

数字资源长期保存:当前进展和最佳实践

——2007 年数字资源长期保存国际会议(iPRES2007)综述

吴振新¹ 刘建华^{1,2} 张玫^{1,2} 赵琦^{1,2} 向菁^{1,2}

¹(中国科学院国家科学图书馆 北京 100080) ²(中国科学院研究生院 北京 100080)

【摘要】 系统而全面地回顾 iPRES2007 数字资源长期保存国际会议,从数字资源长期保存的战略计划与基础设施、相关管理问题、技术研究与实践、认证与评估、教育与培训 5 个方面介绍研究和实践的进展情况,深入分析并总结已有的经验和教训,并就面临的问题和下一步发展进行探讨。

【关键词】 数字资源 长期保存 进展 最佳实践 **【分类号】** G253

Digital Preservation: Sustainable Programs and Best Practices ——A Comprehensive Review of iPRES2007

Wu Zhenxin¹ Liu Jianhua^{1,2} Zhang Mei^{1,2} Zhao Qi^{1,2} Xiang Jing^{1,2}

¹(National Science Library, Chinese Academy of Sciences, Beijing 100080, China)

²(Graduate School of the Chinese Academy of Sciences, Beijing 100049, China)

【Abstract】 This paper reviews iPRES2007 digital preservation international conference comprehensively, introduces the development of policies, strategies, planning and infrastructures of digital preservation, related management issues, technology researches and practices, certification and assessment, etc. It also analyzes and summarizes experiences and lessons existing in practice, discusses problems that we have and brings forward the important concerns for the next phase of digital preservation.

【Keywords】 Digital resource Digital preservation Sustainable programs Best practices

继 2004 年北京、2005 年德国、2006 年北美成功举办之后,2007 年 10 月 11 - 12 日,数字资源长期保存领域的主流国际性系列会议—iPRES(International Conference on Preservation of Digital Objects)再次来到北京,由国家科技图书文献中心主办,中国科学院国家科学图书馆承办,联合其他单位组成会议组委会^[1],为国内外近 200 名专家和相关研究人员提供了一个盛大的交流和学习平台。

近年来,数字资源长期保存领域经历了从基础理论研究到个体实验再到最佳实践的发展过程,iPRES 的研讨内容不断拓展和深入,吸引了越来越多的机构和学者的关注。2004 年的首次国际对话,主要商讨了数字资源长期保存的基本问题,在对长期保存的内涵和外延达成一定共识的基础上开始了长期保存的实验和国际合作。此后,各个实验项目在不同的处理过程中遭遇了相应的阶段性问题,这些问题在 2005、2006 年会议中得到了关注,如 workflow、元数据、知识技术、认证、知识库、保存服务

和项目管理等。这些阶段性的研讨促进了实践思路的日益明晰化,研究体系的不断系统化和具体化^[2-4]。

2007 年的会议围绕“数字资源长期保存:项目进展和最佳实践”,结合具体项目和实践经验,分别从数字资源长期保存的战略计划与基础设施、相关管理问题、技术研究与实践、认证与评估、教育与培训 5 个方面进行了介绍,并就面临的问题和下一步发展进行了交流。

1 数字资源长期保存战略计划与基础设施

目前开展的大量数字资源长期保存实践,大多属于机构或组织内部的独立活动,然而这容易导致长期保存投入与管理的资源配置效率不佳。参与本届 iPRES 会议的代表在本国或本机构长期保存的具体实践基础上,提出了在合作的基础上进行长期保存的规划,内容主要包括保存网络的建立、职责分配及政策制定等。

1.1 建立保存网络

数字长期保存是一项费用高、耗时长、多层面、跨领

域的工程,现有的实验系统基本都采用各自为政、自行建设的发展策略,因此资源重复建设,机构负担较重的问题凸现。鉴于此,部分机构开始考虑构建更大范围的长期保存网络,即一个由多个异地分布的长期保存系统紧密耦合而成的虚拟组织,以促进实践中的资源共享、职责与费用分摊以及交流等。目前,国家范围的保存网络已逐渐兴起,如美国的 NDIIPP^[5]、德国的 nestor^[6]、加拿大的 Multicultural Canada^[7] 及我国拟建的国外科技文献数字保存网络^[8]等;组织间的跨国保存网络也在不断发展,如 Planets^[9]、InterPARES^[10] 及筹建中的 DARIAH^[11]等;此外,少数国家范围的保存网络如荷兰的 e-Depot^[12],也开始向跨国保存网络发展。这种网络型的组织结构与单个保存系统间存在着较大差异,由此引发的对相关的法律、经费、管理、技术等问题的探讨将进一步拓展长期保存领域的研究内容。

