

Harald Sack

Hybride Künstliche Intelligenz in der automatisierten Inhaltserschließung

1 Einleitung

Effizienter (Online-)Zugang zu Bibliotheks- und Archivmaterialien erfordert eine qualitativ hinreichende inhaltliche Erschließung dieser Dokumente. Die passgenaue Verschlagwortung und Kategorisierung dieser unstrukturierten Dokumente ermöglichen einen strukturell gegliederten Zugang sowohl in der analogen als auch in der digitalen Welt. Darüber hinaus erweitert eine vollständige Transkription der Dokumente den Zugang über die Möglichkeiten der Volltextsuche. Angesichts der in jüngster Zeit erzielten spektakulären Erfolge der Künstlichen Intelligenz liegt die Schlussfolgerung nahe, dass auch das Problem der automatisierten Inhaltserschließung für Bibliotheken und Archive als mehr oder weniger gelöst anzusehen wäre. Allerdings lassen sich die oftmals nur in thematisch engen Teilbereichen erzielten Erfolge nicht immer problemlos verallgemeinern oder in einen neuen Kontext übertragen. Das Ziel der vorliegenden Darstellung liegt in der Diskussion des aktuellen Stands der Technik der automatisierten inhaltlichen Erschließung anhand ausgewählter Beispiele sowie möglicher Fortschritte und Prognosen basierend auf aktuellen Entwicklungen des maschinellen Lernens und der Künstlichen Intelligenz einschließlich deren Kritik.

2 Der Siegeszug des maschinellen Lernens

Wenn heute von Künstlicher Intelligenz (KI) die Rede ist, wird damit maschinelles Lernen und im Speziellen meist Deep Learning als spezifischer, aber höchst erfolgreicher Teilaspekt dieser Disziplin der Informatik adressiert. Allerdings geht mit Deep-Learning-Technologien allgemein eine besondere Faszination einher, orientieren diese sich doch am Vorbild von Denk- und Lernprozessen im menschlichen Gehirn. Von Anfang an war die technologische Entwicklung Künstlicher Intelligenz mit überzogenen Erwartungen verbunden. Bereits 1943 legten Warren McCulloch und Walter Pitts das einfache mathematische Modell eines künstlichen Neurons als grundlegendes Schaltelement im Gehirn vor und zeigten, dass sich über die Vernetzung dieser Neuronen beliebige (Turing-

berechenbare) Funktionen implementieren ließen (McCulloch und Pitts 1943). Nachdem der kanadische Psychologe Donald O. Hebb die neurophysiologischen Grundlagen des menschlichen Lernens anhand der Veränderungen der synaptischen Übertragung zwischen Neuronen im Gehirn erklären konnte, lag die Modellierung dieser Lernaktivitäten in einem mathematischen Modell auf der Hand (Hebb 1949). Frank Rosenblatt, ebenfalls Psychologe, formulierte 1958 als erster das so genannte Perzeptron-Lernmodell, das bis heute die Grundlage aller künstlicher neuronaler Netzwerke bildet: Verbindungsgewichte zwischen Eingabe und Schaltelement werden anhand der Differenz von Soll- und Istwerten schrittweise gelernt (Rosenblatt 1958). Dieses einfache Modell ermöglichte bereits Ende der 1950er Jahre die automatische Erkennung handgeschriebener Postleitzahlen. Die aus den ersten Erfolgen des maschinellen Lernens resultierenden euphorischen Erwartungen, dass Computer binnen weniger Jahre die Stufe echter menschlicher Intelligenz erreichen und auch übertreffen würden, stießen jedoch schnell an ihre Grenzen. Marvin Minsky und Seymour Papert wiesen 1969 nach, dass das Perzeptron-Modell nicht in der Lage ist eine einfache binäre XOR-Funktion (exklusives Oder) zu berechnen (Minsky und Papert 1969). Dies führte zu einem Stillstand in der Forschung der künstlichen neuronalen Netze und resultierte im ersten sogenannten „AI-Winter“, der bis zur Mitte der 1980er Jahre andauerte. 1986 legten Rummelhart, Hinton und Williams (Rummelhardt et al. 1986) sowie unabhängig davon David B. Parker (Parker 1986) und Yann LeCun (LeCun 1985) dar, wie ein mehrlagiges neuronales Netzwerk mit Hilfe des Backpropagation-Algorithmus auch sehr komplexe Funktionen erlernen kann. Die Grundidee besteht darin, den Fehler, den das künstliche neuronale Netzwerk bei der Berechnung einer Ausgabe macht, schrittweise rückwärts von der Ausgabe- zur Eingabeschicht weiterzureichen und zur Gewichtsveränderung zu verwenden. 1990 wurden von Jeffrey L. Elman erstmals rückgekoppelte neuronale Netzwerke vorgeschlagen, mit denen sich besonders gut sequenzielle Daten wie z. B. Zeitreihen erlernen lassen (Elman 1990). Eine weitere Verbesserung bestand 1997 in der Einführung einer Art Kurzzeitgedächtnis über sogenannte Long-Short-Term-Memory-Module durch Sepp Hochreiter und Jürgen Schmidhuber (Hochreiter und Schmidhuber 1997). Schließlich wurden 1998 von Yann LeCun die heute dominierenden *Convolutional Neural Networks* vorgeschlagen, die es vor allem in der Bildverarbeitung ermöglichten, das vormals aufwendige und oftmals manuelle Interventionen erfordernde Feature Engineering direkt in den Lernalgorithmus zu integrieren (LeCun 1998). Die künstlichen neuronalen Netzwerke wurden seit Anfang der 2000er Jahre stets komplexer, verbunden mit zahlreichen neu auftretenden Problemen. Erst mit zunehmender Verfügbarkeit billiger, hochparalleler Rechnerkerne (Graphical Processing Units, GPU), wie sie hauptsächlich in der Compu-

tergrafik zum Einsatz kamen, wurde das Deep Learning, also die Verarbeitung vielschichtiger und komplexer künstlicher neuronaler Netzwerke, handhabbar. Weitere Erfolgsfaktoren bestanden in der jetzt vorliegenden Verfügbarkeit extrem großer Trainingsdatensätze, wie sie z. B. über das World Wide Web und seine Social-Media-Plattformen gesammelt werden konnten, sowie die Wiederverwendbarkeit und Anpassbarkeit bereits vortrainierter künstlicher neuronaler Netzwerke an neue Problemstellungen (*Transfer Learning*). Besondere Aufmerksamkeit erlangte im März 2016 das Deep-Learning-basierte System AlphaGo von Google DeepMind, dem es gelang, den südkoreanischen Profispieler Lee Sedol unter Turnierbedingungen im Brettspiel Go zu schlagen. Der Folgeversion AlphaGo Zero, mit keinerlei Vorwissen über das Spiel ausgestattet, sondern ausschließlich mit den Spielregeln und durch Spiele gegen sich selbst trainiert, gelang es nach nur drei Tagen Training die AlphaGo-Version, die Lee Sedol besiegen konnte, mit dem Ergebnis 100:1 zu schlagen (Silver et al. 2017).

