

# 古籍的数字赋能与增值利用

## ——“数智时代的古籍活化与利用”高端论坛述评

### Digital Empowerment and Value-added Utilization of Ancient Books: Review of High-end Forum on Activation and Utilization of Ancient Books in the Digital Intelligence Era

林通<sup>1</sup> 郑翔<sup>1,2</sup> 李明杰<sup>1,2</sup>  
LIN Tong ZHENG Xiang LI Mingjie

(1. 武汉大学信息管理学院, 武汉, 430072; 2. 武汉大学文化遗产智能计算实验室, 武汉, 430072 / 1. School of Information Management, Wuhan University, Wuhan, 430072; 2. Cultural Heritage Intelligent Computing Laboratory of Wuhan University, Wuhan, 430072)

中图分类号: G25 DOI: 10.13366/j.dik.2024.02.081

引用本文: 林通, 郑翔, 李明杰. 古籍的数字赋能与增值利用——“数智时代的古籍活化与利用”高端论坛述评[J]. 图书情报知识, 2024, 41(2): 81-86. (Lin Tong, Zheng Xiang, Li Mingjie. Digital Empowerment and Value-added Utilization of Ancient Books: Review of High-end Forum on Activation and Utilization of Ancient Books in the Digital Intelligence Era[J]. Documentation, Information & Knowledge, 2024, 41(2): 81-86.)

2023年4月21-23日, 武汉大学文化遗产智能计算实验室联合武汉大学大数据研究院、文学院、古籍所、信息管理学院, 以及新闻出版署语义出版与知识服务实验室等单位举办的“数智时代的古籍活化与利用”高端论坛在武汉大学隆重召开。来自国内高校、图书馆和出版社等领域的专家学者齐聚一堂。本次论坛旨在贯彻落实中共中央办公厅、国务院办公厅《关于推进新时代古籍工作的意见》的文件精神, 推进古籍抢救性保护、整理研究、编辑出版、普及推广和人才培养等工作。围绕数智时代的古籍整理与研究、古籍活化利用的前沿理论与先进技术、数智时代的古籍出版与再造三大主题, 与会专家进行了深入的对话和交流。

## 1 数智时代的古籍整理与研究

古籍文献作为传承中华古代文明的重要载体, 其整理研究的传统由来已久, 且已累积了相当丰富的理论经验和学术成果。近年来, 数智化浪潮不可逆转地席卷了整个人类文明, 同时也为传统的古籍整理研究领域带来了新的发展机遇。站在数智时代的风口浪尖, 今后的古籍整理研究工作究竟该何去何从, 已经成为摆在广大学人面前的一个重大课题。

### 1.1 古籍整理研究的历史经验与文献学的数字化转向

首都师范大学文学院南江涛副教授回顾了新中国成立以来古籍编目工作的三次高潮: 第一次是1956年在“向科学进军”的号召下, 全国积极组织编纂古籍书目; 第二次是1978年为贯彻周恩来总理“要尽快把全国善本书目编出来”的指示而开始编纂《中国古籍善本书目》; 第三次是在2007年国务院办公厅发布《关于进一步加强古籍保护工作的意见》后, 为筹备编纂《中华古籍总目》所开展的系列工作。他对将古籍的年代下限刻板地限定于1912年、域外汉籍整理质量整体偏低等问题进行了反思, 提出应当从作者生平和作品内容形式着手界定古籍, 通过提高域外汉籍的整理标准来提升古籍整理质量。而针对影印、点校和注释等不同整理方式, 他给出了分层、分级开展古籍整理的思路, 并对行业内古籍影印不规范的现象提出了中肯的改进建议。最后, 南江涛副教授明确了不同类型读者在古籍工作中的实际角色: 专业读者是古籍保护和整理的受益者, 可能转换为参与者和生产者, 是双向利用古籍的主体力量, 而普通读者则是蕴含在古籍中的传统文化的接受者。