1.2 构建责任体系

构建一个保存网络,其成员在研究领域、工作部门、机构性质等方面必然存在差异,如何合理分配成员职责显得尤为重要,因为这会极大地影响到今后的认证实施和长期保存的有序发展。此次会上,美国 NDIIPP、荷兰 e-Depot、DARIAH 及我国国外科技文献数字保存网络等项目的代表均对该问题进行了探讨。

专家们认为,构建保存网络时,首先需要对参加机构的资格进行审查,确立审核标准、审核主体是一个重要问题;当保存网络形成后,中心组织与成员组织的权责就成为考虑的核心。由于中心组织将负责网络的总体管理与协调,因此“中心组织由谁来担任,具备什么权限”等就成为关键问题;而对其他成员组织之间进行保存资源种类、运行费用等的分配,设计“灾害发生时的行动措施”等也将对整个网络的有效运行产生重要影响。

1.3 制定相关政策

作为长期保存的有力保障,制定合理完善的长期保存政策非常重要。从已经或正在制定的政策来看,其主要内容都围绕以下几点展开:

- (1) 在机构范围内赋予长期保存活动一定的地位,使之成为组织日常工作的一部分;
- (2) 确定长期保存的目标;
- (3) 明确相关人员或组织的职责;
- (4) 制定长期保存的规范流程;
- (5) 规定所需遵守的技术标准或最佳实践;
- (6) 建立对实施效果评估的机制;
- (7) 制定长期保存延续性计划等。

2 数字长期保存中的相关管理问题

长期保存活动是一项复杂的系统工程,需要合理有效的管理来保障它正常、正确地实施。尤其在构建保存网络更增添了它的复杂性和难度,因此,荷兰的 e-Depot、德国的 nestor^[13]、大英图书馆的 Planets^[14] 以及我国国外科技文献数字保存网络^[15] 等项目都针对长期保存系统中的资源层、系统实施层和应用层的有效管理提出了自己的见解。

2.1 资源层的管理

长期保存领域中涉及的资源包括需要被保存的数字对象、保存实施过程中的资金、技术、系统和参与的人员、机构等。

数字对象作为保存的主体,继续得到多方关注。鉴于数字对象存在形式、依附载体以及拥有者的多样性特点,如何根据各自的特长、需求、服务对象特点以及经费等因素确定需要保存的数字对象和阶段性的保存目标,如何从多种渠道,如科学界、图书馆、档案馆、出版商、商业机构、Google/Yahoo 等这类机构处获得数字对象,是至关重要的环节。

在资金管理方面,资金的来源和成本预算都是管理的主要内容。稳定而长期的资金是长期保存活动持续发展的基础,目前各项目资金主要来源于国家或政府机构专项投资、图书馆业务经费分配、出版商赞助、用户赞助等。围绕着资金的管理与有效应用,大英图书馆“数字保存周期成本计算”(LIFE: Costing the Digital Preservation Lifecycle)项目针对保存生命周期,分别从数字对象、技术开发等方面开展成本预算研究,基于关键保存活动和关键影响因素提出了计算公式。DPE 的 DRAMBORA (Digital Repository Audit Method Based on Risk Assessment) 也提出了如何利用机构自审计的方法来对资金预算进行风险管理。

2.2 系统层的管理

本次会议上, workflow 再度成为系统层管理所关注的重点,但研究范围已经从 workflow 设计和 workflow 自动化进一步发展成为 workflow 管理。定义清晰、权责明确的 workflow 管理对降低成本、及时发现错误、保障长期保存系统的正常运转提供了有力的保障, workflow 管理贯穿了数字对象的保存生命周期。Uwe Roseman^[16] 提出: DOI 的注册也应该纳入 workflow 的管理中,借助 DOI 实现对原始数据集的长期保存。在 workflow 管理运行过程中,相关的标准规范是各模块内部和模块间友好交互的基础,无论是 OAIS 模型还是 Planets,无论是元数据标准还是 OAI -

