

# Navigating Babel: A Mid-Level Ontology for Cross-Domain Cultural Heritage Discovery

Mary Ann Tan<sup>1,2,\*</sup>, Genet Asefa Gesese<sup>1,2</sup> and Harald Sack<sup>1,2</sup>

<sup>1</sup>FIZ Karlsruhe, Leibniz Institute for Information Infrastructure, Eggenstein-Leopoldshafen, Germany

<sup>2</sup>Karlsruhe Institute of Technology, Karlsruhe, Germany

## Abstract

The Deutsche Digitale Bibliothek (DDB) aggregates objects from libraries, archives, museums, audiovisual archives, monument conservation, and research institutions into a structurally uniform container, the Europeana Data Model (EDM), yet cross-domain discovery is structurally unsupported. Like Borges' *Library of Babel*, the DDB has order without orientation. We identify three structural parallels and propose MOCHO (Mid-level Ontology for Cross-domain Cultural Heritage Objects), a reusable alignment layer between EDM and GLAM-spanning domain ontologies. MOCHO defines WEMI<sup>1</sup>-typed entities and four alignment patterns. An empirical snapshot of 115k records quantifies the heterogeneity MOCHO addresses.

## Keywords

Ontology, Digital Library, Digital Humanities, Semantic Web, Ontology Alignment, Knowledge Graphs

## 1. Introduction

The Library of Babel is a short story written by the Argentinian author and librarian, Jorge Luis Borges. Borges was the director of the National Library of Argentina between 1955 and 1973. He imagines the library of Babel, hereafter called “the Library”, to be as vast as a universe in which every possible book exists. Yet no one can navigate it: there is order, but without orientation. Every hexagonal room looks like every other, every volume follows the same physical format, and the catalog, if it exists, is itself lost somewhere in the stacks [1]. This tension between completeness and the absence of navigational cues also appears in real-world cultural heritage aggregation.

The Deutsche Digitale Bibliothek (DDB) aggregates approximately 65 million objects from German libraries, archives, museums, audiovisual archives, monument conservation (*Denkmal*), and research institutions into a single portal. DDB objects are represented using the Europeana Data Model (EDM), a generic domain-agnostic container class intended to represent any cultural heritage object (CHO).

An extension of EDM, the DDB-EDM [2], succeeds at what it sets out to do, namely, structural interoperability: cross-domain objects share a common index and interface. However, structural uniformity without ontological differentiation prevents cross-domain discovery. The differentiation that is missing is the WEMI hierarchy of the IFLA FRBR model [3]: a *Work* (abstract intellectual creation), *Expression* (a specific realization), *Manifestation* (a physical or digital embodiment), and *Item* (an individual copy). In EDM, a printed book edition, a film still, and the opera recording inspired by the same literary work are all typed identically as `edm:ProvidedCHO`. There is no node to express that they are all realizations of the same *Work*.

We use the Library as a *partial* structuring metaphor. Borges intended the Library to be self-contradictory, such that its very completeness renders it effectively useless [4]; therefore, we do not extend the metaphor to the complete elements of the short story. Instead, we focus on three structural parallels between the Library and the DDB that highlight why structural uniformity alone is insufficient and what kind of semantic layer would restore orientation.

<sup>1</sup>Work, Expression, Manifestation, Item. See Section 2 for definitions and illustrations.

*Third International Workshop of Semantic Digital Humanities (SemDH 2026), ESWC 2026, Dubrovnik, Croatia*

\*Corresponding author.

✉ ann.tan@fiz-karlsruhe.de (M. A. Tan); genet-asefa.gesese (G. A. Gesese); harald.sack@fiz-karlsruhe.de (H. Sack)

ORCID 0000-0003-3634-3550 (M. A. Tan); 0000-0003-3807-7145 (G. A. Gesese); 0000-0001-7069-9804 (H. Sack)



© 2026 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

This paper makes three contributions:

1. A conceptual framework articulating challenges in cultural heritage aggregation through three structural parallels between Borges' Library and the DDB: the hexagonal rooms and EDM (Section 4), the librarians and facet-based search (Section 5), and the decipherers and domain ontologies (Section 6).
2. MOCHO<sup>2</sup>, a Mid-level Ontology for Cross-domain Heritage Objects, a reusable alignment layer between EDM and GLAM<sup>3</sup>-spanning domain ontologies, defined via four alignment patterns (Section 6.3). MOCHO is designed to be applicable beyond the DDB: any cross-domain aggregator that maps heterogeneous metadata to EDM faces the same structural problem.
3. An empirical snapshot of DDB records drawn from a cross-domain Goethe-*Faust* corpus<sup>4</sup> (Appendix A), henceforth called as "the corpus", quantifying the heterogeneity MOCHO addresses (Section 4).

**Paper structure.** Section 2 reviews related work. This is followed by the introduction of the partial metaphor in Section 3. To reinforce this metaphor, three structural parallels are developed in turn in Sections 4–6, with Section 6 presenting MOCHO in full. Section 7 discusses what MOCHO achieves and what it cannot fix. Finally, Section 8 concludes with mentions of plans for future work.

## 2. Related Work

The EDM is designed to ingest heterogeneous metadata while remaining institutionally neutral. With the re-use of CIDOC-CRM [5], it allows event-centric modeling of CHOs. This modeling framework links contextual descriptions (*who?*, *where?*, *when?*) through intermediary `edm:Event` instances, as opposed to object-centric modeling, where `dc:`<sup>5</sup> and `dcterms:`<sup>6</sup> properties are directly defined with `edm:ProvidedCHO` instances [6].

Rühle et al. [2] describe the DDB-EDM model and document how the transition from CIDOC-CRM to EDM simplified cross-institutional mapping at the cost of semantic precision. Peroni et al. [7] identify three structural limitations of EDM that inhibits cross-domain search: media type ambiguity between CHO and digital representation, multi-layer subject collapse, and role conflation in `dc:creator` and `dc:contributor` that obscures specific agent function, that of the people and the organizations responsible for its existence. Europeana's EDM profiles<sup>7</sup> address domain-specific needs by extending EDM within its existing flat structure; MOCHO differs by introducing WEMI-typed entities as a separate alignment layer above EDM, enabling cross-domain traversal that profiles cannot express. Zapounidou et al. [8] propose a BibFrame-to-EDM crosswalk for the library domain; MOCHO differs in scope (all GLAM domains, not only libraries) and in mechanism (See Section 6.3.1).

In FRBR, bibliographic identity is represented by four conceptual levels: a *Work* is the abstract intellectual creation: Goethe's *Faust* exists as a Work regardless of how many times it has been realized. An *Expression* is a specific realization of that Work in the same medium: a French prose translation and a scholarly critical edition of the German verse are two distinct Expressions of the same Work. Cross-media adaptations such as Gounod's operatic setting and Murnau's silent film are not Expressions of Goethe's *Faust* but derivative Works in their own right, linked to the source Work via `mocho:adaptationOf`, following the treatment in IFLA LRM and [9]. A *Manifestation* is a physical or digital embodiment of an Expression: the 1826 Leipzig printed edition and a digitized PDF of that text are two Manifestations of the German verse Expression. An *Item* is an individual exemplar of a Manifestation: a specific copy of

<sup>2</sup><https://ise-fizkarlsruhe.github.io/ddbkg/mocho/>

<sup>3</sup>Galleries, Libraries, Archives, Museums

<sup>4</sup><https://github.com/ISE-FIZKarlsruhe/ddbkg/tree/main/goethe-faust>

<sup>5</sup><http://purl.org/dc/elements/1.1/>

<sup>6</sup><http://purl.org/dc/terms/type>

<sup>7</sup><https://pro.europeana.eu/page/edm-profiles>

the 1826 edition held at the Staatsbibliothek zu Berlin, or a film still from Murnau's print held at the Bundesarchiv.

FaBiO [10] extended it as an OWL ontology for scholarly publishing domain. Tan et al. [11] mapped DDB-EDM object type terms to FaBiO classes for the library sector, and found no reliable 1-to-1 correspondence, since single `dc:type` or `dcterms:type` value can represent document type, document structure, production process, or even subject heading, requiring manual curation that does not generalize across providers or sectors. MOCHO addresses this by typing at the WEMI level rather than to specific subclasses, and by generalizing from the bibliographic domain to the full GLAM spectrum.

