

---

# 빅데이터 분석 분야 Stack 통합 Test 결과보고서 [ iReport Designer ]

---

2013. 06.

# 목 차

I. Stack 통합 테스트 개요 .....	1
1. 목적 .....	1
II. 빅데이터 .....	2
1. 빅데이터 개요 .....	2
2. 빅데이터 출현 배경 .....	5
3. 빅데이터 중요성 .....	7
4. 빅데이터 3대 핵심 요소 .....	12
III. 빅데이터 구조 .....	15
1. 빅데이터 분석 .....	15
2. 빅데이터 분석 분야 주요 공개SW .....	16
IV.테스트 대상 소개 .....	17
1. Jaspersoft의 iReport Deisgner 소개 .....	17
V. Stack 통합 테스트 .....	19
1. 테스트 환경 .....	19
2. 주요 테스트 방법 .....	20
3. 기능 테스트 수행 결과 .....	21
4. 성능 테스트 수행 결과 .....	22
VI. 종합 .....	28
※ 참고자료 .....	29
[별첨1] 공개SW Jaspersoft 선정지표 테스트 결과	
[별첨2] Jaspersoft 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

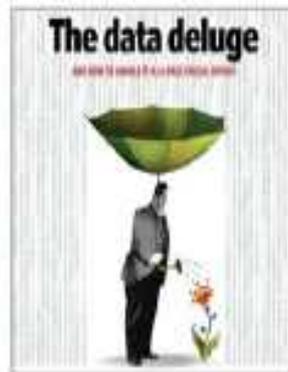
---

---

## II. 빅데이터

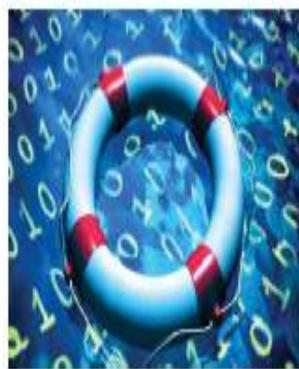
### 1. 빅 데이터 개요

- 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, 공공행정 등 5대 분야에서 6천억불 이상의 가치 창출 예상(출처: McKinsey, 2011.05)

---

---

## 2. 빅 데이터 출현 배경

□ 기업의 고객 데이터 트래킹/분석행위 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

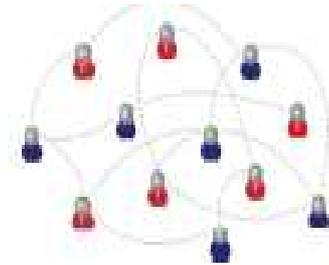
---

---

## 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

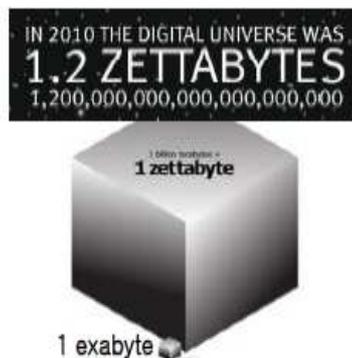
### □ M2M 확산에 따른 센서 저변 확대



- 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, 10<sup>21</sup>) 시대에 살고 있으며 Big Data 추세는 스마트 단말, M2M 센서 확대보급 등으로 더욱 가속화될 전망

