

대용량 과학데이터 아카이브 시스템 설계를 위한 리포지터리 시스템 분석

Analysis of Repository Systems for Designing a Archive System of Large Science Data

임종태*, 서인덕**, 송희섭**, 유승훈*, 정재윤*, 조중권**, Aniruddha Paul**, 고건식**, 김병훈**, 박윤정**, 송진우*, 이서희**, 전현욱*, 최민웅**, 노연우*, 최도잔*, 김연우*, 복경수*, 이정훈***, 이상환***, 유재수†

충북대학교 정보통신공학부*,
충북대학교 빅데이터학과**,
한국과학기술정보연구원 과학데이터연구센터***

Jongtae Lim*, Indeok Seo**, Heesub Song**, Seunghun Yoo*, Jaeyun Jeong*, Jungkwon Cho**, Aniruddha Paul**, Geonsik Ko**, Byounghoon Kim**, Yunjeong Park**, Jinwoo Song*, Seohee Lee**, Hyeonwook Jeon*, Minwoong Choi**, Yeonwoo Noh*, Dojin Choi*, Yeonwoo Kim*, Kyoungsoo Bok*, Jeonghoon Lee***, Sanghwon Lee***, Jaesoo Yoo*†

School of Information and Communication Engineering,
Chungbuk National University*,
Department of Big Data, Chungbuk National University**,
Scientific Data Research Center, Korea Institute of
Science and Technology Information***

요약

본 논문에서는 대용량 과학데이터 아카이브 시스템 설계를 위해 기존 리포지터리 시스템을 분석한다. 대용량 과학데이터를 효율적으로 수집하고 저장하기 위한 아카이브 시스템 아키텍처 설계를 위하여 현재 서비스되고 있는 다양한 과학데이터 리포지터리 시스템을 분석한다. 분석한 내용을 바탕으로 대용량 과학데이터 아카이브 시스템 아키텍처를 설계하기 위한 기술적인 요구사항을 도출한다.

I. 서론

최근 각종 실험 장비의 발전에 따라 천문, 항공, 우주, 유전자 등 첨단 분야에서 생산되는 과학데이터가 빅 데이터화 되고 있는 추세이다. 빅 데이터란 기존 운영 시스템으로 저장, 관리 및 처리가 불가능할 정도의 다양하고 방대한 데이터를 의미한다. 대용량의 과학데이터는 기존의 데이터 수집 및 관리 시스템으로 처리하는데 한계가 존재하기 때문에 분산 저장 관리 등 대용량 빅 데이터 처리를 활용한 접근이 요구된다. 또한 대용량의 과학데이터를 수집, 관리하고 효율적으로 공유, 활용하기 위한 과학데이터 아카이브 기술이 매우 중요하게 요구된다.

본 논문에서는 대용량 과학데이터 아카이브 시스템 설계를 위하여 기존 리포지터리(Repository) 시스템을 분석한다. 대용량 과학데이터를 효율적으로 수집, 저장, 공유, 활용하기 위한 아카이브 시스템 아키텍처 설계를 수행하기 위하여 현재 서비스되고 있는 다양한 과학데이터 리

포지터리시스템을 분석한다. 분석한 내용을 바탕으로 대용량 과학데이터 아카이브 시스템 아키텍처를 설계하기 위한 기술적인 요구사항을 도출한다.

II. 리포지터리 시스템 분석

본 논문에서는 현재 서비스되고 있는 다양한 과학데이터 리포지터리 시스템을 분석한다. 기존의 리포지터리 시스템을 대상으로 해당 리포지터리 시스템의 특징, 시스템 구조, 데이터 수집, 분산 저장 관리, 부하 분산 관리, 그리고 데이터 전송 측면으로 분석을 수행했다.

