

I Einleitung

„Curiously enough, lexicalization has never been used by historical linguistics for the purpose of dating, although its study is extremely rewarding.“¹

Mario ALINEI

Die Datierung von Texten hat in den vergangenen Jahrhunderten Forscher:innen unterschiedlicher Wissenschaften und Kulturkomplexe beschäftigt. Im Vordergrund stehen dabei oft die Exegese und Authentizitätsforschung. Die linguistische Textdatierung stützt sich in aller Regel auf die Beobachtung von und das Vorwissen über Sprachwandel, traditionell durch die Betrachtung bestimmter – einzelner – sprachlicher Phänomene, Wörter, Wortformen, Zeichen oder grammatikalischer Strukturen. Vor allem der Wortschatz jeder aktiv genutzten Sprache, sei es Schrift- oder Umgangssprache, befindet sich in permanentem Wandel. Neue Wörter kommen hinzu, andere fallen nach und nach aus dem Sprachgebrauch heraus.²

Von der Autorschaft unabhängige Textdatierung als Aufgabenfeld der Computerlinguistik ist ein sehr junges Forschungsgebiet, das durch einen Aufsatz von DE JONG, RODE und HIEMSTRA (2005) ins Leben gerufen wurde.³ Die darin vorgestellte Methodik nutzt *Bag of Words* (BoW)-Sprachmodelle und statistische Ähnlichkeitsmaße. Texte können so auf ihre Ähnlichkeit zu anderen Texten oder zu Teilen diachroner Vergleichskorpora geprüft und entsprechend zugeordnet werden. Ein wichtiges Ziel der Datierung ist dabei – im Unterschied zur traditionellen Forschung – die Einordnung von Dokumenten nach Relevanz, z. B. für die Sortierung der Ergebnisse von Suchmaschinen. Vergleichbare Ansätze wurden inzwischen für verschiedene europäische Sprachen mit großem Erfolg angewandt.⁴ Für das Chinesische liegen bisher aber kaum vergleichbare Veröffentlichungen vor.⁵

- 1 Mario ALINEI 2004: „The Problem of Dating in Linguistics“. In: *Quaderni di semantica* 25.2, S. 211–232, S. 225.
- 2 *Zhi* 之 sei an dieser Stelle als für Sinolog:innen anschauliches Beispiel für solche diachronen Unterschiede genannt. In klassischen bzw. schriftsprachlichen Texten (s. u.) wird *zhi* u. a. als Subordinationspartikel eingesetzt. In einer vergleichbaren Funktion wird in der zeitgenössischen Umgangssprache *de* 的 verwendet.
- 3 Franciska M. G. DE JONG, Henning RODE und Djoerd HIEMSTRA 2005: „Temporal Language Models for the Disclosure of Historical Text“. In: *Humanities, computers and cultural heritage: Proceedings of the XVIth International Conference of the Association for History and Computing (AHC 2005)*. Amsterdam: Koninklijke Nederlandse Academie van Wetenschappen, S. 161–168.
- 4 Siehe v. a. Kapitel 3.1, ab S. 40. Vgl. u. a. David BAMMAN et al. 2017: „Estimating the Date of First Publication in a Large-Scale Digital Library“. In: *Proceedings of ACM/IEEE Joint Conference on Digital Libraries, Toronto, Canada, June 2017 (JCDL '17)*, S. 1–10. DOI: 10.475/1234; Filip GRALIŃSKI et al. 2017: „The RetroC Challenge: How to Guess the Publication Year of a Text?“. In: *Proceedings of the 2nd International Conference on Digital Access to Textual Cultural Heritage. DATeCH2017*. Göttingen: ACM, S. 29–34. DOI: 10.1145/3078081.3078095.
- 5 Ausnahmen bilden ein Aufsatz von YAMADA Takahito 山田崇仁 2004: „N-gram moderu o riyōshite senshin bunken no seisho jiki o saguru: ‚Sonshi‘ jūsanhen o jirei toshite N-gram モデルを利用して先秦文献の成書時期を探る — 『孫子』十三篇を事例として —, n-Gramm Ansatz zur Datierung von chinesischen-Qin Texten am Beispiel der 13 *Sunzi*-Kapitel“. In: *Tōkyōdaigaku tōyō bunka kenkyūjo fuzoku tōyō-gaku kenkyū jōhō sentā, Ajia kenkyū jōhō* 東京大学東洋文化研究所附属東洋学研究情報センター, アジア研究情報 (Tokyo University Research & Information Center for Asian Studies, Gateway to Asian Studies in Japan), siehe auch Kapitel 3.1, S. 40, sowie ein kurzer Aufsatz von YU Xuejin und WEI Huangfu 2019:

1 Einleitung

Auf dem Gebiet der Stilometrie wurden in jüngster Zeit auch für schriftsprachliche chinesische Texte⁶ spannende Ergebnisse präsentiert.⁷ Darin deutet sich an, dass chinesischsprachige Textstile über einen Zeitraum von vielen Jahrhunderten eine hohe Rigidität aufweisen können, so dass Unterschiede zwischen Gattungen bzw. Sprachstilen mit statistischen Methoden viel klarer messbar bzw. differenzierbar sind als eine temporale Dimension.⁸

Darüber hinaus liefern orthographische Veränderungen in Sprachen mit alphabetischen Schriftsystemen – sei es bedingt durch tatsächliche phonologische Veränderungen, oder durch Rechtschreibreformen bzw. -konventionen – wertvolle Hinweise zur chronologischen Einordnung von Texten.⁹ Die chinesische Schrift(sprache) hat sich diesbezüglich jedoch über einen langen Zeitraum hinweg verhältnismäßig wenig verändert.¹⁰ Andererseits ermöglicht die fast lückenlose Texttradition, dank der uns Textzeugnisse aus dem Zeitraum von ca. 1000 v. u. Z. bis ins 21. Jh. zur Verfügung stehen, im Falle des Chinesischen eine nahezu unvergleichliche zeitliche Tiefe in diachronen sprachwissenschaftlichen Untersuchungen.

