



SACHSEN-ANHALT

Landesarchiv

Bestandserhaltungsplanung von digitalem Archivgut mittels Verknüpfung von semantischen Technologien und KI

Björn Steffenhagen M. A. mult.

Landesarchiv Sachsen-Anhalt, Abt. 1 Zentrale Dienste

29. Archivwissenschaftliches Kolloquium

20.05.2025 in Marburg



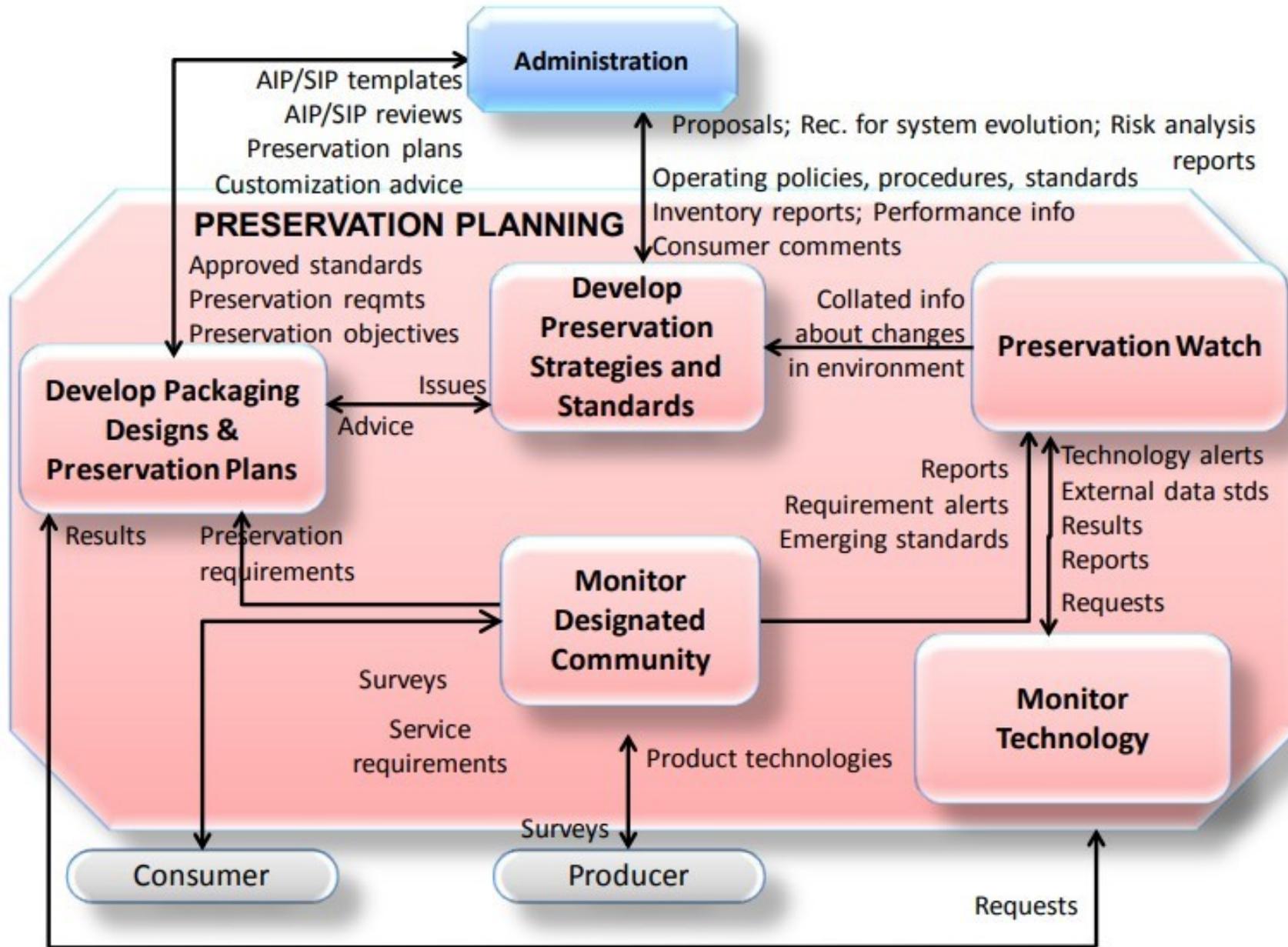
Hintergrund

1. Herausforderung dig. Bestandserhaltung in Archiven
2. KI: Erwartungen und Probleme
3. Semantische Technologien
4. Projekt PROMETHEUS
5. Schlussfolgerungen

Ziele der digitalen Bestandserhaltung im Archivwesen

Sind u. a.:

