Korrespondenzen der Frühromantik: Ein kontrolliertes Vokabular zur Analyse von Kommunikation und Wissenstransfer für das Semantic Web

Suárez Cronauer, Elena

Elena.SuarezCronauer@adwmainz.de Akademie der Wissenschaften und der Literatur I Mainz, Deutschland

Fath, Laura

Ifath@uni-mainz.de Johannes Gutenberg-Universität Mainz

Deicke, Aline

aline.deicke@adwmainz.de Akademie der Wissenschaften und der Literatur | Mainz, Deutschland

Strobel, Jochen

jochen.strobel@uni-marburg.de Johannes Gutenberg-Universität Mainz

Weyand, Sandra

weyands@uni-trier.de Universität Trier

Burch, Thomas

burch@uni-trier.de Universität Trier

Das Projekt "Korrespondenzen der Frühromantik"

Die Jenaer (und Berliner) Frühromantik gilt als die herausragende intellektuelle Revolution junger deutscher Autor*innen und Gelehrter an der Epochenschwelle um 1800. Die Gruppe agierte öffentlichkeitswirksam und nachhaltig, dispers und zugleich netzwerkbildend; sie reflektierte und praktizierte "Geselligkeit" beispielsweise auch mittels der Kommunikationsform "Brief". Die (auch quantitative) Auswertung dieser epistolaren Kommunikationsprozesse zwischen den Frühromantiker*innen einschließlich einer Untersuchung des dabei erfolgen-

den Wissenstransfers ist eines der großen Desiderate der Romantikforschung, dem das DFG-Projekt "Korrespondenzen der Frühromantik. Edition – Annotation – Netzwerkforschung" begegnen möchte. Ein grundlegender Schritt auf dieses Ziel hin ist die Erstellung kontrollierter Vokabulare in Gestalt normierter Meta- und Registerdaten, die in einen Knowledge Graphen auf Grundlage einer domänenspezifischen Ontologie eingebunden werden sollen. Einen Aspekt dieser Modellierungsprozesse präsentiert dieser Beitrag.

Die Datengrundlage bilden die Briefe der wichtigsten Protagonist*innen der Frühromantik (wie z. B. Friedrich, Dorothea, August Wilhelm und Caroline Schlegel, Novalis u.a.) untereinander und mit ihren weiteren Korrespondenzpartner*innen zwischen 1790 und 1802, also von den diversen 'Vorgeschichten' bis zum Zerfall des Jenaer Kreises (Schanze 2018, 18). Die Daten werden systematisch und vollständig erfasst, digital in Open Access publiziert und literaturwissenschaftlich wie netzwerktheoretisch ausgewertet. Darunter fällt auch die semantische Annotation von Aussagen in diesen Brieftexten, bei der dort erfasste Registerentitäten, wie Personen, Werke, Körperschaften oder Periodika, über eine zweiteilige Prädikatstruktur aus Illokution (dem eigentlichen Verb) und Proposition (einer Aussage, die dieses Verb näher spezifiziert) miteinander verknüpft werden. Zusammen mit den Register- und Metadaten der Briefe werden diese Annotationen in einen Knowledge Graphen überführt, in dem die Daten weiter angereichert werden. Auf dieser Basis erfolgen schließlich Auswertungen mittels quantitativer netzwerkanalytischer sowie qualitativer Ansätze. Die digitale Bereitstellung und Erschließung philologisch zuverlässiger Briefdigitalisate, die Annotation der Briefe sowie die parallele und abschließende graphenund netzwerktheoretische Auswertung stellen demnach die drei Kernbereiche des Projekts dar.

Der Vortrag stellt die zwei Arbeitsphasen vor, in denen die kontrollierten Vokabulare der Illokutionen und Propositionen im Zusammenspiel von digitaler Edition und Annotation entwickelt werden: Zunächst wird ein festes Begriffsset aus den Aussagen der Briefe destilliert und definiert, das danach in die Strukturen eines kontrollierten Vokabulars für die Verwendung als Linked Open Data überführt wird. Anschließend wird die Einbindung dieser Vokabulare in das Datenmodell des Knowledge Graphen präsentiert.