Zentrale Fragestellungen und Ziele

Ziel dieser Arbeit ist es, die inhaltsbasierte zeitliche Einordnung bzw. Datierung schriftsprachlicher chinesischer Texte mit computerlinguistischen Methoden zu ermöglichen. Die erwähnten statistischen Sprachmodelle sollen zu diesem Zweck erstmals für chinesisches Textmaterial adaptiert und angewandt werden.¹¹ Einschränkungen ergeben sich dabei aus der stilistischen Rigidität einiger schriftsprachlicher Textgattungen und der Beständigkeit des chinesischen Schriftsystems. Die Nutzung statistischer Sprachmodelle erfordert überdies ein diachrones Trainingskorpus, das den gesamten Zeitraum abdeckt, aus dem Texte datiert werden sollen.

Da ein geringer syntaktischer und ein eher wenig in der Schrift manifestierter phonologischer Wandel die Möglichkeiten statistischer Methoden begrenzen, soll hier der lexikalische Wandel stärker in den Fokus rücken. Bereits DE JONG, RODE und HIEMSTRA „foresee a role for parsed entries from historical dictionaries in this context [...]“,¹² ohne diesen Ansatz

„A Machine Learning Model for the Dating of Ancient Chinese Texts“. In: *International Conference on Asian Language Processing, IALP 2019, Shanghai, China, November 15-17, 2019*. Hrsg. von LAN Man et al. IEEE, S. 115–120. DOI: 10.1109/IALP48816.2019.9037653; Auch Ryan NICHOLS, Edward SLINGERLAND und Kristoffer NIELBO scheinen sich im Rahmen ihrer Beschäftigung mit *machine learning* und *topic modelling* indirekt mit temporaler Klassifizierung antiker chinesischer Texte zu befassen, legen aber bislang keine Veröffentlichung dazu vor. Siehe Ryan NICHOLS et al. 2018: „Modeling the Contested Relationship between *Analects*, *Mencius*, and *Xunzi*: Preliminary Evidence from a Machine-Learning Approach“. In: *Journal of Asian Studies* 77.1, S. 19–57, S. 23, siehe auch Kapitel 4.1, ab S. 60.

6 Zum Begriff *schriftsprachliches Chinesisch* siehe S. 5.

7 Siehe dazu z. B. die Arbeit von Paul VIERTHALER, der unterschiedliche Gattungen von Geschichtstexten untersucht. Siehe Paul VIERTHALER 2016a: „Fiction and History: Polarity and Stylistic Gradience in Late Imperial Chinese Literature“. In: *Journal of Cultural Analytics*, S. 1–32. DOI: 10.22148/16.003, ausführlicher dazu siehe Kapitel 2.3, ab S. 20 und Kapitel 4.1, ab S. 60.

8 Vgl. dazu auch Tilman SCHALMEY 2021: „Thoughts on »Reliable« Learner’s Vocabularies for Classical and Literary Chinese“. In: *Teaching Classical Chinese | Zum Unterricht des Klassischen Chinesischen | Wenyan wen jiaoxue 文言文教学 (Proceedings of the International Symposium on the Teaching of Classical Chinese, December 14–16, 2018)*. Hrsg. von Li Wen 李文 und Ralph KAUF. Gossenberg: Ostasien Verlag, S. 251–261, S. 254–255.

9 Siehe Anne GARCIA-FERNANDEZ et al. 2011: „When Was It Written? Automatically Determining Publication Dates“. In: *String Processing and Information Retrieval*. Hrsg. von Roberto GROSSI, Fabrizio SEBASTIANI und Fabrizio SILVESTRI. Berlin, Heidelberg: Springer Berlin Heidelberg, S. 221–236, S. 7; siehe auch GRALIŃSKI et al. 2017, S. 32.

10 Vgl. z. B. CHOU Ya-Min 周亞民 und HUANG Chu-ren 黃居仁 2010: „Hantology: conceptual system discovery based on orthographic convention“. In: *Ontology and the Lexicon*. Hrsg. von HUANG Chu-ren 黃居仁 et al. Studies in Natural Language Processing. Cambridge & New York: Cambridge University Press, S. 122–143, S. 133. Ausführlicher dazu siehe Kapitel 2.2, ab S. 16.

