herkunft nachgewiesen sein sollten. [3]

Was dies in der praktischen Umsetzung konkret bedeutet, soll im Folgenden erläutert werden. Eine Auswahl von Standards und Leitlinien, die bei der Beschäftigung mit Metadatenqualität Orientierung und Hilfe bieten, finden Sie unter Punkt [4] in den Literaturhinweisen.

## **3. ANALYSE DES STATUS QUO**

Die Ergebnisse von Suchanfragen in der DDB sind noch nicht völlig zufriedenstellend, das Durchführen einer Suche kann zudem recht aufwendig sein. Als Hauptkritikpunkte seien hier genannt: Die Ergebnismengen für einfache Freitextsuchen sind noch zu wenig vollständig; relevante Dokumente können in der Ergebnisliste oft nur schwer, teils

auch gar nicht identifiziert werden; das Filtern von Suchergebnissen ist oft mit Informationsverlust verbunden; der Zugang zum Objekt ist nicht zuverlässig gegeben oder erfordert gegebenenfalls nicht-intuitive Umwege. Zu einem gewissen Teil liegt das auch an den Suchinstrumenten der DDB, die ständig weiterentwickelt werden um den Nutzerbedürfnissen besser zu entsprechen. Aber auch die Qualität der Daten enthält Optimierungspotenzial. Dabei können in allen Bearbeitungsschritten von der Datenerfassung über Mappings und Transformationen bis zur Dokumentpräsentation Fehler und Qualitätsmängel an den Daten entstehen. Daher muss die Datenqualität von Anfang an und in jedem Schritt im Aggregationsprozess überprüft und gesichert werden – angefangen bei der Erschließung der Bestände in den einzelnen Institutionen. Von dort kommen in nicht kleinem Maße Daten, die formal oder sachlich fehlerhaft sind, Metadatenelemente falsch interpretieren, in denen entscheidende Metadaten fehlen, nicht terminologisch kontrolliert sind oder proprietäre Vokabulare benutzen, unregelmäßig oder gar nicht indexiert sind. Ursächlich dafür ist ein individuelles Zusammenspiel folgender Aspekte: Welche zeitlichen und personellen Kapazitäten für Erfassung und Pflege stehen zur Verfügung, werden Regelwerke verwendet, welche individuell unterschiedlichen Ausbildungen haben die Erfasser, welche spezifischen Möglichkeiten und Grenzen bieten die verschiedenen lokalen Anwendungen, wie viel IT-Personal steht zur Verfügung, sind die Schnittstellen korrekt, wie wird Datenqualität in der jeweiligen Einrichtung definiert?

## **4. USE CASE 1: EINFACHE SUCHE**

Beim Durchführen einer einfachen Suche erwartet der Nutzer, alle relevanten – und möglichst nur die relevanten – Ergebnisse für seine Suchanfrage zu finden. Fehlende Indexterme können jedoch zu Informationsverlust oder -ballast führen. Ein Beispiel: Die Suche nach „Schloss“

ergibt 67.740 Ergebnisse, eine Suche nach „Schloß“ führt zu 47.750 Ergebnissen, die Vereinigungsmenge offeriert 112.539 Treffer, „schlösser OR schloss OR schloß“ führt zu 136.075 Ergebnissen. (s. Abb. 1) Auch die fehlende Rückführung auf Grundformen mindert den Recall, zum Beispiel beim Numerus: Die Suche mit der Singularform "Musikinstrument" ergibt 7.587 Treffer, die Pluralform "Musikinstrumente" hat 3.955 Treffer, die Vereinigung liefert 11.454 Ergebnisse. In beiden Beispielen werden die Bezeichnungen nicht normalisiert und führen somit zu Informationsverlust, dieser entsteht ebenfalls durch fehlende Synonymkontrollen (z. B. „Burg“ oder „Kastell“). Durch die fehlende Disambiguierung (z. B. zu „Schloss“ als Vorrichtung zum Verschließen) entsteht wiederum Informationsballast. Weiterhin problematisch sind Schreibfehler sowie historische oder alternative Schreibformen.