Typische Einsatzgebiete von Deep-Learning-Technologien sind heute vor allem die Bilderkennung (*Visual Analysis*) sowie die Verarbeitung natürlicher Sprache (*Natural Language Processing – NLP*). Die meisten Probleme lassen sich dabei auf eine Klassifikation zurückführen, d. h. das künstliche neuronale Netzwerk entscheidet, ob eine Eingabe zu einer bestimmten Klasse gehört oder nicht. Übertragen auf die Inhaltserschließung könnte beispielsweise entschieden werden, ob ein bestimmtes Schlüsselwort für ein Dokument vergeben werden soll oder nicht. Neben diesen traditionellen Klassifikationsaufgaben gewannen seit 2015 insbesondere künstlich generierte wirklichkeitsgetreue Bild-, Video-, Musik- oder Textdokumente große öffentliche Popularität. 2014 wurden von Ian Goodfellow et al. *Generative Adversarial Networks* (GANs) eingeführt, die über vergleichendes Lernen trainiert werden (Goodfellow et al. 2014). GANs setzen sich aus einem generativen Netzwerk und einem diskriminativen Netzwerk zusammen. Während das generative Netzwerk neue künstliche Dokumente erzeugt, lernt das diskriminative Netzwerk, die künstlich generierten Dokumente von realen Dokumenten zu unterscheiden. Das Ziel besteht dann darin, dass das generative Netzwerk künstliche Dokumente erzeugen kann, die das diskriminative Netzwerk nicht mehr von originalen Dokumenten unterscheiden kann.

Besonders beeindruckt auch die Leistungsfähigkeit aktueller, auf Deep-Learning-Technologien beruhender statistischer Sprachmodelle. Statistische Sprachmodelle spiegeln die kontextabhängige Häufigkeitsverteilung bestimmter Wortfolgen im Gebrauch einer Sprache wider und können sowohl zur Sprachanalyse als auch zur Generierung von Textinhalten eingesetzt werden. Die von OpenAI entwickelten generativen Sprachmodelle GPT-2 und GPT-3 sind heute in der Lage, natürlichsprachliche Texte zu vorgegebenen Themen von

solcher Qualität zu generieren, dass deren Unterscheidung von manuell erstellten Texten kaum mehr möglich ist (Brown et al. 2020).

3 Symbolische und subsymbolische Wissensrepräsentation

Wie bereits erwähnt, stellen die heute so populären Deep-Learning-Technologien lediglich ein auf statistischem Lernen basierendes Teilgebiet der Künstlichen Intelligenz dar. Das in einem künstlichen neuronalen Netzwerk gespeicherte Wissen liegt in Form von Kantengewichten zwischen den einzelnen Neuronen vor. Diese implizite, subsymbolische Form der Repräsentation von Wissen lässt sich nur sehr schwer wieder in eine explizite, symbolische und damit nachvollziehbare Repräsentation umwandeln. Symbolische Wissensrepräsentationen dagegen setzen auf ein Kalkül, d. h. ein formales System von Regeln, mit denen sich aus gegebenen Aussagen (Axiomen) weitere Aussagen ableiten lassen. Verbunden mit einer formalen Interpretation lässt sich die Bedeutung (Semantik) dieser Aussagen ableiten. Semantic-Web-Technologien (Berners-Lee et al. 2001) stellen eine aktuelle Form der Wissensrepräsentation auf der formalen Basis von Beschreibungslogiken dar. Als Austauschformat fungiert dabei das Resource Description Format (RDF), das Aussagen in Form einfacher Tripel (Subjekt, Prädikat, Objekt) kodiert und Entitäten über webbasierte Uniform Resource Identifier (URIs) adressiert und identifiziert. Komplexere Semantik lässt sich deskriptiv über die Web Ontology Language (OWL) oder über logische Regeln abbilden. In ihrer semantisch leichtgewichtigen Variante kommen diese Technologien heute im Zuge der Inhaltserschließung als Linked (Open) Data zum Einsatz (Bizer et al. 2009), d. h. Metadaten zu Bibliotheks- und Archivressourcen werden im Web in RDF-kodierter Form zur Verfügung gestellt. Diese lassen sich zur inhaltlichen Verknüpfung unterschiedlicher Datenquellen und damit zur Anreicherung der eigenen Datenbestände ausnutzen. Größere Popularität konnten semantische Technologien erstmals 2012 mit dem Google Knowledge Graph erzielen, der als Grundlage der Websuchmaschine Google zum Einsatz kommt und die Qualität der erzielten Suchergebnisse signifikant verbesserte (Singhal 2012).

Die effiziente Verknüpfung und gemeinsame Nutzung symbolischer und subsymbolischer Wissensrepräsentationen als hybride KI stellt ein aktuelles Forschungsproblem dar. Deep-Learning-Verfahren kommen häufig in der Analyse unstrukturierter Informationen wie z. B. bei Text- und Bilddokumenten

zum Einsatz. Named Entity Recognition und Named Entity Linking erkennt bedeutungstragende Entitäten in natürlichsprachlichen Dokumenten und verknüpft diese mit ihren korrespondierenden Repräsentationen in einem Wissensgraphen (*Knowledge Graph*) bzw. einer Wissensbasis. Moderne Bildklassifikation und Objekterkennung identifiziert in Bilddokumenten dargestellte Objekte und verknüpft diese ebenfalls mit Entitäten oder Klassen eines Wissensgraphen. Verfahren zur Relationsextraktion ermitteln die Verknüpfung bereits identifizierter Entitäten untereinander und können derart zum Aufbau von Wissensgraphen genutzt werden. Von besonderem Interesse sind aktuell Technologien, die eine Abbildung formaler Wissensrepräsentationen in einen Vektorraum ermöglichen, wobei die Semantik der abgebildeten Entitäten mit deren Positionen im Vektorraum korrespondiert (Ristoski und Paulheim 2017). Mit Hilfe dieser Vektorraumeinbettungen (*Embeddings*) lassen sich semantisch ähnliche oder verwandte Entitäten über den Abstand ihrer korrespondierenden Vektoren bzw. über die Anwendung einfacher Vektorraumarithmetik sehr effizient handhaben. Dadurch erschließen sich neue Wege der semantischen und explorativen Suche in großen Dokumentenbeständen sowie die Realisierung einfacher inhaltsbasierter Empfehlungssysteme. Im folgenden Kapitel werden einige ausgewählte Anwendungen symbolischer und subsymbolischer Wissensrepräsentation im Kontext der Inhaltserschließung vorgestellt.

4 Status Quo der Inhaltserschließung an ausgewählten Beispielen

4.1 Automatisierte Verschlagwortung und Klassifizierung

Grundlegend für die Inhaltserschließung unstrukturierter Text- oder Bilddokumente ist die Verschlagwortung oder Klassifizierung. Während bei der Verschlagwortung ein oder mehrere inhaltlich korrespondierende deskriptive Schlagwörter vornehmlich aus einem kontrollierten Vokabular ausgewählt und als Metadaten dem Dokument zugewiesen werden, wird das Dokument in der Klassifizierung einer inhaltlich korrespondierenden Klasse oder Kategorie zugeordnet. Dabei können die Kategorien in hierarchischer Beziehung zueinander stehen.