清华大学人文学院院长刘石教授作了题为“文献学的数字化转向与‘中国古典知识库’建设”的报告。

[通讯作者] 李明杰 (ORCID: 0000-0002-1876-9040), 博士, 教授, 研究方向: 文献学、中国图书文化史, Email: limingjie@whu.edu.cn. (Correspondence should be addressed to LI Mingjie, Email: limingjie@whu.edu.cn, ORCID: 0000-0002-1876-9040)

[作者简介] 林通 (ORCID: 0000-0001-6798-6570), 博士研究生, 研究方向: 古籍整理与保护, Email: qingmls@qq.com; 郑翔 (ORCID: 0000-0001-7933-2822), 博士研究生, 研究方向: 古籍整理与数字人文, Email: zhengxiang059@whu.edu.cn.

他特别强调传统文献与现代数据之间的内在关联,并从文献生产的创革、文本形态的新变、知识获取的拓展等方面对大数据技术推动传统文献学现代转型的作用进行分析:传统文献学较重视经验与思辨,大数据研究更依赖工具与技术,突破了传统文献的生产方式,实现了文献形态的再发现与再生产;文献计量单位从部、册、卷、篇、页、段、行、句等,向电子文本、文本集、数据库、知识库、系统平台等新文献形态转变,向基本储存单元、扩展存储单元等转变,打破传统文献的线性平面结构,完成空间化与可视化,实现跨文本异质同构;传统文献通过篇章划分、页码标记等方式规定阅读顺序,大数据则突破了文本细读和个案研究等传统方式,在知识关联与计量、主题模型提取等方面具有优势。另外,他从国家文化发展战略的高度提出古籍数字化向知识化转变的“中国古典知识库”的构想,即运用计算语言学和自然语言处理前沿技术,使用人工智能、大数据等技术工具和手段,借助过往一切古典学研究成果,周密地设定主题词表,专业地提取各种实体,多维度地构建不同实体间的关系,并通过这些实体及相互关系,在保障古籍文献内容完整性及内部逻辑性的基础上突破文献原有结构,对文献进行深层组织和知识管理,将现存所有古籍建构成关系性、结构化、知识再生型的超大知识库。未了,刘石教授还展示了清华大学在数字人文方面的研究和教学实践经验。

## 1.2 古籍整理研究数智化转型的方向与路径

无论是传统学术范式的人文领域,还是近年来新兴的数字人文领域,古籍整理与研究始终是古籍工作者关注的重要内容。教育部全国高等院校古籍整理研究工作委员会副秘书长、北京大学中文系吴国武副教授认为,古籍数字化首先需要关注数据的质量,经过整理研究、包含优质标注的古籍数据可以提升算法能力,促进古籍数字化的高质量发展。他总结了当前古籍数字化的发展阶段,即从古籍资源的数字化及智能开发,向古籍收藏保护、整理研究、出版利用等全过程的数智化转型,同时也与新兴的数字人文深度融合,合力推进中国自主知识体系、学术范式的构建。此外,吴国武副教授还从历史经验和内在逻辑、功能性质与发展历程、古籍数字化新使命和数字人文研究等角度展开分析,提出古籍数字化的重心应当回归“深度整理”和“综合研究”的命题。他在回顾当前古籍数字化诸多

实践成果的基础之上,分析了古籍数字化整理与研究存在的问题:其一,古籍整理研究的主体作用有所削弱;其二,古籍整理研究专家团队较少参与,更起不到主导作用;其三,现有古籍整理研究成果没有在世界范围内得到共享和有效利用;其四,古籍工作全过程衔接与数字人文研究融合远远不够。同时,他提出了“整理研究”数智化转型的几个重要方向:一是古籍数据资源和数据标注的优化;二是以优质数据提升算法水平,不断生成新的优质数据;三是转向中国自主知识体系的建构和中华文明的全景式呈现。最后,吴国武副教授还特别呼吁数智时代古籍整理研究与数字人文的深度融合。