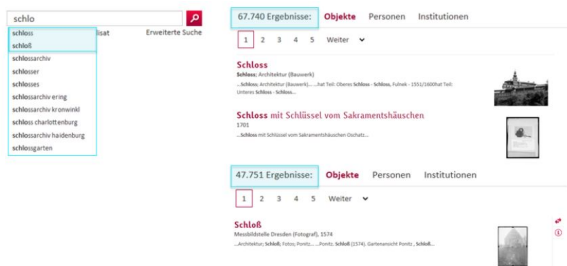


Abb. 1: Screenshotcollage DDB-Portal zur einfachen Suchanfrage nach „Schloss“ und „Schloß“

Eine solche uneinheitliche Erschließung mindert nicht nur die Qualität der Suchergebnisse, sondern erschwert auch automatische Korrekturen. Die Anwendung von Indexierungsregeln und die Kontrolle von Indextermen, gegebenenfalls in Verbindung mit hauseigenen Regeln, ist daher dringend zu empfehlen. Ferner sollte bei der Verfassung ein für alle Mitarbeiter verbindliches Regelwerk verwendet werden, welches detailliert definiert, welche Informationen in welcher Form in welchen Feldern zu erfassen sind. [5] Empfehlenswert ist weiterhin die Nutzung von Datenbanksystemen, die eine systemseitige Kontrolle der Erschließung durch vorgegebene Wertelisten, durch

Vorschlagslisten aus Live-Indices und durch eine automatische, regelbasierte Syntaxprüfung erlauben. Hierbei sollte darauf Wert gelegt werden, dass die Indexfelder individualisiert definierbar, ihre Terme also in verschiedenen Stufen kontrollier- und steuerbar sind.

Zur Sicherung eines einheitlichen Vokabulars zur sinnvollen Filterung der Daten prüft die Datenbank darauf basierend zum einen, ob die einzelnen Terme in den Feldern erlaubt sind und zum anderen, ob sie richtig geschrieben sind, in beiden Fällen erscheint anderenfalls eine Fehlermeldung. Über eine automatische, regelbasierte Syntaxprüfung kann z. B. die standardgerechte Erfassung von numerischen Informationen (Archivnummern, Zeitangaben, u.s.w.) gesichert werden.

## 5. USE CASE 2: FACETTENSUCHE

Normdaten und kontrollierte Vokabulare spielen für gute Retrievalergebnisse im Sinne von Vollständigkeit und Genauigkeit eine entscheidende Rolle. Für eine facettierte und semantische Suchfunktion ist die Nutzung kontrollierter Vokabulare und standardisierter Werte sogar eine notwendige Bedingung, denn nur so kann der Nutzererwartung, dass sich Suchergebnisse über die Facettensuche sinnvoll einschränken lassen, entsprochen werden. Uneinheitliche Ansetzungen und Schreibweisen (s. Abb. 2) führen nicht nur zu Unübersichtlichkeit, sondern auch zu Unklarheit über die Identität von z. B. Personen.

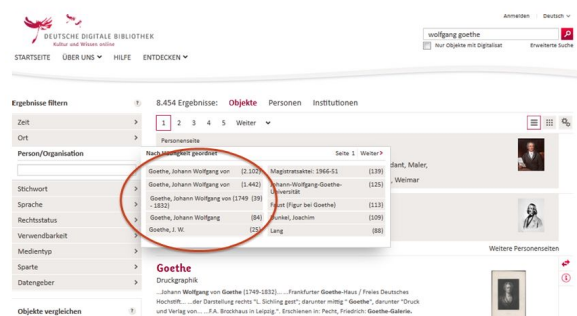


Abb. 2: Screenshot DDB-Portal zur einfachen Suchanfrage nach „wolfgang goethe“

Die Verwendung von Normvokabularen – hier z. B. die GND – würde diese Probleme vermeiden, außerdem gewährleistet die Verwendung eines Normdatensatzes mit URI die Einbeziehung von Synonymen und abweichenden Schreibarten bereits in der Suche: So würden sämtliche Schreibarten des Komponisten Tschaikowski in der Suche gefunden und berücksichtigt, selbiges gilt für Pseudonyme.

Auch die Stichwortfacette der DDB ist im Moment nur bedingt nützlich für die Suche, da sie kein gemeinsames Vokabular für sinnvolle Filterung abbildet, also keine eindeutigen und disjunkten Werte enthält, stattdessen in hoher Zahl Terme aus proprietären oder gar nicht kontrollierten Vokabularen – noch dazu auf unterschiedlichen generischen Niveaus. (s. Abb.3)