The domain ontologies that supply WEMI semantics differ by institutional sector. Dijkshoorn et al. [12] demonstrate CIDOC CRM's practical value for the Rijksmuseum collection. RiC-O [13] addresses the archival sector through fonds hierarchies and provenance chains; the Music Ontology [14] and Audio Commons Ontology [15] cover audio resources. Where prior work has typically aligned to a single one of these families, MOCHO aligns to all of them simultaneously through `owl:unionOf` constraints, enabling cross-domain queries without re-ingestion or schema migration.

Aligning EDM records to domain ontologies is non-trivial. Recent matchers apply transformer embeddings [16] or LLM prompting [17] to match classes and properties across ontologies, achieving competitive results in zero-shot settings, but both rely on lexical similarity in label strings and produce decisions not directly recoverable from the ontology structure alone. Where prior work has typically aligned to a single domain the ontology, MOCHO aligns to all of them simultaneously through several alignment patterns presented in Section 6.

Retrieval-augmented generation [18] and Knowledge Graph(KG)-grounded variants [19] further complement structured KG access by grounding language model capabilities in explicit knowledge; for cultural heritage aggregation, however, LLMs reasoning over raw EDM literals must infer domain- and modality<sup>8</sup>-specific semantics that MOCHO makes explicit in the data.

### 3. The Library of Babel as a Partial Metaphor

Borges' Library contains every possible 410-page book: every combination of 25 orthographic symbols across 1,312,000 characters per volume [1]. It is therefore *total and complete*, holding not only every book that has been written but every book that could be written, including every erroneous version, translation, and commentary on every other book. But completeness is not the same as usefulness: in a collection where every possible statement coexists with every possible contradiction, finding meaningful text is a herculean task.

The librarians have not accepted this situation passively. Some sought *the catalog of catalogs*, others purged volumes they deemed meaningless, and others spent an eternity scouring hexagon after hexagon for canonical volumes. All of these strategies failed. The problem is structural: the Library's physical organization, identical hexagonal rooms and uniform volumes, provides no orientation toward meaning. As Borges' short story illustrates, order does not lead to orientation.

We identify three structural features of the Library that correspond to three structural features of the DDB, which together produce the same discovery failure.

**Parallel 1 (§4)** The hexagonal rooms impose a uniform physical structure on incomprehensible diversity. In the DDB, `edm:ProvidedCHO` plays an analogous role: a structurally uniform category that obscures semantic distinctions of every object regardless of type or domain.

**Parallel 2 (§5)** The librarians' search strategies of seeking the catalog of catalogs, scanning hexagon after hexagon, correspond to keyword search and lowest-common-denominator facets. Knowledge graph navigation is the escape: a shift from scanning uniform containers to traversing typed entities and relations.

---

<sup>8</sup>known as *Media Type* in EDM

**Parallel 3 (§6)** The decipherers who can distinguish nonsense from meaningful text correspond to domain ontologies. Alignment patterns make that principle operational, by restoring semantic orientation without dismantling the existing structure.

To ground these structural observations empirically, we assembled a cross-domain corpus of DDB objects by querying the portal for “Goethe” and “Faust” and retrieving the full EDM records for each result. A detailed description of the corpus is provided in Appendix A.

## 4. Parallel 1: The Hexagonal Rooms and EDM

In Borges’ Library, the hexagonal rooms impose a uniform physical structure on an incomprehensible diversity of content. Every book has the same number of pages, the same character set, the same dimensions, regardless of what it contains. The DDB-EDM plays an analogous role. Its `edm:ProvidedCHO` class is the root type regardless of domain, object classification, modality, or institutional origin.

EDM, as a structural container, is designed for breadth rather than depth. Its deliberate underspecification enables cross-institutional aggregation at a scale that no richer, more prescriptive model could achieve. However, this design choice carries three systematic costs that directly impact cross-domain discovery [7].

1. *Property conflation.* EDM maps a wide range of domain-specific properties onto a small set of Dublin Core properties. For instance, distinction in aboutness such as `crm:P67F_refers_to`, `crm:P129F_is_about`, and `crm:P62F_depicts` are collapsed into a single `dc:terms:subject` property [2].
2. *Weak normalization.* EDM tolerates heterogeneity in literal values rather than resolving at the entity level. Appendix A (Figure 5d) illustrates this: Goethe’s name appears in four variations as `dc:subject` values across the corpus. Without canonical entity identifiers at the aggregation layer, semantically identical references remain syntactically distinct. Discovery therefore depends on string matching rather than entity recognition, limiting the possibility of systematic cross-domain linkage.
3. *Information loss from source schemas.* Transformation from rich source formats to EDM is a one-way reduction. For instance, expert-mediated decisions embedded in MARC<sup>9</sup> subfield codes are not recoverable from the resulting EDM record.

Table 1 summarizes the high-value facets that are lost or degraded per institutional domain.

<sup>9</sup><https://www.loc.gov/marc/bibliographic/ecbdcntf.html>

**Table 1**

Domain-level semantics that cannot be recovered from EDM alone.

Domain	High-value facets	EDM limitation
Libraries	Agent role, edition, subject heading	<code>dc:creator contributor</code> conflates all agent roles; <code>dc:subject</code> uncontrolled
Museums	Material, technique, production event, attribution	<code>dc:type</code> uncontrolled; only creation LIDO event
Photographs	Caption, depicted subject, occasion	No dedicated image or subject semantics
Audio/Visual	Cast, director, broadcast date, format	No AV-specific roles or event properties

## 5. Parallel 2: The Librarians and the Facets

Borges' librarians actively sought solutions. Some journeyed in quest of *the catalog of catalogs*, a book that would describe every other volume and provide its coordinates. Generations of librarians scoured hexagon after hexagon in search of meaningful sequences. Others purged volumes they judged meaningless. None succeeded because scanning for patterns within uniform structure cannot compensate for absent semantic differentiation.

A similar situation emerges in the DDB where a user searching for *Faust*-related objects has two options: keyword search over literal strings, and facet filtering over a set of lowest-common-denominator fields, such as `edm:type` for modality, agent names, names of places, `edm:provider`, `dc:type`. Scanning volumes for meaningful sequences is akin to string-based facet search, matching tokens instead of concepts. The strategies for search and exploration is limited by the representational architecture on which it operates.

For instance, searching for *Faust*<sup>10</sup> returns tens of thousands of objects. DDB-EDM provides no principled mechanisms:

- To distinguish a *Work* from one of its *Manifestations*, such as this digital manifestation of the 1830 German edition<sup>11</sup>,
- To differentiate between two *Expressions*, the realization of the one above and this 1826 French edition<sup>12</sup>.
- To identify that this sculpture of a fist<sup>13</sup> is not a manifestation of *Faust*, or
- To traverse the chain of adaptations from Goethe's text to a theater performance in Weimar in 1893<sup>14</sup> to Gounod's opera<sup>15</sup> to Murnau's silent film in 1925/26 to stills from this film<sup>16</sup>.

A filter on `edm:type = TEXT` recovers the book but excludes the playbill and the illustration; a filter on `dc:type = Fotografie` picks up the film still but also thousands of unrelated photographs. These share one thing that no facet can express: they are all Manifestations or Items realizing the same Work — Goethe's *Faust*, identified in the Gemeinsame Normdatei (GND) as entity `gnd:4099197-0`<sup>17</sup>.

KG-enabled search shifts the search from keyword-based filtering to navigating entities and relations. Once objects are represented as nodes in a typed graph the search patterns described above emerge naturally. WEMI-typed entities connected by relational properties, as shown in Figure 3 allow queries to move up and down the abstraction hierarchy. A query starting at the `:work` node for Goethe's *Faust* can traverse `:hasRealization` links to reach all `:Expressions`: the German text, Gounod's aria, Murnau's silent film, then descend via `:hasManifestation` to all publications of the German text.

However, a knowledge graph does not automatically resolve metadata inaccuracies or transformation loss caused by aggregation. In addition, retrieval performance relies heavily on good quality metadata. A graph propagates wrong metadata with the same fidelity as correct metadata. There will be uncertainties in assigning specific agent roles since this information is neither available nor can they be reconstructed from the digitized manifestation of this 1826 French edition<sup>12</sup>. Even with the availability of task-specific state-of-the-art models, 100% accuracy is not assured.

