
빅데이터 수집 분야 Stack 통합 Test 결과보고서 [Chukwa]

2013. 04.

목 차

I. Stack 통합 테스트 개요	1
1. 목적	1
II. 빅데이터	2
1. 빅데이터 개요	2
2. 빅데이터 출현 배경	4
3. 빅데이터 중요성	6
4. 빅데이터 3대 핵심 요소	11
III. 빅데이터 구조	13
1. 빅데이터 수집	13
2. 빅데이터 수집 분야 주요 공개SW	14
IV.테스트 대상 소개	15
1. Chukwa 소개	15
V. Stack 통합 테스트	17
1. 테스트 환경	17
2. 주요 테스트 방법	18
3. 기능 테스트 수행 결과	19
4. 성능 테스트 수행 결과	20
VI. 종합	32
※ 참고자료	33

[별첨1] 공개SW Chukwa 선정지표 테스트 결과

[별첨2] Chukwa 테스트 케이스

I. Stack 통합 테스트 개요

공개SW Stack 통합테스트는 여러 공개SW들의 조합으로 시스템 Stack을 구성한 후 Stack을 구성하는 공개SW의 상호운용성에 중점을 두고 기능 및 성능테스트 시나리오를 개발하여 테스트를 진행한다.

본 통합테스트를 통해 안정된 Stack 정보를 제공하여 민간 및 공공 정보시스템 도입 시 활용될 수 있도록 한다.

1. 목적

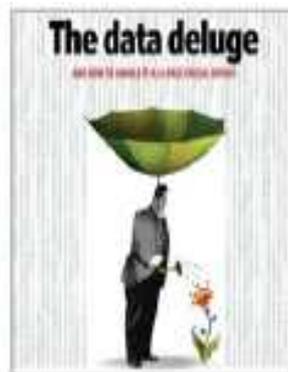
□ 공개SW Stack 통합 테스트 수행 목적

- 공개SW로 구성된 Stack이 유기적으로 잘 동작함을 확인
- 다양한 Stack 구성에 기반을 둔 테스트를 통해 안정된 Stack 조합 규명
- 공개SW 시스템 도입을 위한 Stack 참조모델의 신뢰성 정보로 활용
- 공개SW의 신뢰성과 범용성에 대한 사용자 인식 제고

II. 빅데이터

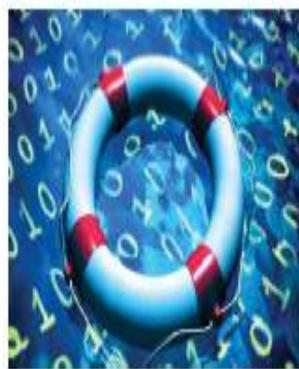
1. 빅 데이터 개요

- 빅 데이터(Big Data, BD)란 기존의 방식으로 저장/관리/분석하기 어려울 정도로 큰 규모의 데이터를 의미
 - DB의 규모에 초점을 맞춘 정의(McKinsey, 2011)
 - 일반적인 데이터베이스 SW가 저장, 관리, 분석할 수 있는 범위를 초과하는 규모의 데이터
 - DB가 아니라 업무수행에 초점을 맞춘 정의 (IDC, 2011)
 - Big Data는 다양한 종류의 대규모 데이터로부터 저렴한 비용으로 가치를 추출하고 (데이터의) 초고속 수집, 발굴 그리고 분석을 지원하도록 고안된 차세대 기술 및 아키텍처
- 최근 글로벌 경제전문지, 컨설팅 그룹이 'Big Data' 관련 특집을 잇따라 출간하며 비중 있게 보도, 분석



-
-
- SNS와 M2M 센서 등을 통해 도처에 존재하는 데이터의 효과적 분석으로 전 세계가 직면한 환경, 에너지, 식량, 의료 문제에 대한 해결책을 제시(출처: Economist, 2010.05)

- ※ SNS(Social Network Service): 특정한 관심이나 활동을 공유하는 사람들 사이의 관계망을 구축해 주는 온라인 서비스인 SNS는 최근 페이스북(Facebook)과 트위터(Twitter) 등의 폭발적 성장에 따라 사회적·학문적인 관심의 대상으로 부상함. SNS는 컴퓨터 네트워크의 역사와 같이 할 만큼 역사가 오래되었지만, 현대적인 SNS는 1990년대 이후 월드와이드웹(WWW) 발전의 산물임
- ※ M2M(Machine to Machine): 모든 사물에 센서 및 통신 기능을 결합해 지능적으로 정보를 수집하고 상호 전달하는 네트워크를 지칭함



- 데이터는 21세기 원유이며 데이터가 미래 경쟁 우위를 좌우함. 또한 기업들은 다가온 데이터 경쟁 시대를 이해하고 정보 공유를 늘려 Information silo를 극복해야 함(출처: Gartner, 2011.03)
- 빅 데이터의 활용에 따라 기업과 공공분야의 경쟁력 확보와 생산성 개선, 사업혁신/신규사업 발굴이 가능하며, 특히 의료, 공공 행정 등 5대 분야에서 6천억불 이상의 가치 창출 예상(출처: McKinsey, 2011.05)

2. 빅 데이터 출현 배경

□ 기업의 고객 데이터 트래킹/수집 행위 증가



- 기업들은 온라인/오프라인 사용자 정보, 소비자 행태에 대한 정보수집에 적극적
- 고객관련 정보 수집의 증가로 더 많은 데이터 스토리지와 정교한 분석 능력을 필요

※ Tesco는 매달 15억 건 이상의 (고객) 데이터를 수집

□ 멀티미디어 콘텐츠와 콘텐츠 사용에 관한 정보의 증가

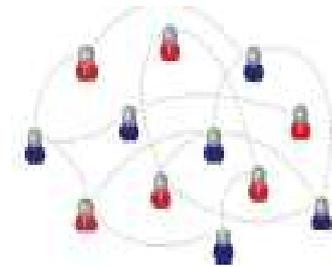


- CT 스캔, CC카메라 등 다양한 부분에서 대용량 멀티미디어 콘텐츠 생산 증가
- 고객관련 정보 수집의 증가로 더 많은 데이터 스토리지와 정교한

분석 능력을 필요

- 오리지널 콘텐츠뿐 아니라 콘텐츠 소비에 관한 정보도 대량 생산 (사용자정보, 선호 등)

□ SNS의 급격한 확산과 비정형 데이터의 폭증



- SNS는 스마트폰의 확산과 더불어 젊은 층에서 중장년 층으로까지 확산
- Facebook에서만 매일 한 이용자당 평균 90개 이상의 콘텐츠를 업로드
- YouTube에서는 1분 마다 24시간 분량의 비디오가 업로드 → SNS 미디어 데이터 폭증

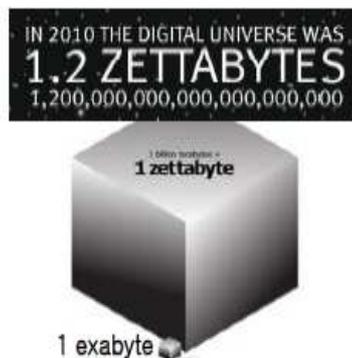
□ M2M 확산에 따른 센서 저변 확대



- 2012년 현재 3천만 개 이상의 사물인터넷 센서가 설치 (향후 5년 동안 CAGR 35% 증가)
- 원격 헬스 모니터링을 통한 헬스케어, RFID를 이용한 소매업, 스마트 미터 기술을 활용한 유틸리티 사업에서도 데이터 발생량이 증가할 것으로 전망
- YouTube에서는 1분 마다 24시간 분량의 비디오가 업로드 → SNS 미디어 데이터 폭증

3. 빅 데이터 중요성

- 우리는 이미 제타(zettabyte, 10²¹) 시대에 살고 있으며 Big Data 추세는 스마트 단말, M2M 센서 확대보급 등으로 더욱 가속화될 전망

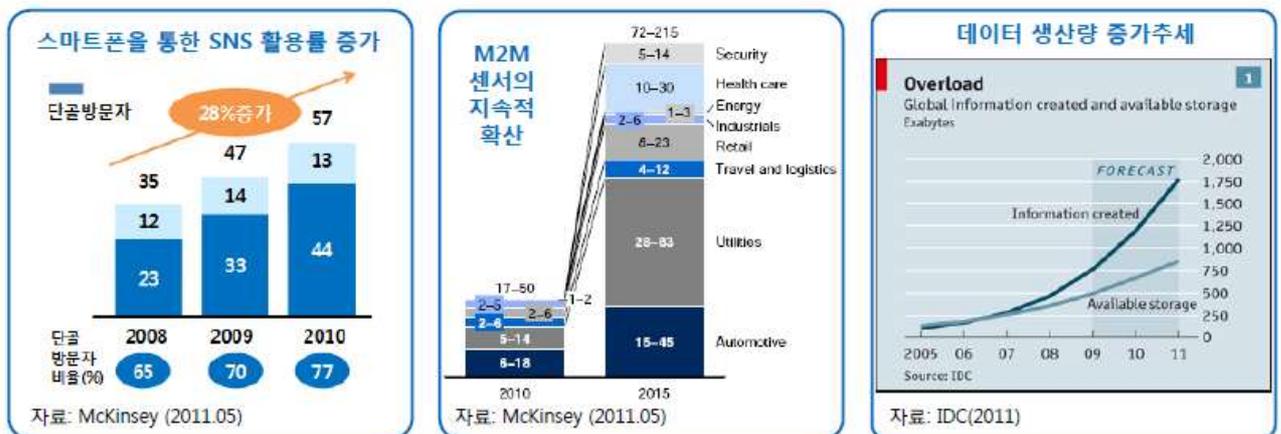


- 세계는 2010년 zettabyte 시대에 돌입(1.2 zettabyte의 정보 생산)
 - 유사 이래 2003년까지 생산된 모든 정보의 합 = 5 exabyte
 - 1 Zettabyte는 美의회도서관 저장정보(235 terrabyte, 11/4 기준)의