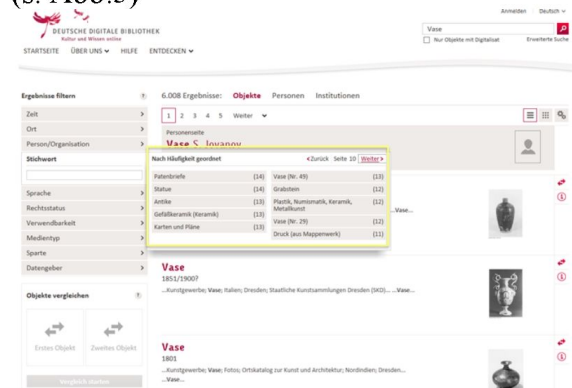


Abb. 3: Screenshot DDB-Portal, Werte der Stichwortfacette zum Suchbegriff „Vase“

Es finden sich in der Stichwortfacette bunt nebeneinander Stichworte aus Objektbezeichnungen, thematische Schlagworte, Klassifikations- oder Gattungsangaben. Es besteht somit dringender Handlungsbedarf, die Stichwortfacette nutzungsfreundlicher zu befüllen. Dafür ist es jedoch nötig, dass die Metadatenelemente korrekt befüllt werden: Häufig werden in den Daten Objektbezeichnung und Themenschlagwörter, Titel und Objektbezeichnung oder Objektbezeichnung und Klassifikation verwechselt oder vermischt.

Abb. 4 hingegen zeigt am positiven Gegenbeispiel der Facette Rechtsstatus, wie komfortabel und übersichtlich sich die

Einschränkung über Filter mit normierten Vokabularen gestalten lässt.



Abb. 4: Screenshot DDB-Portal, Werte der Facette Rechtsstatus

Wichtig ist in jedem Fall, dass für die einzelnen Ebenen und Arten der Informationen separate Metadatenfelder benutzt werden. Obwohl diese Forderung logisch und banal erscheinen mag, zeigt die Praxis, dass diese entscheidende Grundregel beim Erfassen nicht immer beachtet wird. So zeigt Abb. 5 eine gemischte Präsentation textlicher Objektbeschreibungselemente mit Notationen und entsprechenden literalen Werten des Vokabulars Iconclass. [6]

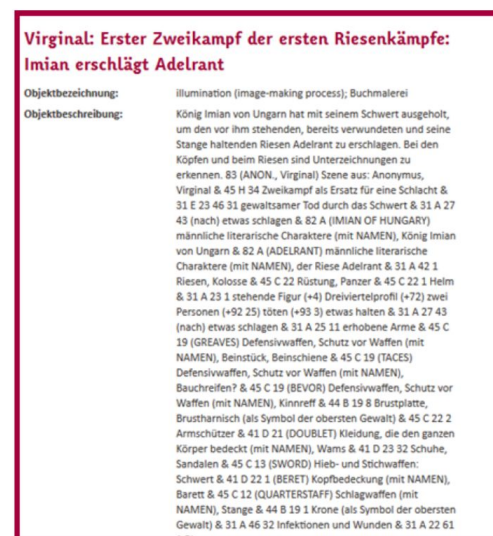


Abb. 5: Screenshot DDB-Portal, Ausschnitt aus einer Detailsansicht, Fokus „Objektbeschreibung“

Gravierender ist die in Abb. 6 deutlich sichtbare undifferenzierte Auflistung und Vermischung von Informationen zu Maßstab, Datierung, Personen, Beschreibungen in einem Feld („inhalt information“). In dieser Form können die Einzelinformationen nicht mehr separiert werden und daher nur als unübersichtliches

und vor allem nicht facettiert durchsuchbares Konglomerat – da die Zuweisung zu konkreten Einzelfeldern nicht möglich ist – für Präsentationszwecke weitergegeben und verwertet werden.

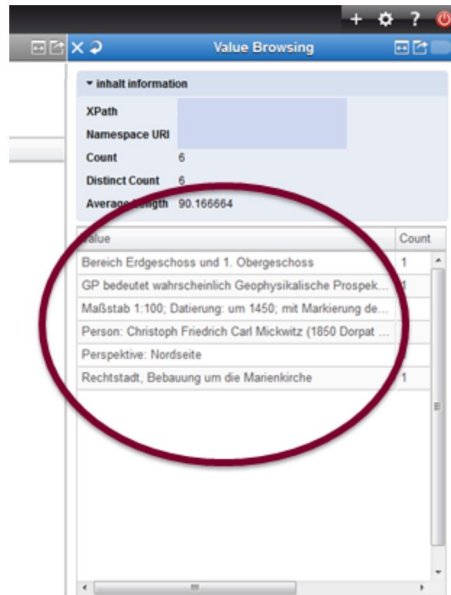


Abb. 6: Screenshot MINT (Mappingtool), statistische Auswertung des Metadatenfelds „inhalt information“