武汉大学信息管理学院李明杰教授回应了吴国武副教授关于古籍数字化向“深度整理”和“综合研究”回归的倡议。他首先将古籍整理划分为实体保存性整理、文本复原性整理、语义阐释性整理、内容组织性整理和知识数据化整理等五个层次,然后从古籍整理五层次的各个环节入手,逐一剖析数字人文技术在古籍研究、古籍推广和古籍整理三大领域的具体应用,特别举出运用不同的数据统计方法研究《红楼梦》前八十回和后四十回的作者所得出的结论存在较大偏差的实例,说明目前的数字人文还无法完全取代传统的人文学术研究,进而强调应将数字人文应用的重点回归到古籍整理领域。李明杰教授指出,传统古籍整理领域的数字人文有两大实现路径:一是古籍整理思想与方法的数据化再现。即将数字人文技术引入版本鉴定、校勘、辨伪、辑佚、标点、注释、翻译、编目、编纂等传统的古籍整理方法中去,同时应在学术范式上融入古典文献学在长期的学术实践中形成和积淀下来的优良学术传统;二是古籍整理成果的数据化再利用。前代遗留下来的丰富的古籍整理成果,如书目、题跋、索引及类书、政书、辞典、方志、谱牒等各种工具书,本身就是为了帮助读者研究和利用古籍,它们在数据化之后更容易实现古籍知识的增值利用。最后,他以亲身科研经历为例,介绍了梅尧臣《宛陵集》版本源流的可视化呈现、明代古籍版刻地理信息系统的构建等实践案例。

大数据、人工智能等新兴技术逐渐被应用于古籍整理与研究的具体实践之中,为校勘、辨伪、标点、注释、编目等传统古籍整理方法的数据化再现提供了智能化的解决思路。但同时,计算机技术水平的提升也对古籍基础数据质量、算法水平和行业标准等相关条

件提出了更高要求。数智时代背景下的古籍整理与研究,一方面要顺应新兴技术发展的浪潮,牢牢把握时代机遇,提升古籍智能整理与研究水平;另一方面还必须关注古籍数字化的现实问题,充分发挥行业和学界的主观能动性,建立精细的行业标准与学术规范,进而全面提升古籍整理研究的整体质量。

## 2 古籍活化利用的前沿理论与先进技术

数智时代对古籍文献的开发利用,既要不断丰富和完善基础理论,又要致力于对数字技术的改良与提升,通过理论知识和智能化手段的有机结合,达到让古籍真正“活起来”的效果。

### 2.1 古籍活化利用的理论反思与设想

武汉大学图书馆古籍保护中心主任周荣教授从市场、公益展示、地方文化宣传、文化普及等角度探讨了当前古籍活化利用的多重维度,对当前古籍“活化利用”语境中传统学术研究被遮掩的现状进行反思,从三个方面重申了传统学术研究在古籍活化利用中的基础地位:其一,“学术”(道)是古籍生命体和中华文明生命力的根源所在;其二,传统学术为古籍活化利用提供了学理保障,表现在音韵、文字、训诂等小学为古籍活化利用“赋形”,注疏、版本、目录、校勘等古典文献学为古籍活化利用“结体”;经史义理、文学艺术、哲学宗教等新旧学术为古籍活化利用“凝神”;其三,古籍活化并不是一个新的学术命题,它一直是传统学术研究追求的目标之一。此外,他结合当前学界对古籍活化利用的探索经验及本人工作实践,谈了对利用数智新技术推进传统学术深化、创新和转型的看法。

四川大学文学与新闻学院王兆鹏教授以文学类古籍特别是古诗文的活化利用为例,指出文学古籍有三种活化形式:一是将古籍由静态变为动态,让单一的文本形态变成可复制、可重排、可检索的“文献典籍活化”;二是利用元宇宙、VR、AI绘画等技术,让文学作品实现可视化呈现的“文学场景活化”;三是将文化资源进行创造性转化、创新性发展,促进文学资源转化为文旅产品的“文学资源活化”。另外,他还提出了文学古籍的“五种融化”:一是打通古诗、新诗界限,探寻不同时空诗歌根源的“融化古今”;二是汇集外诗中译、

