

Otmane Azeroual

Untersuchungen zur Datenqualität und Nutzerakzeptanz von Forschungs- informationssystemen

Framework zur Überwachung
und Verbesserung der Qualität von
Forschungsinformationen

MOREMEDIA



Springer Vieweg

Untersuchungen zur Datenqualität und Nutzerakzeptanz von Forschungsinformationssystemen

Otmane Azeroual

Untersuchungen zur Datenqualität und Nutzerakzeptanz von Forschungsinformationssystemen

Framework zur Überwachung und
Verbesserung der Qualität von
Forschungsinformationen

 Springer Vieweg

Otmane Azeroual 
Deutsches Zentrum für Hochschul- und
Wissenschaftsforschung (DZHW)
Berlin, Deutschland

ISBN 978-3-658-36701-5 ISBN 978-3-658-36702-2 (eBook)
<https://doi.org/10.1007/978-3-658-36702-2>

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de> abrufbar.

© Der/die Herausgeber bzw. der/die Autor(en), exklusiv lizenziert durch Springer Fachmedien Wiesbaden GmbH, ein Teil von Springer Nature 2022

Das Werk einschließlich aller seiner Teile ist urheberrechtlich geschützt. Jede Verwertung, die nicht ausdrücklich vom Urheberrechtsgesetz zugelassen ist, bedarf der vorherigen Zustimmung des Verlags. Das gilt insbesondere für Vervielfältigungen, Bearbeitungen, Übersetzungen, Mikroverfilmungen und die Einspeicherung und Verarbeitung in elektronischen Systemen.

Die Wiedergabe von allgemein beschreibenden Bezeichnungen, Marken, Unternehmensnamen etc. in diesem Werk bedeutet nicht, dass diese frei durch jedermann benutzt werden dürfen. Die Berechtigung zur Benutzung unterliegt, auch ohne gesonderten Hinweis hierzu, den Regeln des Markenrechts. Die Rechte des jeweiligen Zeicheninhabers sind zu beachten.

Der Verlag, die Autoren und die Herausgeber gehen davon aus, dass die Angaben und Informationen in diesem Werk zum Zeitpunkt der Veröffentlichung vollständig und korrekt sind. Weder der Verlag noch die Autoren oder die Herausgeber übernehmen, ausdrücklich oder implizit, Gewähr für den Inhalt des Werkes, etwaige Fehler oder Äußerungen. Der Verlag bleibt im Hinblick auf geografische Zuordnungen und Gebietsbezeichnungen in veröffentlichten Karten und Institutionsadressen neutral.

Planung/Lektorat: Stefanie Eggert

Springer Vieweg ist ein Imprint der eingetragenen Gesellschaft Springer Fachmedien Wiesbaden GmbH und ist ein Teil von Springer Nature.

Die Anschrift der Gesellschaft ist: Abraham-Lincoln-Str. 46, 65189 Wiesbaden, Germany

Vorwort

Die vorliegende Dissertation ist im Rahmen meiner Tätigkeit als wissenschaftlicher Mitarbeiter am Deutschen Zentrum für Hochschul- und Wissenschaftsforschung GmbH (DZHW)¹ zwischen 2016 und 2021 entstanden. Maßgeblich für diese Arbeit war das Projekt „Helpdesk für die Einführung des Kerndatensatzes Forschung (KDSF)“², das vom Bundesministerium für Bildung und Forschung (BMBF) und von allen Bundesländern gefördert wurde. Im KDSF-Projekt wurde ein freiwilliger Standard für die Erhebung, Vorhaltung und den Austausch von Forschungsinformationen entwickelt und sowohl für Dateneigner (z. B. Hochschulen und außeruniversitäre Forschungseinrichtungen (AUFs)) als auch für Datenabfrager (z. B. Forschungsförderer, Behörden usw.) im deutschen Wissenschaftsraum bereitgestellt. Der KDSF beschreibt die Datengruppierung nach Inhalten und die Formate zur Implementierung in Forschungsinformationssysteme (FIS), weil heutzutage die Vorhaltung und der Austausch von Forschungsinformationen an einer zunehmenden Anzahl an Hochschulen und außeruniversitären Forschungseinrichtungen über ein FIS erfolgen. Ein FIS als Datenbanksystem soll mittel- und langfristig interne sowie externe Berichtspflichten und -legungen für die Hochschulverwaltung und politische Entscheidungsträger erleichtern. Es existieren nationale und internationale Standards (z. B. KDSF-Datenmodell und Common European Research Information Format (CERIF)) zur Unterstützung von Forschungsinformationssystemen, und diese ermöglichen Kompatibilität und Interoperabilität zwischen verschiedenen Systemen, um die Forschungsbereiche zu repräsentieren. Um diese systematisch aufbauenden Bereiche effizient nutzen zu können, wird eine gute Datenqualität verlangt, da ein wichtiger Erfolgsfaktor