PMH 协议,都对系统实施的效率影响深远。

尤其让人注目的是本次会议上关于工作流的研究已经超出了保存仓储系统工作流的研究,延伸到保存研究的其它方面,如:维也纳科技大学的 Stephan Strodl^[17]提出了保存计划的工作流模型,中国科学院国家科学图书馆长期保存团队提出了贯穿整个长期保存活动的可信赖工作流管理框架。

2.3 权益管理

长期保存系统的实践中,系统与服务的可信度和使用过程中的权益管理,关系到整个保存体系的安全和可信程度。权益一直是使用管理中的关键问题,如何确定用户的权限等级并提供相应的服务,是保证权益的一个方面,而数字对象的生产者或创造者则又是权益的另一个关键点。本次会议上对此虽未作深入探讨,但在讨论中众多会者不约而同地涉及了这个话题。

3 技术研究与实践

长期保存技术一直是长期保存研究领域的重点。本次会议中,荷兰的 e-Depot、DARIAH 及我国国外科技文献数字保存网络^[18]等项目分别与参会者分享了各自在机构仓储领域、长期保存存档系统、多种保存对象类型等方面的实践经验和研究成果。

3.1 机构仓储领域

机构仓储作为研究机构对自主知识产权统一组织、管理的系统,是数字资源长期保存的重要基础组成部分。本次会议上,欧洲的 DRIVER^[19]项目吸引了不少与会者的目光。该项目致力于构建一个泛欧洲的、交互、可信任的长期知识库网络架构,经过第一阶段的实践,已建成了一个由5个国家知识库提供者组成的常设网络。而为了适应欧洲知识库网络,项目还形成了一套构建本地知识库的指导方针,并实现了对欧洲51个机构知识库数据的再利用。

国内科研领域也在机构仓储的建设方面作出了大量努力。由国家科学图书馆兰州分馆^[20]牵头的中国科学院机构仓储建设已初步完成了试点建设工作,预计将建成覆盖全院近百个研究所的机构仓储网络,同时也在积极探求与欧洲机构仓储网络的合作。清华大学图书馆^[21]也在本次会议上介绍了利用开源的 DSpace 系统构建的本校学生优秀作品数据库。

3.2 存档系统

存档系统作为长期保存的基础组成部分,在各国保存活动中得以充分发展。作为拟建的 NSTL 数字信息资源长期保存网络的一部分,中国科学院国家科学图书馆

采用 Fedora 作为底层存储系统,构建了电子期刊长期保存系统(CAS E-Journal Archiving System),选择一定规模的电子期刊资源作为试验资源,通过数据预处理、保存管理、数据访问服务3个主要系统功能,初步实现了数字对象的长期保存功能。该系统在 OAIS 模型的基础上,遵循开放的技术标准和规范,具有开放性和可扩展性,能够对多种类型的数字对象进行保存管理及有效的生命周期管理。

SDSC 的 Chronopolis 项目^[22]作为一个基于网格的概念性长期保存框架,与前者有明显的差别。Chronopolis 的数据网格由地理上分散的3个站点组成,每个站点根据资源的不同扮演不同的角色,数据完整性和安全性的检查分别由各自连接的主要分段格处理,处理完成后数据资源将独立地推送入每一个存档系统,即每份数据资源存在3个独立管理的复本,同时增加的说明层为数据管理和数据完整性检验提供了额外的安全保证,系统、数据等之间的异构性由 SRB 及 iROD 进行处理。这种存储框架在环境的扩展性、安全性等方面体现了较好的效果。

