

Data Science with Apache Zeppelin

DEVIEW
2015

이문수
NFLabs



contents

DEVIEW
2015

1. Data science lifecycle
2. Apache Zeppelin
3. Zeppelin in your team
4. Helium

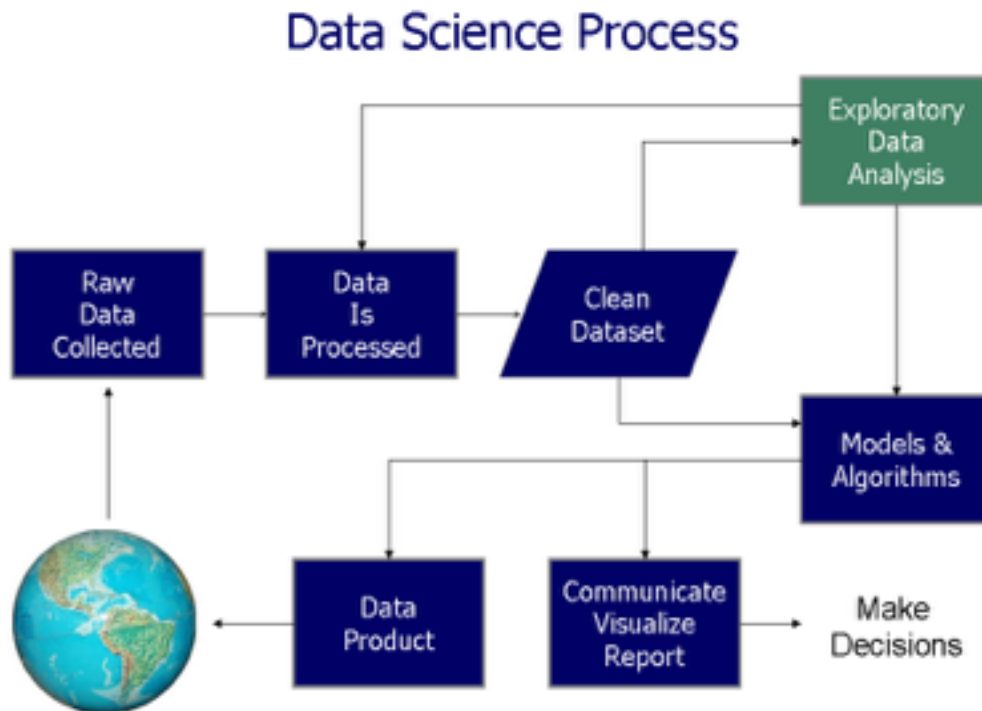
1.