¹ <https://www.dzhw.eu>

² <https://kerndatensatz-forschung.de>

für die Implementierung des FIS in Hochschulen und AUFs die Sicherstellung der Datenqualität darstellt. Allerdings wird dieser Punkt bisher kaum in der Wissenschaft untersucht, obwohl das Thema für Hochschulen, AUFs und deren Forscher relevant und durchaus aufschlussreich ist und in internationalen wissenschaftlichen Publikationen erwähnt wurde. Bislang fehlen die erforderliche Beachtung der Bedeutung der Datenqualität, die Dimensionen zur Messung und Prüfung der Datenqualität in einem FIS, neue Methoden bzw. Techniken zur Verhinderung und Verbesserung von Ursachen mangelnder Qualität in einem FIS und die Verstärkung des Vertrauens bzw. der Akzeptanz in das FIS. Ohne eine ausreichende Datenqualität kann das FIS sein Nutzenpotenzial als Lieferant entscheidungsrelevanter Daten allerdings nicht ausschöpfen. Grundsätzlich gilt, je höher und sicherer die Qualität der Daten in Forschungsinformationssystemen sind, desto größer ist die Nutzerakzeptanz.

Die Entstehung der Dissertation haben zahlreiche Personen ermöglicht und unterstützt. Bei ihnen möchte ich mich an dieser Stelle herzlichst bedanken.

An erster Stelle gilt dies ganz besonders meinem Doktorvater, Herrn Prof. Dr. Gunter Saake. Durch seine persönliche, fortwährende Betreuung, die freundliche Hilfe, die Unterstützung während des Entstehungsprozesses dieser Dissertation, durch die Schaffung von hervorragenden Rahmenbedingungen für ein praxisorientiertes Forschen im Bereich der Datenbanken und Informationssysteme und durch die Freiheiten, die er mir für die Dissertation gewährte, hat er maßgeblich zum erfolgreichen Abschluss des Forschungsvorhabens beigetragen.

Ein besonderer Dank gilt auch Herrn Dr.-Ing. Eike Schallehn für seine hilfreiche und aufmerksame Unterstützung, Ermutigung und stete Motivation.

Zudem habe ich Herrn Prof. Dr.-Ing. Mohammad Abuosba für die Übernahme des Zweitgutachtens sowie für seine ständige, motivierende Beratung, Unterstützung und sein wertvolles Feedback zu diversen Fragen rund um die Dissertation zu danken. Die zahlreichen Gespräche an der Otto-von-Guericke-Universität Magdeburg werden mir stets eine kostbare und gute Erinnerung sein.

Dem Deutschen Zentrum für Hochschul- und Wissenschaftsforschung (DZHW) danke ich vielmals für die finanzielle Unterstützung dieser Dissertation, insbesondere in Form von Förderungen von Konferenz- und Workshopbesuchen, die das Vorhaben wesentlich bereichert haben. Ebenfalls danke ich meinem KDSF-Team am DZHW Berlin für die kollegiale Arbeit und interessante Zeit, die mir in guter Erinnerung bleiben wird.

Verbunden bin ich auch Herrn Dr. Joachim Schöpfel für seine Unterstützung, die konstruktiven Anregungen und die ebenso außerordentliche wie freundschaftliche Zusammenarbeit. Ein zusätzlicher Dank geht an meinem weiteren Gutachter Herrn Prof. Dr. Sören Auer für die zuverlässige und hilfreiche Unterstützung.

Ebenfalls gilt auch ein großer Dank allen Teilnehmenden meiner Befragung für ihre Auskunftsbereitschaft, interessanten Beiträge und Antworten auf meine Fragen.

Außerdem möchte ich Herrn Dr. Volker Manz für das Korrekturlesen meiner Dissertation danken.

Abschließend möchte ich mich von ganzem Herzen bei meinen lieben Eltern Jamila Morjane und Mustapha Azeroual bedanken, die mir das Studium überhaupt erst ermöglicht haben und auf deren Unterstützung ich immer zählen kann und konnte. Sie standen mir während meiner gesamten Ausbildungszeit äußerst liebevoll und mit vollen Kräften zur Seite. Ein besonderer und persönlicher Dank gilt meiner Frau Virginia Azeroual für ihre ganze Bandbreite an Geduld, Motivation und unglaublich hilfreicher Unterstützung während der Anfertigung meiner Dissertation.

Meinen allerherzlichsten Dank an Allah!