- Erhalt der Integrität und Authentizität
- Nutz-, Verfüg-/und Verstehbarkeit erhalten
- Dokumentation des Verwaltungshandelns (Provenienz und Archiv)
 - >Vertrauenswürdigkeit herstellen
 - >evidenzbasierte Auswertungen ermöglichen





Herausforderungen der digitalen Bestandserhaltung

- Masse an Primär-/und Metadaten nur noch (teil-)automatisiert verarbeitbar
- Eingriffe durch automatisierte Erhaltungsmaßnahmen bergen Fehlerrisiko
- fragmentierte und heterogene Informationen erschweren konsistente Erhaltungsentscheidungen
- fehlende archivische Bewertungskriterien für Erhaltungsbedarf und Maßnahmenpriorisierung -> technischer und fachlicher Kontext eines Informationsobjekts



Herausforderungen der digitalen Bestandserhaltung - Beispiele

- Qualität der PRONOM-Datenbank der National Archives UK zur Dateiformaterkennung
- fragmentiertes Wissen: PRONOM, Wikidata for DP, COPTR, LoC, KOST, Empfehlungen staatlicher Archive etc.
- Umgang mit Ergebnissen der Dateiformatvalidierung
- Migration im Kontext: Ist eine JPEG-Datei in einer E-Akte gleich wie in einer Webseite?
- Priorisierung: Fehlende Bewertungskriterien bei Bestandserhaltungsmaßnahmen



Probleme der KI bzw. LLM

Automatisierung ist verlockend, Probleme aber u.a.:

- Halluzinationen
 - unzuverlässige Entscheidungen
 - Black-Box-Problematik – mangelnde Transparenz
 - fehlendes, domänenspezifisches Wissen
- ➔ Archive brauchen jederzeit reproduzierbare, verlässliche und steuerbare Methoden
- ➔ Diese durch den alleinigen Einsatz von KI (noch) nicht erreichbar



Semantische Technologien

- Sollen maschinenlesbare Modelle mit Bedeutung schaffen, um Daten automatisiert zu verknüpfen, anzureichern und qualitativ zu bewerten.
- Explizite, formale Repräsentation von Konzepten + Relationen
-> Web of (meaningful) Data
- 3 wesentliche Technologien:
 - Ontologie
 - Wissensgraph
 - SPARQL

premis-3-0-0 (https://id.loc.gov/ontologies/premis-3-0-0)

dcterms:Policy > Preservation policy > Significant properties

Active ontology x Entities x Individuals by class x DL Query x

Classes Object properties Data properties Annotation properties Datatypes Individuals

Class hierarchy: Significant properties

- owl:Thing
 - Action
 - Activity
 - Event
 - Agent (prov:Agent)
 - Agent
 - Organization (prov:Organization)
 - Person (prov:Person)
 - SoftwareAgent
 - dcterms:FileFormat
 - dcterms:Policy
 - Preservation policy
 - Significant properties
 - Dependency
 - Entity
 - Bundle
 - Collection
 - Object
 - Plan
 - Environment characteristic
 - Fixity
 - foaf:Agent
 - Agent
 - foaf:Organization
 - Organization
 - foaf:Person
 - Person
 - Identifier
 - Influence
 - ActivityInfluence
 - AgentInfluence
 - EntityInfluence
 - Inhibitor

Asserted

Significant properties — http://www.loc.gov/premis/rdf/v3/SignificantProperties

Annotations Usage

Annotations: Significant properties

Annotations +

rdfs:label [language: en]
Significant properties

rdfs:comment [language: en]
Characteristics of a particular object subjectively determined to be important to maintain through preservation actions.

rdfs:isDefinedBy
<http://www.loc.gov/premis/rdf/v3/>

Description: Significant properties

Equivalent To +

SubClass Of +

- 'Preservation policy'

General class axioms +

SubClass Of (Anonymous Ancestor)