Data Science Lifecycle

Data Science

DEVIEW
2015

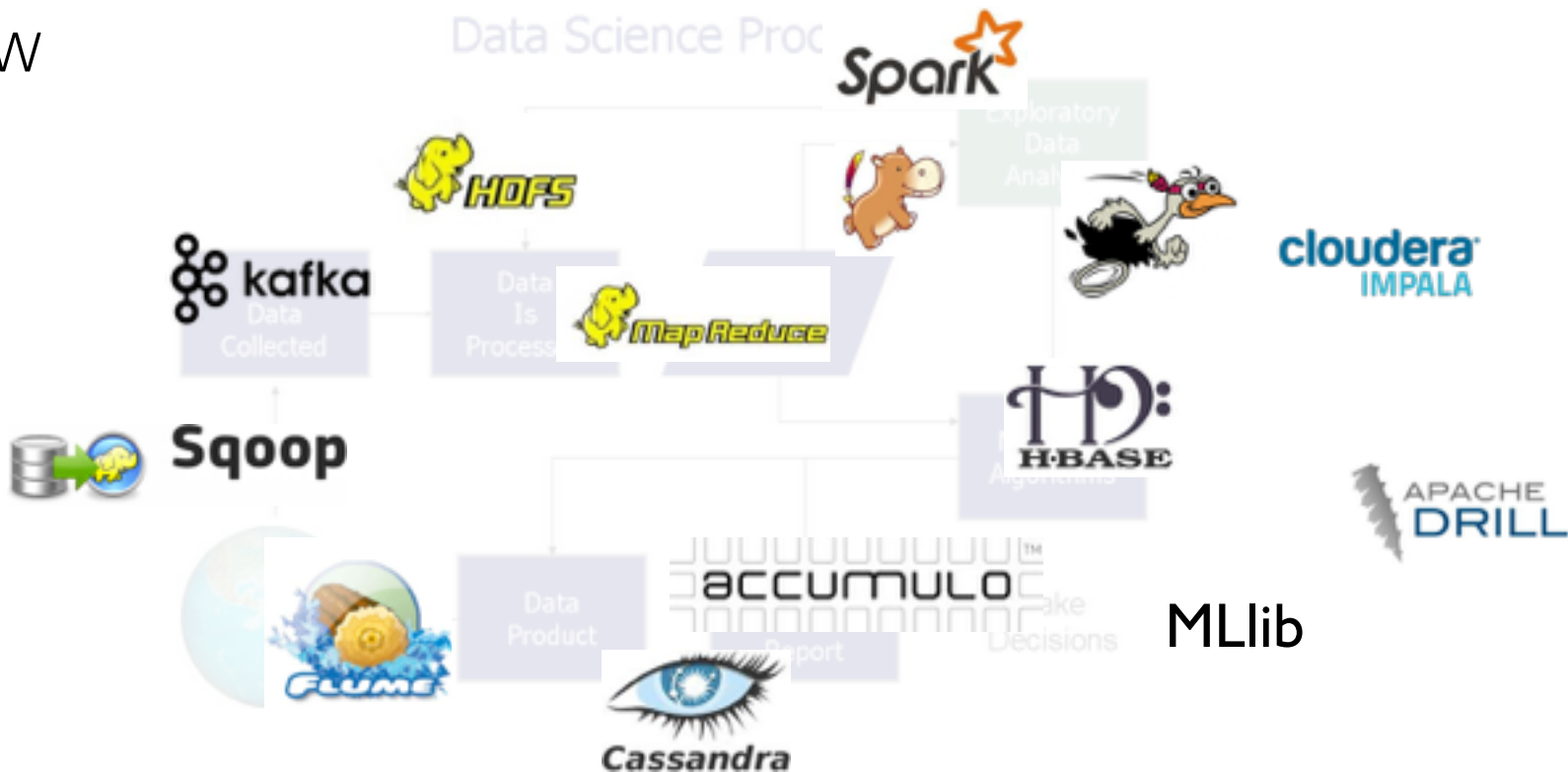
다양한 과정



Data Science

DEVIEW
2015

다양한 SW