4백만 배에 해당사이래 2003년까지 생산된 모든 정보의 합 = 5 exabyte

- 16GB iPad를 축구장 크기로 쌓아도 대기권 2배 높이에 도달

○ 데이터 생산량은 스마트폰의 확산, SNS 사용 확대, M2M 센서 구축 등으로 향후에도 급속히 증가할 전망



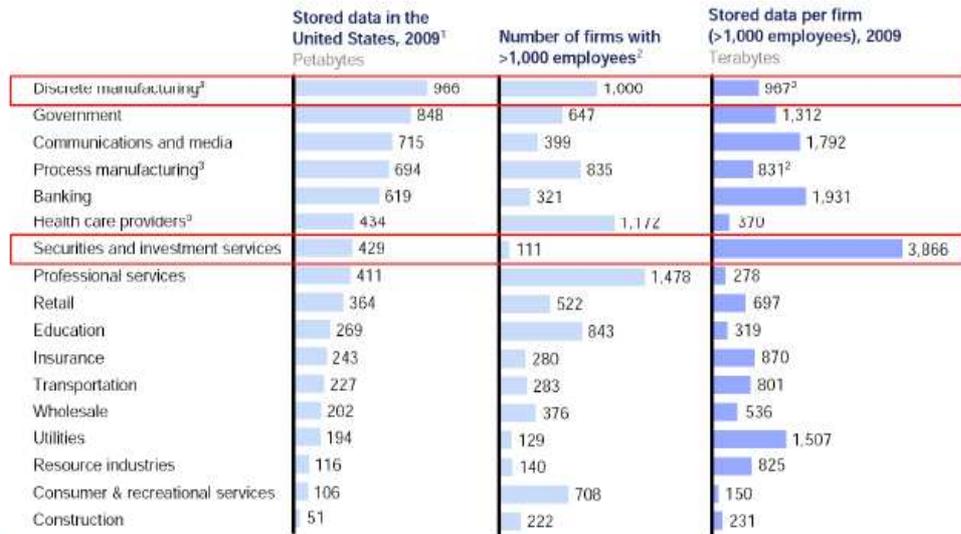
□ 모든 기업이 보유한 Big Data가 '거대한 가치 추출이 가능할 만큼' 충분한 규모에 도달해 누가 먼저 그 가치를 추출해 내느냐가 향후 기업의 성패를 가늠할 상황에 직면

○ Big Data 현상은 거의 모든 산업 부분에서 진행되어 옴

- (산업 부문별 총합으로 보면) 제조업 부분이 보유한 데이터 양이 가장 많고, (1천명 이상 직원 보유 기업 별로 보면) 증권/ 투자 서비스업 부분의 기업들이 가장 많은 정보 보유

○ 각 기업의 Big Data 보유 규모는 '거대한 가치를 창출할 정도의

정보'를 응축하고 있는 수준에 도달



자료: McKinsey (2011.05)

- (미국) 거의 모든 기업이 100 terrabyte 이상의 정보를 보유 중이며, 상당수는 1 petabyte 이상 보유

□ Big Data의 '양적 거대함'은 많은 분야에서 불가능을 가능으로 전환함. Google의 Big Data 솔루션이 빚어낸 Magic - IBM의 실패 프로젝트를 성공으로 변신

- IBM과 Google은 자동 번역 프로그램을 개발하기 위해 기존의 방식과 다른 접근법을 채택

- 40여년 동안 과학자들은 컴퓨터에게 명사, 동사와 같은 구조와 음운을 이해시키려고 노력

- IBM과 Google은 기존 방법과 달리 전문가가 번역한 문건을 DB화해서 비슷한 문장과 어구를 대응 시키는 통계적 기법을

활용하여 번역 문제를 해소하려고 시도

- 매칭에 참고하는 DB 차이가 두 기업의 자동번역 프로젝트의 성패를 좌우



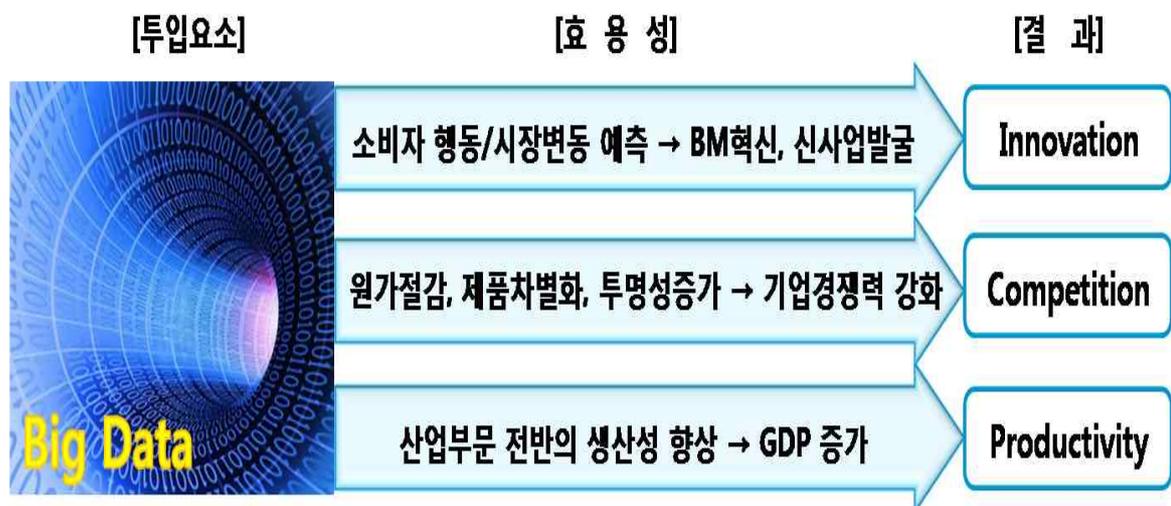
- IBM은 캐나다 의회의 '수백만 건'의 문서를 활용, 영어-불어 자동번역 시스템 개발을 시도했으나 실패
- 반면 Google은 동일방식이지만 '수억 건'의 자료를 활용, 50개 언어 간의 자동번역 시스템 개발 성공
- ※ Google Big Data 구축 방법: 20여개의 언어로 번역된 EC의 문건을 검색을 활용, copy 한 후에 Book스캐닝 프로젝트에서 수천만 권의 전문 번역 DB 구축함
- ※ Google은 Big Data방식을 Spell-check와 음성인식 분야에도 적용하고 있는데, Microsoft가 장기간 대규모 투자로 만들어 낸 스펠링 교정보다 우수한 프로그램을 매일 3억 건씩 발생하는 '검색창의 오타 입력과 수정정보'를 활용하여 개발해 냄. 음성인식 능력의 향상도 반복되는 사용자 자옴교정 정보를 feedback 해서 Big Data를 만들고 이를 활용하여 개선

□ Big Data는 모바일 스마트 혁명의 핵심 자원으로 산업혁명에서의 철과 석탄의 역할을 하며 제 4의 경영자원으로서 혁신과 경쟁력 강화, 생산성 향상을 촉진

○ 산업혁명에서는 철과 석탄이, IT 혁명에서는 인터넷이 세계 경제 변화를 지탱하는 핵심 요소였듯이 다가올 모바일 스마트 혁명에서는 Big Data가 경제 변화의 핵심 자원 역할을 할 것



○ Big Data 는 제 4의 경영자원으로서 혁신과 경쟁력 강화 생산성 향상을 촉진



4. 빅 데이터 3대 핵심요소

□ 클라우드 컴퓨팅

- 클라우드 컴퓨팅은(기존의 IT 환경에 비해) 신속성과 유연성 그리고 규모의 경제를 제공
 - 2020년에는 생산되는 데이터의 약 35%가 클라우드에 있거나 클라우드를 거쳐 갈 전망이며, 클라우드 공급자는 Big Data와 관련된 모든 영역에서 중요한 역할을 할 것
- Big Data는 저장, 보관, 처리 속도 및 비용 측면에서 기업들에게 새로운 도전이 될 것
 - Big Data는 기존방식으로 처리하기엔 데이터 규모가 크고 컴퓨팅 파워가 부족하기 때문에 Hadoop, MapReduce 같은 클라우드 기반 솔루션들의 적용이 본격화 되고 있음

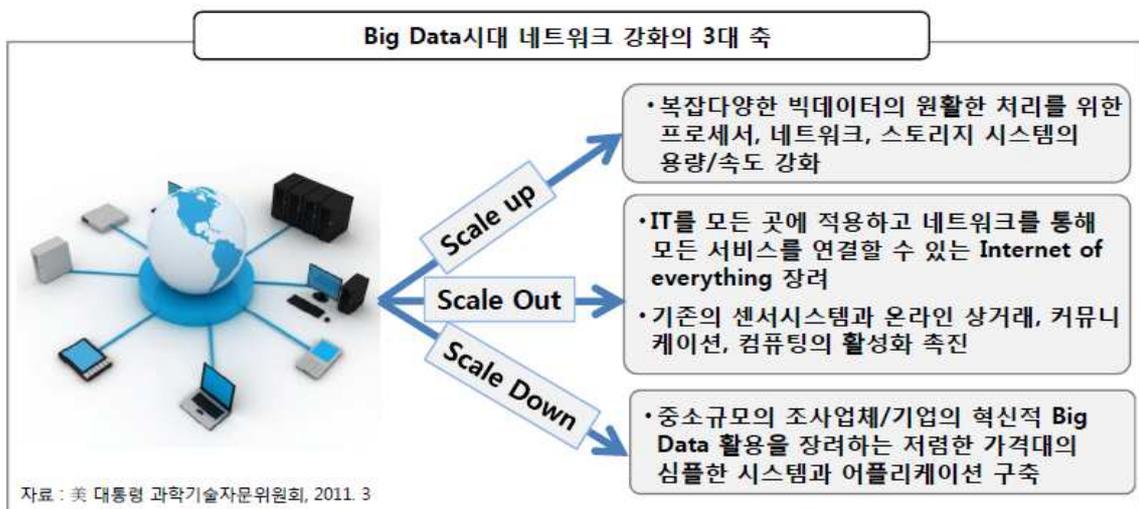
□ 실시간 분석

- SNS, 인터넷 게시판 등의 Big Data를 실시간 분석, 활용은 커다란 가치 창출 기회를 제공하고 있음
 - Big Data를 통해 프로세스 개선, 실행 가능 정보 및 고객만족 이슈 도출로 빠른 next best offer 제공 가능
 - 특히, 이용자들의 好惡가 빠르게 반영되는 SNS Data를 서비스 개선에 실시간으로 적용하는 기업이 증가