Ein aktuelles Projekt, in dem neben der Schaffung themenbasierter Zugänge auch eine automatisierte Klassifizierung von Archivdokumenten durchgeführt wird, ist das DFG-geförderte Projekt *Aufbau einer Infrastruktur zur*

*Implementierung sachthematischer Zugänge im Archivportal-D am Beispiel des Themenkomplexes Weimarer Republik (2018–2021).*¹ Datengrundlage des Projekts stellen digitalisierte Archivalien aus Ministerien, öffentlichen Einrichtungen, Körperschaften und Nachlässen im Bestand des Bundesarchivs und des Landesarchivs Baden-Württemberg aus der Zeit der Weimarer Republik, der ersten deutschen Demokratie, dar, die die Themenbereiche Politik, Wirtschaft, Gesellschaft und Alltag in Deutschland zwischen 1918 und 1933 einschließen. Insgesamt umfasst die Sammlung aktuell 21042 digitalisierte Dokumente, die anhand eines 881 Schlagwörter beinhaltenden Klassifikationssystems (Systematik) klassifiziert werden müssen. Diese Systematik unterteilt sich in 12 Kategorien und 121 Unterkategorien, denen die Schlagwörter zugeordnet werden. Die Klassierung erfolgt manuell, wird aber durch automatisierte Verfahren über die Empfehlung von zu verwendenden Schlagwörtern unterstützt. Traditionelle automatisierte Klassifikationsverfahren basieren auf Algorithmen des überwachten maschinellen Lernens und benötigen eine hinreichend große Menge an Trainingsdaten für jede einzelne Klasse bzw. jedes Schlagwort. Zudem kann jedes Dokument mehreren Klassen gleichzeitig zugeordnet werden.

Der in diesem Projekt verfolgte Ansatz zur automatisierten Klassifizierung macht sich die semantische Ähnlichkeit zwischen den verfügbaren Metadaten (Titel, Kurzbeschreibungen, Angaben zur Ablageorganisation im Archiv) und den Bezeichnungen der Klassen zu Nutze. Erschwerend kommt hinzu, dass Kurzbeschreibungen kaum vorhanden und Titel nicht immer aussagekräftig sind. Auch existiert keine Transkription der oftmals handschriftlichen Dokumente. Als Lösung des Problems wurden die verfügbaren Textdaten über *Word Embeddings* kodiert und die semantische Ähnlichkeit zwischen Dokumenten und Schlagwörtern über den Winkelabstand der jeweiligen korrespondierenden Vektoren bestimmt, die die Zuordnung von Schlagwörtern, gesteuert durch einen zuvor festgelegten Schwellwert, gestatten (Hoppe et al. 2020). Die mit diesem Verfahren erreichte Genauigkeit der Verschlagwortung ist leider unbefriedigend, erspart oftmals aber als Vorschlagsmechanismus in einem manuellen Verfahren Zeit bei der Suche nach dem passenden Schlagwort. In diesem Zusammenhang werden aktuell Möglichkeiten des sogenannten *Zero-Shot Learning* und des unüberwachten Lernens evaluiert, die unter Einbeziehung externer Wissensquellen eine Verbesserung der erzielten Qualität der Verfahren versprechen.

Ausschlaggebend für die Effektivität eines Klassifikationsverfahrens unter den geschilderten Rahmenbedingungen ist die Herstellung von Kontext. Ist die

¹ Projektwebseite *Sachthematische Zugänge zum Archivportal-D*, <https://www.landearchiv-bw.de/de/landearchiv/projekte/sachthematische-zugaenge-im-archivportal-d/63525> (13.1.2021).

ser Kontext in den vorhandenen Daten nur begrenzt verfügbar, müssen externe Wissensressourcen herangezogen werden, vergleichbar mit Hintergrundwissen oder Expertenwissen, das eine:n erfahrene:n Archivar:in oder Bibliothekar:in in die Lage versetzt, eine treffende Klassifizierung zu finden. Dieses externe Wissen wird im vorliegenden Projekt z. B. über die vortrainierten Word Embeddings eingebracht, die auf der Grundlage der Texte der Online-Enzyklopädie Wikipedia trainiert wurden und damit Teile des darin enthaltenen Wissens reflektieren. Um die Qualität dieser Klassifizierungsverfahren weiter zu steigern, können weitere externe Wissensressourcen wie z. B. der der Wikipedia zugrundeliegende Hyperlink-Graph oder alternative Wissensgraphen wie z. B. DBpedia oder Wikidata verwendet werden (Türker 2020).

4.2 Wissensrepräsentationen und zeitliche Dynamik

Das im vorangehenden Abschnitt beschriebene Projekt zur Schaffung sachbezogener Zugänge zu einem Archiv beinhaltet die Entwicklung einer bestandsbezogenen Systematik zur Klassifizierung der die Sammlung umfassenden Archivdokumente. Im Zuge des Projektfortschritts unterlag diese Systematik zahlreichen Veränderungen. Weitere Veränderungen der Systematik sind mit der Erweiterung des damit verwalteten Dokumentenbestands zu erwarten. Zu diesem Zweck müssen automatisierte Workflows entwickelt werden, die dieser Dynamik Rechnung tragen und eventuell obsoletere Klassifizierungen an die neue Systematik anpassen. Zu diesem Zweck wurde die Archive Dynamics Ontology (ArDO) entwickelt, die genau diese dynamischen Veränderungen innerhalb einer hierarchisch strukturierten Systematik, verbunden mit Archivdokumenten, die ihrerseits einem anderen hierarchisch organisierten Ablagesystem folgen, abbilden und nachhalten kann (Vsesviatska et al. 2021). Jede Veränderung der zur Klassifizierung notwendigen Systematik wird in einem den Archivbestand repräsentierenden Knowledge Graph entsprechend der ArDO-Ontologie nachgehalten und kann zur automatischen Reklassifizierung gemäß der neuen, veränderten Systematik verwendet werden. Zudem lassen sich eventuell vorhandene Kopien der Bestände, die noch mit einer alten Version der Systematik klassifiziert wurden, auf den aktuellen Stand bringen.