Entwicklung von Vokabularen zur semantischen Annotation von Aussagen

Der ursprüngliche Plan zur Erstellung des Begriffssets orientierte sich an Tripelstrukturen wie aus dem Semantic Web bekannt (Subjekt-Prädikat-Objekt). Insbesondere das Prädikat wurde, wie bereits angesprochen, im Laufe des Projekts jedoch differenziert und in zwei Kategorien – Illokution und Proposition – aufgesplittet. Diese können zu Analyse- und Publikationszwecken (z.B. als kontrolliertes RDF-Vokabular) wieder zusammengeführt werden, werden jedoch zur Gewährleistung größtmöglicher Flexibilität in der datenhaltenden Schicht zu-

nächst in dieser ausführlicheren Variante vorgehalten. Ein den Semantic-Web-Prinzipien ähnlicher Ansatz findet sich auch in der *quantitative narrative analysis* (QNA) (Franzosi 2010): Nach der sozialwissenschaftlichen Methode bilden semantische Tripel die grundlegenden Einheiten einer Erzählung (Sudhahar u.a. 2015, 2). Von uns wird diese Vorgehensweise erstmals für die Annotation von Briefen genutzt. Durch die zweiteiligen Prädikate werden Meta- und Registerdaten der Briefe (Subjekt und Objekt) semantisch verknüpft. Subjekte können in diesem spezifischen Fall nur Personen sein, Objekte Personen, Körperschaften, Werke und Periodika.

Ziel dieser Auszeichnungen ist es, die Prozesse intellektuellen Austauschs und die Spuren kollaborativen Schaffens in den Brieftexten für eine formale quantitative Analyse zu öffnen. Auf Basis der kontrollierten Vokabulare der Illokutionen und Propositionen (sowie der schon in Gestalt der Meta- und Registerdaten vorgegebenen Registereinträge) werden die konkreten Akte, aus denen sich Strukturen von Kommunikation und Wissenstransfer ergeben, formal expliziert, womit sich weiterführende datengestützte Auswertungsperspektiven im Sinne der Forschungsfrage ergeben.

Die Vokabulare entstehen teils bottom up, teils top down: mehrere Bearbeiter*innen annotieren parallel ein kleines Sample von Briefen. Aus dem im Team diskutierten Abgleich der gewählten Formulierungen ergibt sich eine Liste von zusammengesetzten Prädikaten im weiteren Sinne, die sich aus einem die jeweilige Kommunikationsfunktion bezeichnenden und eine kommunikative Handlung vollziehenden illokutionären Verb (also etwa "behaupten", "erbitten", "grüßen", "positiv bewerten") und fakultativ einer "welthaltigen" Proposition, der Aussage, die hier auf den Punkt zu bringen ist, zusammensetzt, also etwa: "Publikation"; "Buchsendung"; "Arbeitsplan". Die Proposition spezifiziert die abstrakte illokutionäre Aussage und beantwortet Fragen wie: "Was wird positiv bewertet?" oder "Wozu wird jemand (das Objekt) aufgefordert?" Bei der Erstellung dieser Liste(n) stellt der historische Bedeutungswandel der Sprache ein grundlegendes Problem dar. Wir versuchen diesem Problem zu begegnen, indem wir Anachronismen und begriffliche Überschneidungen mit romantischen Konzepten zu vermeiden versuchen. In einzelnen Fällen wird das nicht möglich sein. Hier können wir lediglich darauf verweisen, dass unser modernes Verständnis von "Kritik" nicht identisch ist mit Friedrich Schlegels Begriff. Das Korpus der Propositionen soll 200 bis 300 nicht überschreiten, die Zahl der illokutionären Verben wird etwa 80 bis 90 erreichen. Mit diesen Termini schließen wir an sprachwissenschaftliche Standards, nämlich die Sprechakttheorie, an, die zwischen illokutionärem Akt und Proposition unterscheidet. Dabei sind nicht alle Aussagen eines Briefes zu annotieren. Die Auswahl ergibt sich aus folgenden Fragen: Sind identifizierte/identifizierbare Akteur*innen, Werke, Periodika, Körperschaften beteiligt? Ist die Aussage mit möglichst generischen Formulierungen zu erfassen? Kann die Aussage Fragen zu unseren Forschungsschwerpunkten beantworten? Wollte man auf diese Einschränkungen verzichten und eine komplette maschinenlesbare Paraphrase von 6.000 teils umfangreichen und thematisch oft äußerst heterogenen Briefen leisten, würde der dafür notwendige Zeitaufwand den eines finanzierbaren Forschungsprojekt überschreiten.