<sup>10</sup><https://www.deutsche-digitale-bibliothek.de/searchresults?query=faust>

<sup>11</sup><https://www.ddb.de/item/6N46KJTHTPHEWQI6JRB5U5PKFDHAN5LT>

<sup>12</sup><https://www.ddb.de/item/U73FFILJ4XEH2QITBYJ3AE4Q3754WE3Y>

<sup>13</sup><https://www.ddb.de/item/7NBEMRVL15TBCFODES5X22QRYHRTO6AI>

<sup>14</sup><https://www.ddb.de/item/Q5LA52X7GWYK65FB6LMCFLSORI3BNSPR>

<sup>15</sup><https://www.ddb.de/item/GJAFP4CEYQC6MSM54ENC67MHRC44JRLW>

<sup>16</sup><https://www.ddb.de/item/RW6KFZP37RO2JG3Y4DUMZK6A5IV2BEDT>

<sup>17</sup><https://d-nb.info/gnd/4099197-0>

## 6. Parallel 3: The Decipherers and Domain Ontologies

In Borges' Library, most books are impenetrable: their pages contain nothing but incomprehensible sequences of symbols. The decipherers offer rare respite by even obscure languages of such volumes [1]. Domain ontologies play the decipherer's role for the DDB: they partially restore the semantic vocabulary in which each object's metadata is expressed, recovering the orientation that EDM's design trade-offs removed.

MOCHO (Mid-level Ontology for Cross-domain Heritage Objects) makes this operational through an alignment method and patterns applicable to existing records without re-ingestion. Table 2 summarizes the effect across key dimensions.

**Table 2**  
EDM only vs. EDM + domain ontologies across key dimensions.

Dimension	EDM only	EDM + domain ontologies
Search strategy	Facet filtering	Entity + graph traversal
Semantics	Literal strings, semi-controlled	Typed entities, controlled vocabularies
Relevance	String matching	Semantic proximity
Explainability	Field-value match	Typed relation path

### 6.1. WEMI as Pivot Entities

With emphasis on reuse, eight domain- and modality-specific ontologies act as EDM's "decipherers". Their adaptation in the DDB shares one structural commitment which is the distinction of the WEMI entities amongst each other. This commitment originates in the IFLA FRBR Model [3], and is expressed in domain-specific vocabularies that MOCHO aligns. WEMI therefore serves as the structural pivot within MOCHO. `edm:ProvidedCHO` cannot play this role, as it conflates all four WEMI levels into a single class. In doing so, even a fully materialized knowledge graph over EDM triples has no stable node to serve as a *Work* level anchor, onto which all manifestations of the same intellectual creation can be attached.

The *Faust* running example illustrates this limitation. Three DDB objects: the 1830 printed edition<sup>11</sup>, the 1893 theater playbill from Weimar<sup>14</sup>, and Gounod's "Jewel Song" are *Manifestations* or *Items* realizing the same *Work*: Goethe's *Faust*. However, the sculpture of a fist<sup>13</sup> is, in itself, a distinct *Work*.

<sup>‡</sup>Class-Sharing between WEMI Entities.

<sup>\*</sup>Partial match (restricted to a subtype).

<sup>\*\*</sup>Level collapse where a property bridges two levels in one step.

**Table 3**  
Domain- and modality-specific ontologies and their FRBR WEMI-level equivalents.

Domain/Modality	Ontology	Work	Expression	Manifestation	Item
Library	RDA	<code>rdac:C10001</code>	<code>rdac:C10006</code>	<code>rdac:C10007</code>	<code>rdac:C10003</code>
Museums (event-centric)	CIDOC CRM	<code>E89_Propositional_Object</code>	<code>E73_Information_Object</code> <sup>‡[E/M]</sup>		<code>E24_Physical_Human-Made_Thing</code>
Visual Arts (object-centric)	VRA Core	<code>vra:Work</code> <code>vra:Collection</code>	<code>vra:Inscription</code> <sup>*</sup>	<code>vra:Image</code> <sup>**</sup>	—
Archive	RiC-O	<code>rico:RecordResource</code> <code>rico:RecordSet</code> <code>rico:Record</code> <sup>‡[W/E]</sup>		<code>rico:Instantiation</code> <sup>‡[M/I]</sup>	
Image	LIO	<code>lio:Image</code>	—	—	—
Music	Music Ontology	<code>mo:MusicalWork</code>	<code>mo:MusicalExpression</code>	<code>mo:MusicalManifestation</code> <code>mo:Record</code>	<code>mo:MusicalItem</code>
Audio	ACO	—	<code>aco:AudioExpression</code>	<code>aco:AudioManifestation</code>	<code>aco:AudioItem</code>
Audiovisual	EBUCorePlus	<code>ec:EditorialWork</code>	— <sup>**</sup>	<code>ec:MediaResource</code>	—

## 6.2. The Domain Ontologies

No single domain ontology covers the full semantic range represented in the DDB's collection. MOCHO aligns to eight domain- and modality-specific ontologies, each addressing a distinct institutional model: **RDA** (bibliographic), **CIDOC CRM** (museum), **VRA Core** (visual arts), **RiC-O** (archival), **LIO** (images), **MO** (music), **ACO** (audio), and **EBUCorePlus** (audiovisual). Table 4 additionally lists EDM, DDB-EDM, and FRBR-family ontologies (FRBR, FaBiO, DocO) that serve as structural stepping stones in the BFS traversal (Section 6.3.1). The complete list of ontologies, their namespaces, and source RDF definitions is enumerated in Table 4. Table 3 summarizes the WEMI-level equivalents. Each is described in turn in the succeeding paragraphs.

Ontologies were selected against two criteria: domain coverage matching the DDB's six institutional sectors (libraries, museums, archives, audiovisual media, images, and music/audio), and existing WEMI alignment – either explicit `rdfs:subClassOf` declarations to FRBR Core classes, or documented mappings – to support BFS traversal. Candidates were surveyed using Linked Open Vocabularies<sup>17</sup> and cross-checked against sector standards. For the bibliographic domain, RDA was preferred over BIBFRAME, which does not map directly to WEMI levels, because RDA implements LRM directly and is the current DNB standard [20]. VRA Core and RiC-O were the only RDF-based standards in their respective domains (visual arts and archives) with a WEMI-level structure; EAD (archives) and CDWALite (visual arts) are XML schemas without OWL formalizations. MO provides FRBR-mapped music classes [14]; ACO [15] extends that model to non-musical audio that carries no Work-level concept, making the two ontologies complementary rather than overlapping.

**RDA** (*Resource Description and Access*). RDA is MOCHO's primary bibliographic target ontology. Its classes and properties directly represent the Library Reference Model (LRM) [21]. LRM supersedes FRBR and its variants with a unified conceptual model, covering all types of bibliographic resources, which makes it the most pragmatic entry point for library-sector alignment. The National Library of Germany (*Deutsche Nationalbibliothek*, DNB) adopted RDA in 2015 [20].

**CIDOC CRM**. Museum objects are naturally modeled as products of events rather than as documents: a portrait of Goethe is the outcome of a production event carried out by a named artist at a known place and time. LRMoo harmonizes this event model with the FRBR/LRM conceptual hierarchy [22].

**VRA Core**. Developed by the Visual Resources Association (VRA) and hosted by the MARC Standards Office of the Library of Congress, this standard describes visual arts collections using a tripartite schema: `vra:Work` (the intellectual or physical artifact), `vra:Image` (a visual surrogate or digital reproduction), and `vra:Collection` (an aggregation of works). `vra:Work` and `vra:Collection` align to the Work level; `vra:Image` aligns to Manifestation with a level-collapse annotation, since VRA links Work to Image directly, without an intermediate Expression.

**RiC-O** (*Records-in-Context Ontology*). Archival description follows neither the bibliographic document model nor the museum event model: its fundamental concept is the record hierarchy. Fonds, series, file, item, along with the provenance chain that explains why records exist together [13]. RiC-O formalizes these archival concepts and is the alignment target for the archive sector in MOCHO.

**LIO** (*Lightweight Image Ontology*). Visual and pictorial objects, such as photographs, illustrations, or art prints, require vocabulary for iconographic subject, composition, and physical medium that bibliographic ontologies do not provide. LIO [23] links the CHO with its digital representations with depiction semantics. Its central class `lio:Image`, a digital image, is a subclass of `foaf:Image`, which aligns to the Work level, since in LIO a painting or photograph is the primary object of description.