Data Science

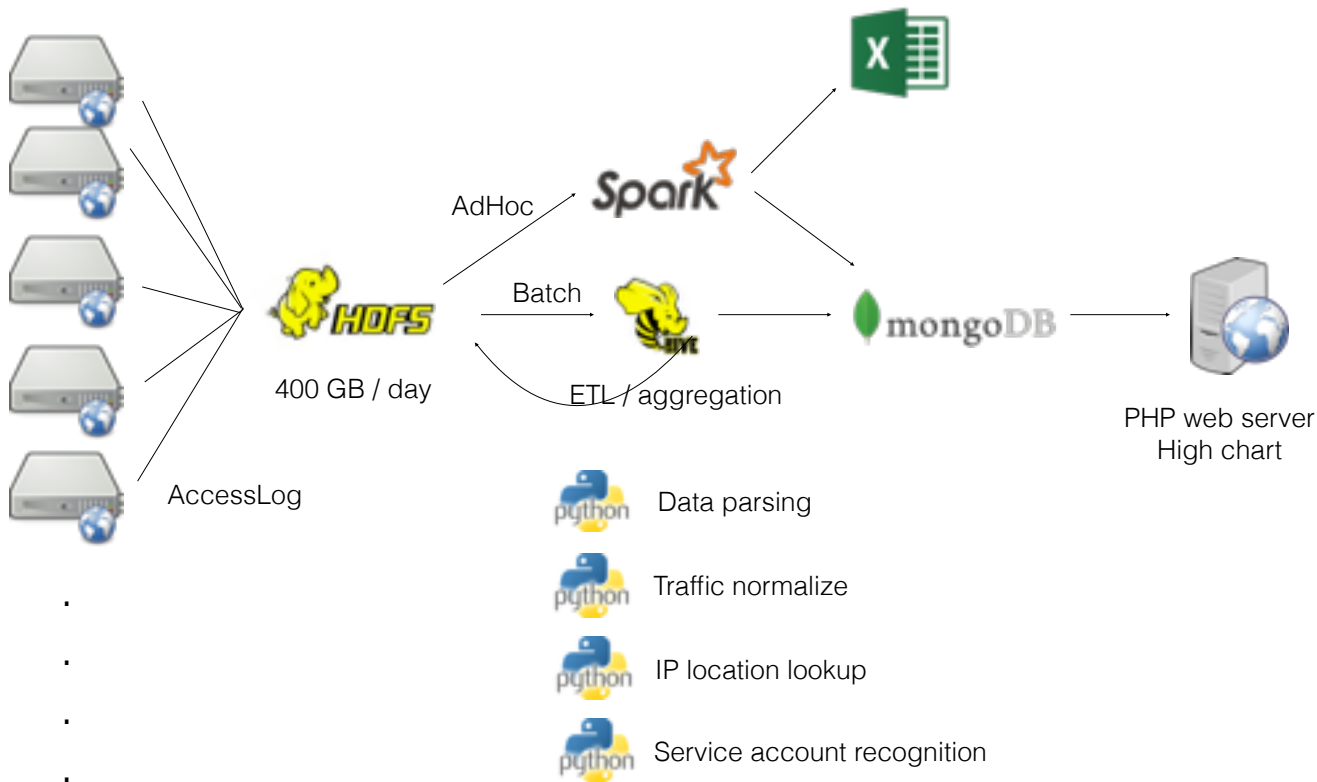
DEVIEW
2015

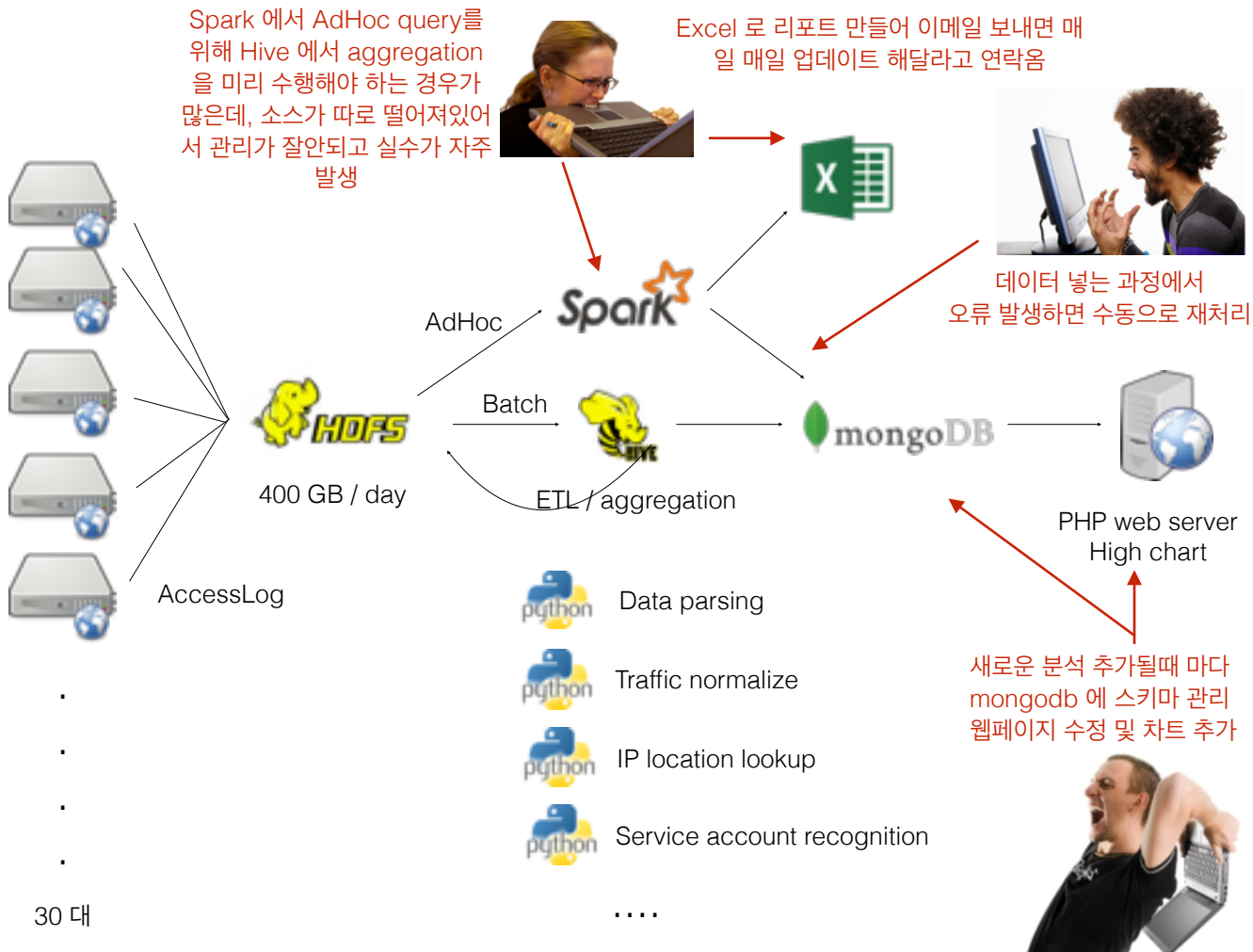
다양한 사람



주







분석에 집중할 수 있을까?