Otmane Azeroual

Inhaltsverzeichnis

1	Einleitung	1
1.1	Motivation des Forschungsvorhabens	2
1.2	Problemstellung	3
1.3	Forschungslücken und Forschungsmethodik	4
1.4	Zielsetzung und Forschungsfragen	6
1.5	Vorgehen und Aufbau der Dissertation	7
2	Konzeptionelle Grundlagen	9
2.1	Hochschulen und außeruniversitäre Forschungseinrichtungen (AUFs)	9
2.2	Forschungsinformationen	11
2.3	Forschungsinformationssysteme (FIS)	18
2.4	KDSF und CERIF als Standard für Forschungsinformationen	26
2.5	Ziele und Herausforderungen für die Anwendung von FIS	34
2.6	Die verschiedenen Sichtweisen von FIS für Stakeholder	35
2.7	FIS-Marktübersicht	38
2.8	FIS-Produkte	40
3	Untersuchung der Datenqualität in FIS	49
3.1	Die Besonderheiten von FIS	49
3.2	Begriffsdefinition der Datenqualität	54
3.3	Dimensionen der Datenqualität	57
3.4	Datenqualitätsprobleme im Vergleich mit anderen Informationssystemen	60
3.5	Datenqualitätsprobleme in FIS	64
3.6	Ursachen für Datenqualitätsprobleme in FIS	73

3.7	Messung der Datenqualität in FIS	75
3.7.1	Vollständigkeit	78
3.7.2	Korrektheit	81
3.7.3	Konsistenz	84
3.7.4	Aktualität	87
3.7.5	Zusammenfassung der Ergebnisse	92
3.7.6	Einfluss von Datenqualitätsdimensionen auf die Gesamtqualität in FIS	93
3.8	Analyse der Datenqualität in FIS	98
3.8.1	Data Profiling	100
3.8.2	Zusammenfassung der Ergebnisse	111
3.9	Verbesserung und Steigerung der Datenqualität in FIS	116
3.9.1	Data Cleansing	117
3.9.2	Datenqualitätsmaßnahmen	123
3.9.3	Zusammenfassung der Ergebnisse	126
3.10	Ergebnisse der Datenqualitätsuntersuchung	127
3.11	Text Data Mining Methoden in FIS	132
3.11.1	Natural Language Processing	134
3.11.2	Informationsextraktion	138
3.11.3	Clustering	141
3.11.4	Zusammenfassung der Ergebnisse	148
4	Ermittlung der Nutzerakzeptanz von FIS	149
4.1	Akzeptanzbegriff und Akzeptanzmodelle	149
4.2	Adaption des Akzeptanzmodells für das FIS	152
4.3	Analyse der Nutzerakzeptanz von FIS	154
4.4	Abhängigkeit zwischen Datenqualität und Nutzerakzeptanz ...	160
4.5	Zusammenfassung der Ergebnisse	169
5	Proof-of-Concept	171
5.1	Methodik	171
5.2	Ergebnisse der Evaluation	173
6	Zusammenfassung und Ausblick	181
	Literaturverzeichnis	189

Abkürzungsverzeichnis

AUFs	Außeruniversitäre Forschungseinrichtungen
BI	Business Intelligence
BMBF	Bundesministerium für Bildung und Forschung
CASRAI	Consortia Advancing Standards in Research Administration Information
CERIF	Common European Research Information Format
CMS	Content Management System
CMS	Campus Management System
CORDIS	Community Research and Development Information Service
CRF	Conditional Random Field
CRIS	Current Research Information Systems
CRM	Customer Relationship Management
DBMS	Datenbankmanagementsystem
DCMI	Dublin Core Metadata Initiative
DFG	Deutsche Forschungsgemeinschaft
DINI	Deutsche Initiative für Netzwerkinformation e. V. – Arbeitsgruppe
DNB	Deutsche Nationalbibliothek
DOI	Digital Object Identifier
DPMA	Deutsches Patent- und Markenamt
DQ	Datenqualität
DQD	Datenqualitätsdimensionen
DQM	Datenqualitätsmanagement
DWH	Data Warehouse
DZHW	Deutsches Zentrum für Hochschul- und Wissenschaftsforschung
eG	eingetragene Genossenschaft
EOSC	European Open Science Cloud

EPO	European Patent Office
ERM	Entity Relationship Modell
ERP	Enterprise Resource Planning
ETL	Extraktion, Transformation, Laden
EU	Europäische Union
EUNIS	European University Information Systems
e. V.	eingetragene Verein
FDM	Full Data Model
FhG	Fraunhofer-Gesellschaft
FI	Forschungsinformationen
FIS	Forschungsinformationssysteme
GEPRIS	Geförderte Projekte Informationssystem
GLM	General Linear Model
GmbH	Gesellschaft mit beschränkter Haftung
HAC	Hierarchical Agglomerative Clustering
HGF	Helmholtz-Gemeinschaft Deutscher Forschungszentren
HTML	Hypertext Markup Language
IDM	Identity Management
Ids	Identifikationen
IE	Informationsextraktion
iFQ	Institut für Forschungsinformation und Qualitätssicherung e. V.
IT	Informationstechnologie
KDSF	Kerndatensatz Forschung
KMO	Kaiser-Meyer-Olkin-Kriterium
MDM	Master Data Management
MPG	Max-Planck-Gesellschaft
NER	Named Entity Recognition
NLP	Natural Language Processing
OCLC	Online Computer Library Center
ODI	Oracle Data Integrator
OPAC	Online Public Access Catalogue
ORCID	Open Researcher and Contributor ID
PAISY	Personal-Abrechnungs-und-Informations-System
PCA	Principal Components Analysis
PDI	Pentaho Data Integration
PLS	Partial Least Squares
POS	Part-of-Speech
RCD	Research Core Dataset
RDBMS	Relational Database Management System