**MO** and **ACO** (*Music Ontology, Audio Core Ontology*). Music Ontology [14] models music production workflows with four superclasses corresponding to FRBR WEMI classes: `mo:MusicalWork`, `mo:MusicalExpression`, `mo:MusicalManifestation`, `mo:MusicalItem`. The Audio Commons Ontology (ACO) [15] generalizes MO for

---

<sup>17</sup><https://lov.linkeddata.es>

**Table 4**  
Domain ontologies and their namespaces.

Ontology	Namespace	Version
EDM	<a href="http://www.europeana.eu/schemas/edm/">http://www.europeana.eu/schemas/edm/</a> Source: <a href="https://github.com/hugomanguinhas/europeana_ld/tree/master/ld-edm/src/main/resources/etc/owl">https://github.com/hugomanguinhas/europeana_ld/tree/master/ld-edm/src/main/resources/etc/owl</a>	5.2.4
DDB-EDM	<a href="https://ise-fizkarlsruhe.github.io/ddbkg/ddbedm/">https://ise-fizkarlsruhe.github.io/ddbkg/ddbedm/</a> Source: <a href="https://ise-fizkarlsruhe.github.io/ddbkg/ddbedm/source/ddbedm-full.owl">https://ise-fizkarlsruhe.github.io/ddbkg/ddbedm/source/ddbedm-full.owl</a>	1.0
FRBR	<a href="https://sparontologies.github.io/frbr/current/frbr.html">https://sparontologies.github.io/frbr/current/frbr.html</a> Source: <a href="https://github.com/SPAROntologies/frbr/tree/master/docs/2018-03-29">https://github.com/SPAROntologies/frbr/tree/master/docs/2018-03-29</a>	1.0.1
FaBiO	<a href="https://sparontologies.github.io/fabio/current/fabio.html">https://sparontologies.github.io/fabio/current/fabio.html</a> Source: <a href="https://github.com/SPAROntologies/fabio/tree/master/docs/2023-05-09">https://github.com/SPAROntologies/fabio/tree/master/docs/2023-05-09</a>	2.1
DocO	<a href="https://sparontologies.github.io/doco/current/doco.html">https://sparontologies.github.io/doco/current/doco.html</a> Source: <a href="https://github.com/SPAROntologies/doco/tree/master/docs/2015-07-03">https://github.com/SPAROntologies/doco/tree/master/docs/2015-07-03</a>	1.3
RDA	<a href="http://rdaregistry.info/Elements/*">http://rdaregistry.info/Elements/*</a> Source: <a href="https://github.com/RDARegistry/RDA-Vocabularies/releases/tag/v5.4.9">https://github.com/RDARegistry/RDA-Vocabularies/releases/tag/v5.4.9</a>	5.4.9
RiC-O	<a href="https://www.ica.org/standards/RiC/ontology/">https://www.ica.org/standards/RiC/ontology/</a> Source: <a href="https://github.com/ICA-EGAD/RiC-O/blob/master/ontology/current-version/RiC-O_1-1.rdf">https://github.com/ICA-EGAD/RiC-O/blob/master/ontology/current-version/RiC-O_1-1.rdf</a>	1.1
CIDOC CRM	<a href="http://www.cidoc-crm.org/rdfs/cidoc-crm#">http://www.cidoc-crm.org/rdfs/cidoc-crm#</a> Source: <a href="https://github.com/erlangen-crm/ecrm/blob/master/ecrm_240307.owl">https://github.com/erlangen-crm/ecrm/blob/master/ecrm_240307.owl</a>	7.1.3
VRA	<a href="http://purl.org/vra/">http://purl.org/vra/</a> Source: <a href="https://s3.amazonaws.com/VRA/ontology.html">https://s3.amazonaws.com/VRA/ontology.html</a>	4.0
LIO	Not maintained: <a href="http://purl.org/net/lio/">http://purl.org/net/lio/</a> Source: <a href="https://imagesnippets.com/lio/lio.owl">https://imagesnippets.com/lio/lio.owl</a>	1.0
ACO	<a href="https://w3id.org/ac-ontology/aco#">https://w3id.org/ac-ontology/aco#</a> Source: <a href="https://github.com/AudioCommons/ac-ontology/releases/tag/v1.2.3">https://github.com/AudioCommons/ac-ontology/releases/tag/v1.2.3</a>	1.2.3
MO	<a href="http://purl.org/ontology/mo/">http://purl.org/ontology/mo/</a> Source: <a href="https://github.com/motools/musicontology/tree/master/rdf">https://github.com/motools/musicontology/tree/master/rdf</a>	1.0
EBUCore Plus	<a href="http://www.ebu.ch/metadata/ontologies/ebucoreplus">http://www.ebu.ch/metadata/ontologies/ebucoreplus</a> Source: <a href="https://github.com/ebu/ebucoreplus/tree/main/ontology/EBUCorePlus">https://github.com/ebu/ebucoreplus/tree/main/ontology/EBUCorePlus</a>	2.0.0

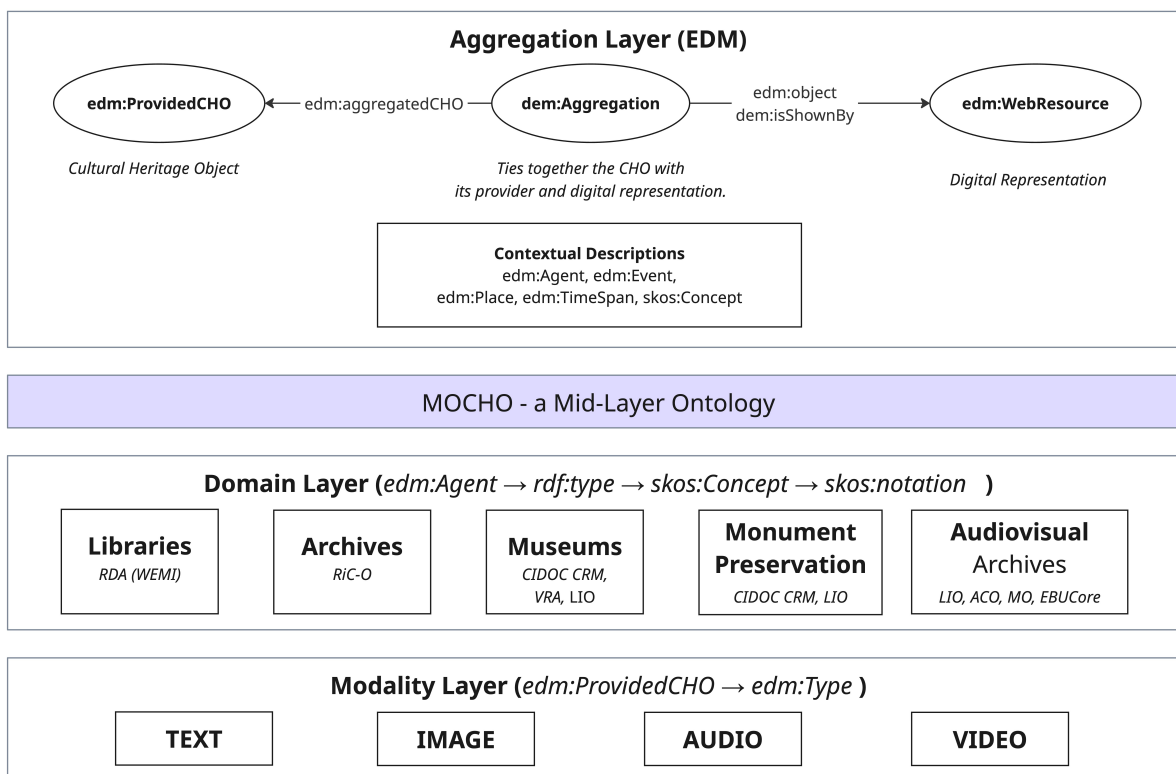
non-musical audio, such as natural sounds, oral traditions, field recordings. In contrast to MO, ACO defines no Work-level class, since not all audio content expresses an intellectual creation. Tan et al. [24] align DDB-EDM to both ontologies using an object-centric approach: event classes are omitted because DDB records carry no performance or production event details, and `edm:ProvidedCHO` is typed as a specialization of `frbr:Item`.