中诗外译等各类文本的“融贯中外”;三是将诗歌意境、场景,动植物的形声、形色立体化呈现的“融汇视听”;四是融合经史子集四部书籍,按需提取不同类别资料的“融化学科”;五是呈现诗歌中饮食、服饰、民俗等文化的“融通文化”。最后,王兆鹏教授提出了融合、贯通文学图谱平台与《汉语大字典》智慧平台的构想。

武汉大学信息管理学院徐雷副教授则从跨学科的视角分享了他对古籍研究领域的认知。他认为,古籍研究不同于一般的自然科学和社会科学研究,属于事实型的自然发现和基于历史证据的认知观点,且研究周期较长,学术成果多以专著的形式呈现。因此,应当加强古籍相关学术成果的利用。当作为本体的原始古籍获取困难时,古籍学术文献便可以成为它的有效替代。他以科技文献的智能化处理为例,从Science IE 软件的知识识别抽取等方面,阐述了古籍学术文献的知识挖掘思路,并特别介绍了元数据化、领域术语词汇抽取、领域实体及关系识别、篇章结构识别、语义组织与活化利用等系列处置办法。最后,徐雷副教授从三个方面总结了古籍学术文献智能处理面临的挑战:一是图书存在版权问题和非一手资料的质量问题;二是学术论文零散不成体系,不易活化;三是资源建设成本偏高,数据评审和交互体验存在实际困难。

### 2.2 古籍活化利用技术的最新探索

现代信息技术的不断升级迭代,是推动古籍数字化开发和利用进程的强大助力。华南理工大学电子与信息学院金连文教授介绍了其团队面对古籍文档版式复杂多样、逻辑版面解析困难、双列夹注、图文混排、大小字密集排版等端到端古籍文档智能识别及结构化理解的挑战,从数据预处理、文本行检测、端到端连接、文本行识别、文字阅读顺序理解等方面,设计了古籍文档图像智能结构化识别解决方案。该方案在2022年大湾区算法算例国际大赛——古籍OCR赛道中取得了综合排名第一的好成绩。

南京农业大学信息管理学院王东波教授提出基于《四库全书》语料预训练的GPT模型——SikuGPT。该模型在以《二十四史全译》为语料的机器翻译下游任务中取得了最佳的性能,且测评效果优于ChatGPT(zero-shot)等其他GPT类模型。此外,他还介绍了其团队构建的“中国古代典籍跨语言知识库平台”,该平台不仅

可以实现在线平行语料应用、典籍智能处理与呈现的功能,而且能提供语料申请与模型下载的服务。

构思精巧的高质量数据库(平台)能为后续的古籍整理研究和活化利用提供更多的可能性。“国学网”的创始人——首都师范大学电子文献研究所所长尹小林副教授提出“元古籍”的新概念,并从古籍创作、编刻与流通年代的不统一,古籍文字形态的不固定,以及古籍体例的特殊性等情况入手,着重分析了建设保留古籍原始形态的元古籍数据库的必要性,指出以海量元数据作为基础,可以通过分类实现快速搭建各种专题数据库的效果,并演示了元古籍数据库在查询字词出现的源流始末、查询若干字词出现的特定位置、字词统计等方面的学术应用场景。

古籍数据化、知识化目标的最终实现离不开对古籍内涵的智能化挖掘、深层次揭示与可视化呈现等具体步骤,而图像识别技术是进行古籍数据化的重要基础,预训练模型的应用则进一步实现了古籍语义信息的提取与表征,并帮助提升古籍自动识别、自动翻译等任务的效率,同时也为准确、高效地推进古籍文本深层次分析等工作提供了有力保障。在此基础上,借助可视化、可交互的技术,融合多源数据和古籍背景知识,能够将古籍深层内涵塑造为更生动形象且易于传播、理解的多样知识形态。

### 3 数智时代的古籍出版与再造

古籍的出版再造肩负着传承中华优秀传统文化的历史使命与社会责任。数智时代的到来,在保留传统纸质图书出版的基础上,又拓展出在线专题数据库的全新出版形式。数字化已成为实现古籍出版内容价值增值的重要路径,同时也为传统的古籍出版提供了新的启示。