Instances +

Target for Key +

Disjoint With +

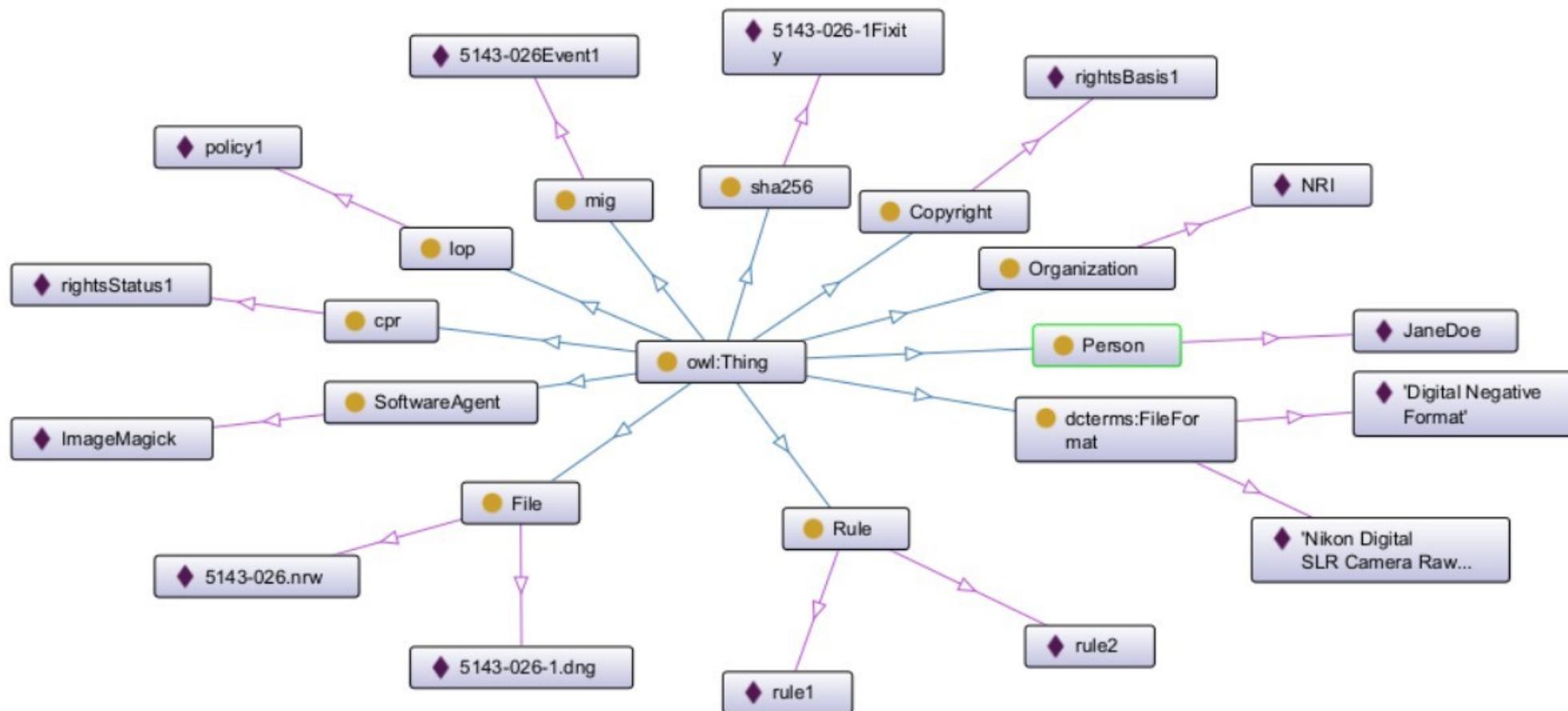
Disjoint Union Of +

Search:

contains

Search

Clear



INHALT

hiv





Probleme der semantischen Technologien

- streng formeller Modellierungsprozess
- Pflegeaufwand
- Expertensysteme für Fachdomänen
- Abfragemöglichkeiten via SPARQL unkomfortabel
- in Kultureinrichtungen (noch) ungenügend etabliert
- Beispiel: Archivführer Deutsche Kolonialgeschichte:
<https://archivfuehrer-kolonialzeit.de/index.php/>