**EBUCorePlus**<sup>18</sup>. Developed by the European Broadcasting Union (EBU), it is designed to describe audiovisual resources including audio/video files and associated broadcast content. In EBUCore, `ec:EditorialWork` and `ec:Programme` align to the Work level, while `ec:MediaResource` and `ec:Essence` align to the Manifestation level. The Expression tier is absent; EBUCorePlus links Work directly to Manifestation via `ec:isInstantiatedBy`, which is an approximate alignment with a flagged level-collapse criterion rather than asserting an Expression node the source data does not support.

### 6.3. MOCHO: A Mid-level Ontology for Cross-domain Heritage Objects

MOCHO mediates between EDM and the eight domain and modality ontologies (Figure 1). The key design decision is to type objects at the *WEMI level* rather than at the level of fine-grained FaBiO classes [11]. This typing is much coarser, but robust enough without per-provider curation.

<sup>18</sup><https://tech.ebu.ch/metadata/ebucore>



**Figure 1:** MOCHO mediates between EDM and the eight domain and modality ontologies.

### 6.3.1. Alignment Method

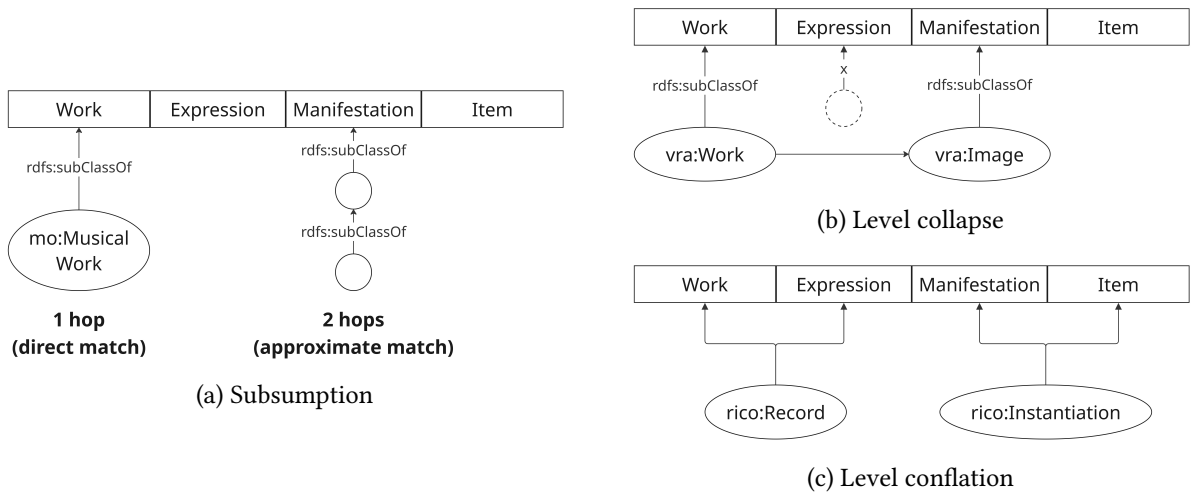
As shown in Table 3, most domain ontologies include at least one class anchored to the FRBR Core Ontology<sup>19</sup> classes, which maintain a 1-to-1 correspondence with the WEMI levels. This WEMI-based structural parallel forms the basis of the alignment procedure. For each class, breadth-first traversal (BFS) of the subclass chain (`rdfs:subClassOf`) terminates at the first FRBR anchor found, identifying the WEMI level. The corresponding FRBR  $\rightarrow$  RDA mapping then yields the RDA target class.

To make this concrete, in Music Ontology, `mo:MusicalWork` declares “`rdfs:subClassOf frbr:work`”. The anchor is therefore reached in one hop, and the FRBR  $\rightarrow$  RDA mapping for FRBR is direct, so `mo:MusicalWork` receives a *direct* match to `rdac:C10001 "work"`. Where no `rdfs:subClassOf` path to a FRBR anchor exists, as in VRA Core 4.0, the class is assigned via the hand-curated mapping in Table 3. The procedure is purely structural: no label comparison, embedding distance, or synonym lookup is performed in this study. In addition to the WEMI classes, FRBR’s *Groups 2* and *3* are also used as pivots: *Group 2* for responsible entities (agent, person, organization, family) and *Group 3* for subject heading.

All alignments described in this section and summarized in Table 3 were performed by the authors. When *direct* matches are absent, alignment proceeds through straightforward structural subsumption in one or two hops (Figure 2). Three recurring deviations from a direct alignment require specific handling, each caused by domain inapplicability rather than ontology gaps:

*Level Collapse.* VRA Core and EBUCorePlus omit the Expression tier. VRA Core has no Expression class because in visual arts cataloging, there is no meaningful step between a painting and its image surrogate. EBUCore links its editorial work directly to a media resource. Both reflect a domain convention in which the carrier is the primary object of description, leaving no counterpart for a separate realized-form class.

<sup>19</sup><http://purl.org/vocab/frbr/core>



**Figure 2:** Three recurring alignment cases in MOCHO: subsumption via `rdfs:subClassOf` traversal (left), level collapse where the Expression tier is absent (top right), and level conflation where a single class spans two WEMI levels (bottom right).

*Level Conflation.* RiC-O spans both Work and Expression because archival description is organized around provenance rather than the bibliographic distinction between an abstract text and its realization. MOCHO therefore places Record classes in both the Work and Expression unions, and Instantiation in both the Manifestation and Item unions, accepting the approximation rather than imposing a bibliographic structure that the archival model rejects.

*Partial Coverage.* LIO defines only a Work-level class (`lio:Image`). In LIO an image is itself the primary object, not a carrier of something else. ACO defines no Work-level class, since not all audio content expresses an intellectual creation. Each union is populated only with classes for which the contributing ontology has domain warrant; the gaps are deliberate.

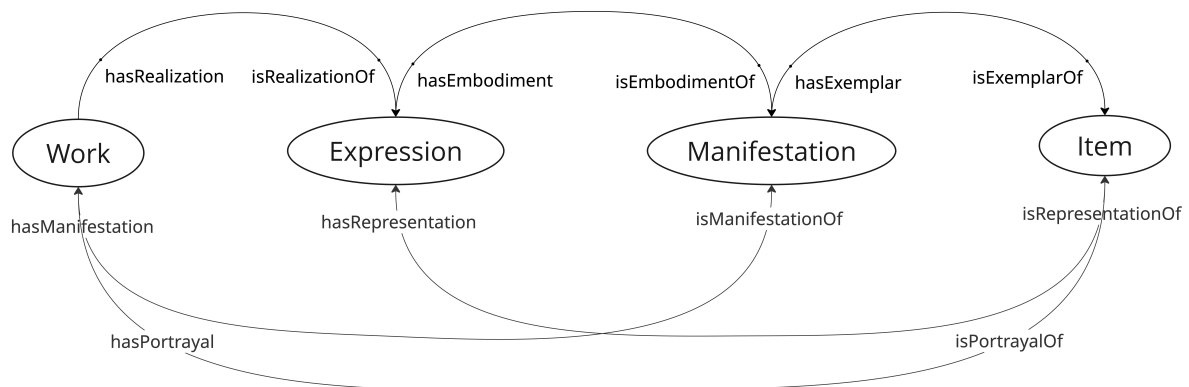
**Property Alignment.** Property alignment follows the same BFS procedure over `rdfs:subPropertyOf` edges. When no path to a FRBR property anchor exists, domain/range pairs identify the closest FRBR bridge property. Links routed through an event class are flagged to mark the *indirection*. CIDOC CRM and Music Ontology both model agent contribution with events (*event-centric*), not of a direct creator-to-object relationship (*object-centric*). Pattern 3 handles the corresponding modeling asymmetry. Structural traversals with some property proved problematic: `frbr:subject` fans out to dozens of subject properties and `frbr:producer` splits into multiple distinct agent roles in RDA.

**Core WEMI Classes and Properties.** MOCHO defines four classes (Work, Expression, Manifestation, and Item) and properties connecting WEMI levels. In addition, MOCHO also defines several FRBR and FaBiO properties connecting adjacent and non-intermediate WEMI levels (Figure 3).

**Pattern 1: Multiple Typing.** A DDB item is typed simultaneously as `edm:ProvidedCHO` and as a MOCHO WEMI class, preserving backward compatibility while adding the typed entity required for graph traversal (Listing 1).