Die Aussagenkette würde den Brief quantitativ um ein Mehrfaches übertreffen.

Bei der Erstellung der Begriffssets ist oberstes Gebot die intersubjektive Prüfbarkeit. Während Subjekt und Objekt positive Gegebenheiten sind, da Kopfdaten des Briefs oder Nennungen in den Registern, müssen Proposition und Illokution erst kontextuell abgeleitet werden. Eine Herausforderung, der man bei der Erstellung der Vokabulare begegnet, ist die Vereinheitlichung des zusammengesetzten Prädikats bei komplexen Aussagen, die auf zwei oder mehr Annotationsketten verteilt werden müssen. Dabei müssen die einzelnen Elemente – Illokution und Proposition – generisch und möglichst abstrakt gehalten werden, um das Vokabular zu begrenzen und eine quantitative Auswertung und potentielle Interoperabilität zu ermöglichen.

Auch die Entscheidung, was als Illokution gelten kann und was nicht, ist nicht immer einfach. Möglichst sollten keine Dopplungen von Begriffen im Propositions- und Illokutionsregister auftauchen wie beispielsweise senden als Illokution und in seiner substantivierten Form Sendung als Proposition. Andererseits können komplexere Aussagen wie "Buchsendung erbitten", also das Phänomen sekundärer Prädikation, nur um den Preis gewisser Unschärfen ausgedrückt werden. So muss zuweilen auf Nuancen verzichtet werden, um die Konsistenz der Annotationspraxis zu gewährleisten. Hier dient das kontrollierte Vokabular auch der Organisation von Informationen (ANSI/NISO, 10).

Die intersubiektive Prüfbarkeit soll durch den Anschluss an sprachwissenschaftliche Standards garantiert werden. So werden Illokutionen und Propositionen jeweils Oberbegriffsklassen zugewiesen. Sie sind die Objektivation des top down-Elements bei der Annotationspraxis, die eben nicht allein auf einer abstrahierenden Paraphrase einzelner Briefpassagen beruht. Im Falle der Illokutionen ist das die Zuordnung zu Illokutionstypen nach Searle (Searle 1976, 1-23). Er teilt Verben in die Gruppen Assertive, Direktive, Kommissive, Expressive und Deklarative ein. Die Oberbegriffe zu den Propositionen sollen abgeglichen werden mit Erkenntnissen der Romantikforschung, der Briefforschung und dem Wissen über die Zeit um 1800 einschließlich ihrer Alltagskultur. Dies ist mehr noch als bei den Illokutionen ein top down-Element, das auf vorgängigen Erkenntnisinteressen und dem Stand der Forschung beruht. Mit diesen in einem Trial-and-Error-Verfahren zu erarbeitenden Begriffsklassen verfolgen wir mehrere Absichten: Wir dokumentieren den Anschluss an vorausgehende wissenschaftliche Überlegungen (wir vermeiden dabei Objektsprache), wir erhöhen den Grad an Disambiguierung in unserem Korpus (oder stellen uns den zwischen den Klassen sicherlich unvermeidlichen Ambiguitäten, wenn etwa Unterbegriffe mehrfach zugeordnet werden können), wir geben Nutzer*innen unserer Editionsplattform als zusätzliche Suchoption Registerelemente an die Hand, erhöhen also die Usability unserer Textsammlung, und fügen schließlich für maschinelle Auswertungen eine zentrale, wenngleich händisch erhobene Datenquelle hinzu. Dieses Angebot soll u.a. die literaturwissenschaftliche Forschung befruchten, aber auch ein Beitrag sein zur Erforschung der Leistung der Kommunikationsform "Brief" in ihrer Textualität. Nicht zuletzt wird im Projekt selbst ein Wechselspiel von quantitativem und qualitativem Arbeiten erprobt, aus dem neue Erkenntnisse und Arbeitsweisen für die digitalen Geisteswissenschaften hervorgehen können.