DEVIEW
2015



편리하고 강력한 분석 언어
인터랙티브 속도
라이브러리
시각화
공유/협업
간편하고 손쉬운 시스템 구성

2. Apache Zeppelin



Apache Zeppelin

DEVIEW
2015



<http://zeppelin.incubator.apache.org/>

2014 년 12월에 ASF 에 incubation 됨

2013 년 8월에 NFLabs 내부 프로젝트로 시작

63 Contributors from worldwide

646 Stars

1 release

Apache 2.0 License

누가 쓰나?

DEVIEW
2015



NAVER



NETFLIX

IBM

Microsoft®



eBay



@WalmartLabs



trulia



BOSCH



Collaboration, Sharing



ZeppelinHub

Zeppelin + Full stack on a cloud

Packaging & Deployment



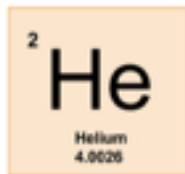
Ambari



Z-Manager



Packages



Backend integration

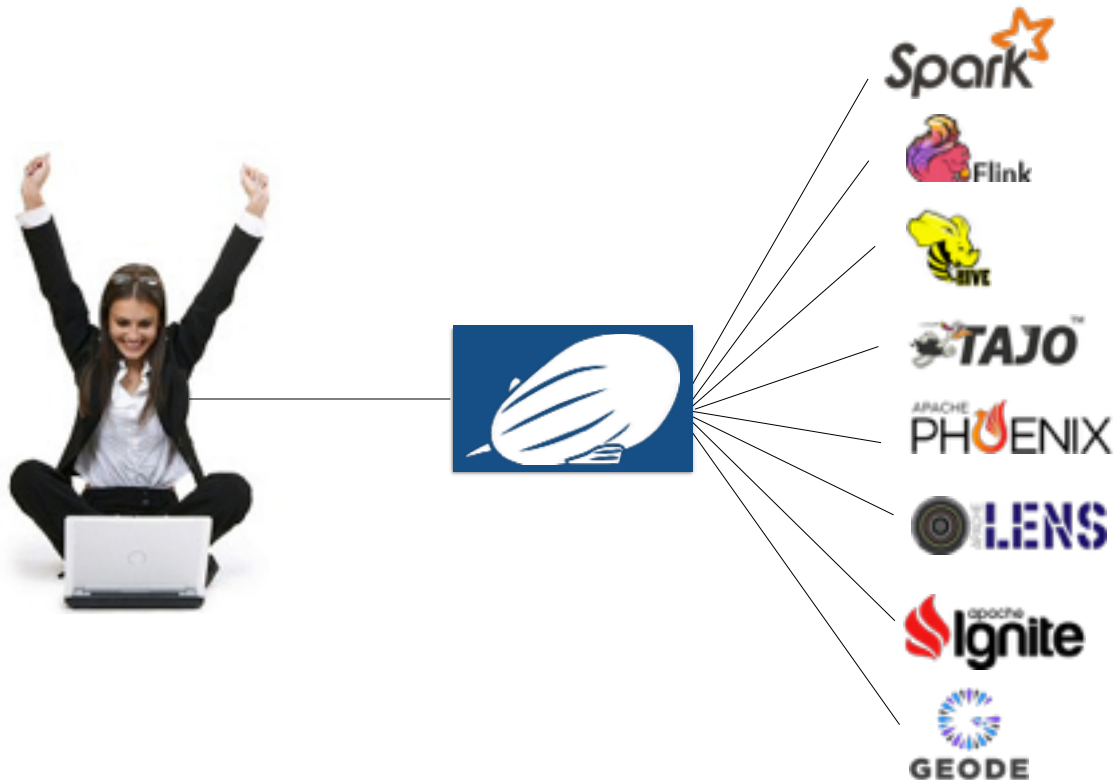


...