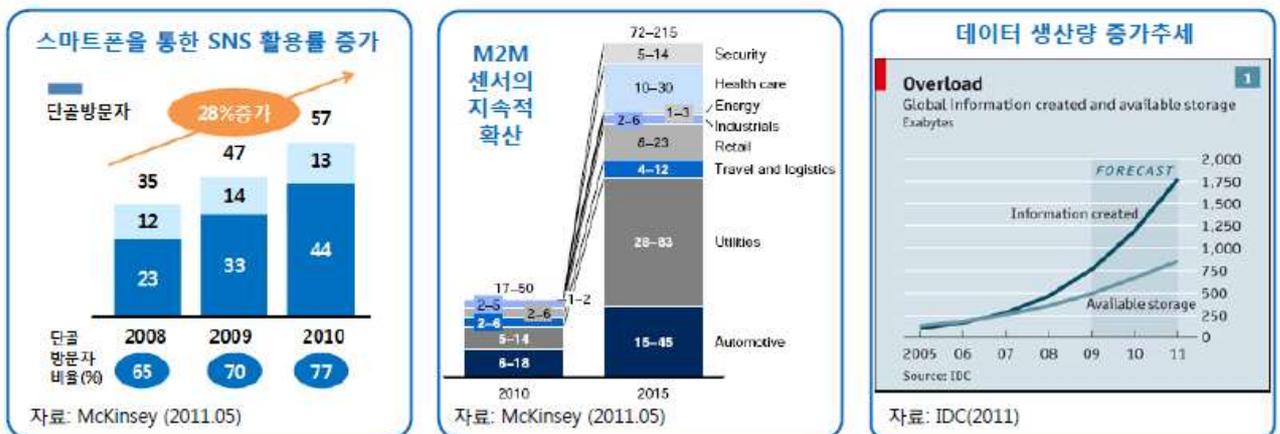


- 세계는 2010년 zettabyte 시대에 돌입(1.2 zettabyte의 정보 생산)
  - 유사 이래 2003년까지 생산된 모든 정보의 합 = 5 exabyte
  - 1 Zettabyte는 美의회도서관 저장정보(235 terrabyte, 11/4 기준)

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



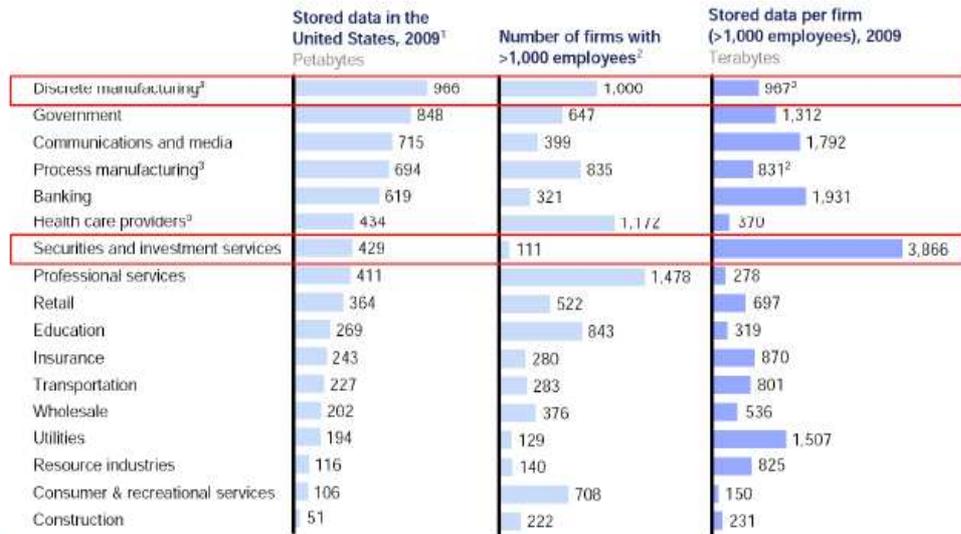
□ 모든 기업이 보유한 Big Data가 '거대한 가치 추출이 가능할 만큼' 충분한 규모에 도달해 누가 먼저 그 가치를 추출해 내느냐가 향후 기업의 성패를 가늠할 상황에 직면

○ 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 terabyte 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



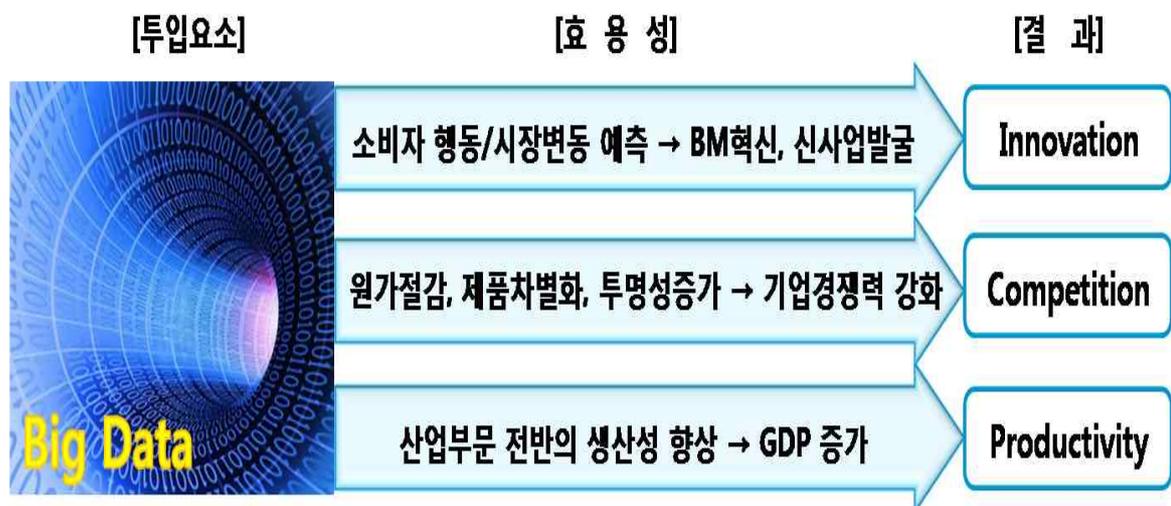
- IBM은 캐나다 의회의 '수백만 건'의 문서를 활용, 영어-불어 자동번역 시스템 개발을 시도했으나 실패
- 반면 Google은 동일방식이지만 '수억 건'의 자료를 활용, 50개 언어 간의 자동번역 시스템 개발 성공
- ※ Google Big Data 구축 방법: 20여개의 언어로 번역된 EC의 문건을 검색을 활용, copy 한 후에 Book스캐닝 프로젝트에서 수천만 권의 전문 번역 DB 구축함
- ※ Google은 Big Data방식을 Spell-check와 음성인식 분야에도 적용하고 있는데, Microsoft가 장기간 대규모 투자로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