3.3 非文献类型资源的保存研究

随着长期保存研究的发展,非文献类型资源的保存进入了许多机构的研究范围。

德国国家科技图书馆 TIB 正致力于建立一套基于 DOI 的科学数据保存体系,启动了原始科学数据的出版和引用项目以实现科学原始数据的访问。该项目借助 TIB 是国际 DOI 机构成员的优势,代理注册非商业化的科学原始数据 DOI,并通过代理中心对数据进行引用与查询。而 DARIAH 则专注于支持访问所有欧洲的数字化人文和文化遗产信息,并实现这些信息长期保存的框架研究。除此之外,上海图书馆^[23]已经实现了馆藏文化遗产的长期保存,系统采用 METS 作存档格式,进行数字资源的 XML 元数据描述和 DOI 标引,实现了数据资源的平台独立性,用户可以通过服务器访问以 TIFF 和 PDF 格式保存在磁带、光盘等介质里的数字资源。

3.4 保存技术策略

受保存内容、保存能力和保存技术3者的影响,保存的技术策略呈现多样化特征。本届会上,荷兰国家图书馆(KB)^[24,25]分别从文件格式的选择方法和仿真实践两个方面交流了经验。KB 受数字保存文献的可持续标准启示,将各个标准按特征分解,赋予各特征不同的权重,量化指标进行文件格式的选择,这为仿真或迁移技术的实施奠定了基础。另外,KB 今年首次展示了他们小组通过构建持久、模块化的虚拟机实现仿真的研究成果,尽管其在效率、数据抽取与插入等方面还存在缺陷,但已然验

证了该方法的有效性。

kopal^[26]项目组介绍了基于 koLibRI 这一开源软件的文件格式迁移保存策略。koLibRI 的迁移管理部分设置事件监听器获取需要迁移的信息,根据存档系统将迁移需求转化为适当的查询式,系统返回标准对象格式,文档分析完成后实现转换工作。这种迁移管理的实现为 koLibRI 的进一步开发奠定了基础。

此外,MIXED^[27]项目组在进行基于迁移的保存技术策略研究中,制定了利用 XML 作为迁移中间格式进行保存的方案,意图通过系列标准和开源软件开发出相关的系统。

4 认证

随着长期保存系统的不断发展,人们愈来愈关注系统及服务可信度的认证。iPRES 在 2006 年首次将该问题纳入讨论范围,今年会议上与认证相关的研究范围逐步从理论研究走向实践,其涉及的内容主要包括以下两点:

4.1 认证指标的构建

有关认证指标的早期研究多从开发一套具备普适性的指标体系入手,忽略长期保存系统的个体差异,而当认证准备进入应用阶段时,本地环境成为一个不容忽视的重要因素。德国 nestor^[28]在一系列国际标准及相关讨论的基础上,构造了一套适用于德国长期保存系统的认证指标体系,现已经进入测试阶段。此外,nestor 还与 RLG 及 OCLC 的相关工作组合作,努力构造一套适用于大多数长期保存系统的认证指标并争取成为 ISO 标准。

在我国的保存研究中,全国范围内的长期保存网络正在筹建中,与之配套的认证机制研究也在同步进行。它拟在 RLG/NARA 制定的《可信赖仓储的审计及认证:指标与列表》(Trustworthy Repositories Audit & Certification: Criteria and Checklist, TRAC)^[29]的基础上,对相关指标加以修正及扩展^[30]。由于保存网络的认证不仅针对单个系统,还包括对整个网络的认证,因此,还对认证执行过程涉及的认证对象、执行主体、认证时间间隔及透明度问题等进行了探讨。

4.2 内部自行审计

长期保存系统在接受外部认证前,往往需要根据认证指标事先自审。针对这种情况,DCC 与 DPE 下属的工作组针对 TRAC 与 nestor 的认证指标,开发出一套“基于风险管理的数字仓储审计方法”(Digital Repository Audit Method Based On Risk Assessment, DRAMBORA)^[31],把长期保存活动中的不确定因素转化为具体的风险,并加以有效地评估、测量、管理及预防,帮助长期保存系统更好

地了解其任务、目标、相关活动及资源等,全面把握自身的优劣势,最终为外部认证作好准备。目前,DRAMBORA 还提供了风险在线注册机制,以促进各保存系统更好地交流。

5 教育与培训

随着数字资源长期保存实践的广泛开展,越来越多的研究和工作人员参与到长期保存活动之中,相应的教育与培训随之兴起。此次会上,DigCCurr^[32]、DPE (Digital Preservation Europe)^[33]及 nestor 的代表^[34]分别介绍了他们在长期保存领域开展教育培训的概况。