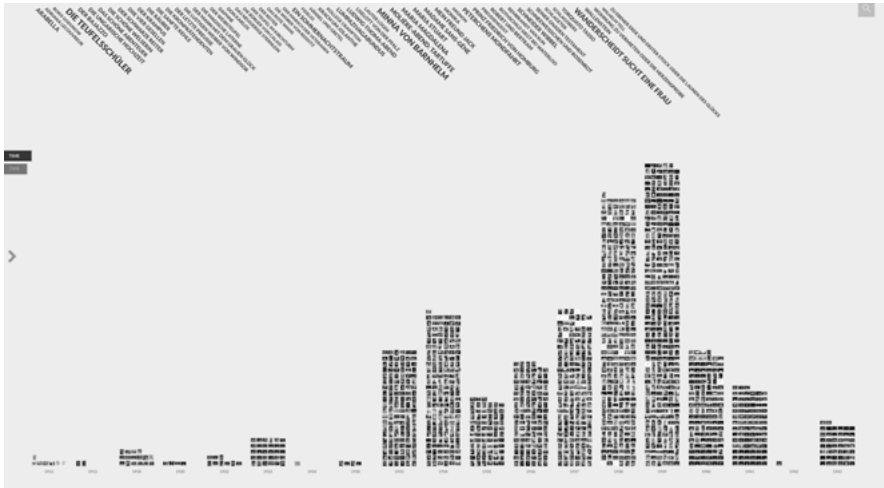


Abb. 1: Visualisierung des gesamten Bilddatenbestands in Linked Stage Graph

Die im vorangegangenen Abschnitt genannten Wissensgraphen bilden bereits vielfach die Grundlage moderner Informations- und Suchsysteme. In ihnen werden Metadaten auf der Basis von Semantic-Web-Technologien miteinander zu einem Graphen verknüpft, der unter anderem eine semantische Suche ermöglicht (siehe Abschnitt 4.1). Beispiel für einen weiteren Wissensgraphen ist der im Rahmen der Cod1ng-Da-V1nc1-Initiative² entstandene Linked Stage Graph (Tietz 2019). Grundlage des Linked Stage Graph sind die Bilddatenbestände *Hof-/Landes-/Staatstheater Stuttgart: Fotos und Graphiken* mit ca. 7 000 historischen Fotografien aus dem Zeitraum von 1896 bis 1942. Die zugehörigen Auführungsmetadaten wurden in einen Wissensgraphen überführt und mit bestehenden externen Wissensressourcen wie Wikidata und der Gemeinsamen Normdatei (GND) der Deutschen Nationalbibliothek verknüpft. Der Linked Stage Graph kann sowohl über einen öffentlich zugänglichen SPARQL-Endpunkt direkt abgefragt werden bzw. wird der Bilddatenbestand in verschiedenen Visualisierungen zugänglich gemacht, wobei unterschiedliche strukturelle Merkmale und Ähnlichkeiten die jeweilige Visualisierung steuern (vgl. Abb. 1).

² Cod1ng-Da-V1nc1-Webseite: <https://codingdavinci.de/> (14.1.2021).

4.3 Automatisiertes Transkribieren historischer Dokumente

Traditionell bestimmen analoge Dokumente die Bestände von Archiven und Bibliotheken. Deren automatisierte inhaltliche Erschließung geht oft mit der Digitalisierung der Bestände einher. Allerdings erfordert die inhaltliche Erschließung eine Aufbereitung der in den Digitalisaten enthaltenen Informationen, d. h. Textdokumente müssen via OCR (*Optical Character Recognition*) transkribiert werden, Textinhalte in Bildern oder audiovisuellen Dokumenten müssen ebenfalls via OCR erkannt und transkribiert werden, Tondokumente bzw. die Tonspuren audiovisueller Dokumente werden via Speech2Text transkribiert. Multimediale Dokumente können über zahlreiche weitere automatische Analyseverfahren inhaltlich erschlossen werden. In diesem Abschnitt soll lediglich ein Projekt vorgestellt werden, das unter anderem das Transkribieren historischer Textdokumente beinhaltet.

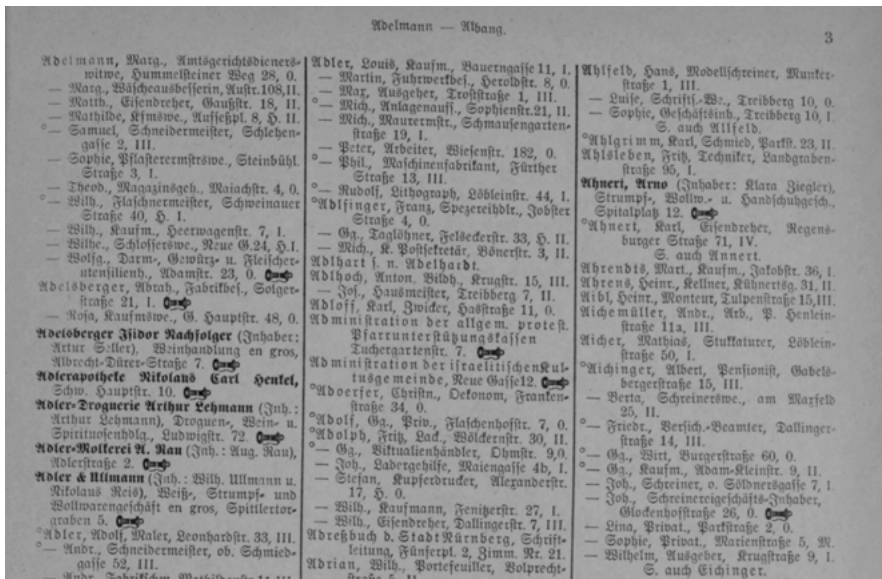


Abb. 2: Beispielseite aus dem Adressbuch für Nürnberg von 1892

Das Projekt TOPORAZ (Nürnberger Topographie in Raum und Zeit) verknüpft ein wissenschaftlich fundiertes 3D-Modell des Nürnberger Hauptmarktes in unterschiedlichen Zeitstufen unmittelbar mit verfügbaren Quellen zu Personen, Ereignissen und Bauten wie z. B. historische Fotos, Zeichnungen, Grafiken und Pläne sowie historische Schriftquellen wie Urkunden, Adressbücher und

Chroniken (Razum et al. 2020). Das Folgeprojekt TRANSRAZ greift die konzeptionellen Vorarbeiten aus TOPORAZ auf und erweitert das Untersuchungsgebiet auf die gesamte Nürnberger Altstadt innerhalb der letzten Stadtumwallung (Tietz et al. 2021). Das Ziel ist die Verknüpfung heterogener historischer Quellen mit dem 3D-Stadtmodell über einen Wissensgraphen. Problematisch ist dabei das automatische Transkribieren handschriftlicher Quellen bzw. Quellen mit historischer Typografie und ausgefallener Strukturierung. Abbildung 2 zeigt eine Beispielseite des Nürnberger Adressbuchs von 1892. Die automatisierte Erfassung der Adressdaten umfasst als erstes die korrekte Auflösung einzelner zusammengehöriger Segmente, dann deren Transkribierung und abschließend eine aufwendige Nachbereitung. Diese besteht in einem Abgleich der Daten mit Referenzlisten von Eigennamen, Straßennamen und historischen Berufsbezeichnungen, einschließlich der Auflösung historisch üblicher Abkürzungsvarianten sowie der Bestimmung von grafischen Sonderzeichen und typografischen Auszeichnungen, die einer eigenen Semantik unterliegen. Weitere historische Quellen erfordern zusätzlich die Analyse natürlichsprachlicher Texte einschließlich der Extraktion von Personendaten, topografischer Daten sowie zeitbezogener Ereignisdaten. Diese Informationen werden semantisch annotiert, in einem manuell erweiterbaren Wissensgraphen verfügbar gemacht und mit externen Wissensressourcen verknüpft. Alle Informationen des TRANSRAZ Wissensgraphen werden direkt mit den topografischen Informationen des 3D-Stadtmodells verknüpft und über dessen grafische Schnittstelle zugänglich gemacht (vgl. Abb. 3).