RDF	Resource Description Framework
RIS	Research Information Systems
RMS	Root Mean Square
SAP	Systemanalyse Programmentwicklung
SEM	Structural Equation Modeling
SSOAR	Social Science Open Access Repository
TAM	Technology Acceptance Model
TDM	Text Data Mining
WGL	Wissenschaftsgemeinschaft Gottfried Wilhelm Leibniz
WR	Wissenschaftsrat
XML	Extensible Markup Language

Abbildungsverzeichnis

Abbildung 1.1	Architektur des Konzepts bzw. Frameworks	5
Abbildung 2.1	Verarbeitung von Forschungsinformationen in das FIS	16
Abbildung 2.2	FIS-Architektur	21
Abbildung 2.3	Nutzung von FIS in deutschen Hochschulen und AUFs (N=160)	22
Abbildung 2.4	Warum von Hochschulen und AUFs kein FIS genutzt wird? (N=48)	23
Abbildung 2.5	Die am meisten verwendeten Metadaten in FIS (N=51)	24
Abbildung 2.6	Integration der unterschiedlichen Datenquellen in FIS (N=51)	25
Abbildung 2.7	Das europäische CERIF-Datenmodell	29
Abbildung 2.8	Das deutsche KDSF-Datenmodell	32
Abbildung 2.9	Die verwendeten Tools von FIS (N=51)	39
Abbildung 3.1	Aspekte zur Beschreibung der Datenqualität im Kontext von FIS (N = 68)	56
Abbildung 3.2	Datenqualitätsprobleme in FIS (N = 68)	64
Abbildung 3.3	Interne und externe Datenquellen	66
Abbildung 3.4	Falsche Erfassung von Umlauten und Sonderzeichen	67
Abbildung 3.5	Unterschiedliche Namensformen von gleichen Autoren	68
Abbildung 3.6	Falsche Erfassung und Reihenfolge von Institutionsangaben	69
Abbildung 3.7	Doppelte Erfassung der DOIs	70

Abbildung 3.8	Fehlerhafte Mehrfacherfassung von Institutionsangaben	71
Abbildung 3.9	Datenqualitätsprüfung in FIS (N = 68)	72
Abbildung 3.10	Datenqualitätsdimensionen in FIS (N = 68)	76
Abbildung 3.11	Klassifizierung von Datenqualitätsdimensionen in FIS	77
Abbildung 3.12	Beispiel für die Messung der Korrektheit	83
Abbildung 3.13	Messung der Korrektheit mithilfe von Levenshtein-Distanz	83
Abbildung 3.14	Beispiel für inkonsistente Daten in einer Publikationsliste	85
Abbildung 3.15	Messung der Konsistenz mithilfe von Levenshtein-Distanz	87
Abbildung 3.16	Framework zur Qualitätsmessung und -verbesserung in FIS	92
Abbildung 3.17	Strukturgleichungsmodell für die Datenqualitätsdimensionen in FIS	95
Abbildung 3.18	Auszug der durchgeführten Umfrage über Datenqualitätsdimensionen	96
Abbildung 3.19	Datenintegration in das FIS	101
Abbildung 3.20	Data-Profiling-Analysearten	103
Abbildung 3.21	Analyse der Publikationsdaten mithilfe des DataCleaner-Tools	110
Abbildung 3.22	Datentyp-Analyse	112
Abbildung 3.23	Null-Werte/Duplikaten-Analyse	113
Abbildung 3.24	Muster- und Domänen-Analyse	114
Abbildung 3.25	Meta-Prozessablauf	115
Abbildung 3.26	Parsing der Daten	118
Abbildung 3.27	Berichtigung und Standardisierung der Daten	119
Abbildung 3.28	Anpassung/Ableichung der Daten	120
Abbildung 3.29	Zusammenführung der Daten	121
Abbildung 3.30	Anreicherung der Daten	122
Abbildung 3.31	Anwendungsfalldiagramm zur Verbesserung der Datenqualität in FIS	123
Abbildung 3.32	Maßnahmenportfolio nach verschiedenen Kriterien	124
Abbildung 3.33	Maßnahmen zur Qualitätsverbesserung in FIS (N = 68)	128