다양한 Backend 동시에 사용

DEVVIEW
2015



다양한 Backend 동시에 사용

Shell 명령 이용해 데이터 카피



Hive 이용해 데이터 transformation



Spark 의 MLlib 으로 분석



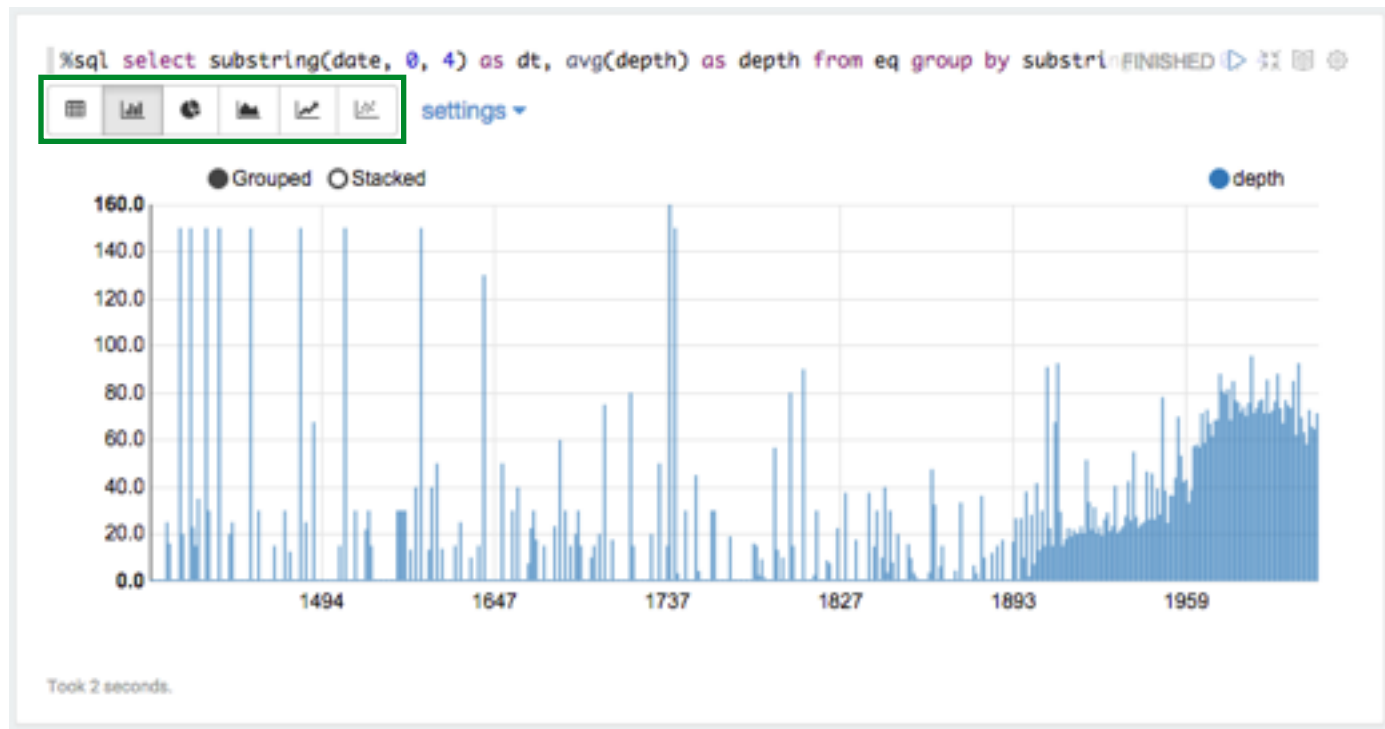
python 이용하여 시각화

하나의 노트북에서 순차적으로 작업을 처리



Visualization

DEVVIEW
2015



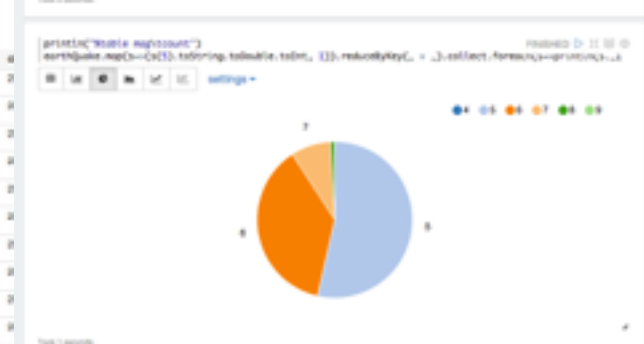
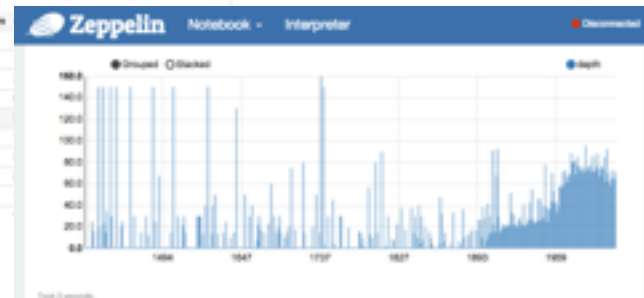
Apache Zeppelin