Listing 1: Multiple Typing in MOCHO.

```
ddb:GJAFP4CEYQC6MSM54ENC67MHRC44JRLW
a edm:ProvidedCHO ;           # aggregation
a mocho:Work ;                # bibliographic anchor
a mo:MusicalWork ;           # domain-specific type
dc:title "Faust_:Jewel_song/_(Gounod)"@en ;
dc:alternative "Faust_(Einheitssachtitel)"@de ;
...
```



**Figure 3:** Core WEMI classes and properties from FRBR and FaBiO defined in mocho.

**Pattern 2: Subproperties.** Domain-specific properties are asserted as `rdfs:subPropertyOf` of the generic Dublin Core terms they shadow, preserving backward compatibility while exposing precise role distinctions. A query on `dc:creator` returns all creators via sub-property inference; `rdaa:P50190` retrieves only the composer role (Listing 2).

Listing 2: Definition of domain-specific properties using `rdfs:subPropertyOf`.

```

@prefix rdaa: <http://rdaregistry.info/Elements/a/> .
@prefix rdaw: <http://rdaregistry.info/Elements/w/> .
@prefix rdae: <http://rdaregistry.info/Elements/e/> .
@prefix dc: <http://purl.org/dc/elements/1.1/> .

# has editor agent of
rdaa:P50450 rdfs:subPropertyOf dc:contributor .

# has creator agent of work
rdaw:P10065 rdfs:subPropertyOf dc:creator .

# has creator agent of expression
rdae:P20053 rdfs:subPropertyOf dc:creator .

# has composer agent of work"
rdaw:P10053 rdfs:subPropertyOf dc:creator .

```

**Pattern 3: Event- or Object-Centric.** Museum and music ontologies model attribution through production events: a portrait is the outcome of `crm:E12_Production`; a musical work is the outcome of `mo:Composition`. Bibliographic ontologies attach creator properties directly to the object. mocho supports both conventions, bridging them through `edm:Event` and the CIDOC CRM event hierarchy, so that cross-domain queries retrieve objects from either modeling convention transparently.

**Pattern 4: Union-of Inheritance.** The `owl:unionOf` necessary condition ensures that a mocho WEMI typing is valid only if the node belongs to at least one domain-specific class. Instances without domain membership are flagged as under-specified by a reasoner. Listing 3 shows the definition of `mocho:Work` according to Table 3.

The union contains domain-specific classes rather than `frbr:Work` directly. Most members already declare `rdfs:subClassOf frbr:Work` (e.g., `mo:MusicalWork`), so FRBR ancestry follows transitively.

Listing 3: MOCHO Work class as shown in Table 3.

```
mocho:Work rdfs:subClassOf [
  owl:unionOf ( rdac:C10001 crm:E89_Propositional_Object
                vra:Work vra:Collection
                lio:Image rico:RecordResource
                rico:RecordSet rico:Record
                mo:MusicalWork ec:EditorialWork) ] .
```

## 6.4. Alignment Strategies for Existing EDM Records

The four patterns can be applied to existing DDB-EDM records without modifying the underlying records, re-running ingest pipelines, or requiring provider cooperation. Two complementary strategies drive the enrichment.

**Heuristic typing.** In addition to provider information, `dc:type`, `edm:type`, and `ddb:hierarchyType` field values can be mapped to MOCHO WEMI subclasses by rule. For example, metadata provided by libraries with `edm:type = "TEXT"` and `dc:type = "Buch"` or `ddb:hierarchyType = "Monografie"` are considered instances `mocho:Manifestation`. While records with `edm:type = "IMAGE"` combined with a LIDO source schema assigns a museum Item or Manifestation.

**Statistical and NER-based enrichment.** Heuristic rules address modality and structural type but not entity identity. Fine-grained named entity recognition over literal values (titles, contributor strings, date fields) can surface entity candidates that are then linked to authority file URIs (GND, Wikidata, Getty vocabularies) and mapped to MOCHO properties. A contributor string “hrsg. von”<sup>20</sup> identifies an editorial role suitable for `rdaa:P50450`; a place name in a `dc:spatial` field can be linked to a TGN or Wikidata URI suitable for `edm:Place`. Where entity linking confidence is low, the enriched triples carry a provenance annotation rather than being asserted as ground truth.

## 7. Discussion

**What alignment achieves and what it cannot fix.** MOCHO adds semantic structure to EDM records but does not correct metadata itself. This alignment exercise is a structural solution, not a data quality solution. Ingestion from MARC, LIDO, or EAD to EDM is a one-way reduction: expert-mediated decisions encoded in MARC subfield codes are silently dropped [2]. MOCHO can recover some of this precision by recognizing editorial markers in `dc:contributor` literals and asserting `rdaa:P50450` (*is editor agent of*). However, objects with no such indication in the EDM record are unrecoverable by alignment alone. Agent distinction for such objects can only be recovered with manual intervention from the metadata provider.

**The LLM objection.** LLMs can perform named entity recognition (NER) and surface-level inference over raw EDM literals without ontological scaffolding. What they cannot reliably determine is the WEMI level of a record, whether it represents a Work, a Manifestation, or an Item. LLMs and MOCHO are complementary: WEMI typing provides language models with a stable entity anchor, improving retrieval precision beyond string-level inference. In this way, LLMs and KGs complement each other.

**Broader applicability.** MOCHO is designed for reuse across aggregators: the four alignment patterns are broadly applicable, and the `owl:unionOf` constraint on each WEMI class is agnostic to which domain ontology a given institution uses.

## 8. Conclusion

Three structural parallels connect Borges’ Library to the DDB. EDM defines *what*: its `edm:ProvidedCHO` container imposes structural uniformity across every contributing institution, but carries no semantic

<sup>20</sup>*Herausgeben von* is the standard German phrase for “edited by”

priority. KG-enabled search defines *how objects relate*: entity-first traversal replaces field-value scanning with typed relation paths. Domain ontologies define *what type of thing it is*: MOCHO makes this operational through four alignment patterns that restore semantic orientation without dismantling the existing aggregation structure.

MOCHO mediates between EDM's flat model and the expressive depth of modality and GLAM-spanning domain ontologies. Its four WEMI classes, defined via `owl:unionOf` across eight domain ontologies, act as cross-domain pivot entities. Together the four patterns enable queries that EDM alone cannot support: cross-domain Work-level faceting, WEMI-level traversal from intellectual creation to physical item, and agent role disambiguation across the DDB's institutional sectors. MOCHO is formalized and published as a standalone OWL ontology with a versioned namespace, enabling independent adoption and formal validation.

Several directions remain open. The search improvement claim requires evaluation against retrieval metrics, in addition to domain-expert-developed competency questions. A concrete next step is to evaluate MOCHO's alignment coverage and WEMI-level typing precision on DDB's digitized subset, using it as a held-out benchmark for both recall of correctly typed objects and precision of heuristic rules. Metadata quality remains unaddressed by alignment alone. Heuristics and statistical-based refinement are to be considered for future planning.

Aggregation without orientation is the Library without a catalog. Without a semantic layer, every object looks like every other: the traveler who enters in search of *Faust* emerges with thousands of records and no way to ask which of them is the Work. MOCHO is the orientation the Library lacked, it makes "*the catalog of catalogs*" possible.

## Declaration on Generative AI

This work employed AI tools during ideation (ChatGPT) and data preparation scripting (Claude Code). All content was reviewed and edited by the authors.