Abbildung 3.34	Grad der Datenqualität von FIS in Hochschulen und AUFs (N = 68)	129
Abbildung 3.35	Richtlinie für die Qualitätsüberwachung und -verbesserung in FIS	131
Abbildung 3.36	Workflow zur Text- und Dokumentenanalyse in FIS	135
Abbildung 3.37	Anwendung von NLP-Funktionen	137
Abbildung 3.38	IE-Prozess	138
Abbildung 3.39	Entitätenextraktion aus dem Beispiel des Publikationstextes	140
Abbildung 3.40	Bildung von Clustern	142
Abbildung 3.41	Funktion des <i>k</i> -means Algorithmus	143
Abbildung 3.42	Beispielberechnung mit <i>k</i> -means Clustering	144
Abbildung 3.43	Beispielberechnung einer hierarchischen agglomerativen Clusterbildung	146
Abbildung 3.44	Beispiel für das Clustering von Dokumenten	147
Abbildung 4.1	FIS-Akzeptanzmodell	153
Abbildung 4.2	Gestellte Fragen zur Nutzerakzeptanz von FIS	156
Abbildung 4.3	Ergebnisse der Akzeptanz von 51 FIS-Nutzern	157
Abbildung 4.4	Korrelationsanalyse „Streudiagramm“	160
Abbildung 4.5	SEM für die Abhängigkeit zwischen Datenqualität und Nutzerakzeptanz	164
Abbildung 4.6	Regressionsanalyse	168
Abbildung 5.1	Analyse der Datenqualitätsprobleme mithilfe des DataCleaner-Tools	174
Abbildung 5.2	Optimierung der Datenqualitätsprobleme mithilfe des DataCleaner-Tools	175
Abbildung 5.3	Datenqualitätsprobleme in Relation zur Qualitätsanalyse	178

Tabellenverzeichnis

Tabelle 2.1	Überblick über die vier wichtigsten AUFs in Deutschland	11
Tabelle 2.2	Definitionen des FIS-Begriffes	19
Tabelle 2.3	Die Top fünf Anbieter von FIS-Tools	38
Tabelle 2.4	Kriterienkatalog zur Bewertung von FIS-Produkten	44
Tabelle 3.1	Definitionen von Datenqualität	55
Tabelle 3.2	Objektive Datenqualitätsdimensionen	58
Tabelle 3.3	Subjektive Datenqualitätsdimensionen	59
Tabelle 3.4	Ursachen für Datenqualitätsprobleme in FIS	74
Tabelle 3.5	Beispiel für die Messung der Vollständigkeit	79
Tabelle 3.6	Beispiel für die Berechnung der metrischen Aktualität ...	89
Tabelle 3.7	Ergebnis des Cronbachs-Alpha für Datenqualitätsdimensionen	97
Tabelle 3.8	Ergebnis der Hauptkomponentenanalyse für Datenqualitätsdimensionen	97
Tabelle 3.9	Beispiel für die funktionale Abhängigkeit	105
Tabelle 3.10	Beispiel für die Referenzanalyse	106
Tabelle 3.11	Kriterienkatalog zur Bewertung von Datenqualitätstools	108
Tabelle 3.12	Datenqualitätsmaßnahmen	125
Tabelle 4.1	Definitionen des Akzeptanzbegriffs	150
Tabelle 4.2	Akzeptanzmodelle	151
Tabelle 4.3	Ergebnisse der deskriptiven Analyse	159
Tabelle 4.4	Ergebnisse der Reliabilität und Validität von Nutzerakzeptanzindikatoren	165
Tabelle 5.1	Datenqualitätsprobleme vor und nach der Anwendung ...	176

Formelverzeichnis

3.1	Definition der Metrik einer Datenqualitätsdimension [LPFW06]	75
3.2	Metrik zur Datenqualitätsdimension Vollständigkeit [LPFW06]	79
3.3	Metrik zur Datenqualitätsdimension Korrektheit [LPFW06]	81
3.4	Metrik zur Datenqualitätsdimension Konsistenz [LPFW06]	85
3.5	Messergebnis der Konsistenz [ASW18]	86
3.6	Metrik zur Datenqualitätsdimension Aktualität [HK09]	87
3.7	Aktualitätsdimension basierend auf dem Beispiel (Verfall(A))=0 [HK09]	88
3.8	Aktualitätsdimension basierend auf dem Beispiel (Alter(w,A))=0 [HK09]	88
3.9	Aggregierte Aktualität des Beispiels [KKH07]	89
3.10	Ähnlichkeit zwischen zwei Zahlen „x“ und „y“ [CMA+12]	142
3.11	Mittelwert Ähnlichkeit [CMA+12]	143
3.12	Quadratisches Mittel Ähnlichkeit [CMA+12]	143
3.13	Peak Ähnlichkeit [CMA+12]	143
3.14	Berechnung des Durchschnittsvektors [Wol05]	144