Implementierung in den Knowledge Graphen

Für die Implementierung in den Knowledge Graphen dient folglich zunächst ein festes Set aus Begriffen als Grundlage, das die Briefaussagen, die im Zusammenhang mit Kommunikationsprozessen und dem daraus folgenden Wissenstransfer stehen, als die Verknüpfung von Registerentitäten durch Illokution und Proposition mit zugehörigen Oberklassen abbildet. Dieses Begriffsset wird bei der Modellierung einer Ontologie der "Korrespondenzen der Frühromantik" berücksichtigt und anschließend in einer RDF-Serialisierung in ein kontroliertes Vokabular im Sinne des Semantic Webs überführt.

Eine Ontologie wird als eine formale und explizite Spezifikation einer gemeinsamen Konzeptualisierung von Wissen definiert (Gruber 1993, 199). Ein kontrolliertes Vokabular als eine "Zusammenstellung von Bezeichnern (URIs) mit klar definierter Bedeutung" (Hitzler u.a. 2008, 48) stellt zunächst Informationen dar, identifiziert diese Informationen darüber hinaus aber auch nochmals eindeutig in einer maschinenlesbaren Form und bildet somit die Voraussetzung für die Beschreibung von semantischen Beziehungen innerhalb einer Ontologie. Hier werden diese Begriffe dann in komplexe, ggf. hierarchische Beziehungen gesetzt.

Herausfordernd ist in vorliegendem Fall die Transformation der oben beschriebenen Aussagen in Konzepte: Die Aussagen in den Briefe spiegeln deren natürliche Sprache; sie werden in den Vokabularen aus Illokutionsund Propositionsklassen formalisiert und in einen semantischen Zusammenhang gestellt. Diese Arbeit wird schließlich im Knowledge Graphen fortgesetzt, indem diese Aussagen als Konzepte innerhalb einer Ontologie abgebildet werden. Der Begriff 'Konzept' wird hier als eine aus der Wahrnehmung abstrahierte Vorstellung, also als eine mentale, wortähnliche Repräsentation von Dingen verstanden und somit eben nicht als ein spezifischer, in einem Brief beschriebener Sachverhalt (Margolis 2022). Von vornherein werden nicht lediglich die Aussagen der Briefe rekonstruiert bzw. formalisiert. Vielmehr wird das übergreifende und somit für weitere semantische Verarbeitung relevante Konzept hinter dieser Aussage bereits in den generischen Elementen des Begriffssets sichtbar. So ist beispielsweise nicht relevant, dass Schlegel in seinem Brief Schleiermacher um seine Kritik an einer bestimmten von ihm übersetzten Textpassage bittet. Relevant sind vielmehr die Begriffe "Kritik" und "erbitten" bzw. "Bitte" (also: Proposition und illokutionäres Verb) im Verhältnis zu Subjekt und Objekt, da durch diese Konzepte sowohl die Selbstreferentialität der Kommunikation (IIlokutionen) als auch die relevanten Themen der Kommunikationsprozesse (insbesondere der Wissenstransfer) angesprochen werden. Wie diese Aussagen in das kontrollierte Vokabular übernommen werden - z.B. ob man "Kritik erbitten" als ein Konzept wertet oder nicht hat Auswirkungen auf das weitere Datenmodell, also die Ontologie, und ihre Fähigkeit, generische Aussagen der frühromantischen Wissensdomäne abzubilden.