11 Siehe v. a. Kapitel 6.1, ab S. 156.

12 DE JONG, RODE und HIEMSTRA 2005, S. 2.

aber weiter zu verfolgen. Die Idee, Informationen über das Entstehen und Verschwinden von Wörtern aus Wörterbüchern zur Textdatierung zu nutzen, findet sich erneut in der Arbeit von GARCIA-FERNANDEZ et al. (2011) für die Datierung französischsprachiger Texte.¹³ Dennoch ist die computerlinguistische Nutzung lexikographischer Quellen zur Textdatierung bisher ausgeblieben. Sie beklagen:

[...] there is no pre-compiled list of words with their year of appearance or disappearance. This type of information is sometimes included in dictionaries, but depends on the availability of these resources.¹⁴

Für das Chinesische steht mit dem *Hanyu da cidian* 漢語大詞典 (*Großes Wörterbuch der chinesischen Sprache*)¹⁵ (HYDCD) jedoch eine sehr umfangreiche Quelle digital zur Verfügung,¹⁶ aus der zumindest Informationen über das Erscheinen von Wörtern bzw. Lexemen extrahiert werden können. Zu diesem Zweck wird die *Plain Text* Fassung des HYDCD in eine relationale Datenbank umgeformt, in der Angaben über die frühesten Belege von Lexemen abrufbar sind. Auf dieser Grundlage entwickle ich alternative lexikographische Methoden zur Schätzung der Entstehungszeit schriftsprachlicher Texte.¹⁷ Es ergeben sich zwei zentrale Fragestellungen:

1. Lassen sich die erwähnten statistischen Methoden erfolgreich für die Verwendung mit chinesischem Textmaterial adaptieren?
2. Können die beschriebenen Limitationen mit lexikographischen Vorgehensweisen reduziert werden?

Vorbereitend werden grundsätzliche Fragen zu geeigneten Ressourcen und dem computerlinguistischen Umgang mit schriftsprachlichem Chinesisch erörtert. Welche diachronen Trainings- und Testkorpora sind geeignet? Ist eine verlässliche Segmentierung schriftsprachlicher Texte möglich? Wie eingangs erwähnt besteht überdies eine enge Verbindung zwischen linguistischer Datierung und Sprachwandel. Sekundär sollen daher anhand des verwendeten Materials Beobachtungen zum Sprachwandel, insbesondere dem lexikalischen Wandel des Chinesischen, ermöglicht und diskutiert werden.

Die Verwendung der untersuchten Methoden bleibt auf digitales *Plain Text*-Material beschränkt. Sie eignen sich nicht zur Datierung gescannter Drucke oder gar von Handschriften, die nur als Bilddaten vorliegen, auch wenn darin sichtbare Merkmale von Originaldokumenten wie Materialität, Schriftstile, Drucktypen oder Zeichenvarianten wichtige Aspekte der Textdatierung darstellen können.¹⁸

¹³ Siehe GARCIA-FERNANDEZ et al. 2011, S. 5.

¹⁴ Ebd.

¹⁵ LUO Zhufeng 羅竹風, Hrsg. 1986–1994: *Hanyu da cidian* 漢語大詞典 (*Großes Wörterbuch der chinesischen Sprache*). Bd. 1–13. Shanghai 上海: Cishu chubanshe 辭書出版社 (im Folgenden zit. als HYDCD).

¹⁶ LUO Zhufeng 羅竹風, Hrsg. 2005: *Hanyu da cidian* 漢語大詞典 UTF-8 (*Großes Wörterbuch der chinesischen Sprache, Unicode-Version*). Shanghai 上海. URL: <http://bbs.gxsd.com.cn/forum.php?mod=viewthread&tid=498015> (besucht am 13. 01. 2013) (im Folgenden zit. als DHYDCD).

¹⁷ Siehe v. a. Kapitel 6.2, ab S. 179 und 6.3, ab S. 210.

¹⁸ Neuere archäologische Befunde sprechen zudem dafür, dass die typischerweise linear dargestellte Entwicklung der chinesischen Schrift von Orakel-, Bronze, Siegel- und Kanzleischrift bis hin zu einer standardisierten Schrift, tatsächlich weniger linear verlief als gemeinhin angenommen. Siehe dazu Imre GALAMBOS 2006: *Orthography of Early Chinese Writing: Evidence from Newly Excavated Manuscripts*. Budapest: Department of East Asian Studies, Eötvös Loránd University.

Digital Humanities

Die vorliegende Arbeit verbindet Sinologie und Computerlinguistik und kann übergreifend den *Digital Humanities* (DH, *shuzi renwen* 数字人文) zugerechnet werden.¹⁹ Damit sie für computerlinguistisch nicht vorgebildete Leser:innen verständlich und nachvollziehbar bleibt, werden an geeigneter Stelle Konzepte und Begriffe aus der Computerlinguistik erklärt, bzw. auf die weiterführende Fachliteratur verwiesen.