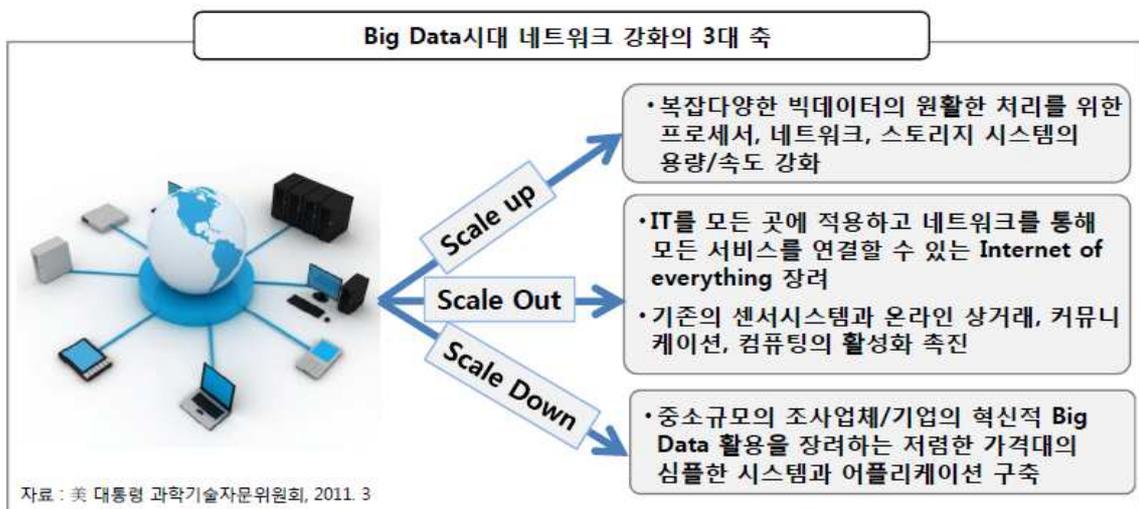
### □ 실시간 분석

- 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

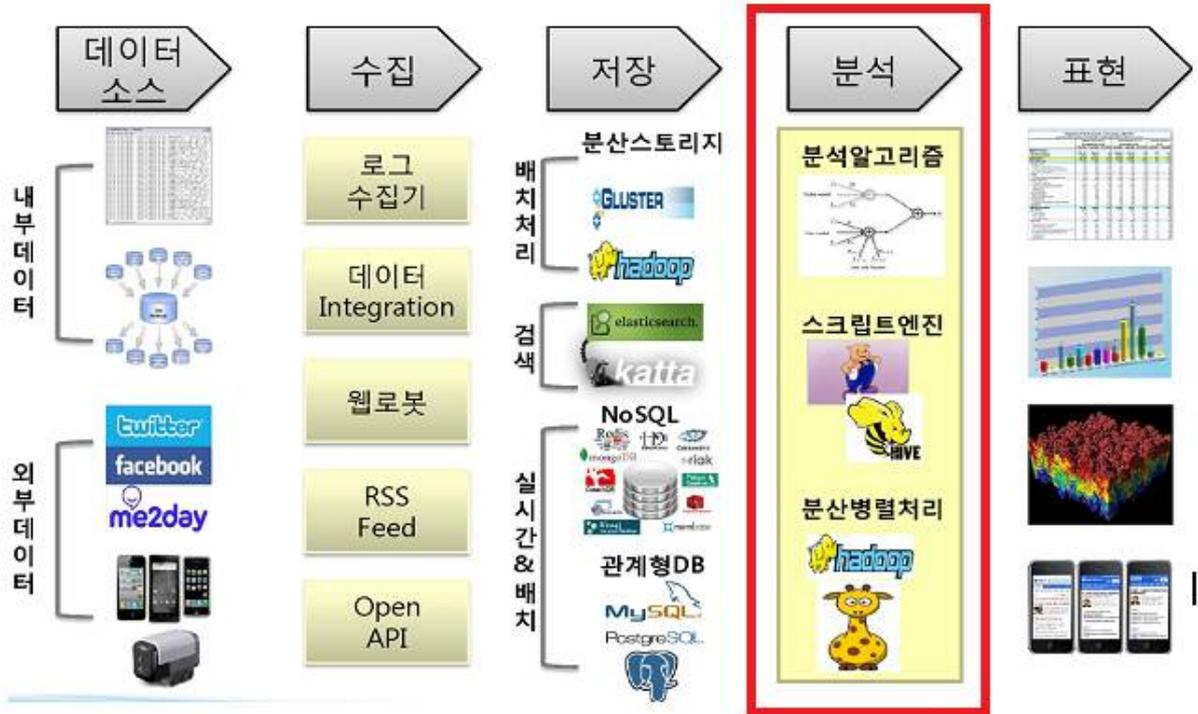
- 폭증하는 실시간 Big Data의 효율적 처리를 위해서는 네트워크 역량 강화도 병행해야 함
  - Big Data는 단순 양의 증가뿐만이 아니라 데이터의 복잡성/다양성의 증가를 의미하기 때문에, 폭증하는 다양한 데이터 처리를 위한 네트워크 강화가 핵심요소로 부상
  - 3가지 측면의 네트워크 강화(scaling of network)를 통해 Big Data를 통한 가치창출의 기반 조성 가능
  - 기존 유무선 네트워크 및 주파수 인프라 관리 또한 복잡/다단한 Big Data시대에 부합하도록 대응 필요



### III. 빅데이터 구조

#### 1. 빅 데이터 분석

빅데이터 처리흐름은 아래의 그림과 같이 5가지 로 표현할 수 있다.



[빅 데이터 처리 흐름]

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

빅데이터의 분석은 데이터의 종류와 목적에 따라 분석 기법이 다르다. 빅데이터의 분석 기법은 텍스트 마이닝, 평판 분석, 소셜네트워크 분석, 클러스터 분석, 통계 분석이 있는데 본 테스트는 비정형 데이터 처리를 위하여 텍스트 마이닝 기법으로 접근한다.

## 2. 빅 데이터 분석 분야 주요 공개SW

빅 데이터 분석 분야 주요 공개SW 중 Pentaho, JasperSoft, Talend를 테스트 SW로 선정하고 공개SW 선정지표를 통하여 가장 점수가 높은 JasperSoft로 테스트를 진행 하였다.

[표 III-1. 분석 분야 주요 공개SW]

제품명	Stack 환경	홈페이지	비고
Pentaho	Linux/Window/Mac	<a href="http://community.pentaho.com/">http://community.pentaho.com/</a>	
JasperSoft	Linux/Window/Mac	<a href="http://community.jaspersoft.com/">http://community.jaspersoft.com/</a>	
Talend	Linux/Window/Mac	<a href="http://www.Talend.com/">http://www.Talend.com/</a>	

[표 III-2. 분석 분야 주요 공개SW 선정지표 점수]

분야	세부분야	대상	항목[배점]				총점 [100]
			Document [25]	Support [25]	Product [30]	Community [20]	
빅데이터	분석	Pentaho	22.0	17.5	25.0	5.0	69.5
		JasperSoft	24.6	17.5	26.3	6.7	75.0
		Talend	23	15.8	22.5	8.3	69.6

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

※ [별첨1]공개SW JasperSoft 선정지표 테스트 결과

---

---

## IV. 테스트 대상 소개

### 1. Jaspersoft의 iReport Designer 소개

iReport Designer는 Jaspersoft Corporation에서 만든 오픈소스 리포트 툴로 데이터 분석 및 보고서 작성까지 완료 해주는 공개SW다.