Um dieser Problematik für das Datenmodell adäquat zu begegnen, wurden zunächst verschiedene Ansätze geprüft, die sich mit der Modellierung geisteswissenschaftlicher Daten, insbesondere Briefdaten, beschäf-

tigen. ' Das Datenmodell der "Korrespondenzen der Frühromantik" referenziert auf das Cidoc Conceptual Reference Model (CRM) als Upper Ontology. Auch wenn das Cidoc CRM eigentlich als eine Ontologie für die Museumsdomäne entwickelt wurde, folgt die Orientierung daran der Definition einer Ontologie als einer gemeinsamen Konzeptualisierung von Wissen, da sich zahlreiche geisteswissenschaftliche Projekte an dem Modell orientieren. Zudem kommt die logische Ausrichtung des Cidoc CRMs auf Ereignisse (event driven architecture) den Forschungsfragen des Projekts zugute, da so bspw. die Korrespondenzen als Prozesse modelliert werden und somit Flexibilität innerhalb des Modells garantiert wird. So wird jede Klasse der frühromantischen Domänenontologie als Superklasse einer Klasse des Cidoc CRMs übergeordnet oder es werden, falls möglich, die Klassen des Cidoc CRMs direkt nachgenutzt. Ebenso wird bei den Properties vorgegangen. Darüber hinaus wurden Aspekte der Auszeichnungslogik von correspdesc miteinbezogen, da die bereits publizierten Daten der August-Wilhelm-Schlegel-Edition (Strobel 2014-2020) dieser folgen.

Für die Modellierung der Aussagen bedeutet dies nun Folgendes: Innerhalb des Konzeptes Brief (Klasse Letter mit Superklasse E33 Linguistic Object des Cidoc CRMs) können Personen (Klasse E21 Person), Zeitschriften (Periodical mit Superklasse E33), Werke (Work mit Superklasse E33), Institutionen (Institution mit Superklasse E74 Group), Schlagwörter (Theme mit Superklasse E55 Type) und Aussagen (Statement mit Superklasse E13 Attribute Assingment) miteinbezogen sein. Verbunden sind diese Klassen mit Letter jeweils mit der Property P129 is about des Cidoc CRMs. Das kontrollierte Vokabular wird in den Klassen Illocution und Proposition abgebildet, beide mit Superklasse E55 Type. Illocution ist zudem mit der Gruppe der Illoktutionstypen durch die Klasse Illocutiongroup verbunden, die ebenfalls die Superklasse E55 Type hat. Innerhalb der Klasse Statement als Domain werden die Aussagen nun modelliert, indem das Subjekt durch has_Subject (Subproperty von P140 assigned attribute to) mit E21 Person als Range verbunden wird, das Objekt als Range mit den Klassen E21 Person, Periodical, Work und/oder Institution mit der Property has_Object (Subproptery von P141 assigned) sowie das Prädikat mit has_Predicate (Subproperty von P177 assigned property of type) mit den Klassen Illocution und Proposition (Abb. 1). Die jeweiligen Einheiten einer Aussage werden somit als gesonderte Klassen betrachtet, die erst innerhalb der Aussage selbst wieder zusammengeführt werden, wobei die Properties die Subjekt-Illokution-Proposition-Objekt-Struktur festlegen. Perspektivisch bedeutet dies, dass bspw. die Proposition "Kritik" und die Illokution "erbitten" erst für die Netzwerkanalyse aus zuvor getrennt vorgehaltenen Einheiten der Graphdatenbank zusammengesetzt werden. Hierdurch wird eine flexible Struktur garantiert: Die Komplexität der Konstruktion von Aussagen wird durch das Aufgliedern in ihre einzelnen Bestandteile reduziert, sodass eine Nachnutzung der Daten - unter der Berücksichtigung der Auswahl der Begriffe, die von den angesprochenen Forschungsfragen bestimmt war - auch in anderen Kontexten möglich ist. Es gilt dementsprechend aus den spezifischen und eine konkrete "Sache" betreffenden Aussagen des Quellmaterials generische und allgemeine Konzepte zu entwickeln und diese als Linked Open Data darzustellen, welche dann einerseits für den im Projekt zu entwickelnden Knowledge Graphen bzw. die Ontologie genutzt, andererseits aber auch für andere Forschungsinteressen, welche die Zeit um 1800 oder das Kommunikationsmedium "Brief" betreffen, nachgenutzt werden können. Diese Entwicklung generischer und offener Forschungsdaten aus konkretem Quellenmaterial stellt eine Aufgabe dar, denen viele Forschungsprojekte aus dem Bereich der Diaital Humanities begegnen, die iedoch nicht minder signifikant ist: Nur so werden Forschungsdaten erzeugt, die auch über den Projektkontext hinaus genutzt werden und folglich die Forschungslandschaft ergänzen können. Zudem können durch den Linked Open Data-Ansatz die projektinternen Forschungsdaten an andere Datenkorpora angeschlossen werden, was einerseits zu einer grö-Beren Sichtbarkeit, andererseits zu vielfältigen Anschlussperspektiven führt, beispielsweise für die Historische Netzwerkforschung, die Geschichtswissenschaften (v.a. zur Kultur-, Sozial- und Ideengeschichte, aber auch zu Alltags- und Emotionsgeschichte), die Genderforschung oder die Judaistik.