## References

- [1] Borges, Jorge Luis, *The Library of Babel*, in: *Collected Fictions*, Penguin, New York, 1998, pp. 112–120.
- [2] S. Rühle, F. Schulze, M. Büchner, *Applying a linked data compliant model: The usage of the Europeana Data Model by the Deutsche Digitale Bibliothek*, in: *Proceedings of the International Conference on Dublin Core and Metadata Applications*, volume 2014 of *Dublin Core Conference*, Dublin Core Metadata Initiative, Dublin, OH, USA, 2014. URL: <https://dcpapers-data.dublincore.org/articles/dc-2014/952136503/files/dcmi-952136503.pdf>. doi:10.23106/dcmi.952136503.
- [3] IFLA Study Group (ISG), *Functional Requirements for Bibliographic Records - Final Report*, K. G. Saur, Berlin, Boston, 1998. URL: <https://doi.org/10.1515/9783110962451>. doi:doi:10.1515/9783110962451.
- [4] J. Basile, *Tar for Mortar: "The Library of Babel" and the Dream of Totality*, punctum books, Earth, Milky Way, 2018. doi:10.21983/P3.0196.1.00.
- [5] M. Doerr, *The CIDOC conceptual reference module: an ontological approach to semantic interoperability of metadata*, *AI Mag.* 24 (2003) 75–92.
- [6] Europeana, *Definition of the Europeana Data Model v5.2.7*, Technical Report, Europeana, 2016. URL: [https://pro.europeana.eu/files/Europeana\\_Professional/Share\\_your\\_data/Technical\\_requirements/EDM\\_Documentation/EDM\\_Definition\\_v5.2.7\\_042016.pdf](https://pro.europeana.eu/files/Europeana_Professional/Share_your_data/Technical_requirements/EDM_Documentation/EDM_Definition_v5.2.7_042016.pdf).
- [7] S. Peroni, F. Tomasi, F. Vitali, *Reflecting on the Europeana Data Model*, in: *IRCDL 2012*, 2012, pp. 228–240.
- [8] S. Zapounidou, M. Sfakakis, C. Papatheodorou, *Integrating library and cultural heritage data models: the bibframe - edm case*, in: *Proceedings of the 18th Panhellenic Conference on Informatics*,

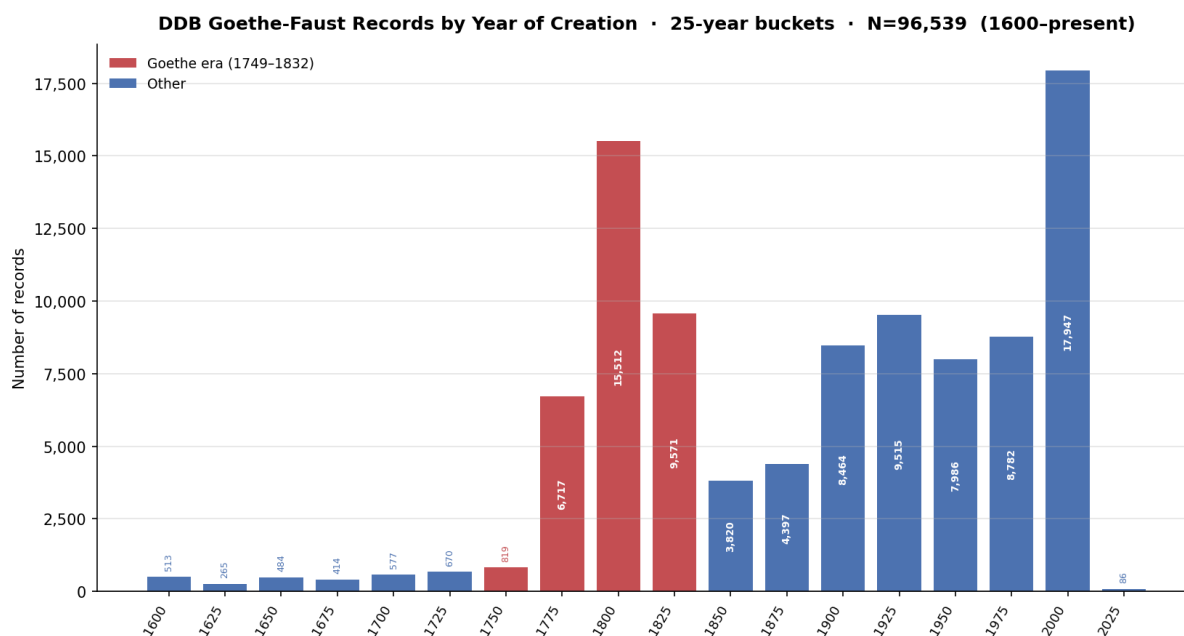
- PCI '14, Association for Computing Machinery, New York, NY, USA, 2014, p. 1–6. URL: <https://doi.org/10.1145/2645791.2645805>. doi:10.1145/2645791.2645805.
- [9] R. Bennett, B. F. Lavoie, E. T. O'Neill, The Concept of a Work in WorldCat: an Application of FRBR, *Library Collections, Acquisitions, and Technical Services* 27 (2003) 45–59. URL: <https://www.sciencedirect.com/science/article/pii/S1464905502003068>. doi:[https://doi.org/10.1016/S1464-9055\(02\)00306-8](https://doi.org/10.1016/S1464-9055(02)00306-8).
- [10] S. Peroni, D. Shotton, FaBiO and CiTO: Ontologies for describing bibliographic resources and citations, *Journal of Web Semantics* 17 (2012) 33–43. doi:10.2139/ssrn.3198992.
- [11] M. A. Tan, T. Tietz, O. Bruns, J. Oppenlaender, D. Dessì, H. Sack, DDB-EDM to FaBiO: The Case of the German Digital Library, in: *Proc. of the 20th Int. Semantic Web Conference - Posters and Demos – ISWC 2021*, volume 2980, CEUR-WS.org, 2021.
- [12] C. Dijkshoorn, L. Aroyo, J. van Ossenbruggen, G. Schreiber, Modeling Cultural Heritage Data for Online Publication, *Appl. Ontology* 13 (2018) 255–271.
- [13] International Council on Archives, Experts Group on Archival Description (ICA-EGAD), Records in Contexts Ontology (RiC-O), Version 1.0, 2023. URL: <https://www.ica.org/standards/RiC/ontology>, ontology IRI: <https://www.ica.org/standards/RiC/ontology>.
- [14] Y. Raimond, S. Abdallah, M. Sandler, F. Giasson, The Music Ontology, in: *Proc. of the 8th Int. Conference on Music Information Retrieval, ISMR 2007*, Vienna, Austria, September 23-27, 2007, Austrian Computer Society, 2007, pp. 417–422.
- [15] M. Ceriani, G. Fazekas, Audio commons ontology: A data model for an audio content ecosystem, in: D. Vrandečić, K. Bontcheva, M. C. Suárez-Figueroa, V. Presutti, I. Celino, M. Sabou, L.-A. Kaffee, E. Simperl (Eds.), *The Semantic Web – ISWC 2018*, Springer Int. Publishing, Cham, 2018, pp. 20–35. doi:10.1007/978-3-030-00668-6\_2.
- [16] Y. He, J. Chen, D. Antonyrajah, I. Horrocks, BERTMap: A BERT-Based Ontology Alignment System, in: *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022*, AAAI Press, 2022, pp. 5692–5700. URL: <https://doi.org>. doi:10.1609/aaai.v36i5.20510.
- [17] S. Hertling, H. Paulheim, OLaLa: Ontology Matching with Large Language Models, in: *Proceedings of the 12th Knowledge Capture Conference 2023, K-CAP '23*, Association for Computing Machinery, New York, NY, USA, 2023, p. 131–139. URL: <https://doi.org/10.1145/3587259.3627571>. doi:10.1145/3587259.3627571.
- [18] P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, S. Riedel, D. Kiela, Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks, in: *Advances in Neural Information Processing Systems*, volume 33, Curran Associates, Inc., 2020, pp. 9459–9474. URL: [https://proceedings.neurips.cc/paper\\_files/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf).
- [19] D. Edge, H. Trinh, N. Cheng, J. Bradley, A. Chao, A. Mody, S. Truitt, D. Metropolitanaky, R. O. Ness, J. Larson, From Local to Global: A Graph RAG Approach to Query-Focused Summarization, 2025. URL: <https://arxiv.org/abs/2404.16130>. arXiv:2404.16130.
- [20] Deutsche National Bibliothek, Jorge Luis Borges' Babel's Library Designed By Various Artists, [https://www.dnb.de/EN/Professionell/Standardisierung/Standards/\\_content/rda.html](https://www.dnb.de/EN/Professionell/Standardisierung/Standards/_content/rda.html), 2022. Accessed: 2026-01-26.
- [21] M. Žumer, IFLA Library Reference Model (IFLA LRM)— Harmonisation of the FRBR Family, *International Society for Knowledge Organization* 45 (2018) 310–318. doi:10.5771/0943-7444-2018-4-310.
- [22] P. R. Aalberg, Trond, M. Žumer, LRMoo: Object-Oriented Definition and Mapping from the IFLA Library Reference Model, Technical Report, International Federation of Library Associations and Institutions, 2024. URL: <https://repository.ifla.org/handle/20.500.14598/3677>.
- [23] P. Hayes, M. Warren, A Lightweight Ontology for Describing Images., in: *Linked Data Meets Artificial Intelligence, Papers from the 2010 AAAI Spring Symposium*, Technical Report SS-10-07, 2010.
- [24] M. A. Tan, E. Posthumus, H. Sack, Audio Ontologies for Intangible Cultural Heritage, in: *Proc. of the 19th European Semantic Web Conference - Posters and Demos – ESWC 2022*, 2022.