Lösungsansatz: Kombination aus semantischen Technologien und KI

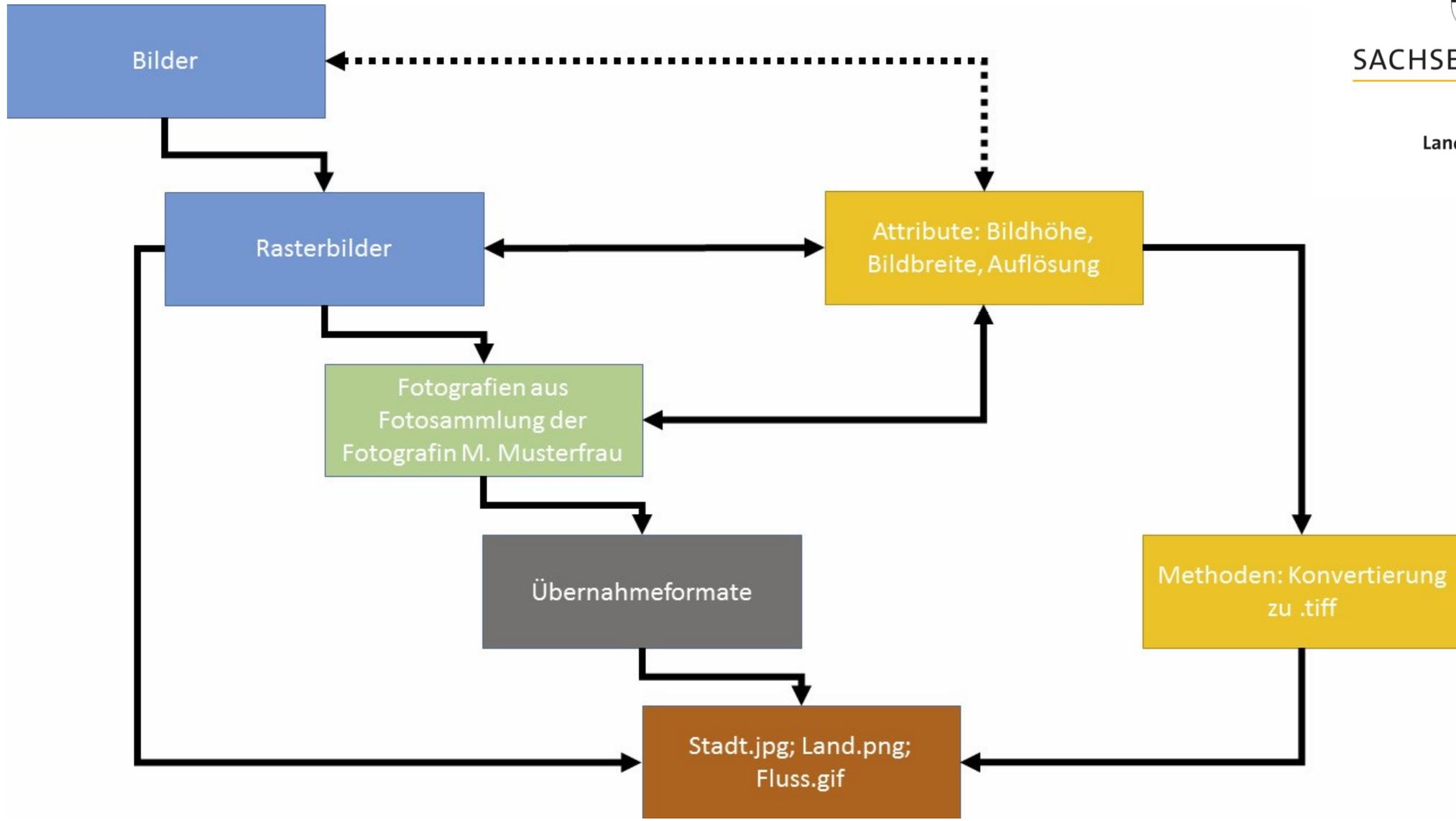
- Ontologie und Graphen als Wissensbasis
- KI-Modelle als
 - Abfragemöglichkeit
 - Möglichkeit der Entscheidungsfindung
 - Erweiterung des Gesamtmodells
- Ergebnis bildet ein RAG-System (Retrieval-Augmented Generation)



Projekt PROMETHEUS

PROMETHEUS – **P**reservation of digital **o**bjects and **m**etadata: encoding and decoding **t**he digital **e**nvironment under the **u**se of ontologies

- Kernbereich Ontologie
 - darauf aufbauend hybrid-reasoning-basierte Entscheidungslogik für risikoorientierte Priorisierung
 - Klassifikation und Priorisierung von Erhaltungsmaßnahmen anhand von technischen und archivischen Merkmalen
- >Einsatz von KI für die Risikoplanung und Preservation Watch
- >Bildung eines einheitlichen Modells zur Bestandserhaltungsplanung





PROMETHEUS – Stand und Ausblick

- Modellierung der Ontologie
- Auswahl von KI-Modellen zur Anbindung an Ontologie
- Fachkonzept

- dauerhafte Verknüpfung mit externen Informationsquellen oder eigenes Vorhalten von Informationen?
- Erweiterung zur interaktiven Ontologieabfrage über einen Chatbot?
- GUI?



Schlussfolgerungen

1. Digitale Bestandserhaltung ist ein Big-Data-Problem
 2. KI im Archiv benötigt eine qualitativ hochwertige Wissensbasis
 3. Ontologien und Wissensgraphen eignen sich dafür
 4. Die Verknüpfung von KI und Ontologien fördert die Stärken beider Ansätze
- >Kein entweder-oder der Technologien, sondern ein sowohl-als-auch
5. Kooperationen zwischen Archiven werden noch relevanter



Kontakt

Björn Steffenhagen

Abteilung 1 - Zentrale Dienste

Landesarchiv Sachsen-Anhalt

bjoern.steffenhagen@sachsen-anhalt.de

<https://orcid.org/0000-0002-5849-7123>