Die Verbindung von Geisteswissenschaft und Informatik ist keineswegs erst ein Trend der vergangenen Jahre. Für den Versuch, die Stücke William SHAKESPEARES chronologisch zu ordnen, wurden bereits Ende des 19. Jhs. – noch ohne den Einsatz von Computern – statistische Methoden für die Beantwortung geisteswissenschaftlicher Fragestellungen eingesetzt.²⁰ 1949 gelang es Roberto BUSA, Thomas WATSON von IBM zu überzeugen, eine elektronische Konkordanz der Texte Thomas von AQUINS für dessen Forschung zu erstellen.²¹ Anhand dieser untersuchte BUSA erfolgreich „nicht bewusste Spracheigentümlichkeiten“ mit wahrscheinlichkeitstheoretischen Ansätzen.²² Eine 1964 veröffentlichte Arbeit über die Verfasserschaft der *Federalist Papers* stützte sich ebenfalls auf die Häufigkeit bestimmter Funktionswörter.²³ Solche stilometrischen Analysen der Autorschaft von Texten gehören also zu den frühesten Forschungsbereichen der DH.²⁴

Toolstack

Für die Programmierung der im Rahmen dieser Arbeit entwickelten Software kommt hauptsächlich *Python 3* zum Einsatz.²⁵ *Python* gehört „zu den beliebtesten Programmiersprachen weltweit“²⁶ und „hat sich in den letzten Jahrzehnten zu einem erstklassigen Tool für wissenschaftliche Berech-

19 Die DH haben sich als „neues Arbeitsfeld etabliert, das an der Schnittstelle zwischen den Geisteswissenschaften und der Informatik angesiedelt ist.“ Fotis JANNIDIS, Hubertus KOHLE und Malte REHBEIN, Hrsg. 2017: *Digital Humanities – Eine Einführung*. Stuttgart: Metzler, S. XI. Unter dem Begriff DH vereint werden sowohl Arbeiten, die vorhandene Tools wie *Voyant* einsetzen, um einer geisteswissenschaftlichen Fragestellung nachzugehen Stéfán SINCLAIR und Geoffrey ROCKWELL 2016: *Voyant Tools*. URL: <https://voyant-tools.org/> (besucht am 13. 04. 2023), (Bsp. unter <https://voyant-tools.org/docs/#!/guide/gallery>); als auch Arbeiten, die neue Werkzeuge zur Beantwortung geisteswissenschaftlicher Fragestellungen entwickeln oder verbessern. Vgl. z. B. Adam KILGARRIFF et al. 2004: „The Sketch Engine“. In: *Proceedings of the 11th EURALEX International Congress*. Hrsg. von Geoffrey WILLIAMS und Sandra VESSIER. Lorient, France: Université des lettres et des sciences humaines, S. 105–115.

20 Siehe Andrew MURPHY 2003: *Shakespeare in Print: A History and Chronology of Shakespeare Publishing*. Cambridge: Cambridge University Press, S. 209–210. Siehe dazu auch Kapitel 3.1, S. 41.

21 Siehe Benjamin MANGRUM 2018: „Aggregation, Public Criticism, and the History of Reading Big Data“. In: *PMLA* 133.5, S. 1207–1224. DOI: 10.1632/pmla.2018.133.5.1207, 1207. Obwohl BUSA zunächst nicht selbst programmierte, wird er als einer der Gründungsväter der Computerlinguistik gefeiert.

22 Siehe auch Manfred THALLER 2017: „Geschichte der Digital Humanities“. In: *Digital Humanities – Eine Einführung*. Hrsg. von Fotis JANNIDIS, Hubertus KOHLE und Malte REHBEIN. Stuttgart: Metzler, S. 3–12, S. 3–4.

23 Siehe Frederick MOSTELLER und David L. WALLACE 1984 [1964]: *Applied Bayesian and Classical Inference – The Case of The Federalist Papers*. 2. Aufl. New York: Springer, z. B. S. 10–14.

24 Weitere Beispiele, sowie ein Abriss über unterschiedliche Schwerpunkte und Trends der vergangenen 50 Jahre findet sich z. B. in THALLER 2017, einige aktuellere Arbeiten werden in Kapitel 3.1 und mit Schwerpunkt auf das Chinesische in 4.1 vorgestellt.

25 Den entscheidenden Impuls für die Wahl von *Python* gaben ein Workshop sowie die Arbeit von Paul VIERTHALER, siehe auch Kapitel 4.1, ab S. 60. Das vorliegende Projekt wurde anfangs in *Python 2.7* entwickelt, jedoch wegen der deutlich verbesserten nativen Unicode-Unterstützung zur Version 3.8 gewechselt. Das vorher umständliche En- und Decodieren der Zeichenrepräsentanzen (siehe auch Kapitel 4.3, S. 69) entfällt seit Version 3 völlig – bzw. geschieht automatisch. Vgl. auch Paul VIERTHALER 2020: „Digital humanities and East Asian studies in 2020“. In: *History Compass* e12628. DOI: 10.1111/hic3.12628, S. II (EN 7).

26 Fotis JANNIDIS 2017a: „Grundbegriffe des Programmierens“. In: *Digital Humanities – Eine Einführung*. Hrsg. von Fotis JANNIDIS, Hubertus KOHLE und Malte REHBEIN. Stuttgart: Metzler, S. 68–95, S. 68.

nungen entwickelt, insbesondere auch für die Analyse und Visualisierung großer Datensätze“,²⁷ z. B. mit der Bibliothek *pandas*²⁸, sowie zur Lösung von Problemstellungen der quantitativen Linguistik.²⁹ Für die Verwendung von *Python* zur Bearbeitung von Fragestellungen der historischen Linguistik des Chinesischen spricht zudem die Verfügbarkeit einer wachsenden Anzahl spezialisierter Bibliotheken.³⁰

Zur Berechnung und Bewertung statistischer Zusammenhänge wird *R* verwendet, eine Programmiersprache, die sich im Bereich der Statistik als wissenschaftlicher Standard etabliert hat. Die *R*-Programmiersprache *ggplot2* wird zudem zur Visualisierung von Daten und Ergebnissen eingesetzt.³¹