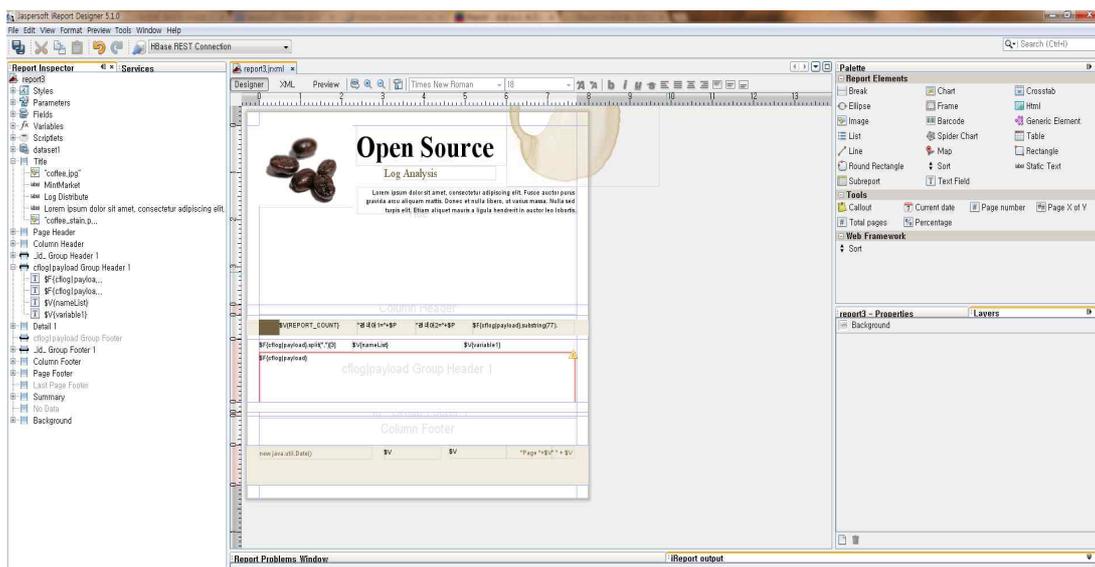
사용되는 내부 언어는 JAVA이며 Netbean기반의 보고서를 디자인 할 수 있다.

UI는 공개SW 개발 툴인 Eclipse 형태를 하고 있으며 각 종 Component를 제공한다.

빅데이터에 대한 커넥터를 별도로 제공하고 있으며 Hadoop, Hive, Hbase, Avro를 플러그인 형태로 지원한다.

보고서는 PDF, XML, CSV, HTML, 오픈오피스 등의 형식을 지원한다.

[그림 IV-1. iReport Designer 구조]



---

---

## □ iReport Inspector구성

### o Parameters

데이터 분석에 사용되는 변수 관리

`$P{변수명}`으로 접근 할 수 있음 (ex `$P{Input_id}.toString()`)

### o Fields

연결된 데이터베이스의 필드, 커스텀 필드를 생성 할 수 있음

`$F{필드명}`으로 접근 할 수 있음

### o Variables

페이징, 카운트 등의 함수 관리 ,커스텀 함수 추가 가능

`$V{함수명}`으로 접근 할 수 있음

### o Style

디자인 스타일 관리

## □ Design

보고서 영역 디자인 관리

## □ XML

보고서 영역의 실제 코드 부분

## □ Preview

분석 결과 미리 보기 영역

※ 추가적인 자세한 정보는 아래의 링크 정보 참조

- <http://community.jaspersoft.com/documentation?version=7114>

- <http://community.jaspersoft.com/system/files/documentation/ireport-ultimate-guide.pdf>

---

---

## V. Stack 통합 테스트

### 1. 테스트 환경

테스트SW 버전

[표 V-1. 테스트 SW ]

SW	Version
HBase	0.94.5
iReport Designer	5.1.0

Stack IP 환경

[표 V-2. Stack 환경]

Stack	OS	IP	OS	IP
A	Ubuntu 12.04	121.162.249.88	CentOS 6.2	121.162.249.18
B	Window	121.162.249.88	CentOS 6.2	121.162.249.18

서버 HW 환경

[표 V-3.서버 HW 환경]

제조사	모델명	CPU	MEM	Disk	NIC
IBM	X3550M2	Intel Xeon(R)CPU 2.40GHz	8GB	320GB	Gigabit 1Port

※ 클라이언트 PC는 동일 환경에서 멀티 부팅 사용하여  
A,B Stack 구성(Ubuntu,Window)

※ IBM 동일 사양의 HW 3대로 Hadoop Stack 구성(CentOS)

---

---

## 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 iReport Designer를 사용하는 사용자들이 사용할 수 있는 항목 중 Parameters, Fields, Variables 에 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 「iReport Designer 테스트 케이스」 문서를 참고한다.