1. Fedora Commons

Fedora Commons(이하 페도라)는 디지털 콘텐츠의 관리 및 보급을 위한 강력한 모듈형, 오픈 소스 저장소 시스템이다[1]. 페도라는 ModeShape와 Infinispan으로 구성된다. ModeShape는 질의, 전체 텍스트 검색, 이벤트, 버전, 참조, 유연한 동적 스키마를 지원하는 분산, 계층, 트랜잭션 데이터 저장소이다. Infinispan은 분산 인메모리(In-memory) Key/Value 데이터 저장소이다. 페도라는 모든 콘텐츠 유형을 지원하며, 관련 메타데이터를 관리함으로써 리포지터리 서비스를 제공한다. 분산 저장 관

† 교신저자 : yjs@chungbuk.ac.kr

본 연구는 미래창조과학부 및 정보통신기술진흥센터의 대학ICT연구센터육성 지원사업(ITP-2016-H8501-16-1013)과 2016년 KISTI의 '과학기술빅데이터 공동 활용 거버넌스 체제 지원' 과제의 재원을 지원받아 수행된 연구임

리는 Infinispan이 제공하는 분산 저장 관리 기법을 이용한다. 페도라는 부하 분산 관리를 위해 AWS(EC2)와 Tomcat7을 이용한 클러스터링을 수행한다. 클러스터의 전 단계에 Apache Server와 mod_jk를 사용하여 부하분산을 수행한다. Apache Server를 활용한 부하분산은 사용자 접근을 위한 웹 서버는 단일하게 유지하면서 실제 작업을 수행하는 어플리케이션 서버를 다중으로 유지하여 실제 작업량에 따라 작업량이 적은 어플리케이션 서버에 작업을 할당하는 방법이다.

2. dSpace

dSpace는 논문, 회의자료, 이미지, 동료간의 리뷰 자료, 기술보고서, 연구중인 자료 등에 대한 저장, 공유, 검색을 위해 개발된 무료 공개 리포지터리 시스템이다[2]. dSpace는 데이터관리, 사용자 및 권한관리, 저장관리, 검색 및 열람 기능을 지원한다. dSpace는 저장 계층, 내부처리 계층, 응용 계층으로 구성된다. 저장 계층은 실제 데이터 저장에 사용되는 운영 시스템으로 구성된다. 내부처리 계층은 dSpace 서비스 지원을 위한 워크플로우(workflow), 콘텐츠 관리, 운영, 검색 및 브라우징 모듈 등을 포함하는 dSpace만의 특징적인 기능으로 구성된다. 응용 계층은 dSpace를 이용하여 실제 서비스를 제공하기 위한 기능들로 구성된다. dSpace는 학술논문, 책, 회의 보고서, 3D, 사진, 필름, 비디오, 연구 자료 세트 등을 포함하는 폭 넓은 데이터 형태를 지원한다. 분산 저장 관리는 함께 운영되는 저장 시스템(PostgreSQL, MySQL 등)의 기능을 사용하며, 부하 분산 관리와 관련해서는 페도라와 같이 톱갯을 이용한 부하 분산 관리 기능을 수행한다.

3. Open Science Date Cloud

Open Science Date Cloud(OSDC)는 프로젝트 단위로 관리되는 클라우드 리포지터리 서비스이다[3]. OSDC는 OpenStack을 기반으로 클라우드 리포지터리 서비스를 제공한다. OpenStack은 일반 서버에서 클라우드 컴퓨팅 서비스를 실행할 수 있도록 해주는 오픈 소스 플랫폼이다. OSDC의 핵심 모듈은 Cloud Controller이며, Cloud Controller는 Nova, Glance, Keystone으로 구성된다. Nova는 계산(Compute) 인스턴스들을 제어 및 관리하기 위한 서비스이다. Glance는 가상 머신 이미지들을 저장, 등록, 관리, 전달하기 위한 서비스이다. Keystone은 모든 서비스의 Identity를 증명하는 역할을 수행한다. OSDC는 Biology, Genomics, Earth science, Social science, Text data, Astronomy, Music, Model reduction 등의 데이터 종류를 지원한다. OSDC는 오브젝트 스토리지, Hadoop 등의 운영 시스템과 연동하여 저장 관리를 수행하기 때문에, 운영시스템의 분산 저장 관리와 부하 분산 기능을 사용한다.