Schriftsprachliches Chinesisch

Da die zeitliche Einordnung von Texten aus einem Zeitraum von den Anfängen des klassischen Schrifttums (erstes Jahrtausend v. u. Z.) bis ins 20. Jh. Gegenstand dieser Arbeit ist, wird – in Anlehnung an die ebenfalls weit gefassten Begriffe *wenyan* (*wen*) 文言 (文) bzw. *Literary Chinese* – hier der Begriff *schriftsprachliches Chinesisch* bzw. *Chinesische Schriftsprache* als Sammelbegriff für alle vormodernen Entwicklungsstufen des Chinesischen verwendet. Dies geschieht in Abgrenzung zur modernen Hoch- bzw. Umgangssprache (*putonghua* 普通話 und *guoyu* 國語 bzw. *kouyu* 口語). Häufig werden – auch in der sinologischen Literatur – die Bezeichnungen *Klassisches Chinesisch* bzw. *Classical Chinese* und *Literary Chinese* synonym verwendet, insbesondere im anglo-amerikanischen und chinesischen Sprachgebrauch. Im Deutschen bezeichnet *Klassisches Chinesisch* im Wesentlichen aber das antike Schrifttum – die aus Textzeugnissen etwa der zweiten Hälfte des ersten Jahrtausends v. u. Z. überlieferte Schriftsprache.³² Im Englischen wird der Begriff auch von Sprachwissenschaftler:innen etwas weiter gefasst: „Classical Chinese is a conventional way of

27 Jake VANDERPLAS 2018: *Data Science mit Python – Das Handbuch für den Einsatz von IPython, Jupyter, NumPy, Pandas, Matplotlib, Scikit-Learn*. Übers. von Knut LORENZEN. 1. Aufl., S. 14.

28 Als Bibliotheken bzw. *libraries* werden fertige Programmpakete oder Funktionen bereitgestellt, die beliebig innerhalb eigener Programme eingesetzt werden können. *pandas* etwa vereinfacht die Analyse und Manipulation von Daten in tabellenähnlichen Objekten. Siehe Wes MCKINNEY 2010: „Data Structures for Statistical Computing in Python“. In: *Proceedings of the 9th Python in Science Conference*. Hrsg. von Stéfan van der WALT und Jarrod MILLMAN, S. 56–61. DOI: 10. 25080/Majora-92bf1922-00a; THE PANDAS DEVELOPMENT TEAM 2020: *pandas 1.5.0*. DOI: 10. 5281/zenodo.3509134.

29 Ralf JÜNGLING und Gabriel ALTMANN 2003: „Python for linguistics?“ In: *Glottometrics* 6, S. 70–82.

30 Einige wichtige Beispiele sind der Tokenizer *Jieba*, siehe SUN Junyi 2018: *Jieba zhongwen fenci* 结巴中文分词 (*Jieba Chinesisch-Tokenizer*). GitHub Repository. URL: <https://github.com/fxsjy/jieba> (besucht am 26.02.2019), ausführlicher siehe Kapitel 4.5, ab S. 83; *sinopy*, das Funktionen zur Konvertierung von Zeichen in versch. Umschriften bereitstellt, wobei nicht nur *Hanyu Pinyin* 漢語拼音, sondern auch die historisierende Umschrift für das mittelchinesische von BAXTER, sowie das *International Phonetic Alphabet (IPA)* zur Verfügung stehen. Siehe Johann-Mattis LIST 2018: *SinoPy: Python Library for quantitative tasks in Chinese historical linguistics*. Jena. URL: <https://pypi.org/project/sinopy/> (besucht am 26.04.2020); sowie *mafan*, das Anwender:innen *mafan* 麻煩 („Unannehmlichkeiten“) ersparen soll. Die Bibliothek enthält Tools zur Umwandlung zwischen und Erkennung von traditionellen und vereinfachten Zeichen, sowie zum Umgang mit Codierungen. Siehe Herman SCHAAF 2017: *Mafan – Toolkit for working with Chinese in Python*. Python module. URL: <https://pypi.org/project/mafan/> (besucht am 26.04.2020).

31 Hadley WICKHAM 2016: *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer. URL: <https://ggplot2.tidyverse.org>.

32 Siehe z. B. Ulrich UNGER 1985a: *Einführung in das Klassische Chinesisch, Teil I: Allgemeines, Chinesische Texte, Indices*. Wiesbaden: Harrassowitz, S. 1; für eine ausführliche Diskussion der Begrifflichkeit siehe z. B. Tilman SCHALMEY 2009: „Überlegungen zur Konzeption eines neuen Lehrbuchs für das Klassische Chinesische“. Magisterarbeit. München: Ludwig-Maximilians-Universität, S. 7–11; Mit dem Problem, dass die Definition einer „Klassischen“ Sprache *flexibel* ist, haben nicht nur Sinolog:innen zu kämpfen, ähnlich schwammige Terminologien existieren auch für das Hebräische. Siehe z. B. Ian YOUNG und Robert REZETKO 2014: *Linguistic Dating of Biblical Texts*. Bd. 1. Routledge, S. 8.