#### □ 테스트 시나리오 현황

[표 V-4. 테스트 시나리오 현황]

기능	테스트 시나리오	테스트 케이스
HBase Connector	1	4
Report Inspector	3	9
시작/종료	1	3
모니터링	1	2
합 계	6	18

#### □ 테스트 결과

기능 테스트 시나리오를 통한 테스트 수행 결과 HBase Connector, Report Inspector 등 시나리오 상의 모든 기능이 예상 결과와 동일하게 동작함을 확인하였다.

[표 V-5. 테스트 결과]

분류		PASS	FAIL	N/A
기능	개수			
HBase Connector	4	4	0	0
Report Inspector	9	9	0	0
시작/종료	3	3	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.

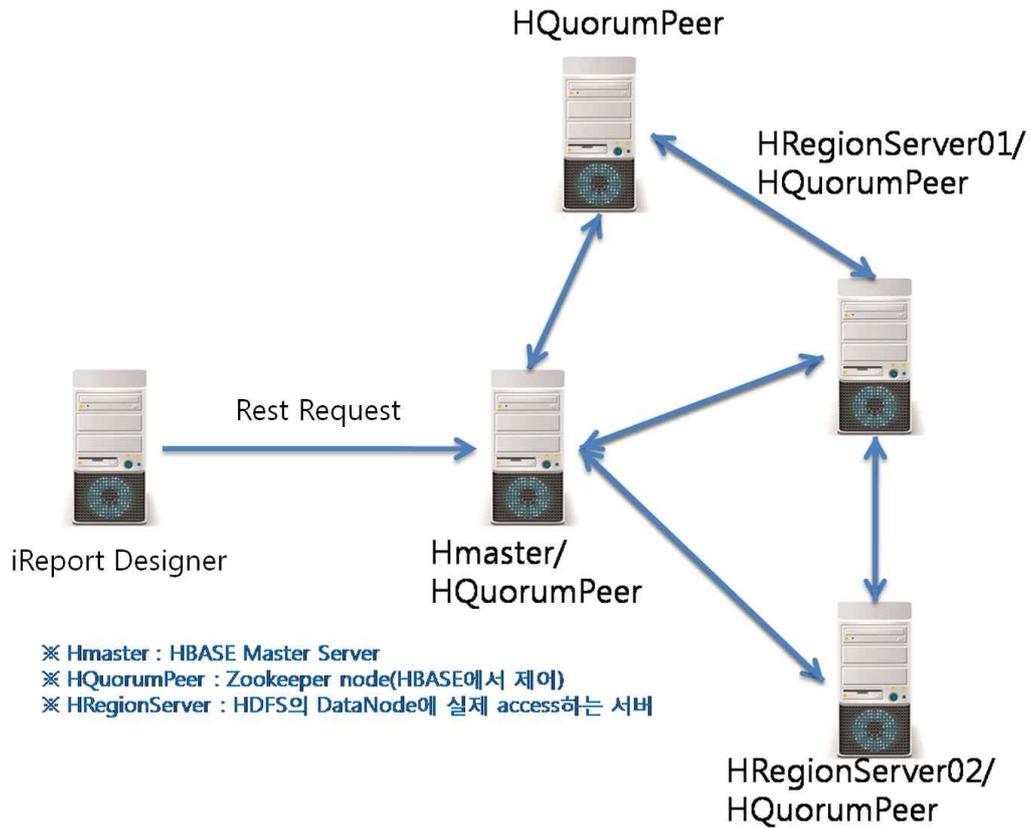
본 성능 테스트는 iReport Designer를 사용하여,100MB의 비정형 데이터(웹로그)를 분석하는 시나리오를 재현하고 성능을 측정한다.