Ausblick

Durch den Aufbau solcher kontrollierter Vokabulare, die Entwicklung des Datenmodells und die Integration in einen Knowledge Graphen eröffnet sich die Möglichkeit tiefergehender Analysen mit Methoden der historisch-geisteswissenschaftlichen Netzwerkforschung, innerhalb derer die Strukturen des Netzwerks der Jenaer (und Berliner) Frühromantik sowie seine Genese und Entwicklung als relationale Phänomene rekonstruiert werden. So wird ein neues Bild der einleitend erwähnten Akteur*innen gezeichnet, werden zentrale Personen, homo- und heterogene Strukturen, Überschreitungen sozialer Barrieren und die Kreuzung sozialer Kreise ausgemacht sowie Dynamiken dieser Strukturen aufgeschlüsselt. Da sich die Netzwerkforschung in Bezug auf die umfassende, auch inhaltliche Erschließung von Korrespondenzen noch in den Anfängen befindet, ergibt sich in der Kombination mit einer umfassenden, hochstrukturierten Datenbasis die Möglichkeit, neue Forschungsansätze in der Synthese literaturwissenschaftlicher und netzwerktheoretischer Arbeitsweisen zu entwickeln und für den Anwendungsfall der Jenaer Frühromantik eingehend zu prüfen.

Abbildung

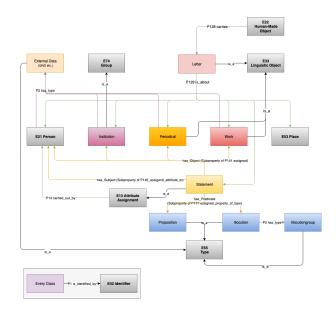


Abb. 1: Ausschnitt aus Datenmodell der "Korrespondezen der Frühromantik"

Fuβnoten

1. Hier sei auf das CIDOC CRM (Bekiari u.a. 2022), das "Letter Model", das im Rahmen der COST Action "Reassembling the Republic of Letters" entwickelt wurde (Jeffries u.a. 2019), das *correspdesc*-Element der TEI Guidelines (TEI Consortium [Hg.] 2021), das Europeana Data Model (Doerr u.a. 2010), die Ontologie "OntoAndalus" (Almeida u.a. 2021) sowie die Ontologie "OntoBellini-Letters" (Cristofaro u.a. 2022) verwiesen.

Bibliographie

Almeida, Bruno, and Rute Costa. "OntoAndalus: an ontology of Islamic artefacts for terminological purposes." Semantic Web 12.2 (2021): 295–311.

ANSI/NISO. 2005. Guidelines for the Construction, Format, and Management of Monolingual Controlled Vocabularies, 10. https://www.niso.org/publications/ansiniso-z3919-2005-r2010 (zugegriffen: 28. Juli 2022).