中华书局古联公司洪涛总经理首先介绍了由其公司开发运营的集各类专业数据库产品、智能化辅助工具、在线众包整理平台、古籍人才培养平台为一体的国家级古籍整理出版资源平台——“籍合网”的基本情况。而后,他重点展示了古联公司与南京农业大学协同开发的方志物产项目。该项目以手抄本《方志物产》为底本,分为图书及数字出版、技术研发与知识库建设、产品融合与市场应用三个阶段目标:其一,整理、

影印南京农业大学的重要馆藏资源,建设“方志物产全文数据库”;其二,开展中国古代物产文献的智能化研究与自动化处理,利用智能技术进一步挖掘古代物产文献;其三,同领域专家合作,将古代物产与现代物种资源对接,服务于农、牧、林、渔等领域,同时在更大范围内与文化、旅游、创意产业对接,促进古籍数字化业务朝向古今融通、产业融合的方向发展。

上海古籍出版社吴长青副总编分别介绍了该社在诸如影印、点校、注释、今译等古籍传统“活”化方式——古籍整理,以及以电子书、有声书、专题数据库等为代表的古籍现代“活”化方式——古籍数字化的主要成果,并重点分享了“汇典·古籍数字服务平台”的建设情况。该平台汇集了上海古籍出版社的核心古籍,并囊括世纪出版集团内外相关出版社的古籍资源,具有阅读和检索的双重功能,能够实现图文对读、原书图片缩放、文本复制和引用、笔记记录、字典查询等多种阅读体验,目前已上线3.4亿字内容,未来总字数预计将突破100亿。

上海外国语大学图书馆欧阳剑研究员对现阶段古籍出版的传统线性模式进行了深刻反思,指出物理形态向数字形态的单纯转变并不能充分挖掘古籍的全部价值。他认为,数智时代要深度融合大数据、人工智能的数据思维和新兴技术应用,古籍数字出版不仅要面向用户,也要面向智能机器的“阅读”和理解,因而必须打破古籍知识固化、单一、封闭的逻辑导向,破除长期以来存在的多重知识区隔,通过重组和聚合形成新的优质内容产品,再借助新技术将古籍数字化内容、知识融入场景中,完成新知识的生产与新场景的构建,最终实现古籍数字出版的价值再造。

武汉大学文学院博士后张亚静女士分享了目前正在建设的集原始文本呈现、外围资料纂集、互动再造为一体的“元杂剧智能化数据平台”的具体情况。该平台由核心、外围、扩展三层构成:核心层是对剧目的不同版本呈现,通过版本比对和算法演练来分析文字差异的背后原因,同时也能提取曲牌、用韵、典故等关键信息,进行跨作家、作品的比对统计和分析,并可视化呈现结果;还可以通过分析主题、人物形象、故事线、审美偏好、文体风格、意象等,来辅助定位佚名作品的作者及时代信息。外围层则是对包括作者生平、行迹、交友等个人信息,文本关联的历史或文学材料、工尺谱、服饰、剧评等剧目资料,以及衍生文本、舞台、影视剧等改

编资料的纂集。扩展层的主旨是互动与再造,读者在浏览前两层和使用工具的基础上,可以构建专属的研究资料集并上传至平台共享,形成主题社区。

古籍的数字出版与再造最大程度调和了古籍实体保护与传统文化传播之间的矛盾,并助力古籍突破纸本载体的形式束缚,最终转化成为灵活的数据形态知识。传统文献与现代信息技术的深度融合,对古籍内容进行整合、挖掘、可视化等系列操作,推动原本固定不变的古籍文本实现内容数量的增长与呈现形态的革新,这样所形成的古籍数字产品,在为专业领域的研究者提供更多使用便利和全新的研究思路的同时,又能够辐射到更加广泛的读者群体,拓展出更为丰富的应用场景。