- 상당 부분의 Data는 폐기되거나 이어지는 데이터에 의해 대체
 - 생산되는 데이터를 모두 저장한다는 것도 이미 불가능 ('07년 생산량이 저장량 증가를 추월)
 - 의료 분야, 입자물리학 실험실 등에서 발생하는 데이터의 90%가 폐기되고 있음

□ 네트워크 역량 강화

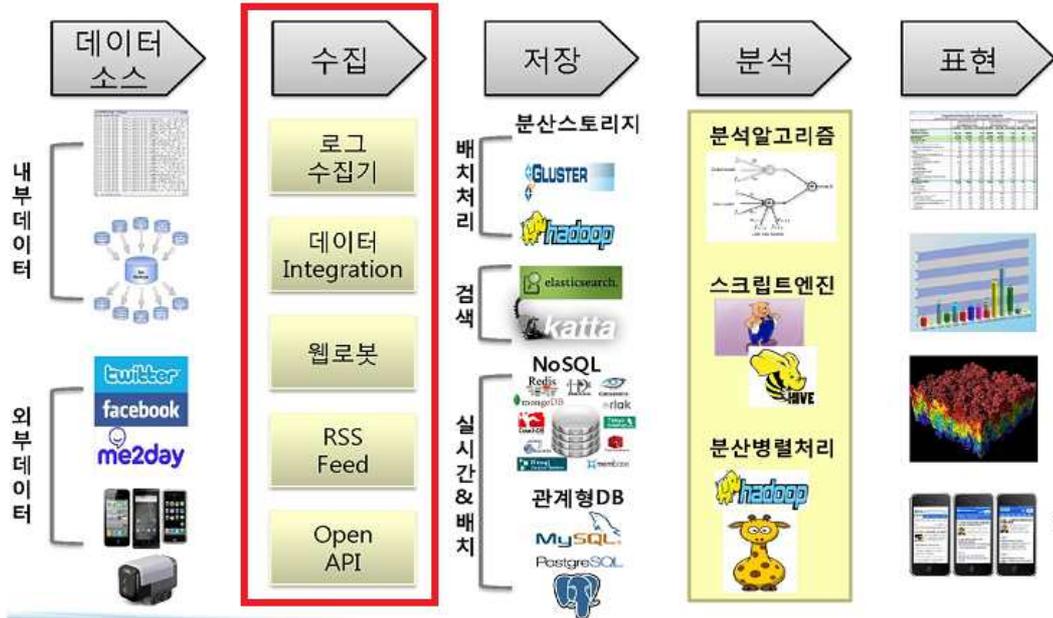
- 폭증하는 실시간 Big Data의 효율적 처리를 위해서는 네트워크 역량 강화도 병행해야 함
 - Big Data는 단순 양의 증가뿐만이 아니라 데이터의 복잡성/다양성의 증가를 의미하기 때문에, 폭증하는 다양한 데이터 처리를 위한 네트워크 강화가 핵심요소로 부상
 - 3가지 측면의 네트워크 강화(scaling of network)를 통해 Big Data를 통한 가치창출의 기반 조성 가능
 - 기존 유무선 네트워크 및 주파수 인프라 관리 또한 복잡/다단한 Big Data시대에 부합하도록 대응 필요



III. 빅데이터 구조

1. 빅 데이터 수집

빅데이터 처리흐름은 아래의 그림과 같이 5가지 로 표현할 수 있다.



[빅 데이터 처리 흐름]

이 다섯 가지 처리 흐름 중 데이터의 수집 부분에 해당하는 공개SW의 테스트를 진행하기로 하였다.

데이터의 수집이란 자료 처리 시스템에 들어갈 자료를 모으는 일 또는 여러 장소에 있는 자료를 한 곳으로 모으는 일을 말한다. 빅 데이터는 데이터의 생성 양·주기·형식 등이 기존 데이터에 비해 너무 크기 때문에, 종래의 방법으로는 수집이 어려운 방대한 데이터를 말한다. 이런 방대한 데이터를 분석하고 활용하기 위해서는 대용량 데이터를 수집할 수 있는 데이터 수집 프레임워크가 필요하다.

2. 빅 데이터 수집 분야 주요 공개SW

빅 데이터 수집 분야 주요 공개SW 중 많이 알려진 Apache군의 Flume, Kafka, Chukwa와 FaceBook의 Scribe를 테스트 SW군으로 선정하고 공개SW 선정지표를 통하여 하둡 시스템과 높은 호환성을 보이는 Apache Chukwa로 테스트를 진행 하였다.

[표 III-1. 수집 분야 주요 공개SW]

제품명	Stack 환경	홈페이지	비고
Flume	Linux	http://flume.apache.org/	
Scribe	Linux	http://github.com/facebook/scribe/	
Chukwa	Linux	http://incubator.apache.org/chukwa	
Kafka	Linux	http://kafka.apache.org/	

[표 III-2. 수집 분야 주요 공개SW 선정지표 점수]

분야	세부분야	대상	항목[배점]				총점 [100]
			Document [25]	Support [25]	Product [30]	Community [20]	
빅데이터	수집	Chukwa	16.2	17.5	18.8	0.0	52.5
		Flume	22.2	17.5	25.0	5.8	70.5
		Kafka	18.6	20.0	17.5	6.7	62.8
		Scribe	20.4	15.0	21.3	6.7	63.3

※ 공개SW 선정지표에 의해 선정된 공개SW가 품질/성능의 우수성을 뜻하는 것은 아님

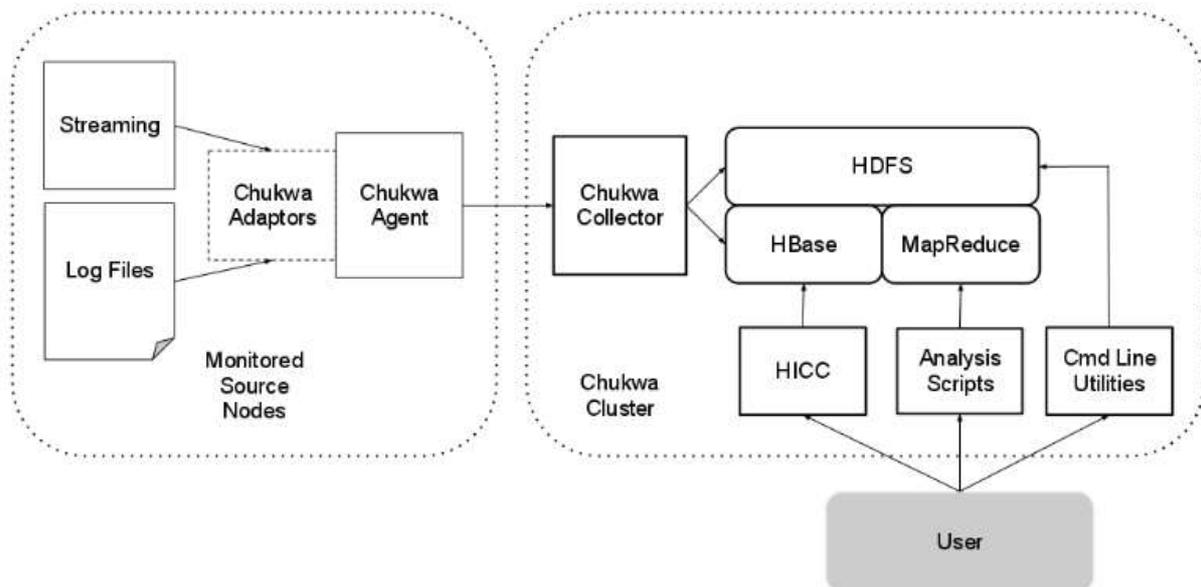
※ [별첨1] 공개SW Chukwa 선정지표 테스트 결과

IV. 테스트 대상 소개

1. Chukwa 소개

Chukwa은 하둡 프로젝트의 서브 프로젝트로 진행 중인 오픈소스 프로젝트로 분산된 서버의 로그를 하둡 파일 시스템에 저장하고 하둡 맵리듀스를 이용해 로그의 분석을 수행하는 솔루션이다. 척와는 범용적인 로그 수집과 로그 관리를 위한 솔루션으로 개발됐지만 하둡 클러스터의 로그와 서버의 상태 정보를 관리하는 기능도 들어있다. 하둡 역시 분산된 환경에서 수십 ~ 수천 대 규모로 운영되기 때문에 하둡 데몬에서 출력하는 로그와 사용자가 수행시킨 작업에서 출력하는 로그는 각 서버에 저장된다. 하둡을 이용해 작업을 수행하다 보면 사용자 로그를 보기위해 여러 노드를 접속해야하는 불편함이 있는데, 척와를 이용하면 쉽게 로그를 수집할 수 있다.