4.4 Automatisierte Erzeugung und Befüllung von Wissensgraphen

Die in den vorangegangenen Abschnitten vorgestellten Projekte setzen alle einen Knowledge Graph als zentrales Bindeglied zwischen Daten, Dokumenten und externen Informationsquellen ein, um auf dieser Basis ein effizientes Informationssystem zur Verfügung zu stellen, das sowohl für den Endnutzenden als auch für einen automatisierten Zugriff zugänglich ist. Zum Aufbau solcher Wissensgraphen müssen unstrukturierte Dokumente und Daten semantisch analysiert werden, um deren Inhalte auf korrespondierende Ontologien und Vokabulare abzubilden. Während in Bilddokumenten über die visuelle Analyse Objekte erkannt werden können, die sich direkt auf semantische Entitäten beziehen und deren räumliche Beziehung zueinander aus den Bilddokumenten ermittelt werden kann (vgl. Abb. 4), müssen Textdokumente zunächst einer sprachlichen Analyse (NLP) unterzogen werden.

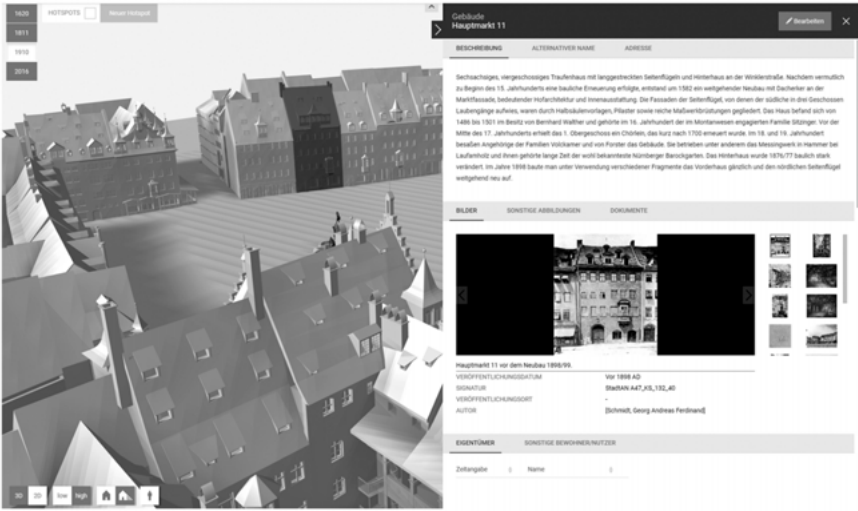


Abb. 3: Ausschnitt des 3D-Stadtmodells des historischen Nürnberger Hauptmarkts mit zusätzlichen topografisch verknüpften Informationen

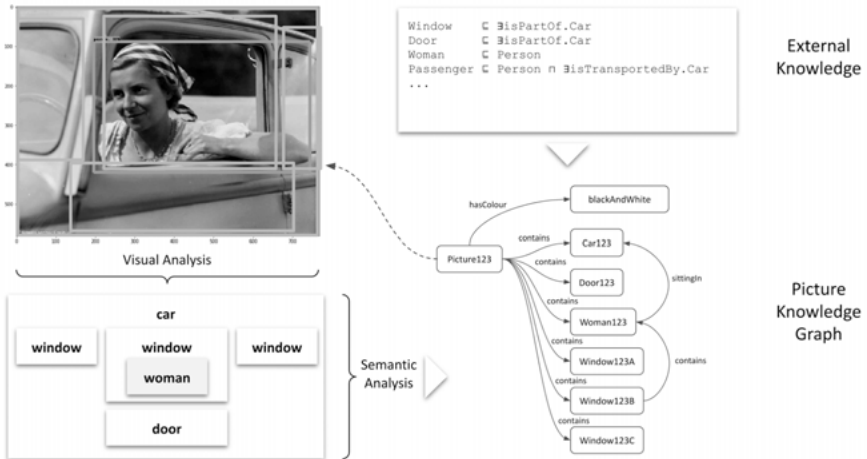


Abb. 4: Visuelle Analyse bestimmt die Bildobjekte und deren Beziehung zueinander im Bild. Mit Hilfe externen Wissens wird daraus ein Knowledge Graph generiert.

Beispiel für die Erzeugung eines Knowledge Graph aus einer Sammlung wissenschaftlicher Arbeiten ist der Artificial Intelligence Knowledge Graph (Dessì et al. 2020). Als Grundlage dienten mehr als 300 000 wissenschaftliche Arbeiten im Bereich der Künstlichen Intelligenz aus dem Microsoft Academic Graph (MAG). Mit Hilfe einer komplexen Verarbeitungspipeline, die statistische und linguistische Deep-Learning-basierte Verfahren zur Identifikation und Extraktion von bedeutungstragenden Entitäten und zugehörigen Relationen kombiniert, konnte ein Knowledge Graph erzeugt werden, der bislang 14 Millionen Fakten über 820 000 extrahierte Entitäten mittels 27 Relationen miteinander und mit weiteren externen Informationsressourcen verknüpft. Zum einen lässt sich der erzeugte Knowledge Graph zur Implementierung einer semantischen Suche oder eines inhaltsbasierten Empfehlungssystems über die darin enthaltenen wissenschaftlichen Arbeiten verwenden. Darüber hinaus lässt sich die angewandte Technologie auch an andere Wissensgebiete zur Erzeugung neuer Wissensgraphen anpassen, wie aktuell z. B. für COVID-19 oder für mathematische wissenschaftliche Literatur.

5 Effiziente Nutzung inhaltserschlossener Dokumente

Symbolische und subsymbolische Wissensrepräsentation lassen sich zur inhaltlichen Erschließung nutzbringend kombinieren. Multimodale Dokumente, d. h. Texte, Bilder oder audiovisuelle Inhalte können mit Hilfe moderner Deep-Learning-Verfahren inhaltlich analysiert werden, sodass deren Inhalte in Form eines Wissensgraphen repräsentiert werden können. Informationssysteme setzen Wissensgraphen zur Verbesserung des Zugangs zu den in ihnen gespeicherten Informationsinhalten ein. Beispielsweise lässt sich eine semantische Suche implementieren, Informationsinhalte und deren Zusammenhänge untereinander können durch geeignete Visualisierungen deutlich gemacht werden, oder inhaltsbasierte Empfehlungssysteme können die Suche in Informationssystemen hin zu einer explorativen Suche unterstützen (Waitelonis und Sack 2012).