DEVIEW
2015

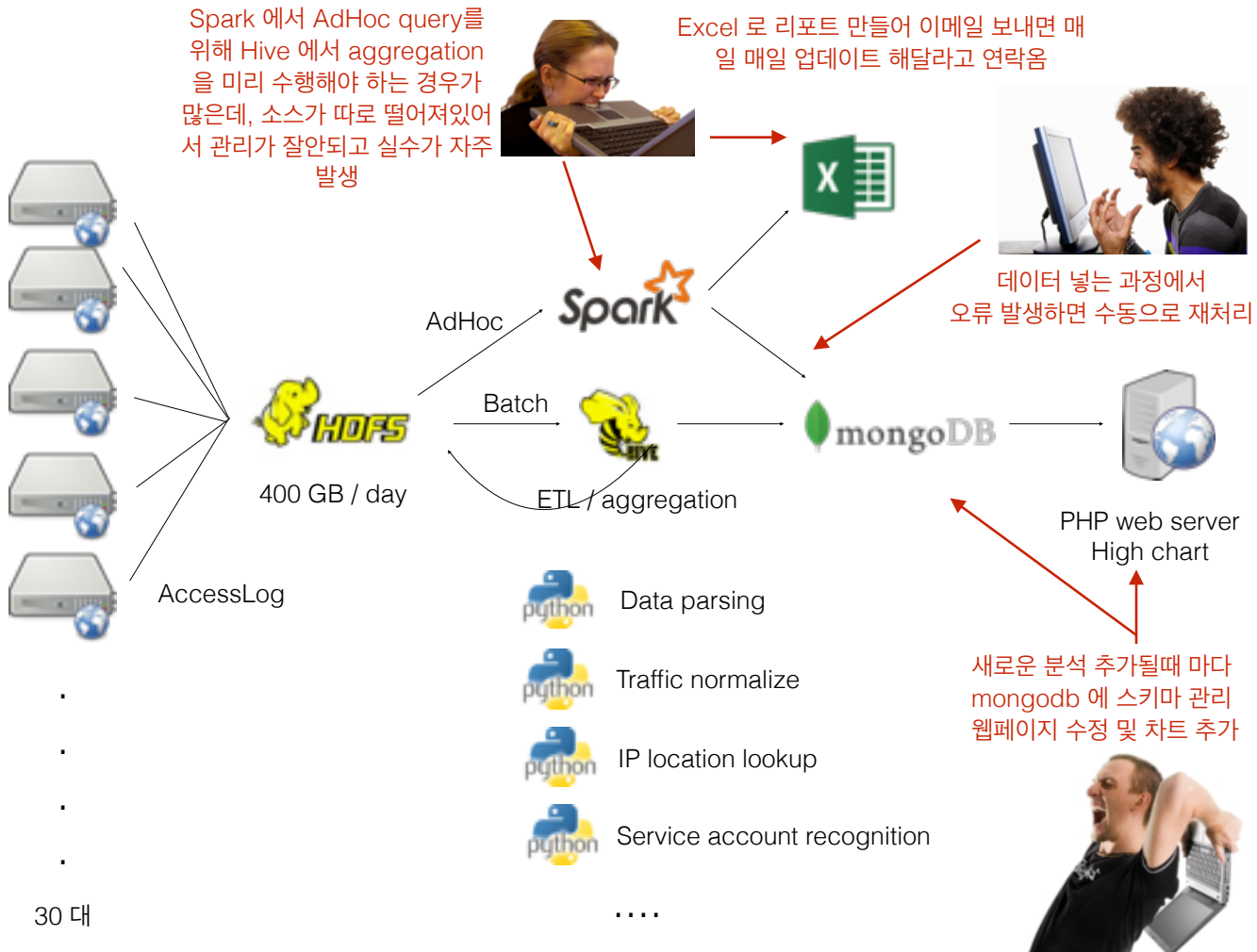
Interactive Notebook



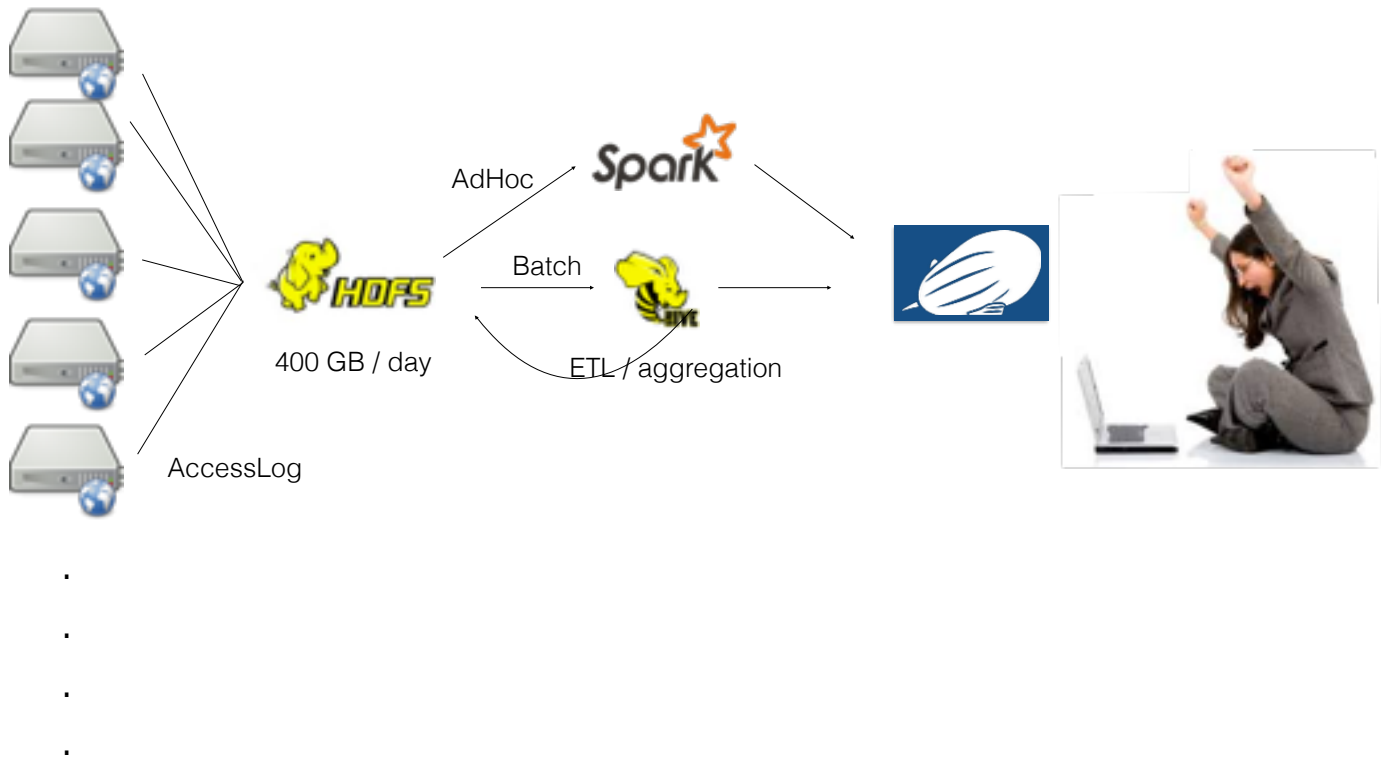
number	state	login	avatar	created_at	id
20	closed	jongjoo		2015-04-01T01:50:36Z	2
21	closed	jongjoo		2015-04-01T01:52:04Z	3
22	closed	PamVandenberg		2015-04-01T08:09:57Z	4
23	closed	jongjoo		2015-04-01T07:50:56Z	5
24	closed	terryky		2015-04-02T01:15:05Z	6
25	closed	jongjoo		2015-04-02T13:58:30Z	7
26	closed	Leemooosoo		2015-04-02T09:02:04Z	8
27	closed	Leemooosoo		2015-04-02T13:48:45Z	9
28	closed	lee		2015-04-07T07:39:56Z	10
29	closed	comaeedong		2015-04-07T08:04:11Z	11
30	open	lee		2015-04-07T09:00:06Z	12