Als Status quo lässt sich konstatieren: Individualnamen für Personen, Körperschaften und Orte sind überwiegend nicht kontrolliert und die Entitäten folglich nicht eindeutig identifiziert. Eine zusätzliche Hürde sind Synonyme, Schreibfehler sowie historische oder syntaktische Schreibweisen, die dem Nutzer oft nicht bekannt sind. Noch schwieriger gestaltet sich ohne terminologische Kontrolle bzw. Nutzung kontrollierter Vokabulare die Harmonisierung ungebundener, freier und undefinierter Sachschlagwörter. Es sei daher ausdrücklich empfohlen, überall wo dies möglich ist, kontrollierte Vokabulare für die Erfassung zu verwenden, denn diese erhöhen sowohl die Vollständigkeit der Suchergebnisse (Recall) als auch die Präzision der Suchergebnisse (Precision). Sie sind Voraussetzung für multilinguale Schnittstellen, ermöglichen eine automatische Suchausweitung, sind Voraussetzung für eine facettenbasierte

Suche und für semantische Suchen und bilden die Grundlage für eine Top-down-Navigation (Klassifikationen).

Kontrollierte Vokabulare sollten möglichst multilingual und reich an Synonymen sein, keine ambigen Bezeichnungen (Homonyme ohne klärenden Zusatz) enthalten, gut hierarchisch strukturiert, verbreitet, verlässlich gepflegt und langfristig nutzbar sein, sowie möglichst als nachnutzbare Linked-Open-Data-Vokabulare vorliegen und einen breiten Abdeckungsgrad haben. Keiner der in der Praxis der Bildarchive verwendeten Thesauri oder Klassifikationen erfüllt alle diese Anforderungen, aber dennoch sind sie aus unterschiedlichen Gründen in verschiedenen Kontexten sehr nützlich und sollten in die Erschließungsarbeit integriert werden. Als Beispiele seien die Thesauri: Gemeinsame Normdatei (GND) der Deutschen Nationalbibliothek; Art & Architecture Thesaurus (AAT); Getty Thesaurus of Geographic Names (TGN) sowie die Klassifikationen Dewey Dezimalklassifikation (DDC) und Iconclass genannt.

Nicht zuletzt sind Normdaten mit persistenten URIs als eindeutige Identifier die wichtigste Voraussetzung für Linked Open Data. Kontrollierte Vokabulare, die ihre Inhalte als LOD zur Verfügung stellen, können schon jetzt zur Anreicherung von Metadaten sowohl mit entsprechenden Normdaten-URIs als auch für semantische Anreicherungen mittels API-Call genutzt werden.

### 6. USE CASE 3: ERGEBNISLISTE

Zur Bewertung einer Ergebnisliste erwartet der Nutzer aussagekräftige Informationen, die ihm die Relevanzbeurteilung der Suchtreffer erlauben. Die Informationen der Trefferliste reichen jedoch oft nicht aus, um einschätzen zu können, ob ein Ergebnis dem Suchwunsch entspricht. Ein Beispiel: Die DDB-Suche nach "Schloss" erbrachte 67.740 Ergebnisse, von denen bereits auf den ersten Ergebnisseiten 143 Dokumente in Folge identische Titel und

Beschreibungen haben. Die Größe und Qualität der Vorschaubilder reicht für eine Identifizierung des jeweiligen Objektes nicht aus.

Um eine bessere Unterscheidbarkeit der Ergebnisse zu ermöglichen, erarbeitet die DDB momentan ein Konzept zur sinnvollen Einbeziehung von Orts- oder Zeitangaben in die Ergebnisanzeigen der Trefferliste, um die bereits jetzt existierenden, spartenspezifisch unterschiedlichen Regeln zur Bildung der „Untertitel“ der Ergebnisse in der Trefferliste zu optimieren. An die datengebenden Institutionen ist jedoch zu appellieren, dass bei der Erschließung kultureller Objekte möglichst aussagekräftige Titel gebildet werden, die den Gegenstand kurz und prägnant beschreiben, selbsterklärend sind – auch wenn der Titel alleine steht. Des Weiteren sollten sie möglichst keine Abkürzungen – es sei denn, sie sind im Original enthalten oder allgemein etabliert – und keine nicht informativen Wörter, wie z. B. Wiederholungen des Objekttyps enthalten. Ein weiterer Aspekt in Bezug auf die Beurteilung der Suchergebnisse ist die Tiefe der erschließenden Angaben, also der Aussagegehalt der mitgegebenen Informationen. Beispiele wie in Abb. 7, wo nur rudimentärste Informationen mitgegeben werden, die lediglich einem Platzhalter gleichkommen, beeinträchtigen die Glaubwürdigkeit in die Qualität von Kulturportalen ebenso wie vermehrte Schreibfehler, falsche oder nicht mehr aktuelle Informationen.