III. 아카이브 시스템 요소 기술 도출

기존 시스템을 분석한 결과, 많은 기존의 리포지터리 시스템이 다양한 운영 시스템과 연동되어 사용되고 있음을 확인할 수 있었다. 경우에 따라 자체적인 분산 저장을 사용하는 시스템도 존재하긴 했지만 대부분의 리포지터리 시스템이 연동되는 운영 시스템의 분산 저장과 부하

분산 기능을 활용하고 있음을 확인할 수 있었다. 본 논문에서는 분석된 내용을 바탕으로 아카이브 기술 요구사항을 도출하였다. 그림 1은 도출된 아카이브 기술 요구사항을 보여준다. 아카이브 기술 요구사항은 수집, 저장관리, 검색, 전송 측면으로 구분된다. 저장 관리를 위한 기능으로는 분산 저장 관리, 부하 분산 관리, 데이터 압축 저장이 존재한다. 분산 저장 관리는 데이터의 특성에 따라서 내용량의 데이터를 나누어 저장하거나 작은 용량의 데이터를 그룹화하여 함께 관리하는 기능이다. 부하 분산 관리는 데이터의 접근이나 용량 측면으로 특정 노드에 부하가 발생했을 때, 데이터 이주, 동적 분할, 데이터 복제 등을 통해 부하의 불균형을 해소하는 기능이다. 검색을 위한 기능으로는 데이터 이력 관리, 메타데이터 관리, 데이터 캐싱, 부분 검색 지원이 존재한다. 대표적으로 데이터 이력 관리는 데이터의 변경이 발생할 경우, 데이터의 신뢰성 향상이나 이력 기반 검색을 지원하기 위하여 데이터가 변화한 내용(이력)을 버전 또는 프로버넌스를 통해 관리하는 기능이다. 이렇게 도출된 아카이브 기술 요구사항은 대용량 과학데이터 아카이브 시스템 아키텍처를 설계하는데 활용된다.

대분류	소분류	구분	역매	관련 데이터 분류 기준
저장 관리	분산 저장 관리	Partitioning	○	데이터 크기(㎉) 데이터 연관성 등
		Categorizing	○	데이터 크기(㎉) 데이터 연관성 등
	부하 분산 관리	작업 부하 분산	△	데이터 변경 등
		데이터 부하 분산	△	데이터 크기(㎉) 등
		데이터 복제	○	-
데이터 압축 저장	데이터 압축	△	데이터 사용 패턴 데이터 변경 등	
검색 (활용)	데이터 이력 관리	Versioning	△	데이터 변경 등
	메타데이터 관리	메타데이터 관리	△	데이터 크기 등
	데이터 캐싱	데이터 캐싱	△	데이터 사용 패턴 데이터 연관성 등
	부분 검색 지원	라이브러리 연동	△	데이터 포맷 데이터 접근성 등
수집 전송	데이터 임출력 인터페이스	데이터 Import	○	-
		데이터 Export	○	-
	데이터 통신	데이터 전송	○	-
	압축 전송	△	-	

▶▶ 그림 1. 아카이브 기술 요구사항

IV. 결론

본 논문에서는 대용량 과학데이터 아카이브 시스템 설계를 위한 리포지터리 시스템 분석을 수행했다. 대표적인 과학데이터 리포지터리 시스템으로 페도라, dSpace, OSDC를 대상으로 시스템 구조, 데이터 수집, 분산 저장 관리, 부하 분산 관리, 그리고 데이터 전송 측면으로 분석을 수행했다. 분석된 내용을 바탕으로 아카이브 기술 요구사항을 도출하였다. 추후 연구로는 도출된 아카이브 기술을 기반으로 대용량 과학데이터 아카이브 시스템 아키텍처를 설계할 것이다.

■ 참고 문헌 ■

- [1] <http://fedora-repository.org/>
- [2] <http://www.dspace.com/>
- [3] <https://www.opensciencedatacloud.org/>