### □ 테스트 시나리오

[표 V-6. 테스트 시나리오]

대분류	시나리오ID	시나리오
F_Result	F_Result_Seoul	웹 로그 중 서울에 관련된 정보를 검색한다.
	F_Result_Inchoen	웹 로그 중 인천에 관련된 정보를 검색한다.
	F_Result_Busan	웹 로그 중 부산에 관련된 정보를 검색한다.
	F_Result_Kim	웹 로그 중 성이 김씨인 사람의 정보를 검색한다.
	F_Result_Lee	웹 로그 중 성이 이씨인 사람의 정보를 검색한다.
	F_Result_Park	웹 로그 중 성이 박씨인 사람의 정보를 검색한다.
	F_Result_login	웹 로그 중 로그인에 관련된 정보를 검색한다.
	F_Result_Exception	웹 로그 중 예외처리에 관련된 정보를 검색한다.
	F_Result_Sql	웹 로그 중 DB에 관련된 정보를 검색한다.

□ 서버 구성



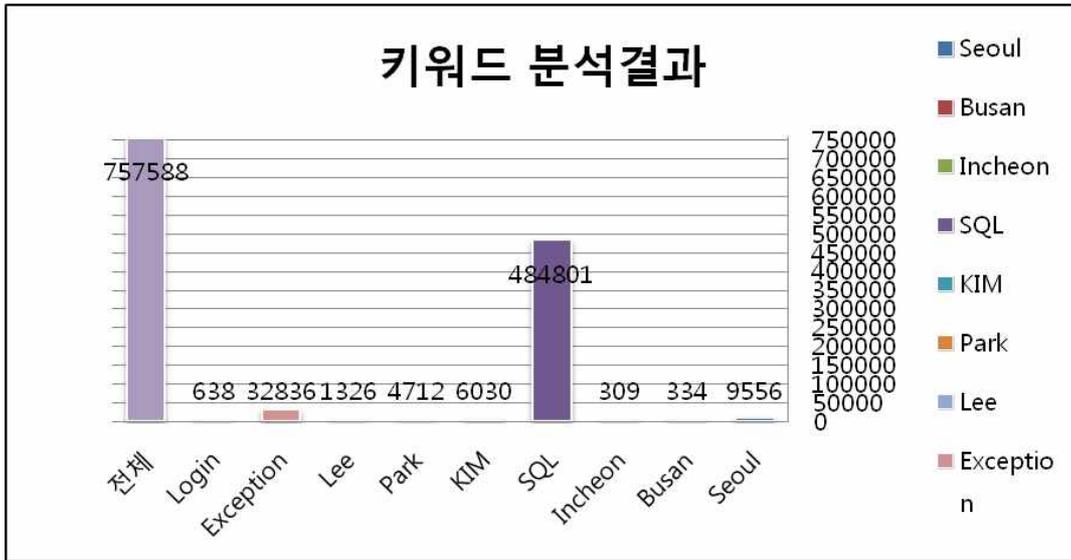
[그림 V-1. 성능 테스트 환경 정보]

□ 측정항목

[표 V-7. 측정항목]

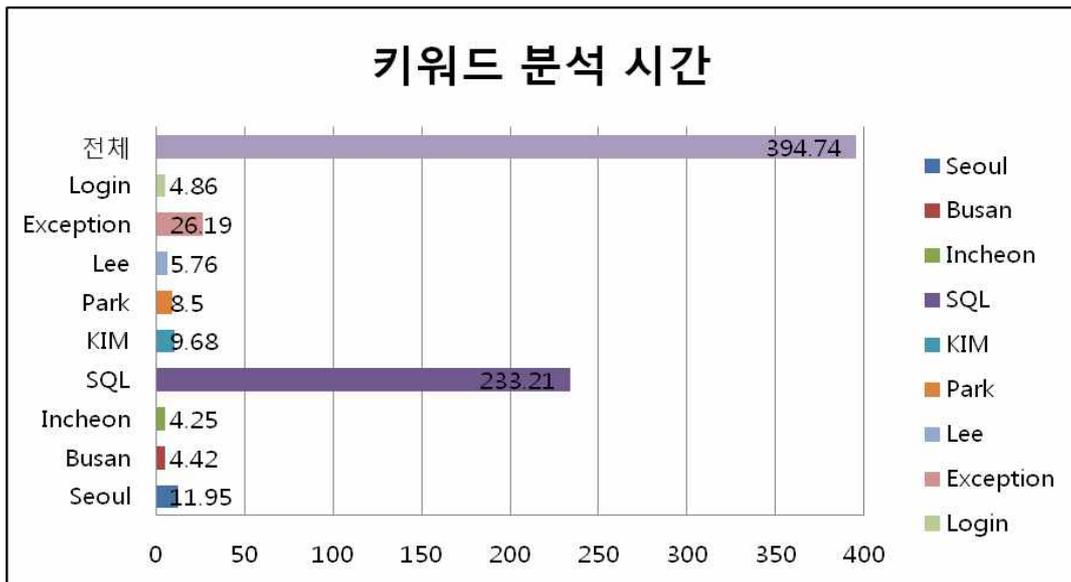
항목	내용
라인	로그파일에서 키워드로 검색된 데이터의 라인 수
초	키워드로 검색 시 데이터 분석 시간