Abb. 7: Screenshot DDB-Portal Detailansicht zu einem Objekt

Um Fotografien oder andere Bildmedien möglichst umfassend und nutzbringend zu beschreiben, sollten die Daten sowohl in Bezug auf die technischen/administrativen

Angaben als auch bei der inhaltlichen Beschreibung angemessen erschlossen werden. Zu den administrativen Angaben gehören Informationen zu ID, Signatur, Eigentümer, Erhaltungszustand, Provenienz, u.s.w. die technischen Informationen umfassen Aussagen zu Technik, Material, Maße/Größe, Medientyp, Ausrichtung, Farbe, etc. Die Erfassung des Bildinhalts sollte neben dem Titel, einer Bildbeschreibung und der Datierung auch Beziehungen zu Personen (z.B. Dargestellte, Urheber eines abgebildeten Kunstwerks, Auftraggeber), Beziehungen zu Körperschaften (z.B. Dargestellte, Verwalter eines Kunstwerks), Beziehungen zu Orten/ Bauwerken, Schlagwörter sowie möglichst normierte Angaben zu Objekttyp und Klassifikation enthalten. Wenn auch nicht für jedes kulturelle Objekt alle genannten Informationen ermittelt und bereitgestellt werden können, so soll diese Liste zumindest eine Richtschnur bieten, welche Bandbreite an erschließenden Metadaten das Ziel der wissenschaftlichen Primäerschließung von Sammlungsbeständen sein sollte.

## 7. WEITERE ANFORDERUNGEN AN DIE METADATEN

Um den Rechtsstatus auf Objektebene eindeutig nachweisen zu können sollten Lizenzen und Rechteauszeichnungen [7] sowohl für die originalen Objekte als auch für die digitalen Derivate und die Metadaten möglichst normiert dokumentiert werden. Abb. 8 zeigt ein Negativbeispiel, in dem Rechteangaben und Verwendungshinweise nicht separiert in einem Feld erfasst werden.

Field	Value
CI_Stadt	14 1 4
CI_Telefonnummer	14 1 22
CI_URL	14 1 18
Copyright_Vermerk	14 8 35,3
Datensamenerweiter...	14 1 3
Datentyp	14 1 4
Datensatz_ID	14 14 34,5
Datum_der_Aufnahm...	13 9 8
Datum_und_Uhrzeit...	14 11 14
Value	Köröchel, Franz/Josef / CCBY-NC-SA 3.0 Dokumentationsst...te Regierungsbunker (Pressefoto) Axel/H / gemeinfrei Bindara, Vanessa Ellgaard, Hölger / CCBY-SA 3.0 Ketschenbach, Sascha / Dokumentationsst...te Regierungsbunker Raymond / CCBY-SA 3.0 Vorderstra...Dirk / CCBY 2.0

Abb. 8: Screenshot MINT, Ausschnitt der statistischen Auswertung des Metadatenfelds "Copyright\_Vermerk"

Aus der Fülle technischer Anforderungen, die an dieser Stelle nicht ausführlich behandelt werden können, seien abschließend einige genannt, auf die hinzuweisen aus Erfahrung lohnt: Inkludierte Weblinks sind häufig nicht aktuell oder fehlerhaft, mittels eines Linkvalidators lässt sich dieses Problem beheben. Für den Datenaustausch sollten die Daten bevorzugt in XML-Dateien zur Verfügung gestellt werden. Hierbei ist darauf zu achten, dass die Dateien sowohl valide als auch wohlgeformt sind. Kostenfreie Tools wie XML Marker (<http://symbolclick.com/>) oder Validome (<http://www.validome.org/>) helfen dabei, diese Grundvoraussetzungen zu erfüllen. Für den Nachweis und/oder die Anzeige der Objekte in der DDB müssen diese auf einer eigenen lokalen Webseite publiziert und über einen stabilen Link erreichbar sein. Besonders bei hierarchischen Beziehungen von Objekten ist ein direkter Zugang (ohne zusätzliches Suchen) zum Quelldokument gewünscht.