[그림 IV-1. Chukwa 구조]



□ Chukwa 구성

o chunk

Chukwa에서 전달되는 데이터의 단일단위, 하둡에서도 사용된다.

o adaptors

이벤트를 발생시키는 대상을 정의

o Agent

로그를 collector서버로 전송하는 기능

o Collector

Agent에서 전송된 로그를 HBASE와 HDFS에 저장하는 기능

※ 추가적인 자세한 정보는 아래의 링크 정보 참조

- <http://incubator.apache.org/chukwa/docs/r0.5.0/admin.html>

V. Stack 통합 테스트

1. 테스트 환경

테스트 SW

[표 V-1. 테스트 SW]

SW	Version
Chukwa	0.5.0
Hadoop	1.1.2
HBase	0.94.5
Zookeeper	3.4.5

Stack 환경

[표 V-2. Stack 환경]

Stack	Agent	IP	Collector	IP
A	Ubuntu 12.04	121.162.249.12	CentOS 6.2	121.162.249.18

HW 환경

[표 V-3. HW 환경]

제조사	모델명	CPU	MEM	Disk	NIC	서버수
HP	DL360G6	Quad-Core 2.40Ghz~2P	8GB	500GB	Gigabit 1Port	X1
IBM	X3550M2	Intel Xeon(R)CPU 2.40GHz	8GB	320GB	Gigabit 1Port	X3
SUN	X4140	2x AMD 2380(2.5GHz)	8GB	140GB	Gigabit 1Port	X1

※ Agent 서버 = HP HW 1대로 Stack 구성(Ubuntu 12.04)

※ Collector 서버 = IBM 3대+ SUN 1대로 Hadoop+Zookeeper+Hbase Stack 구성
(CentOS 6.2)

2. 주요 테스트 방법

□ 시나리오 테스트

시나리오 테스트 기법은 단일 기능에 대한 결함 여부를 확인하는 것이 아니라, 서로 다른 컴포넌트 사이의 상호작용과 간섭으로 발생할 수 있는 결함을 발견하기 위한 기법이다.

본 테스트에서는 사용자 시나리오 테스트 기법을 적용하여 Chukwa를 사용하는 사용자들이 사용할 수 있는 항목 중 Agent와 Collector에 대한 사용자 시나리오를 도출하였다. 각각의 항목에서 도출한 세부 시나리오는 사용자가 일반적으로 수행할 수 있는 시나리오를 추출하여 테스트케이스로 작성하였다.

□ 상호 운용성 기반 테스트

상호 운용성은 서로 다른 기술로 이루어진 제품이나 서비스가 상호작용 상의 오류가 없는지 검증하는 기법으로, 본 테스트에서는 애플리케이션이 지원하는 Stack을 구성하여 애플리케이션과 Stack 환경 사이의 상호작용 상의 동작여부를 검증하였다.

3. 기능 테스트 수행 결과

기능 테스트 수행 관련 세부 시나리오는 별첨 「Chukwa 테스트 케이스」 문서를 참고한다.

□ 테스트 시나리오 현황

[표 V-4. 테스트 시나리오 현황]

기능	테스트 시나리오	테스트 케이스
Agent	11	11
Collector	5	5
시작/종료	1	2
모니터링	1	2
합 계	18	20

□ 테스트 결과

기능 테스트 시나리오를 통한 테스트 수행 결과 Agent, Collector 등 시나리오 상의 모든 기능이 예상 결과와 동일하게 동작함을 확인하였다.

[표 V-5. 테스트 결과]

분류		PASS	FAIL	N/A
기능	개수			
Agent	11	11	0	0
Collector	5	5	0	0
시작/종료	2	2	0	0
모니터링	2	2	0	0

4. 성능 테스트 수행 결과

성능 테스트의 경우 하드웨어 사양뿐 아니라, OS 및 애플리케이션 환경 구성에 따라 성능 측정 결과가 상이하므로, 실제 운영 시스템 환경에 따라 테스트 결과가 다를 수 있다.

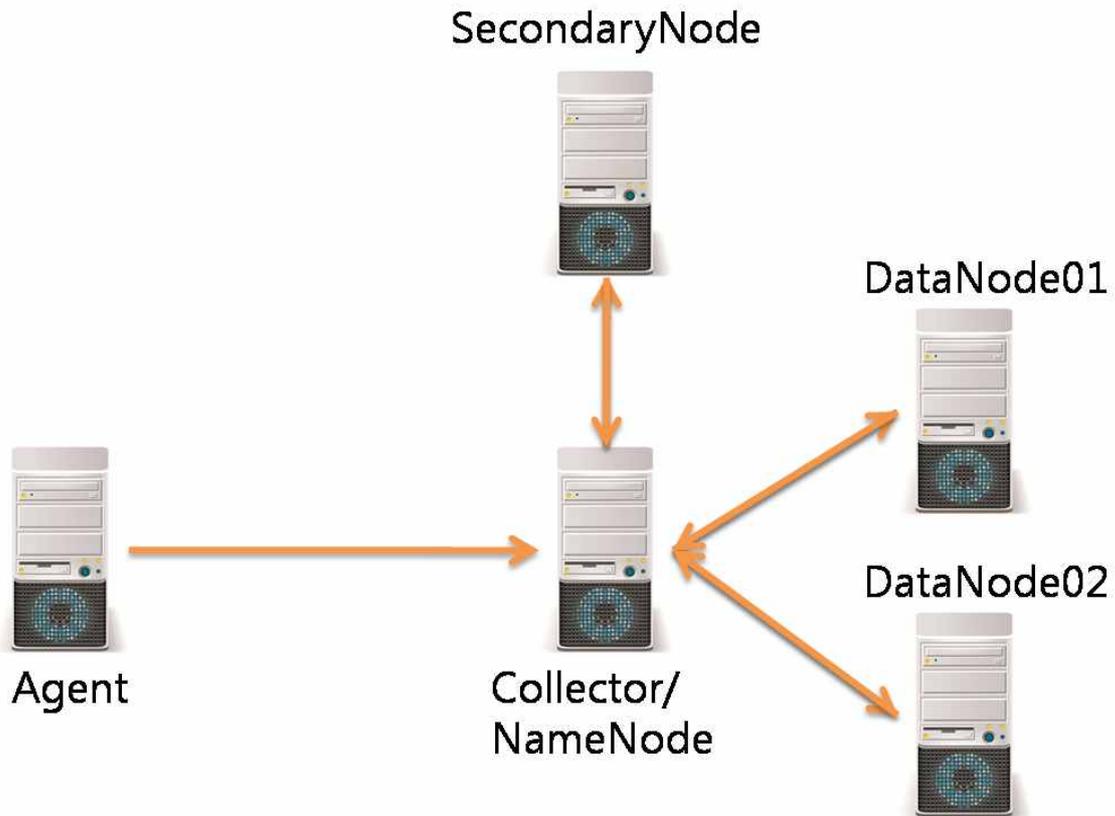
본 성능 테스트는 Chukwa 시스템이 가동되는 상황에서 로그파일 전송 시나리오를 재현하여, Chukwa Agent에서 전달된 정보를 수집하는 Collector Server가 HDFS와 HBASE에 저장 할 때의 자원사용률을 측정한다.

□ 테스트 시나리오

[표 V-6. 테스트 시나리오]

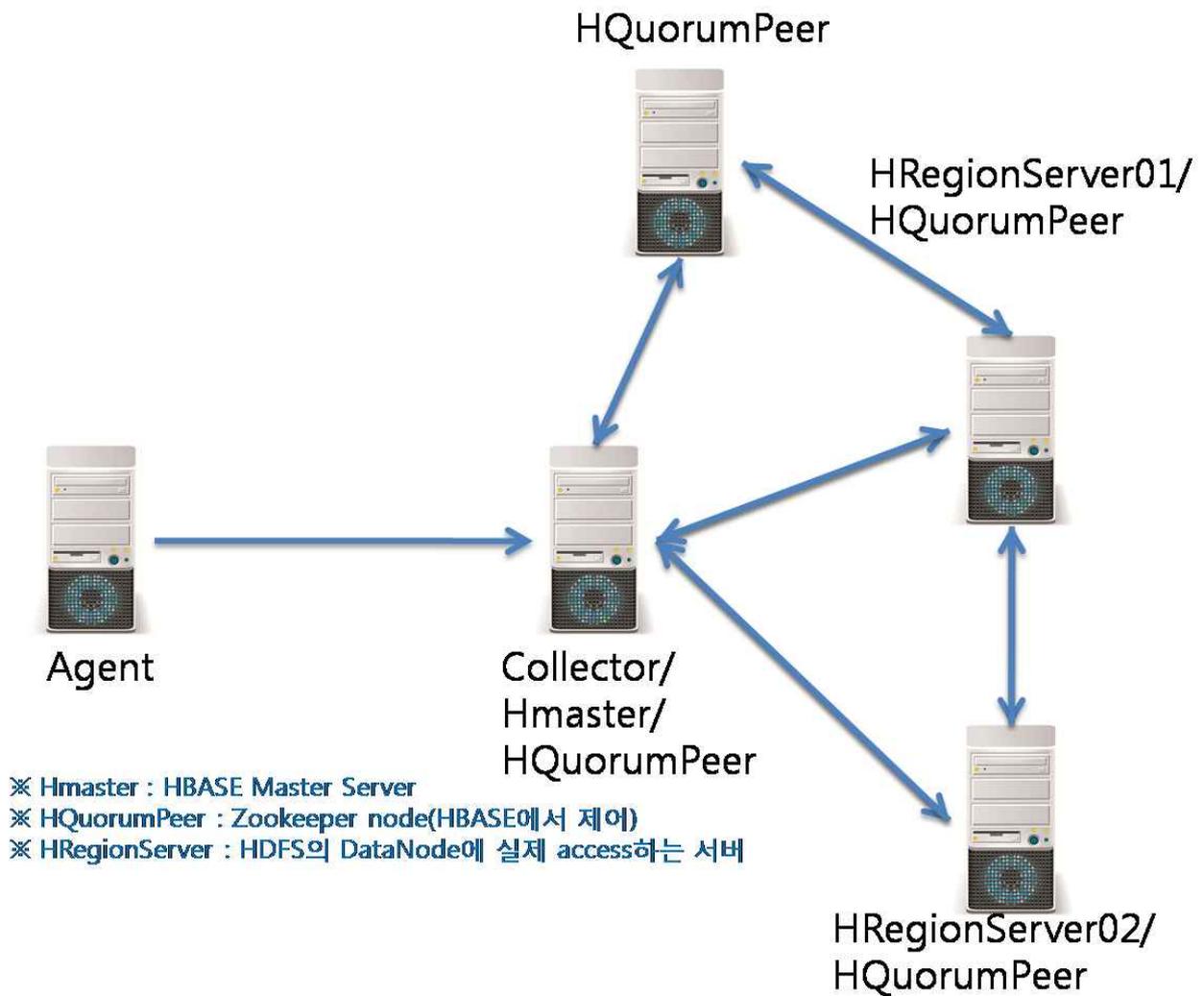
대분류	시나리오ID	시나리오
F_HA	F_HA_U_100	100MB의 단일 파일을 전송하여 Hadoop에 저장(Ubuntu->CentOS)
F_HB	F_HB_U_100	100MB의 단일 파일을 전송하여 HBase에 저장(Ubuntu->CentOS)