5.1 Semantische Suche

Im Unterschied zur textbasierten Suche konzentriert sich die semantische Suche auf Entitäten, d. h. einerseits werden so sprachliche Mehrdeutigkeiten der tradi-

tionellen Textsuche einfach vermieden beziehungsweise aufgelöst, andererseits können Informationsinhalte auch über Synonyme oder Umschreibungen aufgefunden werden. Auf diese Weise wird eine vollständigere und gleichzeitig auch genauere Suche realisiert. Die im Wissensgraph repräsentierten Entitäten sind mit Ontologien verbunden, über die zugehörige Kategorien und Oberbegriffe ermittelt werden können. Mit diesen lässt sich die ursprüngliche Suche erweitern, so dass zusätzliche relevante Suchergebnisse gewonnen werden können. Eine weitere Möglichkeit besteht in der Erzeugung semantischer Suchfacetten, d. h. das erzielte Suchergebnis lässt sich nach inhaltlich relevanten Kriterien filtern und einschränken. Zusätzlich kann der Wissensgraph auch direkt in das verwendete Suchmaschinenmodell integriert werden. Das zahlreichen Suchmaschinen zugrundeliegende Vektorraummodell lässt sich zu einem generalisierten Vektorraummodell erweitern. Dazu werden dessen Basisvektoren, die vormals z. B. Schlüsselwörter repräsentieren, selbst noch einmal in weitere Vektoren zerlegt, so dass Ähnlichkeitsbeziehungen unter den ursprünglich orthogonalen Basisvektoren direkt im generalisierten Vektorraummodell repräsentiert werden. Auf diese Weise lassen sich sehr effizient semantisch ähnliche Suchergebnisse oder aber auch semantisch miteinander in Bezug stehende Suchergebnisse ermitteln, die die ursprünglich erzielten Suchergebnisse ergänzen (Waitelonis et al. 2016).

5.2 Visualisierung

Ein grundsätzliches Problem der Umsetzung einer semantischen Suche besteht in der Erwartungshaltung des:der jeweiligen Benutzer:in. Seit dem Aufkommen der ersten indexbasierten Websuchmaschinen Ende der 1990er Jahre haben sich die Benutzer:innen an die angebotene Visualisierung der Suchergebnisse in Form einer verlinkten linearen Liste gewöhnt.

Einfache Volltextsuche wird der Verwendung von komplexen Suchoperatoren vorgezogen. Allerdings verbesserte sich die Qualität der so erzielten Suchergebnisse fortlaufend durch die Auswertung von Feedback und zunehmende Personalisierung. Die Suchmaschine lernt zusammen mit ihrer:m Benutzer:in, insofern dieser dazu bereit ist, Daten über das eigene Suchverhalten preiszugeben. Grundsätzlich sind für eine semantische Suche auch semantisch erschlossene Dokumente notwendig, d. h. inhaltlich relevante Entitäten müssen in den unstrukturierten Daten als solche annotiert werden, um einen wie im vorangegangenen Abschnitt geschilderten Nutzen zu erzielen. Automatisierte semantische Annotation, wie z. B. Entity Linking zur Identifikation relevanter Entitäten, ist oft fehlerbehaftet und muss durch manuelle Annotation ergänzt werden. Die-

se scheitert aber oft am Design intuitiver und einfach zu benutzender Bedienoberflächen. *refer*³ ist ein als Wordpress Plugin realisiertes Entity-Annotations- und Visualisierungswerkzeug, das eine einfache und intuitive Benutzungsschnittstelle zur semantischen Annotation bietet (vgl. Abb. 5). *refer* erlaubt die automatisierte Annotation mit semantischen Entitäten des DBpedia⁴ Knowledge Graph (Tietz et al. 2016). In gleicher Weise lassen sich Werkzeuge der automatischen Objektidentifikation in Bildarchiven nutzen, um erkannte Objekte mit korrespondierenden Entitäten aus einem Wissensgraphen zu annotieren.



Abb. 5: Typgesteuertes semantisches Annotationswerkzeug refer implementiert als WordPress Plugin zur teilautomatisierten semantischen Annotation am Beispiel des SciHi Weblogs⁵

5.3 Explorative Suche und Empfehlungssysteme

Ein weiteres Problem der Nutzung inhaltlich erschlossener Dokumente liegt oftmals darin, dass die Suchenden nicht notwendigerweise über das benötigte Fachwissen in Form eines bestimmten Vokabulars verfügen, um die gewünschte Suchintention unter den Restriktionen der Suchmaschine eines Informationssystems zielbringend umzusetzen. Oftmals fällt es den Suchenden auch schwer, die gewünschte Suchintention überhaupt auszudrücken oder zu fixieren. Dabei

³ *refer*, Annotations- und Visualisierungswerkzeug, <https://refer.cx/> (21.1.2021).

⁴ DBpedia Knowledge Graph, <https://wiki.dbpedia.org/> (22.1.2021).

⁵ SciHi Weblog, <http://scihi.org/> (21.1.2021).

wäre es hilfreich zu wissen, welche Informationen insgesamt im betreffenden Informationssystem vorhanden sind. Nur zu leicht überfordern aktuelle Informationssysteme dabei die Benutzer:innen durch die immense Fülle der verfügbaren Informationen. Explorative Suche ermöglicht es den Benutzer:innen eben diese Herausforderungen explizit zu adressieren und bietet vielfältige Lösungsmöglichkeiten. Liegt das im Informationssystem verfügbare Wissen in Form eines Knowledge Graph vor, kann das explizit darin repräsentierte Wissen dazu genutzt werden, neue Wege durch den verfügbaren Informationsraum aufzuzeigen, indem semantische Ähnlichkeiten, relationale Beziehungen der Inhalte untereinander sowie inhaltliche Beziehungen zum:r jeweiligen Benutzer:in aufgezeigt und Suchempfehlungen ausgesprochen werden können. Das bereits genannte Annotations- und Visualisierungswerkzeug *refer* ermöglicht auch die Visualisierung inhaltlicher Zusammenhänge zwischen den in einem Dokument referenzierten Entitäten (vgl. Abb. 6). Auf diese Weise lassen sich Informationsinhalte schneller rezipieren sowie inhaltlich verwandte bzw. ähnliche Informationsinhalte direkt auffinden.

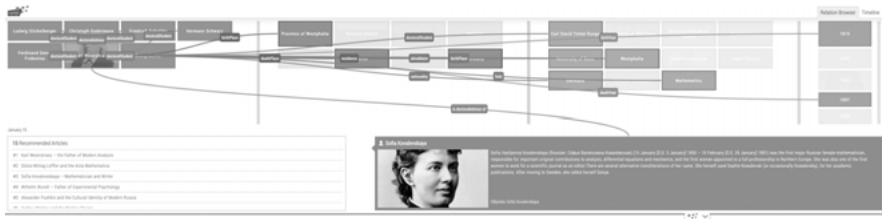


Abb. 6: Visualisierung inhaltlicher Zusammenhänge zwischen Informationsinhalten mit *refer*