1 Einleitung

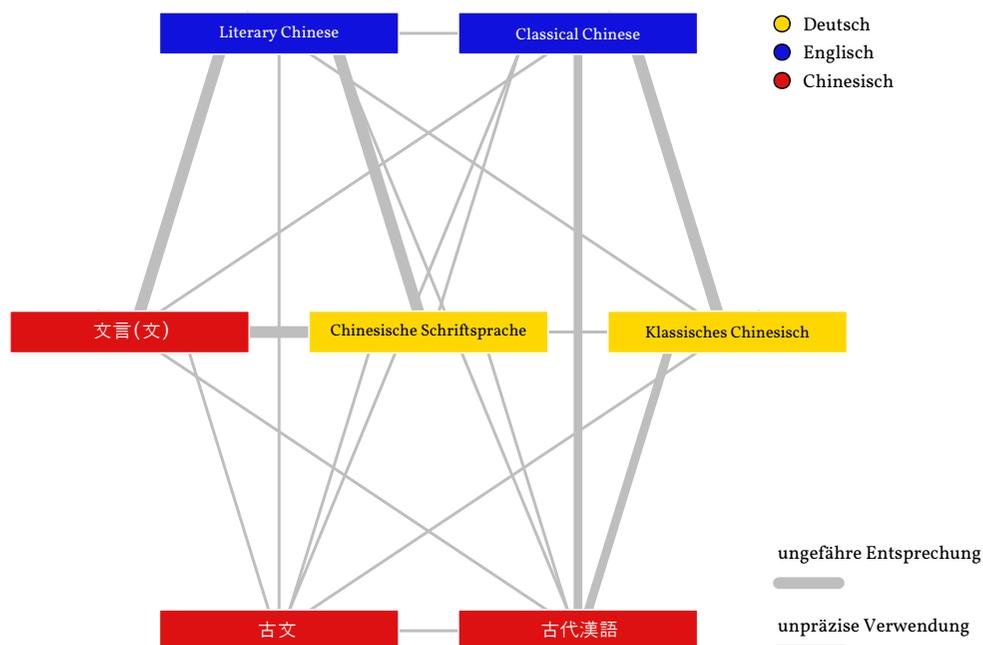


Abbildung 1.1 Klassisches und schriftsprachliches Chinesisch: Übersetz- und Austauschbarkeit der Begriffe

referring to the written form of Old Chinese, the period from the Spring & Autumn period to the end of the Han Dynasty.³³

Die problematische Austauschbarkeit dieser eigentlich unterschiedlichen Bezeichnungen *Literary* und *Classical Chinese* (Abb. 1.1) hat ihren Ursprung sicherlich in der Dehnbarkeit der chinesischen Begriffe *gudai Hanyu* 古代漢語, *wenyan (wen)* 文言(文) und *guwen* 古文. Vor allem die beiden letzteren implizieren – wie auch *schriftsprachliches Chinesisch* – häufig keinerlei epochale Einteilung. Stattdessen beziehen sie sich teils eher auf einen schriftsprachlichen *Stil* – im Kontrast zur (gesprochenen) Umgangssprache (*baihua* 白話)³⁴ der jeweiligen Zeit – oder, wie von der GABELNTZ formuliert, Chinesisch „mit Ausschluss des niederen Stiles und der heutigen Umgangssprache“.³⁵ Abgesehen von diesen temporalen und stilistischen Aspekten ist bereits das Konzept des Klassischen Chinesischen an sich eine „kühne Vereinfachung der sprachlichen Vielfalt, die im alten China existierte“,³⁶ denn bereits für die Zeit der Frühlings- und Herbstannalen (*Chunqiu shidai* 春秋時代, ca. 770–476 v. u. Z.)³⁷ sind dialektale Unterschiede bzw. unterschiedliche Sprachen textlich belegt.³⁸ Über die Dialekte dieser Zeit ist jedoch zu wenig bekannt, um sie in einer Untersuchung computerlinguistischer Datierungsmethoden separat zu betrachten. Zudem ist es „offensichtlich,

33 Jerry NORMAN 1988: *Chinese*. Cambridge: Cambridge University Press, S. 83. NORMANS Definition ist tatsächlich sinnvoller und präziser, da die uns vorliegenden Fassungen klassischer Texte zumeist durch eine frühestenfalls Han-zeitliche (漢, 202 v. u. Z.–220) Redaktion gegangen sind. Siehe z. B. Martin KERN 2004: „Die Anfänge der chinesischen Literatur“. In: *Chinesische Literaturgeschichte*. Hrsg. von Reinhard EMMERICH. Stuttgart: Metzler, S. 46.

34 Für schriftliche Zeugnisse dieser Umgangssprache sind auch die Begriffe *tongsu* 通俗 bzw. im Englischen (*written vernacular*) gebräuchlich.

35 Hans Georg Conon von der GABELNTZ 1881: *Chinesische Grammatik: mit Ausschluss des niederen Stiles und der heutigen Umgangssprache*. Leipzig: T. O. Weigel, S. U1.

36 Kai VOGELSANG 2021: *Introduction to Classical Chinese*. Oxford: Oxford University Press, S. 363, übersetzt durch den Verfasser.

37 WILKINSON 2000, S. 10.

38 Siehe NORMAN 1988, S. 208–209; siehe auch VOGELSANG 2021, S. 363–364.

dass eine Art vornehmer Standard entstanden war“.³⁹ Geographische Aspekte der Sprachentwicklung sollen hier also außer Acht gelassen werden.