□ 서버 구성



- ※ NameNode : 하둡의 ROOT같은 역할
- ※ SecondaryNode : 하둡의 ROOT의 BACKUP역할
- ※ DataNode : 하둡에 저장된 데이터가 복제되는 곳. 현재 replication은 3(NameNode포함)

[그림 V-1. HDFS 성능 테스트 환경 정보]



[그림 V-2. HBASE 성능 테스트 환경 정보]

□ 측정항목

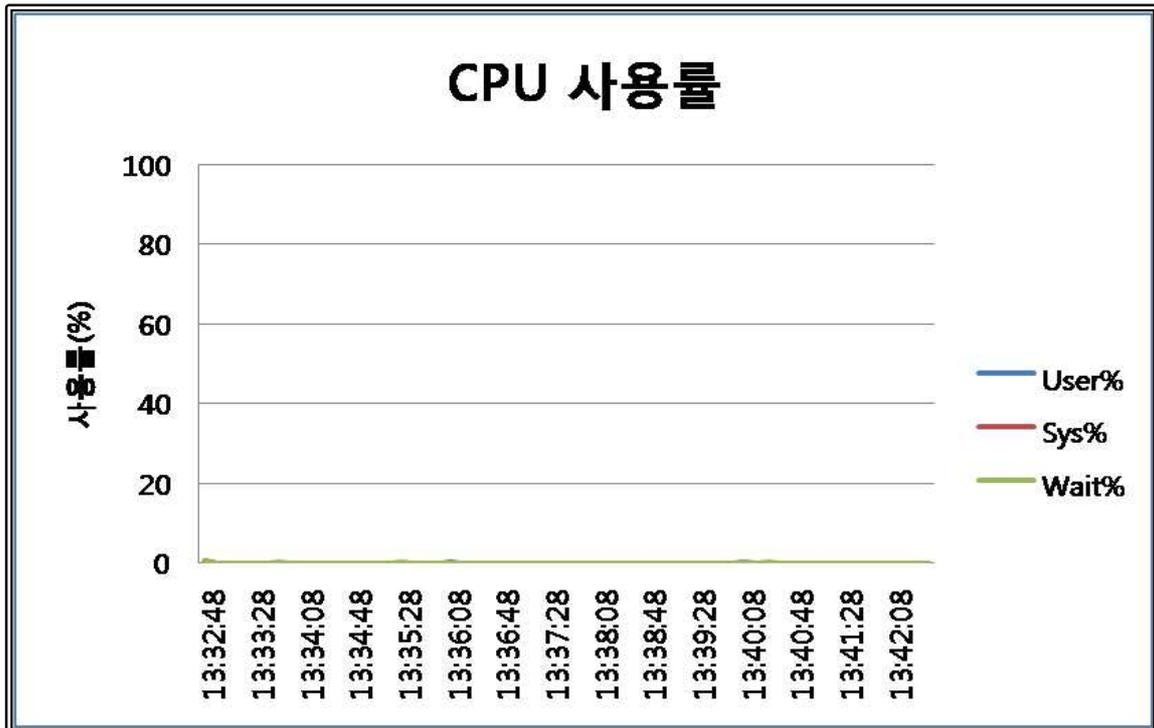
[표 V-7. 측정항목]

항목	내용
CPU 사용률	프로세스에서 CPU(Central Processing Unit)를 사용한 비율(%)
메모리 사용률	Physical 메모리 사용량
네트워크 I/O	네트워크 입출력 작업의 양(K/B)

□ 테스트 결과

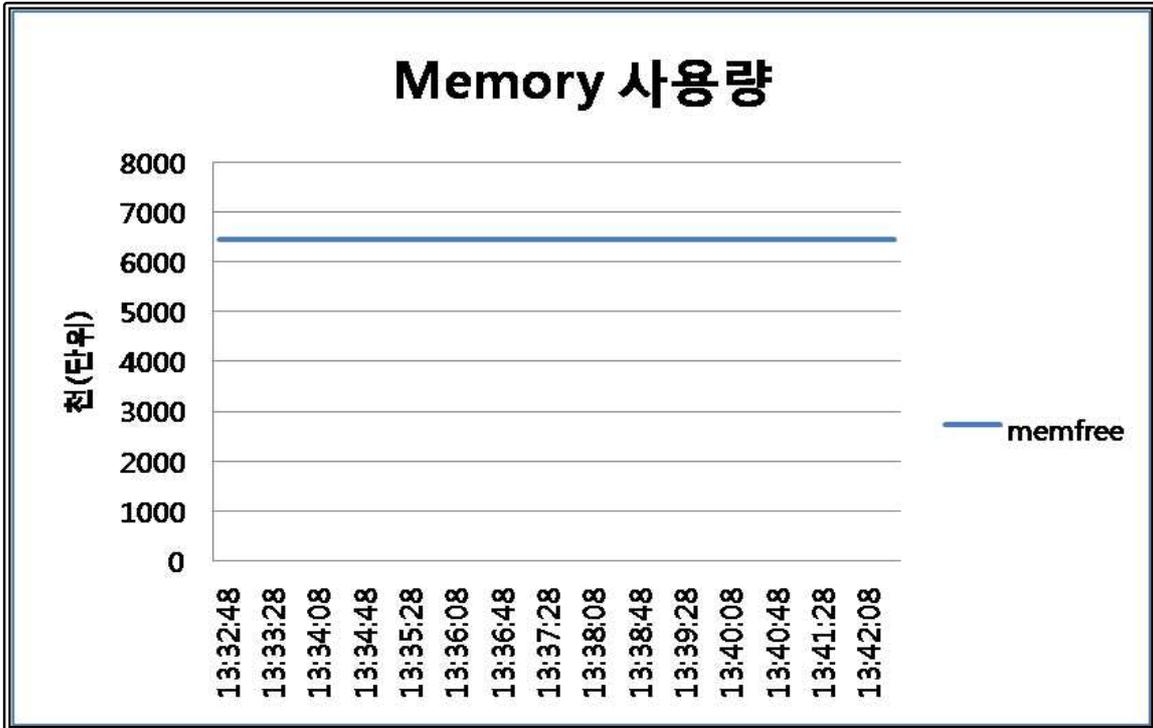
o Chukwa Collector 실행 전 자원사용률 측정

- Master (Hadoop+HBASE_Zookeeper) Server 측정결과



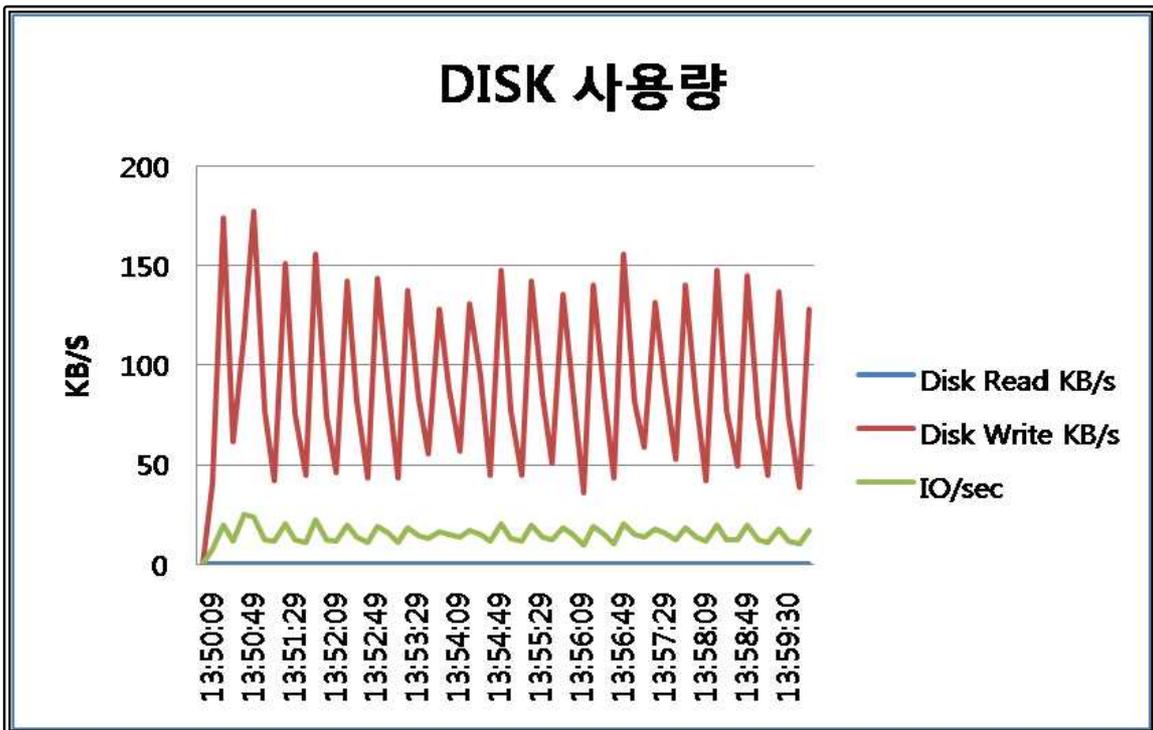
[그림 V-3. CPU 자원사용률]