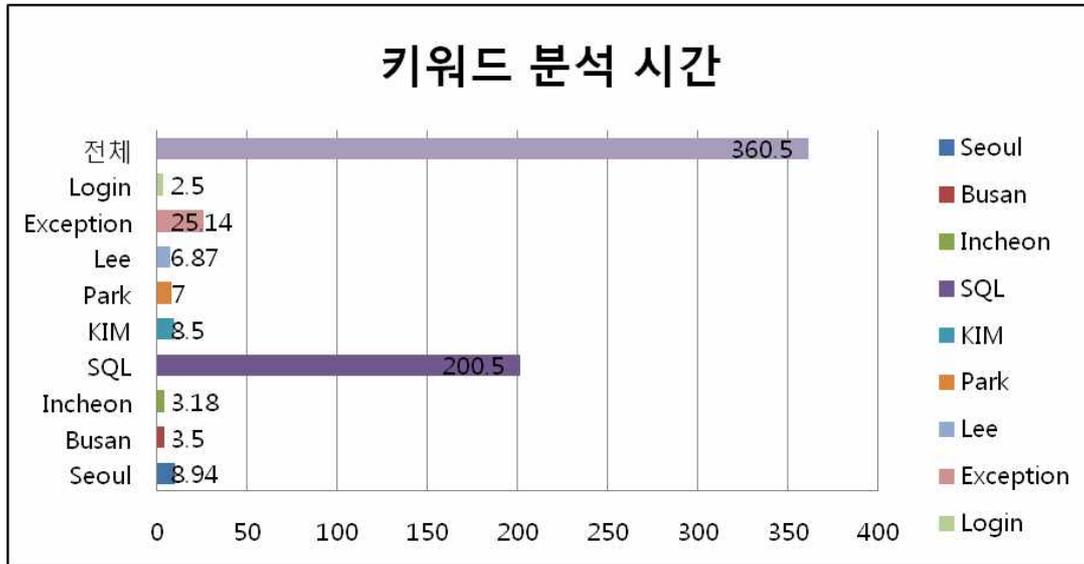


[그림 V-3. 키워드에 해당하는 라인 수 분석 결과]

HBase에 저장된 비정형 데이터에서 키워드 분석을 통해 이용자의 지역과 이용자의 비율을 알아낼 수 있었다.



[그림 V-4. Window용을 사용한 키워드 분석 시간]



[그림 V-5. Ubuntu용을 사용한 키워드 분석 시간]

우분투에서 사용한 iReport Designer가 전체적으로 Window용 iReport Designer보다 분석에 사용된 시간이 짧음을 알 수 있다.

---

---

## VI. 종합

- iReport Designer 기능 테스트 수행 결과 공개SW로 구성된 Stack 상에서 각 기능 시나리오 수행 시 치명적 오류 또는 심각한 장애가 발생하지 않았으며, Stack을 구성하는 각 공개SW(Hadoop, Zookeeper, HBase)가 유기적으로 동작함을 확인하였다.
  
- iReport Designer 성능 테스트 수행 결과 100MB, 총 76만 라인의 비정형 웹 로그에서 각각의 키워드 별로 분석, 가장 많이 사용된 것은 DB(SQL) 검색이었고 두 번째는 예외 처리(Exception) 였다. 또한 Seoul 지역에서 접속한 사람이 가장 많았으며 성이 Kim인 사람이 가장 많이 사용하였다.

---

---

## ※ 참고 자료

- [1] <http://community.jaspersoft.com/>
- [2] <http://hbase.apache.org/>
- [3] ireport designer ultimate guide.pdf
- [4] Hadoop 완벽가이드 - 한빛 출판사
- [5] 클라우드 컴퓨팅 구현 기술 - 에이콘
- [6] IDG\_TechReport\_BigData-20120426.pdf