Listingsverzeichnis

3.1	Vollständigkeitsmessung	80
3.2	Korrektheitsmessung	84
3.3	Konsistenzmessung	86
3.4	Aktualitätsmessung	90



Einleitung

1

Mit zunehmendem Bedarf der Forschungsaktivitäten, den steigenden Anforderungen an die Forschungsberichterstattung im deutschen Wissenschafts- und Hochschulsystem und den wachsenden Datenmengen bei gleichzeitigen Veränderungen der Informationsbedürfnisse von Forschern und seitens der Öffentlichkeit werden alle Organisationen vor neue Herausforderungen gestellt. Langfristig werden sich jene Hochschulen und außeruniversitären Forschungseinrichtungen (AUFs) durchsetzen, die sich auf diese Bedingungen einstellen, die also in der Lage sein werden, flexibel und schnell auf Veränderungen zu reagieren, und gleichzeitig ihre Kosten im Griff haben und auch reduzieren. Hierfür ist jedoch eine genaue Kenntnis der aktuellen Organisationssituation unverzichtbar. Um dies zu gewährleisten und das Management von Forschungsinformationen bei den eigenen Planungs- und Entscheidungstätigkeiten mit den benötigten Daten zu versorgen, werden hoch entwickelte Forschungsinformationssysteme eingesetzt. Diese wurden von einer europäischen Community namens euroCRIS¹ (*Current Research Information Systems (CRIS)*) entwickelt und unterstützt. Mittlerweile hat sich in der Forschung der Begriff Forschungsinformationssystem (FIS) etabliert. Dabei beschreibt ein FIS Ansätze wie das Sammeln, Integrieren, Speichern und Analysieren von Forschungsinformationen einer Organisation und dient als föderierte Datenbank [ASA18b].

In den letzten Jahren hat sich das Thema Forschungsinformationssysteme in den deutschen und internationalen Organisationen stark verbreitet und ist zu einem festen Bestandteil der universitären IT-Landschaften geworden [ASW18], [ASA18a],

¹<https://www.eurocris.org>

Ergänzende Information Die elektronische Version dieses Kapitels enthält Zusatzmaterial, auf das über folgenden Link zugegriffen werden kann https://doi.org/10.1007/978-3-658-36702-2_1.

[ASA18b]. Gleichzeitig arbeiten viele Hochschulen und AUFs derzeit noch an der Implementierung solcher Informationssysteme. Durch den Einsatz von Forschungsinformationssystemen können die Organisationen und die Forscher darin unterstützt werden, ihre Forschungsaktivitäten und Forschungsergebnisse transparent und intelligent zu gestalten. Mit der Hilfe eines FIS können Forschungsinformationen wie *Projekte, Drittmittel, Patente, Kooperationspartner, Preise, Publikationen usw.* gespeichert und verwaltet sowie in den jeweiligen Webauftritt eingebunden werden. Für die Hochschulen und AUFs bilden Forschungsinformationen die Grundlage für die interne und externe Bewertung ihrer Leistungen als selbstständige wissenschaftliche Organisation. Nebenbei haben sie die Möglichkeit, das eigene Forschungsprofil übersichtlich darzustellen; ein FIS dient somit als ein effektives Werkzeug des Informationsmanagements. Die Vereinheitlichung und Erleichterung der Berichterstattung schafft einen Mehrwert sowohl für die Verwaltung und das Management als auch für die Forscher. Forschende und andere Interessengruppen bekommen mit dem FIS ein geeignetes Mittel, um Informationen über ihre Forschungsaktivitäten und aktuelle Trends, das Erstellen von Lebensläufen, die Ergebnisse und Kooperationen in ihren Forschungsfeldern, existierende Projekte, Publikationen und Förderungen sowie den Kontakt zu den Interessenten aus der Wirtschaft zu managen und zu begleiten bzw. den Anforderungen der Berichterstattung gerecht zu werden. Durch all diese Prozesse soll eine Mehrarbeit für die FIS-Nutzer vermieden werden, was z. B. die Reduktion des Zeitaufwandes bei der Erstellung von Berichten oder bei der Außerdarstellung ihrer Forschungsleistung und wissenschaftlichen Expertise bedeutet.