- CPU 자원을 거의 사용하지 않는 것을 확인 할 수 있다(2% 미만)



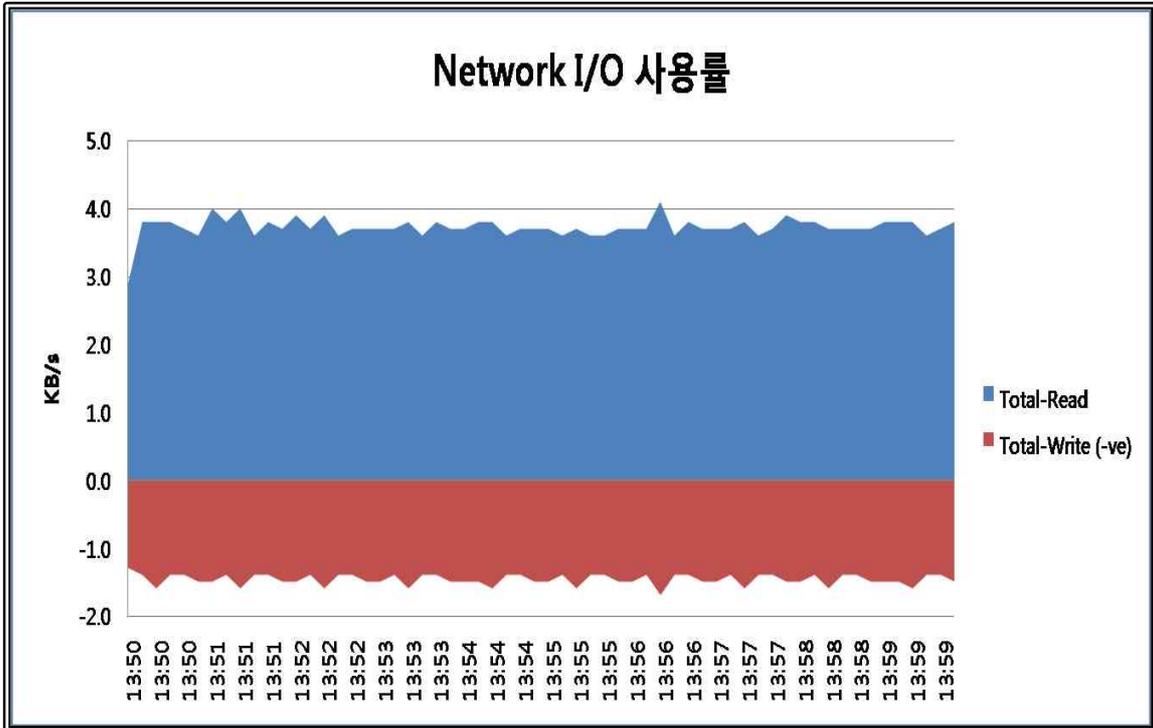
[그림 V-4. Memory 자원사용률]

- 여유 Memory는 평균 6500을 유지하고 있음



[그림 V-5. DISK 자원사용률]

- DISK 사용량은 초당 200KB 이내



[그림 V-6. DISK 자원사용률]

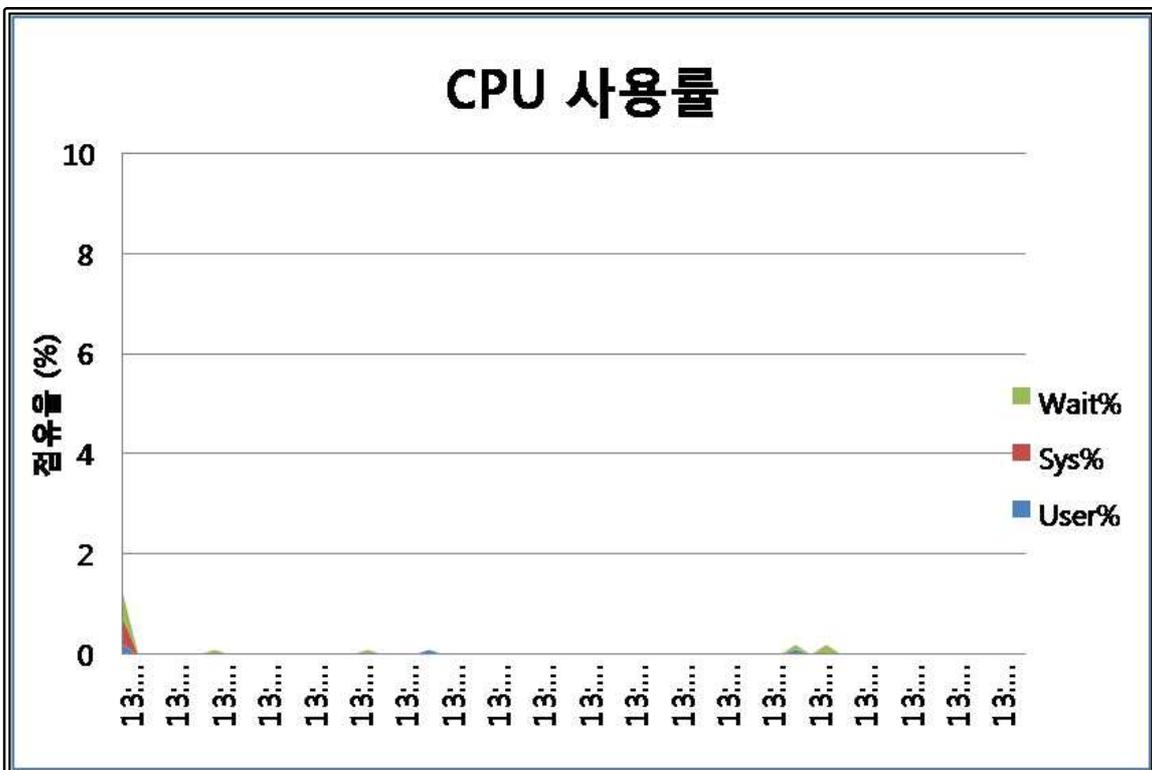
- Total-Read = Agent Server로부터 데이터를 전송 받는 양
- Total-Write = 실제로 저장되는 데이터의 양

○ 대분류 F_HA 성능 테스트 결과

- 수행정보

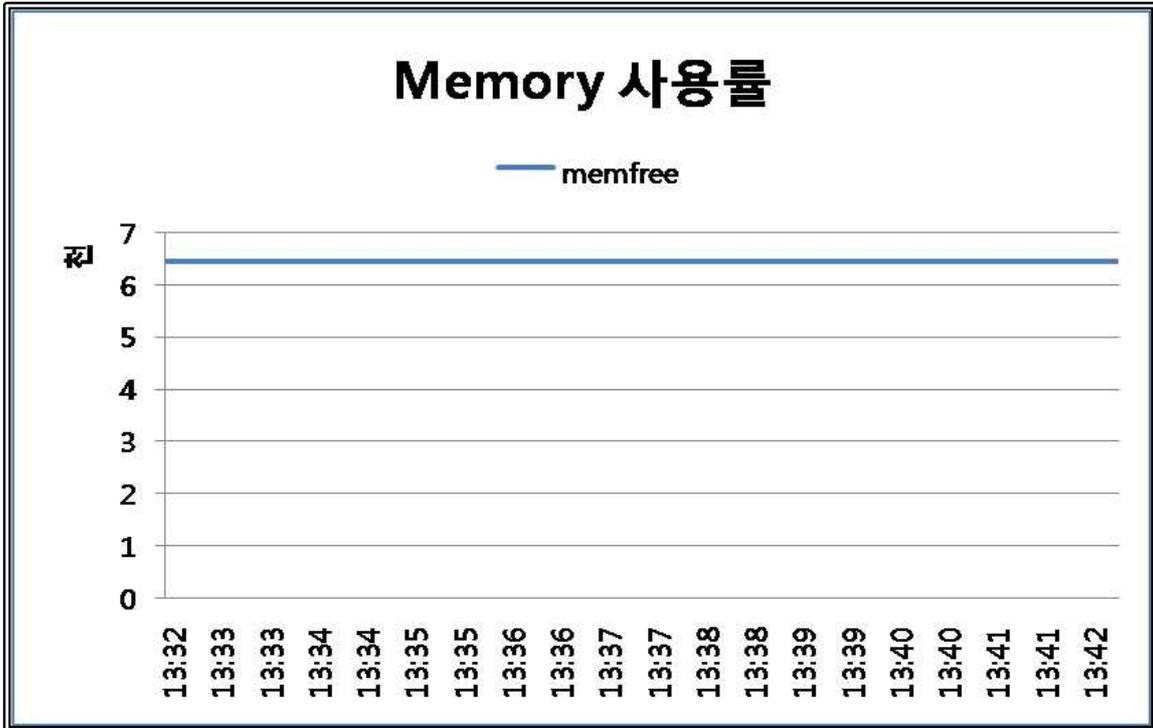
수행 조건	- 100MB의 단일로그파일을 Hadoop에 저장 (Ubuntu에서 CentOS로 전송)
-------	------------------------------------------------------