## 4 讨论与展望

学术报告结束后,与会专家还就古籍定本、古籍年代下限、古籍数字化完成进度、纸本古籍及石刻文献残缺文字的识别、语言模型系统训练、数字古籍总目编纂等一系列相关问题展开了自由而热烈的讨论。有的问题达成了共识,比如对古籍数字资源采取书目控制,以避免数字古籍的重复开发;对民国文献的数字化,采取与古籍同等重视的措施。有的问题尚存在分歧,比如对古籍定本的权威性与多样性的处理,采用计算机模糊识别技术“识别”出的古籍文字的校勘价值,等等。此次会议将传统古籍整理与研究为数智时代的信息技术结合起来,对于古籍活化与利用基础理论、技术场景与产业应用等方面的研究起到了积极的推动作用。展望未来,数智时代古籍活化利用研究有以下几个问题值得重点关注:

第一,古籍文本的识别与数据化转化。文本的识别与数据化是古籍活化利用与出版再造的前提条件。受生产年代、收藏条件等诸多因素的影响,古籍纸张时常伴随有发黄、老化、破损、虫蛀等异常情况,有的文字漫漶不清或缺损不全,无形中增加了计算机自动识别页面进而实现数据化转化的难度。为了更好地解决这一难题,现有研究围绕古籍文档的结构化识别方案进行了积极探索,如将卷积神经网络<sup>[1]</sup>、注意力机制<sup>[2]</sup>等深度学习技术,以及迁移学习<sup>[3]</sup>的执行策略引入古籍文本的识别领域。以此作为基础,后续研究还应积极

探索提升古籍文本识别精度的方案,推进古籍文本大规模自动识别的进度,降低古籍数据化的人力、物力成本。此外,对古籍缺失页面的智能联想与补全也是未来应重点关注的方向。基于古籍文本识别,训练具有古籍语义逻辑提示、上下文文本预测的模型,从概率和推理的角度向文史学者提供缺失页面可能的候选文本,将有助于古籍数据化向纵深发展。

第二,数智时代古籍的版本问题。对古籍数字化的底本进行甄别和优选,已成为大家的共识。不过事实上,多数典籍其实很难找出一种真正善本,所谓的版本优劣仅仅是相对而言,不同版本之间的出入可能蕴藏着有重要研究价值的文献信息,这提示我们在数字化过程中不应仅拘泥于某一种善本,而是应当在条件允许的情况下,尽可能让更多的版本完成数字化流程,最大程度保留不同版本的原始信息,为后续智能提取、多重比对和数据分析等更深层次的研究提供更多发展空间。除了对数字古籍的底本信息和来源进行必要的标注外,还应当著录不同版本的数字化古籍,形成一部实时更新的在线数字古籍总目,这样既可以为业内人士提供后续的重要参考,避免重复开发造成的资源浪费,又能够向读者提供便捷的使用引导,提高数字古籍的传播效率。

第三,推进古籍知识的挖掘与聚合。知识挖掘和关联是古籍活化利用与出版再造的重要环节。以古诗词、引文、史书等为代表的古籍蕴含着大量的主题内容、情感分布、人物关联、语义联系等知识内涵,运用主题识别、情感分析、社会网络分析、知识图谱等方法深入挖掘大规模的古代典籍,识别并关联古籍知识,可视化古籍知识或历史人物,能够从定量的角度帮助用户迅速、准确地理解大规模古籍的内涵。通过对古诗词中情感、景观、意象等的识别与挖掘,可反映古人生活场景和人文风貌<sup>[4]</sup>;通过对古籍引文的条目进行抽取与分析,可呈现引文的分布特征,从而定量地研究古人的引文行为<sup>[5]</sup>;以前期的知识挖掘作为基础,并借助知识关联技术,构建诸如古代百家思想知识图谱<sup>[6]</sup>、古典诗词知识图谱<sup>[7]</sup>、《宋元学案》知识图谱<sup>[8]</sup>等可视化系统。后续研究应充分吸收和借鉴相关领域的先进经验,再进一步拓宽古籍文本知识挖掘与关联的范围。

第四,构建多元化的古籍知识服务数字平台。作为连通古籍数字化成果与现实应用的桥梁,知识服务数字平台充分实现了古籍知识的传播和利用。已有的