## A. Goethe-Faust Corpus

The choice of “Goethe” as one of the search queries <sup>21</sup> used to gather metadata for this corpus reflects Johann Wolfgang von Goethe’s importance in Germany’s cultural identity. A true *Renaissance* man, he was a writer, scientist, politician, and avid collector, making him a natural anchor for a cultural heritage dataset. The keyword search in the DDB portal returned 96,773 objects. A further 25,275 objects were retrieved by searching for “*faust*”, one of Goethe’s most important literary works. After deduplication, the corpus contains 115,432 objects. Generated data analysis, their associated scripts, and the accompanying documentation are hosted in a GitHub repository<sup>4</sup>.

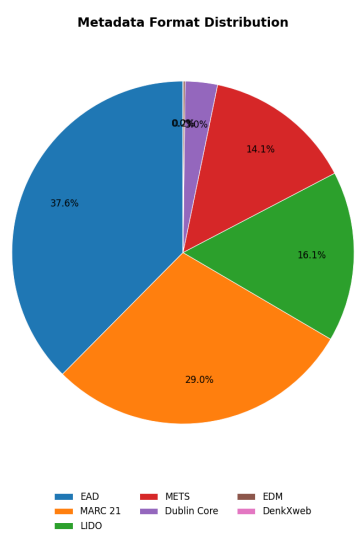
This corpus includes metadata of literary texts, illustrations, screenplays, photographs, audio recordings, and videos from 454 unique providers across various memory institutions, making it a suitable test case for examining cross-domain heterogeneity within a uniformly structured aggregation.

The distribution of approximately 88.8% of objects with date attributions is shown in Figure 4 in 25-year intervals. A small fraction of objects (2.6%) dated before 1600 are excluded from the graph. Objects dated in the first quarter of the 21st century are predominantly theses or dissertations (*Hochschulschrift*) and working papers (*Arbeitspapier*), most likely with Goethe as their subject (See Figure 5d).

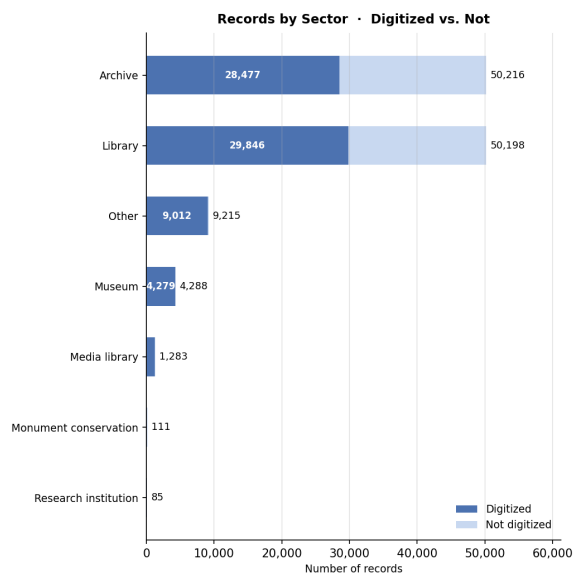


**Figure 4:** DDB Objects associated with the keywords “Goethe” and “*Faust*” by year of creation.

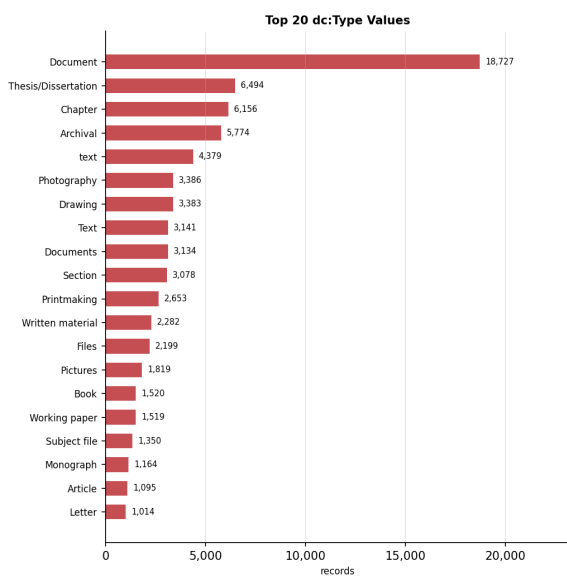
<sup>21</sup><https://api.deutsche-digitale-bibliothek.de/2/search/index/search/select?q=goethe>



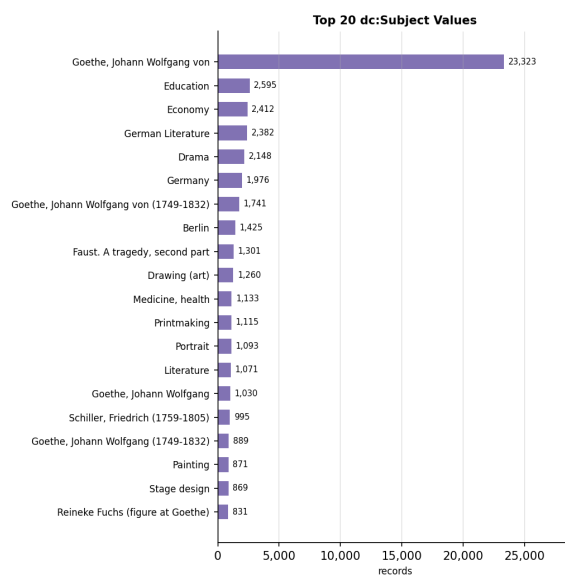
(a) Metadata format distribution



(b) Records by sector and digitization status



(c) Top 20 dc : type values



(d) Top 20 dc : subject values

**Figure 5:** Overview of the DDB Goethe-*Faust* corpus ( $N = 115,432$  records). Subfigure (a) shows the source metadata schema distribution; (b) breaks records down by institutional sector and digitization status; (c) and (d) list the most frequent object type terms and subject terms after German-English translation.

## B. Automatic alignment of `dc:type` to FaBiO Subclasses

From the corpus provided by libraries (`sec_02`, 27%), 359 unique object type terms (`dc:terms`) were first extracted, then looked up against 98 FaBiO and 53 DocO subclass names.

Table 5 shows that a large proportion of the object type terms have no correspondence to the names of FaBiO and DocO subclasses, as the terms refer to non-bibliographic objects such as buildings, medals, artworks, and furniture.

The four-tier matching criteria start from strict matching, both surface form and direct translation using Helsinki NLP Research Group’s `opus-mt-de-en`<sup>22</sup>. Levenshtein edit distances are computed on the remaining unmatched terms. Lookups on the last tier use cosine distance of SBERT embeddings<sup>23</sup>

**Table 5**

Automatic alignment of `dc:type` to FaBiO subclasses from the corpus.

Tier	Method	<i>n</i>	%	Example
1	Strict (original)	7	1.9	<i>Index</i> → <code>fabio:Index</code>
2	Strict (translated)	20	5.6	<i>Abschnitt</i> <sup>“Section”</sup> → <code>doco:Section</code>
3	Edit distance ( $\text{conf} \geq 88\%$ )	1	0.3	<i>Tonträger</i> <sup>“Sound recordings”</sup> → <code>fabio:SoundRecording</code> ( $\text{conf} = 0.93$ )
4	Sentence embeddings	86	24.0	<i>Abschnitt</i> ( <i>Publikation</i> ) <sup>“Section (publication)”</sup> → <code>doco:Section</code> ( $\text{conf} = 0.70$ )
<b>Total matched</b>		<b>114</b>	<b>31.8</b>	
Unmatched		245	68.2	Non-bibliographic objects

<sup>22</sup>The model is purpose-built for short DE→EN noun-phrase style terms, can run locally without requiring GPU access, <https://huggingface.co/Helsinki-NLP/opus-mt-de-en>

<sup>23</sup><https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>