- Hadoop Master Server 측정결과



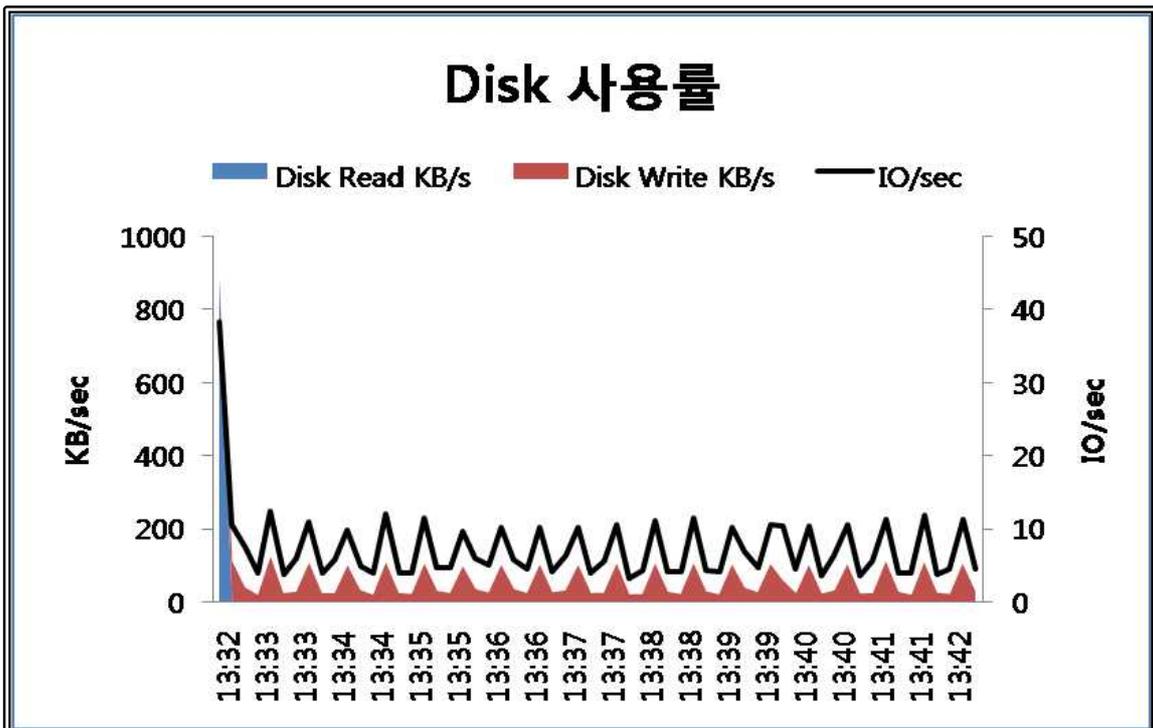
[그림 V-7. CPU 자원사용률]

- CPU 자원 사용률이 낮은 것을 확인 할 수 있다(2% 미만)



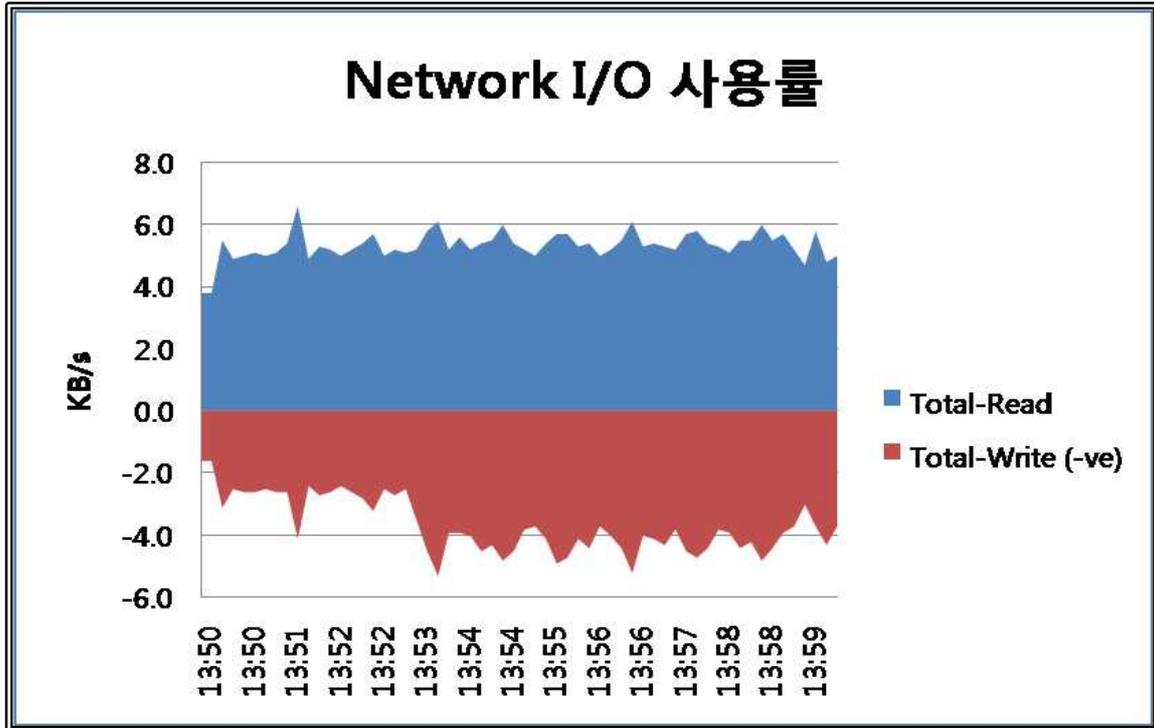
[그림 V-8. Memory 자원사용률]

- 메모리 사용률에 변화가 없음



[그림 V-9. DISK 자원사용률]

- 초당 200KB이내의 사용률이 확인됨



[그림 V-10. Network 자원사용률]

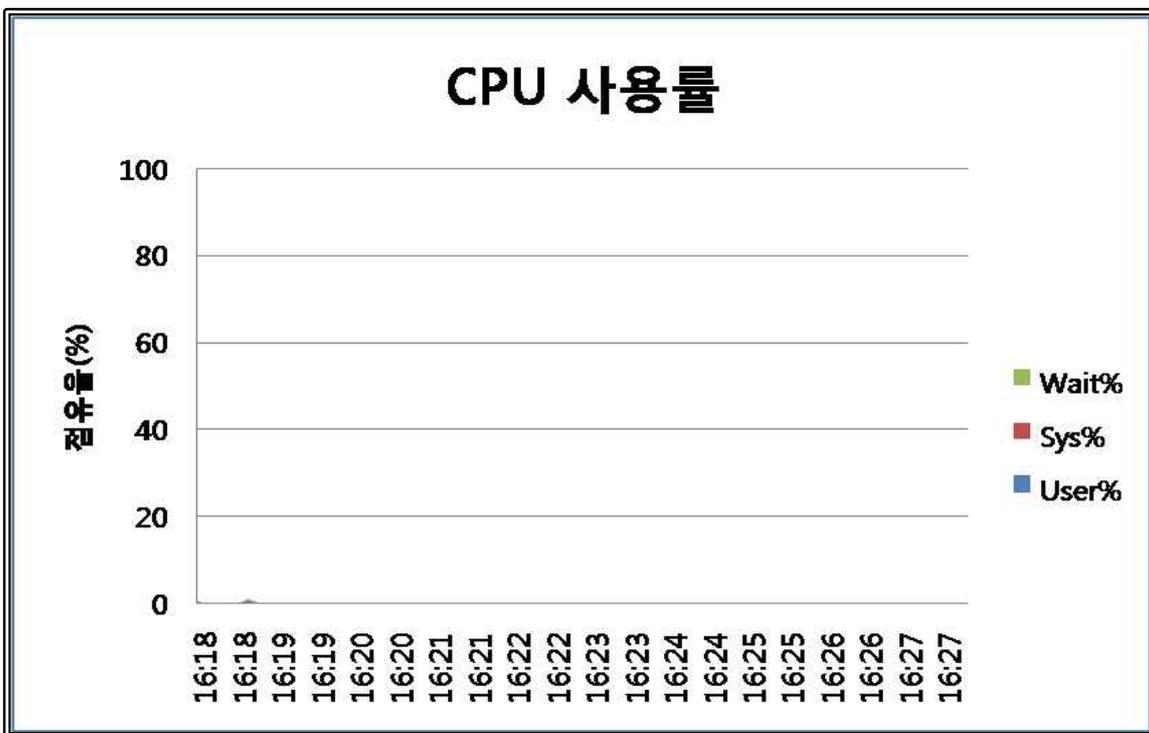
- Total-Read = Agent Server로부터 데이터를 전송 받는 양
- Total-Write = 실제로 저장되는 데이터의 양