1.1 Motivation des Forschungsvorhabens

Obwohl das FIS schon seit Längerem existiert, hat es erst vor wenigen Jahren besonders durch viele verschiedene Konferenzen etwa im Rahmen der europäischen Organisation euroCRIS und der Deutschen Initiative für Netzwerkinformation e. V. (DINI AG-FIS)² mit Experten und Entwicklern an Bekanntheit gewonnen. Hier wurde das Potenzial von FIS entdeckt, und einige deutsche Hochschulen und AUFs begannen damit, solche Systeme zu implementieren. Ein FIS einzuführen, bedeutet, die erforderlichen Informationen über Forschungsaktivitäten und -ergebnisse in gesicherter Qualität zur Verfügung zu stellen. Da die Einführung und der Betrieb von Forschungsinformationssystemen mit erheblichen Kosten und einem hohen Ressourcenbedarf verbunden sind, gewinnen die Qualität der Daten und die Nutzerakzep-

² <https://dini.de/ag/fis/>

tanz zunehmend an Bedeutung. Die Datenqualität gehört neben der Datensicherheit, der Bedienerfreundlichkeit und anderen Variablen zu den wesentlichen Bedingungen der Nutzerakzeptanz eines FIS. Dabei geht es in erster Linie um Vertrauen – Vertrauen in das System, in dessen Anbieter und in die Administration. Einem System, das Datenprobleme nicht zuverlässig identifiziert oder korrigiert oder das selbst eine Quelle für Datenqualitätsmängel ist, kann (und wird) kein Vertrauen entgegengebracht werden. Wahrgenommene oder empfundene Qualitätsprobleme beeinträchtigen die subjektive Leistungserwartung gegenüber dem System. Dieser Sachverhalt ist nicht neu [Dav93], [WT05]; im Fall eines FIS ist eine mangelhafte Datenqualität aber umso problematischer, als es um strategische und zum Teil hochsensible Informationen und Entscheidungshilfen geht, wie personenbezogene oder finanzielle Daten. Für die Nutzbarkeit und Interpretation einrichtungsspezifischer Forschungsinformationen spielt die Datenqualität eine wichtige Rolle; sie entscheidet über Erfolg oder Misserfolg bei Hochschulen und AUFs. Nur mit einer hohen Qualität können die FIS-Nutzer zuverlässige und nutzbringende Forschungsergebnisse liefern und eine fundierte Entscheidungsfindung mit einer guten Präsentation von Forschung ermöglichen. Daher sollte das Thema Datenqualität Gegenstand ständiger Sorgfalt und Aufmerksamkeit sein; ein „Qualitätsmanagement“ ist erforderlich, um die Systemleistung und die Zufriedenheit und Akzeptanz der Benutzer zu gewährleisten. Die vorliegende Dissertation adressiert die Problemstellung der Datenqualität bei FIS nutzenden Hochschulen und AUFs in Deutschland und verschiedenen europäischen Ländern und stellt entsprechende Lösungen bereit, um die Akzeptanz des FIS seitens dieser Einrichtungen zu erhöhen.

1.2 Problemstellung

Das Management von Forschungsinformationen wird zunehmend zu einer wichtigen Aufgabe für die Hochschulen und AUFs [AA18]. Wachsende Datenmengen und die größer werdende Anzahl an internen und externen Datenquellen, so z. B. durch Informationssysteme für Humanressourcen, Finanzhaushalte und Bibliotheken, führen in einem FIS zu mehr Fehlern, z. B. Rechtschreibfehlern, Dubletten, zu einer fehlerhaften Formatierung oder zu fehlenden, inkorrekten und uneinheitlichen Daten, die zunehmend zu einer Herausforderung in Hochschulen und AUFs werden und sich zu einem ernsthaften Problem entwickeln können. Diese unterschiedlichen Datenfehler können bei der Erfassung, der Übertragung sowie der Integration von Forschungsinformationen in das FIS entstehen und sich über verschiedene Bereiche erstrecken sowie schwer auffindbar sein [ASA18b]. Außerdem kann die Qualität der externen Datenquellen (wie z. B. unterschiedliche Publikationsdatenbanken,

Identifiers, Textdateien oder XML-Files usw.) einen ungünstigen Einfluss auf die Qualität der FIS haben. Heutzutage gehört der Umgang mit großen Datenquellen für die Hochschulen und AUFs zum täglichen Betrieb. Dies liegt hauptsächlich an der Tatsache, dass viele Schnittstellen im Forschungsmanagementprozess entwickelt werden müssen, um einen Informationsaustausch zwischen den Quellsystemen zu ermöglichen, da das FIS als integraler Bestandteil einer kompletten Verwaltungssoftware für die Administration von Einrichtungen konzipiert werden kann und eine Plattform für die Datenintegration darstellt [BEL12]. Während der Datenintegration können die Hochschulen und AUFs jederzeit über den Zustand ihrer Forschungsinformationen informiert und ihre Qualitätsfehler schnell identifiziert und korrigiert werden. Daher muss bei der Integration von Forschungsinformationen auf die Datenheterogenität und Datenverteilung geachtet werden. Hierbei ist eine einmalige Bereinigung in den internen und externen heterogenen und verteilten Quellsystemen nicht ausreichend, denn Forschungsinformationen müssen kontinuierlich gepflegt und ihre Qualität muss immer wieder optimiert werden [ASA18a], [ASA18b]. Solange die Nutzer nicht in der Lage sind, auf die am dringendsten benötigten Informationen zuzugreifen und schnelle Entscheidungen zu treffen, sinkt der Wert der verwendeten Forschungsinformationen und das Vertrauen in das FIS und dessen Akzeptanz [ASA18b]. Die Datenqualität bei den Hochschulen und AUFs sollte daher mit hoher Priorität behandelt werden; dies bedingt die Einführung einer Strategie, die als Frühindikator dient, um die Probleme der Datenqualität im FIS anzugehen.