Inhaltsübersicht

Ungeachtet der sprachgeschichtlich unbefriedigenden Ungenauigkeit des Begriffs *schriftsprachliches Chinesisch* unterteilt man die Epochen in der historischen Linguistik selbstverständlich feiner. Kapitel 2 (ab S. 11) ist der Erarbeitung eines Grundverständnisses chinesischer Sprachgeschichte, des Sprachwandels und insbesondere des Wortschatzwandels gewidmet. Es fällt auf, dass in der historischen Sprachwissenschaft – der traditionellen chinesischen wie der westlichen⁴⁰ – besonders der Aspekt des phonologischen Wandels für das Chinesische intensiv untersucht wird, während der lexikalische Wandel eher stiefväterlich behandelt wird.⁴¹ Anhand der offiziellen Dynastiegeschichten (*zhengshi* 正史) werden zudem beispielhaft konkrete Erscheinungsformen des syntaktischen und vor allem des lexikalischen Wandels beobachtet (Kapitel 2.3, ab S. 20). Dieses Korpus ermöglicht solche diachronen Betrachtungen innerhalb einer stilistisch gefestigten Textgattung über einen Zeitraum von mehr als 2.000 Jahren.

Kapitel 3 (ab S. 35) gibt einen Überblick über die für westliche Sprachen zur Verfügung stehenden computerlinguistischen Textdatierungsmethoden. Besondere Aufmerksamkeit wird dabei den eingangs erwähnten statistischen Sprachmodellen zuteil. Die bestehenden Ansätze werden mit Blick auf ihre Eignung für die Anwendung auf chinesischsprachiges – vor allem schriftsprachliches bzw. klassisches Textmaterial beleuchtet.

In Kapitel 4 (ab S. 59) wird geprüft, ob und wie gut sich bestehende computerlinguistische Ressourcen, Methoden und Tools für die Verarbeitung schriftsprachlicher Texte nutzen lassen. Bevor geeignet erscheinende Datierungsmethoden angewandt werden können, müssen passende diachrone Test- und Trainingskorpora festgelegt werden (Kapitel 4.2, ab S. 62). Überdies muss der Einfluss besonderer Eigenschaften der chinesischen Schrift und Schriftsprache auf die Verwendung computerlinguistischer Werkzeuge berücksichtigt werden. Während z. B. das fast vollständige Fehlen von Flexion⁴² von Vorteil sein kann, erschweren Zeichenvarianten (Kapitel 4.3, ab S. 69) und vor allem das Fehlen von Leerzeichen zur Markierung von Wortgrenzen quantitative Analysen. Ein zentraler Aspekt ist daher die zur Erstellung von Worthäufigkeitslisten notwendige Möglichkeit der Segmentierung bzw. Tokenisierung von Texten (Kapitel 4.4–4.6, ab S. 73). Implikationen ergeben sich auch für die Erkennung von Personennamen (Kapitel 4.7, ab S. 97) und *temporal expressions* (Kapitel 4.8, ab S. 103).

In Kapitel 5 (ab S. 107) werden Entstehungsgeschichte und Datenqualität des *HYDCD* untersucht. Neben den rein lexikographischen Informationen, Lexemen und ihren Bedeutungen,

³⁹ NORMAN 1988, S. 209, übersetzt durch den Verfasser.

⁴⁰ Während eine professionelle historische Linguistik im Westen erst eine Erscheinung des späten 18. Jahrhunderts ist (siehe Roger LASS 2014: „Lineage and the Constructive Imagination: The Birth of Historical Linguistics“. In: *The Routledge Handbook of Historical Linguistics*. Hrsg. von Claire BOWERN und Bethwyn EVANS. London & New York: Routledge, S. 45–63, v. a. S. 45–48;), deutet sich mit der in China bereits zur Han-Zeit florierenden Kultur der Kommentarliteratur mit ihren erklärenden Glossen für die Aussprache und Bedeutung von Zeichen in früheren, kanonischen Texten ein frühes Bewusstsein von Sprachwandel an. Siehe z. B. KERN 2004, S. 82–87. Solche Kommentare im Kontext der *xiaoxue* 小學 (Philologie) stellen damit eine frühe Form der historischen Sprachwissenschaft dar.

⁴¹ Siehe James H.-Y. TAI und Marjorie K. M. CHAN 1999: „Some Reflections on the Periodization of the Chinese Language“. In: *Studies in Chinese Historical Syntax and Morphology: Linguistic Essays in Honor of Mei Tsu-lin*. Hrsg. von Alain PEYRAUBE und SUN Chaofen. Paris: École des Hautes Études en Sciences Sociales, S. 223–239, S. 227, siehe auch S. 233.