## 8. DANKSAGUNG

Ich danke meine Kolleginnen Francesca Schule, Stefanie Rühle, Jutta Lindenthal und Angela Kailus, die sich vielseitig mit dem Thema Qualität von Metadaten auseinandersetzen und deren Erkenntnisse in diesen Beitrag mit eingeflossen sind.

## 9. LITERATURHINWEIS

[1] Zu Leitlinien und künftigen Arbeitsschwerpunkten der DDB s. Deutsche Digitale Bibliothek. Kultur und Wissen online (2016): Strategie 2020, URL: [https://www.deutsche-digitale-bibliothek.de/static/files/asset/document/ddb\\_strategie\\_2020\\_download.pdf](https://www.deutsche-digitale-bibliothek.de/static/files/asset/document/ddb_strategie_2020_download.pdf) (Stand 19.10.2016)

[2] Dangerfield, Marie-Claire; Kalshoven, Lisette (2016): Report and Recommendations from the Task Force on Metadata Quality, URL: [http://pro.europeana.eu/files/Europeana\\_Professional/Publications/Metadata%20Quality%20Report.pdf](http://pro.europeana.eu/files/Europeana_Professional/Publications/Metadata%20Quality%20Report.pdf) (Stand: 19.10.2016)

[3] Bruce, Thomas R.; Hillmann, Diane I. (2004): The Continuum of Metadata Quality: Defining, Expressing, Exploiting, URL: <http://hdl.handle.net/1813/7895> (Stand: 19.10.2016)

[4] Standards und Leitlinien zu Metadatenqualität

-Deutsche Forschungsgemeinschaft (2013): DFG-Praxisregeln "Digitalisierung". DFG-Vordruck 12.151 -02/13, URL: [http://www.dfg.de/formulare/12\\_151/12\\_151\\_de.pdf](http://www.dfg.de/formulare/12_151/12_151_de.pdf), S. 25-29, (Stand: 19.10.2016)

-Deutschen Gesellschaft für Informations- und Datenqualität (o. J.): Die 15 Dimensionen der Datenqualität, URL: [http://www.az-direct.ch/fileadmin/pdf/15\\_Dimensionen\\_Datenqualitaet\\_DGIQ.pdf](http://www.az-direct.ch/fileadmin/pdf/15_Dimensionen_Datenqualitaet_DGIQ.pdf) (Stand: 19.10.2016)

-NISO Framework Working Group: A framework of guidance for building good digital collections, A NISO Recommended Practice, 3rd edition, National Information Standards Organization (NISO), Baltimore, 2007

-Calhoun, Karen; Cantrell, Joanne; Gallagher, Peggy; Hawk, Janet: Online catalogs: what users and librarians want. An OCLC Report, OCLC Online Computer Library Center, Inc., Dublin, Ohio, 2009

-Harpring, Patricia (2015): CONA: Subject Access for Art Works, URL: [https://www.getty.edu/research/tools/vocabularies/cona\\_and\\_subject\\_access.pdf](https://www.getty.edu/research/tools/vocabularies/cona_and_subject_access.pdf) (Stand: 19.10.2016)

-Harpring, Patricia; Baca, Murtha (2014): Categories for the Description of Works of Art, URL: [http://www.getty.edu/research/publications/electronic\\_publications/cdwa/index.html](http://www.getty.edu/research/publications/electronic_publications/cdwa/index.html) (Stand 19.10.2016)

[5] Für den Bereich der Bildmedien sei hier beispielhaft das Marburger Informations-, Dokumentations- und Administrations-System erwähnt, ein Regelwerk, das v.a. in kunsthistorischen Bildarchiven und Museen Anwendung findet, online verfügbar unter: [http://archiv.ub.uni-heidelberg.de/artdok/3770/1/Bove\\_Heusinger\\_Kailus\\_MIDAS\\_Handbuch\\_2001.pdf](http://archiv.ub.uni-heidelberg.de/artdok/3770/1/Bove_Heusinger_Kailus_MIDAS_Handbuch_2001.pdf)

[6] Iconclass ist ein ikonographisches Klassifizierungskonzept zur Erfassung und inhaltlichen Erschließung von Bildinhalten im Bereich der Kunstgeschichte, online abrufbar unter: <http://www.iconclass.org/>

[7] Eine Übersicht über alle Lizenzen und Rechteausszeichnungen, die momentan in der DDB verwendet werden, finden Sie unter: <https://www.deutsche-digitale-bibliothek.de/content/lizenzen-und-lizenzhinweise-rechtssicherheit-der-deutschen-digitalen-bibliothek>