6 Conclusion

Neueste Deep-Learning-basierte Verfahren können sowohl zur Klassifikation als auch zum Generieren von Informationsinhalten genutzt werden. Generative Adversarial Networks nutzen vergleichendes Lernen, um z. B. realistische Bilder zu erzeugen. In Verbindung mit Sprachmodellen können sprachliche Beschreibungen und inhaltlich korrespondierende Bilder gemeinsam erlernt werden, sodass generative Modelle in der Lage sind, Bilder zu textuellen Beschreibungen zu generieren (Xu et al. 2018). Während diese Technologie in jüngster Vergangenheit lediglich auf thematisch eng begrenzte Domänen sinnvoll zur Anwendung gelangen konnte, eröffnen neue hochkomplexe Sprachmodelle, wie z. B. aktuell GPT-3, nahezu treffsicher die Erzeugung passender realistischer Bilder

zu beliebigen realen oder auch rein fiktiven Objektbeschreibungen (Ramesh et al. 2021). Angewandt auf die Suche in Bildarchiven eröffnet diese Technologie neue Möglichkeiten, da sich die aktuell bereits eingesetzte visuelle Ähnlichkeitsbestimmung mit der semantischen Ähnlichkeit textueller Beschreibungen verbinden lässt. Damit werden auch vormals unerschlossene Bildbestände inhaltlich durchsuchbar. Verbindet man diese Deep-Learning-Verfahren mit Wissensgraphen, kann ebenso eine explorative Suche implementiert werden, die den Benutzer:innen bei der inhaltlichen Erkundung eines Bild- oder Textarchivs jenseits der Möglichkeiten aktueller Suchtechnologie unterstützt. Diese hybride Anwendungsform der Künstlichen Intelligenz verbindet die traditionell symbolische Wissensrepräsentation mit dem subsymbolisch repräsentierten Wissen der Deep-Learning-Verfahren. Dadurch eröffnen sich neue Möglichkeiten der inhaltlichen Erschließung, die qualitativ einen Quantensprung zu aktuell eingesetzten Technologien bieten und die Suche in Informationssystemen verändern werden. Schwierig bleibt dabei nach wie vor die Erfüllung der Erwartungshaltungen unterschiedlicher Benutzer:innen. Idealvorstellung bleibt eine Suchmaschine, die stets individuell mit ihren Benutzer:innen lernt, ohne dabei persönliche Daten preiszugeben.

Aber wird dadurch die intellektuelle Inhaltserschließung in Zukunft obsolet? Noch scheint es nicht ganz so weit. Auch wenn Deep-Learning-basierte Verfahren in speziellen Anwendungsgebieten die intellektuellen Fähigkeiten des Menschen bereits übertroffen haben, bleibt die allgemeine menschliche Intelligenz mit ihren kognitiven, sensomotorischen, emotionalen und sozialen Komponenten und der Ausbildung eines Bewusstseins bislang unerreicht. Semi-automatische Verfahren zur Inhaltserschließung unterstützen die intellektuelle Erschließung durch intelligente Vorschlagsmechanismen oder eine intellektuelle Bestätigung und eventuelle Korrektur der automatisch gewonnenen Ergebnisse. So können die Geschwindigkeit automatisierter Verfahren mit der höheren Genauigkeit intellektueller Erschließung effizient kombiniert werden. Wann werden also hybride KI-Verfahren die intellektuelle Inhaltserschließung vollständig ersetzen? Die Antwort auf diese Frage hängt nicht allein vom Stand der technischen Entwicklung ab. Es ist vor allen Dingen eine Frage der Bereitschaft, den dazu notwendigen Aufwand zu investieren: vollständige Digitalisierung der Inhalte, vollständige Transkription textbasierter Dokumente, vollständige visuelle und semantische Analyse. Dazu kommt das Problem eventueller fehlender Akzeptanz nicht intellektuell gewonnener Ergebnisse, obwohl unterschiedliche menschliche Bearbeiter:innen auch zu unterschiedlichen Ergebnissen gelangen können, bedingt durch unterschiedliches Fach- und Hintergrundwissen sowie verschiedene Kontexte. Eine preisgünstige universelle und umfassende Out-of-the-Box-Lösung zur automatischen Inhaltserschließung für Bibliotheken und

Archive wird aktuell noch auf sich warten lassen. Erste Ansätze in diese Richtung werden z. B. von dem DFG-geförderten Projekt *EEZU – Einfaches Erschließungs- und Zugriffssystem für kleine und mittlere Archive als Open-Source-Software und gehosteter Dienst*⁶ verfolgt, das semantische Technologien zur Datenverwaltung einsetzt und dazu auch Deep-Learning-gestützte Verfahren zur Inhaltserschließung erprobt und implementiert.

7 Literaturverzeichnis

- Berners-Lee, Tim, James Hender und Ora Lassila: The Semantic Web: a new form of Web content that is meaningful to computers will unleash a revolution of new possibilities. In: *Scientific American* (2001) Bd. 284 Nr. 5. S. 34–43.
- Bizer, Chris, Tom Heath und Tim Berners-Lee: Linked Data – The Story So Far. In: *International Journal on Semantic Web and Information Systems* (2009) Bd. 5 H. 3. S. 1–22.
- Brown, Tom B., Benjamin Mann, Nick Ryder, Melanie Subbiah u. a.: Language Models are Few-Shot Learners, arXiv:2005.14165, 2020. <https://arxiv.org/abs/2005.14165> (1.3.2021).
- Dessì, Danilo, Francesco Osborne, Diego Reforgiato Recupero, Davide Buscaldi, Enrico Motta und Harald Sack: AI-KG: an Automatically Generated Knowledge Graph of Artificial Intelligence. In: *The Semantic Web – ISWC 2020. 19th International Semantic Web Conference, Athens, Greece, November 2–6, 2020, Proceedings, Part II, Resources Track*. Hrsg. v. Jeff Z. Pan, Valentina Tamma, Claudia d’Amato, Krzysztof Janowicz, Bo Fu, Axel Polleres, Osmani Seneviratne und Lalana Kagal. Cham: Springer 2020. S. 127–143. https://doi.org/10.1007/978-3-030-62466-8_9.
- Elman, Jeffrey L.: Finding Structure in Time. In: *Cognitive Science* (1990) Bd. 14 H. 2. S. 179–211. https://doi.org/10.1207/s15516709cog1402_1.
- Goodfellow, Ian, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville und Yoshua Bengio: Generative Adversarial Networks. In: *Proceedings of the 27th International Conference on Neural Information Processing Systems (NIPS 2014)*, Volume 2. Hrsg. v. Z. Ghahramani, M. Welling, C. Cortes, Neil D. Lawrence und K. Q. Weinberger. Cambridge: MIT Press 2014. S. 2672–2680. <https://dl.acm.org/doi/10.5555/2969033.2969125>.
- Hebb, Donald O.: *The organization of behavior. A neuropsychological theory*. New York: Wiley 1949.
- Hochreiter, Sepp und Jürgen Schmidhuber: Long short-term memory. In: *Neural Computation* (1997) Bd. 9 Nr. 8. S. 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>.
- Hoppe, Fabian, Tabea Tietz, Danilo Dessì, Nils Meyer, Mirjam Sprau, Mehwish Alam und Harald Sack: The challenges of German archival document categorization on insufficient labeled data. In: *Proceedings of the Third Workshop on Humanities in the Semantic Web (WHiSe*

⁶ EEZU – Einfaches Erschließungs- und Zugriffssystem für kleine und mittlere Archive als Open-Source-Software und gehosteter Dienst, in *Geförderte Projekte der DFG*, <https://gepris.dfg.de/gepris/projekt/449727012?context=projekt&task=showDetail&id=449727012&> (30.1.2021).