⁴² Siehe Axel SCHUESSLER 2015: „Old Chinese Morphology“. In: *Encyclopedia of Chinese Language and Linguistics Online*. Hrsg. von Rint SYBESMA et al. Leiden: Brill, für eine aktuelle Diskussion über das Vorhandensein von Flexion in früheren Entwicklungsstufen des Chinesischen.

enthält es zahlreiche Textbelegstellen, die – bestenfalls – die früheste Verwendung dieser Wörter dokumentieren. Daraus wird eine diachrone *mySQL*⁴³-Lexemdatenbank aufgebaut. Die so strukturierten Daten können bei geringer Redundanz bequem für unterschiedliche Zwecke abgefragt werden. Die spärlichen bibliographischen Angaben im *HYDCD* werden unter anderem mit Daten aus dem *China Biographical Database Project (CBDB)*⁴⁴ verknüpft, was eine genauere chronologische Einordnung der lexikalisierten Wörter anhand ihrer Belegstellen ermöglicht. Aus den Textbelegen wird zudem ein provisorisches, diachrones Behelfskorpus für den Zeitraum vom 7. Jh. v. u. Z. bis ins 20. Jh. erzeugt (Kapitel 5.6, S. 137).

Die gewonnenen Daten werfen ein neues Licht auf die Machart dieses wichtigen Standardwerks und ermöglichen darüber hinaus Beobachtungen zum Wortschatzwachstum und -wandel des Chinesischen (Kapitel 5.7, ab S. 138).

Kapitel 6 (ab S. 155) widmet sich der Entwicklung und Anwendung von computerlinguistischen Datierungsmethoden für schriftsprachliches Chinesisch. Die statistischen Methoden aus der westlichen Computerlinguistik werden erstmals für die Datierung chinesischsprachiger Texte implementiert und mit verschiedenen (Behelfs-)korpora getestet (Kapitel 6.1, ab S. 156). Temporale Sprachmodelle erweisen sich dabei als für die chronologische Einordnung schriftsprachlicher Texte grundsätzlich geeignet. Versuche mit dem *HYDCD*-Arbeitskorpus zeigen, dass auch über einen langen Betrachtungszeitraum hinweg eine grobe Einordnung möglich ist. Der Erfolg bleibt aber stark vom Stil der zu datierenden Texte und den verfügbaren Trainingsdaten abhängig. Während das Verständnis von Sprachwandel ein Fundament der linguistischen Textdatierung darstellt, können die verwendeten statistischen Modelle andersherum auch genutzt werden, um Sprachwandel zu erkennen.⁴⁵

Wünschenswert wäre eine von Genre und Trainingskorpus unabhängige Datierung schriftsprachlicher Texte. Eine lexikographische Herangehensweise auf Basis der diachronen Lexemdatenbank könnte dies ermöglichen. In Kapitel 6.2 (ab S. 179) wird eine Methode vorgestellt, anhand derer Texte aufgrund der enthaltenen Zeichenkombinationen chronologisch eingeordnet werden. Auch Personennamen und die Erkennung von *temporal expressions* können hierfür eingesetzt werden. Die so erzeugten temporalen Textprofile eignen sich neben einer automatisierten Datierung auch als Ausgangspunkt für eine qualitative Analyse.

Darüber hinaus wird untersucht, ob sich auch eine stark abstrahierte Messgröße, die durchschnittliche Entstehungszeit der in einem Text nachgewiesenen Lexeme (*Average Year of Lexicalization, AYL*) eignet, um Rückschlüsse auf seine Entstehungszeit zu ziehen (Kapitel 6.3, ab S. 210). Experimente offenbaren einen linearen Zusammenhang zwischen *AYL* und Textgenese, der zur Datierung aber nur bedingt herangezogen werden kann.

In Kapitel 6.4 (ab S. 229) wird ein Vergleich der vorgestellten linguistischen Datierungsmethoden angestrebt und auf Vor- und Nachteile der unterschiedlichen Herangehensweisen eingegangen. Zuletzt entwickle ich ein *user interface*, das mit geringen Vorkenntnissen die Verwendung der erarbeiteten Methodik im *Browser* ermöglicht. Benutzer:innen erhalten damit

43 *mySQL* ist eine verbreitete *Structured Query Language (SQL)*. *SQLs* sind Programmiersprachen zur Formulierung von Datenbankabfragen. Die Daten werden in einer tabellenartigen Struktur gehalten und können so mit allen verbreiteten, modernen Programmiersprachen problemlos gelesen, geschrieben, sortiert, durchsucht, miteinander verknüpft und transformiert werden. Siehe z. B. Edwin SCHICKER 2017: *Datenbanken und SQL*. 5. Aufl. Wiesbaden: Springer Vieweg, S. 3–7.

44 Michael A. FULLER 2017: *China Biographical Database Project (CBDB)*. URL: <https://projects.iq.harvard.edu/cbdb> (besucht am 24. 04. 2017) (im Folgenden zit. als *CBDB*).

45 Siehe DE JONG, RODE und HIEMSTRA 2005, S. 1. Diese Perspektive wird von den Autor:innen zwar angesprochen, aber nicht weiter beleuchtet und findet auch in der auf dieser Studie aufbauenden Arbeit anderer Forscher:innen wenig Beachtung. In Kapitel 6.1.4 (ab S. 177) wird diese Idee aufgegriffen.

Anhaltspunkte für die zeitliche Einordnung eines unbekanntes Texts. Die Arbeit einer Philolog:in kann damit nicht ersetzt, aber unterstützt und erleichtert werden.

In Kapitel 7 (ab S. 243) werden wichtige Erkenntnisse dieser Studie zusammengefasst und Limitationen der (computer-)linguistischen Datierung schriftsprachlicher chinesischer Texte diskutiert. Abschließend werden Ideen für künftige Erweiterungen und Verbesserungen lexikographischer und statistischer Datierungsmethoden skizziert.