Bekiari, Chryssoula, George Bruseker, Martin Doerr, Ore Christian-Emil, Stead Stephen, Athanasios Velios, Erin Canning, und Philippe Michon. 2022. Volume A: Definition of the CIDOC Conceptual Reference Model. Version 7.2.1. ICOM/CIDOC Documentation Standards Group/CRM Special Interest Group. http://www.cidoc-crm.org/sites/default/files/CIDOC%20CRM_v.7.0_%2020-6-2020.pdf (zugegriffen: 28. Juli 2022).

Cristofaro, Salvatore, Pietro Sichera und Daria Spampinato. 2022. "An ontology proposal for a corpus of letters of Vincenzo Bellini: formal properties of physical structure and the case of rotated texts". International

Journal of Metadata, Semantics and Ontologies 15, Nr. 4: 269-279.

Dumont, Stefan, Ingo Börner, Dominik Leipold, Jonas Müller-Laackman und Gerlinde Schneider. 2019. "Corresponde Metadata Interchange Format." In Encoding Correspondence. A Manual for Encoding Letters and Postcards in TEI-XML and DTABf, hg. von Stefan Dumont, Susanne Haaf und Sabine Seifert. Berlin. https://encoding-correspondence.bbaw.de/v1/CMIF.html (zugegriffen: 28. Juli 2022).

Franzosi, Roberto. 2010. *Quantitative narrative analysis.* Los Angeles u.a.: SAGE Publications (=Quantitative Applications in the Social Sciences 162).

Gruber, Thomas. 1993. "A Translation Approach to Portable Ontology Specifications." *Knowledge Acquisition* 5, Nr. 2: 199–220.

Hitzler, Pascal, Markus Krötzsch, Sebastian Rudolph und York Sure. 2008. Semantic Web: Grundlagen. Berlin: eXamen.press.

Jefferies, Neil, Howard Hotson, Christoph Kudella, Miranda Lewis, Thomas Stäcker, und Gertjan Filarski. 2019. "Letter Model". In Reassembling the Republic of Letters in the Digital Age. Standards, Systems, Scholarship, hg. von Howard Hotson und Thomas Wallnig, 171–89. Göttingen: Universitätsverlag Göttingen. http://www.univerlag.uni-goettingen.de/handle/3/isbn-978-3-86395-403-1 (zugegriffen: 28. Juli 2022).

Margolis, Eric und Stephen Laurence, "Concepts". In *The Stanford Encyclopedia of Philosophy* (Fall 2022 Edition), hg. von Edward N. Zalta und Uri Nodelman, https://plato.stanford.edu/archives/fall2022/entries/concepts/ (zugegriffen: 07. Dezember 2022).

Miles, Alistair und Dan Brickley. 2005. "SKOS Core Vokabular" https://www.w3.org/TR/2005/WD-swbp-skoscore-spec-20051102/ (zugegriffen: 28. Juli 2022).

Schanze, Helmut. 2018. *Erfindung der Romantik.* Stuttgart: Metzler.

Searle, John R. 1976. "A classification of illocutionary acts." *Language in Society* 5, Nr. 1: 1–23.

Sudhahar, Saatviga, Giuseppe A Veltri und Nello Cristianini. 2015. "Automated Analysis of the US Presidential Elections Using Big Data and Network Analysis." *Big Data & Society 2*, Nr. 1. https://doi.org/10.1177/2053951715572916 (zugegriffen: 28. Juli 2022).

Strobel, Jochen und Claudia Bamberg (Hg.). 2014-2020. August Wilhelm Schlegel: Digitale Edition der Korrespondenz. https://august-wilhelm-schlegel.de (zugegriffen: 07. Dezember 2022).

TEI Consortium (Hg.) 2021. "2.4.6 Correspondence Description." *TEI P5: Guidelines for Electronic Text Encoding and Interchange.* 4.2.1. TEI Consortium. https://www.tei-c.org/release/doc/tei-p5-doc/en/html/HD.html#HD44CD (zugegriffen: 28. Juli 2022).