- 2020) co-located with 15th Extended Semantic Web Conference (ESWC 2020). Hrsg. v. Alessandro Adamou, Enrico Daga und Albert Meroño-Peñuela. CEUR 2020 (Bd. 2695). S. 15–20. <http://ceur-ws.org/Vol-2695/paper2.pdf> (1.3.2021).
- Le Cun, Yann: Learning Processes in an Asymmetric Threshold Network. In: *Disordered Systems and Biological Organization*. Hrsg. v. E. Bienenstock, F. Fogelman-Soulié und G. Weisbuch. Berlin: Springer 1986. S. 233–240. https://doi.org/10.1007/978-3-642-82657-3_24.
- McCulloch, Warren und Walter Pitts: A logical calculus of the ideas immanent in nervous activity. In: *Bulletin of Mathematical Biophysics* (1943) Bd. 5. S. 115–133. <https://doi.org/10.1007/BF02478259>.
- Minsky, Marvin L. und Seymour A. Papert: *Perceptrons*. Cambridge, MA: MIT Press 1969.
- Parker, David B.: *Learning Logic*. Technical Report TR-87. Cambridge, MA: MIT, Center for Computational Research in Economics and Management Science 1985.
- Ramesh, Aditya, Mikhail Pavlov, Gabriel Goh, Scott Gray u. a.: DALL-E: Creating Images from Text. *OpenAI Blog*, January 5, 2021. <https://openai.com/blog/dall-e/> (21.1.2021).
- Razum, Matthias, Sandra Göller, Harald Sack, Tabea Tietz, Oleksandra Vsesviatska, Gerhard Weilandt, Marc Grellert und Torben Scharm: TOPORAZ – Ein digitales Raum-Zeit-Modell für vernetzte Forschung am Beispiel Nürnberg. In: *Information – Wissenschaft & Praxis* (2020) Bd. 71 H. 4. S. 185–194. <https://doi.org/10.1515/iwp-2020-2094>.
- Ristoski, Petar und Heiko Paulheim: RDF2vec: RDF graph embeddings for data mining. In: *The Semantic Web – ISWC 2016*. 15th International Semantic Web Conference, Kobe, Japan, October 17–21, 2016, Proceedings, Part I. Hrsg. v. Paul Groth, Elena Simperl, Alasdair Gray, Marta Sabou, Markus Krötzsch, Freddy Lecue, Fabian Flöck und Yolanda Gil. Cham: Springer 2016. S. 498–514. https://doi.org/10.1007/978-3-319-46523-4_30.
- Rosenblatt, Frank: The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain. In: *Psychological Review* (1958) Bd. 65 Nr. 6. S. 386–408. <https://psycnet.apa.org/doi/10.1037/h0042519> (21.1.2021).
- Rummelhart, David E., Geoffrey E. Hinton und Ronald J. Williams: Learning Representations by Back-Propagating Errors. In: *Nature* (1986) Bd. 323. S. 533–536. <https://doi.org/10.1038/323533a0>.
- Silver, David, Julian Schrittwieser, Karen Simonyan u. a.: Mastering the game of Go without human knowledge. *Nature*, Bd. 550 (2017). S. 354–359.
- Singhal, Amit: Introducing the Knowledge Graph: Things, Not Strings. *Google Official Blog*, 16.05.2012, <http://googleblog.blogspot.com/2012/05/introducing-knowledge-graph-things-not.html> (12.2.2021).
- Tietz, Tabea, Oleksandra Bruns, Sandra Göller, Matthias Razum, Danilo Dessi und Harald Sack: Knowledge Graph enabled Curation and Exploration of Nuremberg’s City Heritage. In: *Proceedings of 2nd Conference on Digital Curation Technologies* (2021).
- Tietz, Tabea, Joscha Jäger, Jörg Waitelonis und Harald Sack: Semantic Annotation and Information Visualization for Blogposts with refer. In: *Visualization and Interaction for Ontologies and Linked Data*. Proceedings of the Second International Workshop on Visualization and Interaction for Ontologies and Linked Data (VOILA 2016) co-located with the 15th International Semantic Web Conference (ISWC 2016). Hrsg. v. Valentina Ivanova, Patrick Lambrix, Steffen Lohmann u. Catia Pesquita. CEUR 2016 (Bd. 1704). S. 28–40. <http://ceur-ws.org/Vol-1704/paper3.pdf> (21.1.2021).
- Tietz, Tabea, Jörg Waitelonis, Kanran Zhou, Paul Felgentreff, Niels Meyer, Andreas Weber und Harald Sack: Linked Stage Graph. In: *Proceedings of the Posters and Demo Track of the*

- 15th International Conference on Semantic Systems co-located with 15th International Conference on Semantic Systems (SEMANTICS 2019). Hrsg. v. Mehwish Alam, Ricardo Usbeck, Tassilo Pellegrini, Harald Sack und York Sure-Vetter. CEUR 2019 (Bd. 2451). <http://ceur-ws.org/Vol-2451/paper-27.pdf>.
- Türker, Rima, Lei Zhang, Mehwish Alam und Harald Sack: Weakly Supervised Short Text Categorization Using World Knowledge. In: *The Semantic Web – ISWC 2020. 19th International Semantic Web Conference, Athens, Greece, November 2–6, 2020, Proceedings, Part I*. Hrsg. v. Jeff Z. Pan, Valentina Tamma, Claudia d’Amato, Krzysztof Janowicz, Bo Fu, Axel Polleres, Oshani Seneviratne und Lalana Kagal. Cham: Springer 2020. S. 584–600. https://doi.org/10.1007/978-3-030-62419-4_33.
- Vsesviatska, Oleksandra, Tabea Tietz, Fabian Hoppe, Mirjam Sprau, Niels Meyer, Danilo Dessi und Harald Sack: ArDO: An Ontology to Describe the Dynamics of Multimedia Archival Records. In: *Proceedings of the 36th ACM/SIGAPP Symposium On Applied Computing*. 2021.
- Waitelonis, Jörg, Claudia Exeler und Harald Sack: Linked Data Enabled Generalized Vector Space Model To Improve Document Retrieval. In: *Proceedings of the Third NLP&DBpedia Workshop (NLP & DBpedia 2015) co-located with the 14th International Semantic Web Conference 2015 (ISWC 2015)*. Hrsg. v. Heiko Paulheim, Marieke van Erp, Agata Filipowska, Pablo N. Mendes und Martin Brümmer. CEUR 2016 (Bd. 1581). S. 33–44. <http://ceur-ws.org/Vol-1581/paper4.pdf> (21.1.2021).
- Waitelonis, Jörg und Harald Sack: Towards exploratory video search using linked data. In: *Multimedia Tools and Applications* (2012) Bd. 59 H. 2. S. 645–672. <https://doi.org/10.1007/s11042-011-0733-1>.
- Xu, Tao, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Gan, Zhe Gan, Xiaolei Huang und Xiaodong He: AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks. In: *Proceedings of 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE 2018. S. 1316–1324. <https://doi.ieeecomputersociety.org/10.1109/CVPR.2018.00143>.