o 대분류 F_HB 성능 테스트 결과

- 수행정보

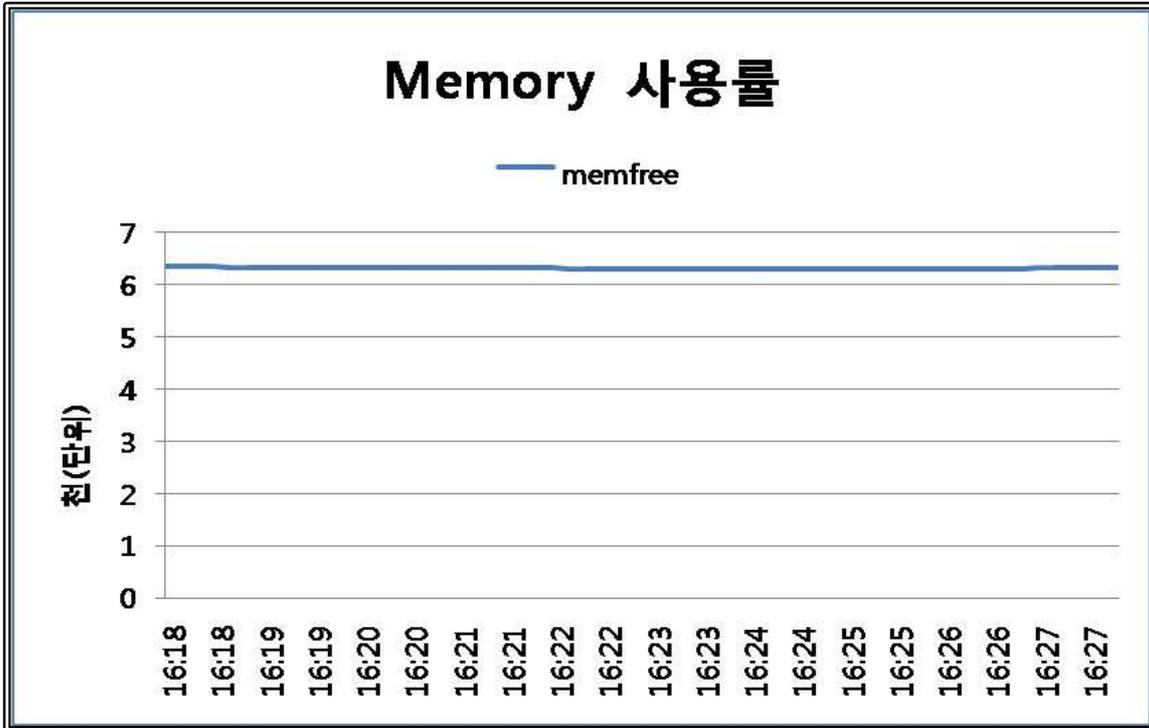
수행 조건	- 100MB의 단일로그파일을 HBase에 저장 (Ubuntu에서 CentOS로 전송)
-------	-----------------------------------------------------

- HBase+Zookeeper Server 측정결과



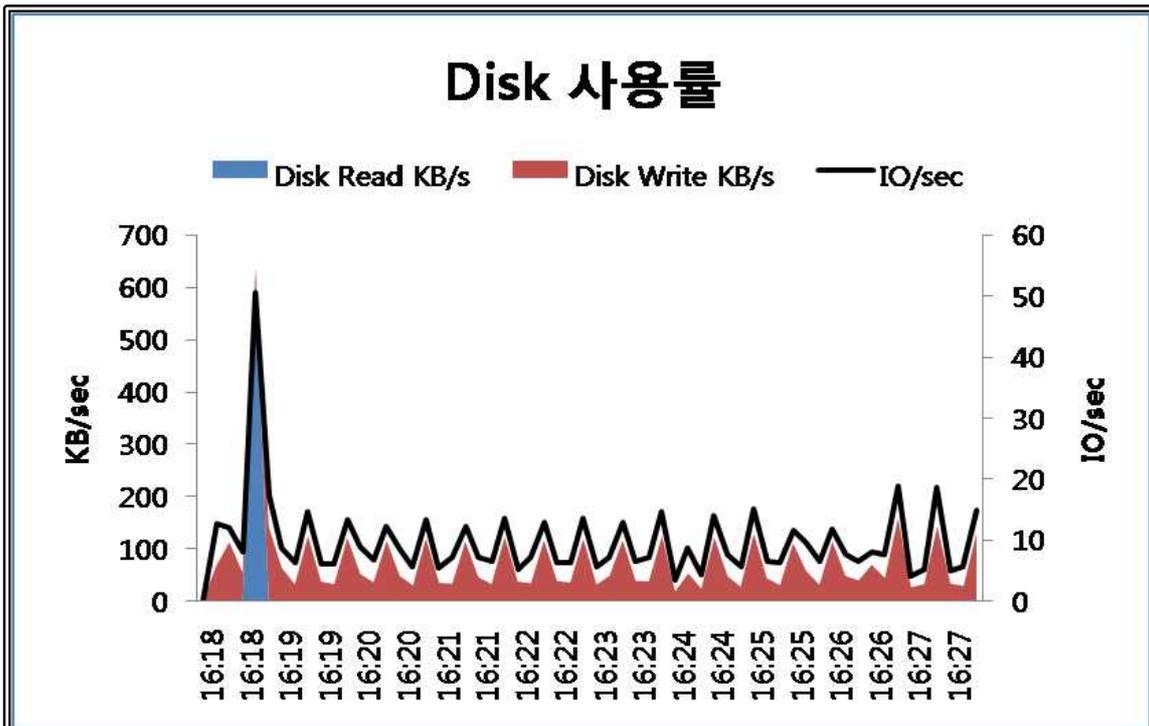
[그림 V-11. CPU 자원사용률]

- CPU 자원 사용률이 낮은 것을 확인 할 수 있다(2%미만)



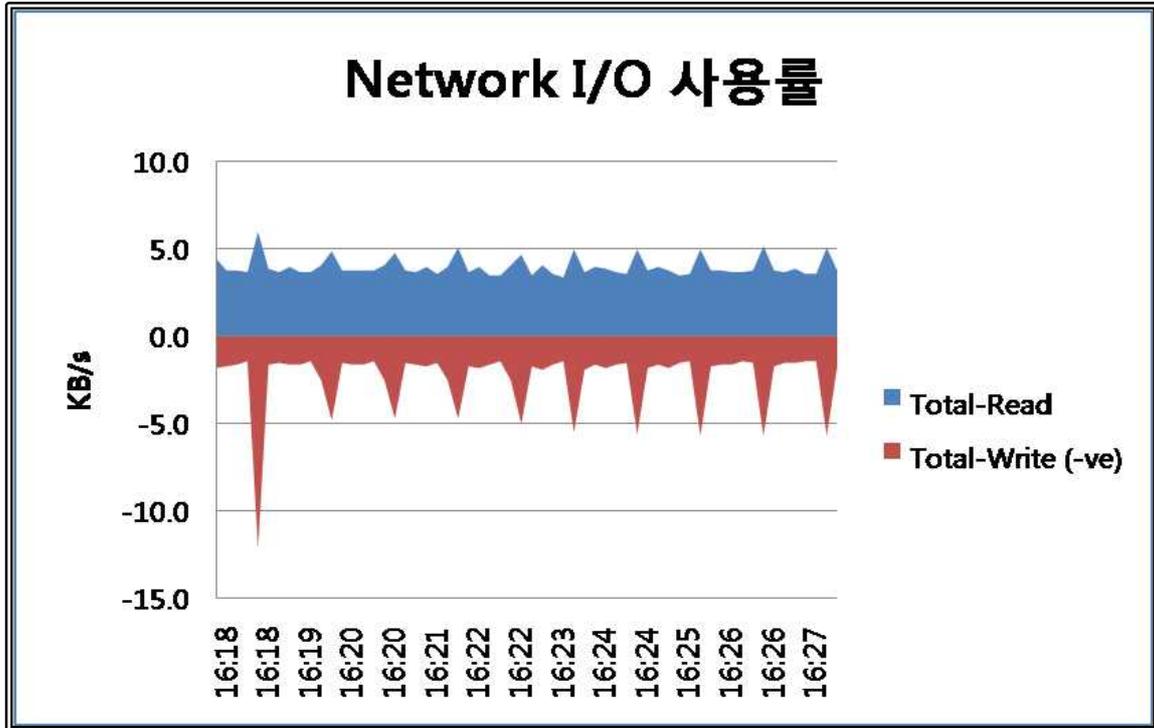
[그림 V-12. Memory 자원사용률]

- 약간의 메모리 사용(100~200MB 이내 사용됨)



[그림 V-13. DISK 자원사용률]

- 초당 200KB이내의 사용률이 확인됨



[그림 V-14. Network 자원사용률]

- Total-Read = Agent Server로부터 데이터를 전송 받는 양
- Total-Write = 실제로 저장되는 데이터의 양

VI. 종합

- Chukwa 기능 테스트 수행 결과 공개SW로 구성된 Stack 상에서 각 기능 시나리오 수행 시 치명적 오류 또는 심각한 장애가 발생하지 않았으며, Stack을 구성하는 각 공개SW가 유기적으로 동작함을 확인하였다.

- Chukwa 성능 테스트 수행 결과 CPU, Memory, Disk, NetWork 사용률은 HBase와 Hadoop 모두 자원 사용률이 매우 낮아 전체적으로 안정적 수치를 보였다. 이번 테스트에서는 시간적인 여건으로 인해 Chukwa의 웹 모니터링 기능인 HICC는 설치하지 않고 콘솔 모드에서 로그를 확인하여서 매우 불편함을 느꼈다. 관리적인 측면에서 HICC기능은 꼭 검토해야 할 것으로 추정된다.

- ※ HICC(Hadoop Infrastructure Care Center)- 하둡 서버군의 자원 활용률 및 로그 데이터의 실시간 확인 등의 기능을 가진 웹 관리 툴, 필수 SW로 Mysql이 설치되고 별도의 설정을 해주어야 사용이 가능함.

※ 참고 자료

- [1] <http://incubator.apache.org/chukwa/>
- [2] 클라우드 컴퓨팅 구현기술 - 에이콘 출판사
- [3] <http://www.bicdata.com/>
- [4] HBase 클러스터 구축과 관리 - 에이콘 출판사