1.3 Forschungslücken und Forschungsmethodik

Je mehr Forschungsinformationssysteme es gibt und je mehr Forschungsinformationen dort gesammelt, gespeichert und verarbeitet werden, desto wichtiger wird die Betrachtung der Qualität, damit die Stakeholder qualitativ hochwertige Ergebnisse erhalten können und die Systemakzeptanz erhöht wird. Ohne eine akzeptable Datenqualität im jeweiligen FIS ist eine Akzeptanz der Nutzer nicht möglich. In der Literatur zu FIS beschäftigen sich nationale und internationale FIS-Experten in ihren wissenschaftlichen Publikationen im Allgemeinen nur mit dem Aufbau, dem Mehrwert und den Herausforderungen eines FIS in Form von Erfahrungsberichten. Mit Blick auf die Datenqualität des FIS ließen sich hingegen kaum Quellen in der Literatur finden, und es gibt bisher keine Methodik, die vorschreibt, welche Qualitätsschritte in welcher Reihenfolge durchgeführt werden müssen, um anschließend ein hohes Qualitätsniveau des FIS langfristig und dauerhaft zu garantieren. Insbesondere mangelt es an Untersuchungen über die Maßnahmen und deren praktischen

Einsatz zur Bewertung und Verbesserung der Datenqualität im FIS. Allerdings wurden diese Punkte in wissenschaftlichen Arbeiten zu anderen, ähnlichen Bereichen der Informationssysteme und des Informationsmanagements behandelt und publiziert.

Immerhin haben Hochschulen, AUFs und Forschende dem Thema Datenqualität in FIS in den ersten Jahren im Rahmen von Konferenzen und Workshops große Aufmerksamkeit geschenkt, und immer öfter ist die Rede von Forschungslücken im Hinblick auf die Sicherstellung und Verbesserung der Datenqualität bei der Integration von Forschungsinformationen mit unterschiedlichen Quellsystemen in das FIS sowie von der Notwendigkeit, die Nutzerakzeptanz von FIS zu steigern [Mün17], [SSC+19]. Vor diesem Hintergrund ist die vorliegende Dissertation im Forschungsbereich der Wirtschaftsinformatik (WI) zum Themengebiet Datenbanken und Software Engineering entstanden und soll dabei unterstützen, FIS effektiv und effizient zu nutzen. Das Ziel besteht darin, im Rahmen einer anwendungsbezogenen Forschung zum einen die untersuchten Datenqualitätsprobleme in FIS aus der Praxis zu lösen. Grundlage hierfür bilden Erfahrungen, die aus dem KDSF-Projekt und organisierten Workshops am DZHW Berlin stammen, sowie die Erkenntnisse aus einer quantitativen Untersuchung (Umfrage) mit FIS nutzenden Hochschulen und AUFs. Zum anderen geht es darum, ein Konzept bzw. Framework für einen kontrollierten Umgang mit der Datenqualität in FIS zu entwickeln.

Abbildung 1.1 stellt die Referenzarchitektur als Konzept- bzw. Framework-Lösung für die vorliegende Dissertation dar. Die einzelnen Komponenten werden in **Kapitel 3** weiter vertieft. Das Thema ist von hoher Praxisrelevanz an praktisch allen

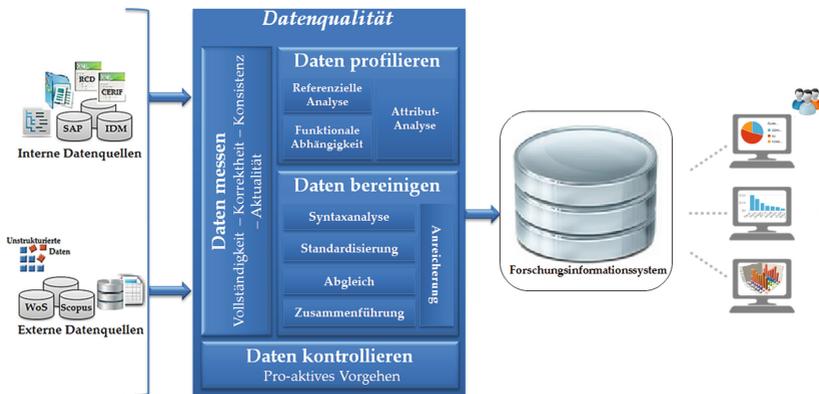


Abbildung 1.1 Architektur des Konzepts bzw. Frameworks