5.1 开展形式

各项目普遍采用了举办专家讲座的方式介绍长期保存领域最新理论、技术研究前沿,除此之外,各项目均发展了其他的特色方法。DigCCurr 作为一项培养数字资源保存领域研究人才的研究生教育项目,设计得较完善,学生课程、专业研讨会、实习 3 大部分将基础教育、深度研究和就业紧密相连,在广度、深度及实践上均有较好的体现。而 nestor 则充分利用春冬季学校、手册、联合培养课程、E-learning 等形式,开发了信息服务平台(德语)、主题交流平台(德语),极大地扩大了学习机会。与前两者相比,DPE 的培训较为短暂,但 DPE 为参会者配套培训练习题,将理论与实践结合,有利于提高培训质量。

5.2 主要内容

作为一个教育项目,DigCCurr 在扩展数字资源生命周期、建立核心知识体系等 5 项规则的指导下建立了课程的六维模型,如原则与标准、专业或组织环境等,每门课程专注模型的某一方面或涵盖更多方面,循序渐进实现教育目标。而 DPE 与 nestor 则更多针对实践工作提供培训,主要涵盖数字资源保存的基础知识、 workflow、元数据、OAI 模型及所涉及的管理、法律问题等。另外,nestor 还结合自身项目,组织专家编写《nestor 手册——数字资源保存的百科全书》^[35](第 1 版),内容涉及 OAI 模型、可信赖仓储、长期保存标准及策略等,及时跟踪前沿发展动态,保持动态更新。

5.3 相关计划

教育与培训是一项长期的工作,DigCCurr 计划今后对相关行业的企业招聘作深度调研,为学生提供就业指导。DPE 则计划通过多次培训,提高培训质量,吸引社会对长期保存的关注。而 nestor 不仅立足于当前,还准备冲破现有网页资源主要为德语的限制,将更多内容以英语呈现。同时,nestor 还跟与其有合作关系的欧洲 10 所大学在联合进行长期保存教育的课程组织、选择标准、联

合培养学位的模式等方面达成共识,并计划与美国相关机构合作,力图在基础教育中取得成绩。

6 结 语

数字资源长期保存已经积累了较丰富的经验和知识,我国的数字资源长期保存研究也在在战略计划、 workflow管理、技术实施和认证等方面都取得了初步进展。但数字资源长期保存也面临很多问题,如:

(1)形成机构内或国家内的长期保存系统时,如何制定阶段性目标?在构建保存网络时,如何建立责任体系?

(2)长期保存经费来源、成本预算、系统模型与标准选择、权益管理等如何有效解决?

(3)随着数字对象存储复杂性的增加,如何长期保证数据的精确可读性(技术的可读性和人的可读性)?如何最小程度地削弱技术过时的影响?

(4)为了降低成本,减少重复开发,选用已有软件或开源软件进行继续开发是一种资源优化的选择,但如何有效规避软件自身存在的硬伤对系统的健壮性和效率等造成的影响?

(5)如何针对特定数据(科学数据、文化遗产资源)的保存制定有效、统一的标准?

(6)数据安全如何保障?如何进一步在保存的基础上利用数据开发服务?系统的可信度和提供的服务、内容如何鉴定?

(7)教育与培训由谁组织?如何实施?受教育和培训对象如何确定?

此外,在大会主旨报告专家 Seamus Ross^[36]对研究前景焦点的描述中,如数据恢复、互操作、自动化、组织环境、存储和实验等方面也存在着很多盲点需要研究,长期保存依旧需要长期深入的研究和实践。

参考文献:

- [1] iPRES2007—“数字资源长期保存:项目进展和最佳实践”[EB/OL]. [2007-10-14]. http://ipres.las.ac.cn/index_cn.jsp.
- [2] Chinese-European Workshop on Digital Preservation[EB/OL]. [2007-10-15]. http://159.226.100.135/meeting/cedp/index_en.html.
- [3] International Conference on Preservation of Digital Objects[EB/OL]. [2007-10-15]. <http://rdd.sub.uni-goettingen.de/conferences/ipres05/>.
- [4] 2006 iPRES Conference[EB/OL]. [2007-10-15]. <http://ipres.library.cornell.edu/>.
- [5] Molly Johnson, Creating a Digital Preservation Network with Shared Stewardship and Cost[EB/OL]. [2007-10-25]. http://ipres.las.ac.cn/pdf/Molly%20Johnson-speech_CreatingDigitalPreservationNetwork_MJohnson_final.ppt.
- [6] Reinhard Altenhöner. The Next Step—Establishing a Long-Term Preservation Infrastructure for Science and Research; the Nestor2 Approach[EB/OL]. [2007-10-25]. http://ipres.las.ac.cn/pdf/Short_info_nestor_english_final.pdf.
- [7] Ian Yiliang Song. Preservation of Digitized Canadian Multicultural Heritage[EB/OL]. [2007-10-25]. http://ipres.las.ac.cn/pdf/Ian%20Song%20%20%20%20iPRES2007_IanSong_final.pdf.
- [8] Xiaolin Zhang. Chinese National Archival Network for Digital STM Material[EB/OL]. [2007-10-25]. <http://ipres.las.ac.cn/pdf/IPRES2007-XiaolinZhang.pdf>.
- [9] Helen Hockx-Yu. A Practical Approach to Digital Preservation: Updates from PLANETS[EB/OL]. [2007-10-25]. http://ipres.las.ac.cn/pdf/Planets_iPRES_HHY_V1.0.pdf.
- [10] Sherry L. Xie. Foundations for Developing Digital Preservation Strategies and Policies; The InterPARES Policy Framework and Guidelines[EB/OL]. [2007-10-25]. http://ipres.las.ac.cn/pdf/Sher-ry%20L.%20Xie%20%20%20%20%20iPRES2007_Beijing-Presentation_SherryXie_Canada_20071011.pdf.
- [11] Ellen Willemsse. Dariah - Digital Research Infrastructure for the Arts and Humanities[EB/OL]. [2007-10-25]. <http://ipres.las.ac.cn/pdf/Ellen.ppt.pdf>.
- [12] Elisabeth van Eijck van Heslinga. SHAPING COURSE; The Development of the Strategy for the e-Depot of the Koninklijke Bibliotheek, National Library of the Netherland[EB/OL]. [2007-10-25]. <http://ipres.las.ac.cn/pdf/Els%20Van%20Eijck%20OUTLINE%20VOOR%20ipres%202007.pdf>.
- [13] Perla Innocenti, Andrew McHugh. Digital Curation Centre (DCC) and Digital Preservation Europe (DPE) Audit Toolkit; DRAMBORA[EB/OL]. [2007-10-25]. http://ipres.las.ac.cn/pdf/inno-centi_ipres07.pdf.
- [14] Paul Wheatley. LIFE; Costing the Digital Preservation Lifecycle[EB/OL]. [2007-10-25]. <http://ipres.las.ac.cn/pdf/Paul%20Wheatley%20%20%20%20LIFE-PaulWheatley.pdf>.
- [15] Zhenxin Wu. Preservation Management in Practice; Trusted Workflow[EB/OL]. [2007-10-25]. <http://ipres.las.ac.cn/pdf/wu%20zhenxin%20-Workflows.pdf>.
- [16] Uwe Rosemann. A System for Digital Preservation of Scientific Data using DOI Names[EB/OL]. [2007-10-25]. http://ipres.las.ac.cn/pdf/Uwe%20Rosemann_A%20system%20for%20digital%20preservation%20of%20scientific%20data%20using%20DOIs.pdf.
- [17] Stephan Strodl. Preservation Planning in the OAIS Model[EB/OL]. [2007-10-25]. http://ipres.las.ac.cn/pdf/strodl-final_3.pdf.
- [18] Zhixiong Zhang. CAS STM Journal Archival[EB/OL]. [2007-10-25]. http://ipres.las.ac.cn/pdf/zhang_zhixiong_Developing%20a%20CAS%20E-Journal%20Archiving%20System_new.pdf.
- [19] Norbert Lossau. DRIVER - Digital Repository Infrastructure Vision for European Research[EB/OL]. [2007-10-25]. <http://ipres.las.ac.cn/pdf/Nobert%20Lossau%20DRIVER-IPRES->