Before



After



3. Zeppelin in your team



Team = Multi user

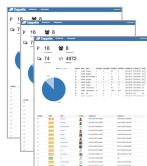
Multi user = ACL

다중 사용자 구성 A

각각 노트북에서 interpreter
setting 선택 하여 사용

사용자별로
Interpreter setting 추가
하여 각각 다른 리소스 할당

하나의 Zeppelin 인스턴스 사
용



50 core, 100GB mem



100 core, 300GB mem

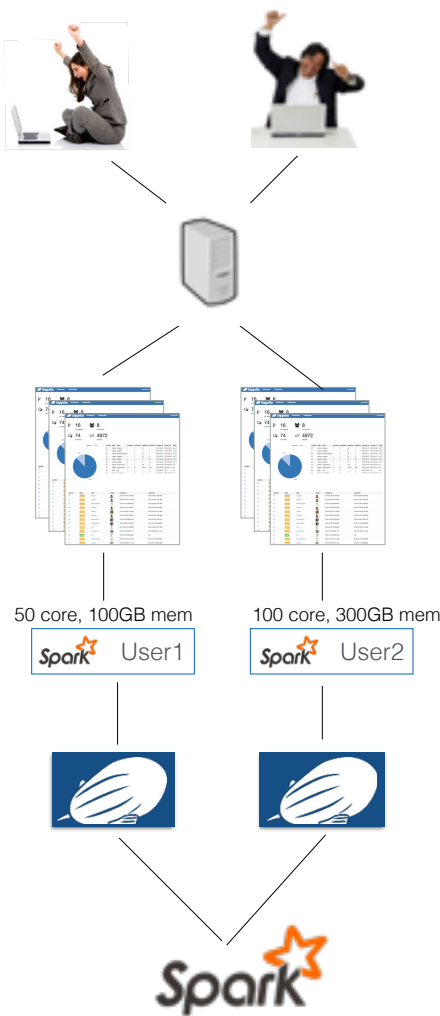


DEVIEW
2015

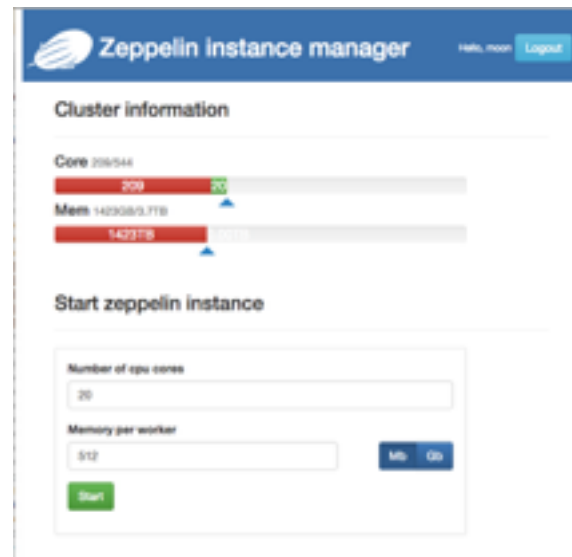
다중 사용자 구성 B

Proxy server 에서 인증 및
Instance 할당 관리

사용자별로 Instance 생성



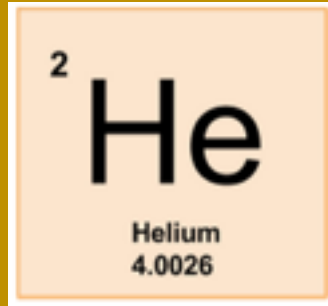
DEVVIEW
2015



<http://nflabs.github.io/z-manager/>

4.

DEVIEW
2015

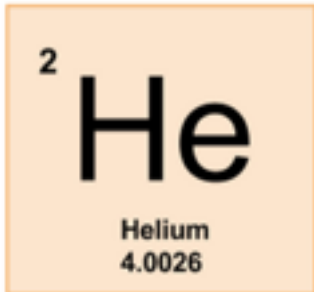


Platform for **Data analytics application**
on top of Apache Zeppelin

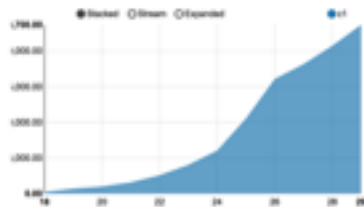
<https://cwiki.apache.org/confluence/display/ZEPPELIN/Helium+proposal>

Helium Application

DEVIEW
2015



=



+

```
/** Configures the OAuth Credentials for accessing Twitter API */
def configureTwitterCredentials(apiKey: String, apiSecret: String, accessToken: String) {
  val config = new OAuthConfig(apiKey, apiSecret, accessToken)
  "apiKey" => apiKey, "apiSecret" => apiSecret, "accessToken" => accessToken
}

println("Configuring Twitter OAuth")
config.foreach { case (key, value) => {
  if (value.trim.isEmpty) {
    throw new Exception("Error setting authentication - value for " + key + " is empty")
  }
  val fullKey = "twitter.oauth." + key.replace("api", "consumer")
  System.setProperty(fullKey, value.trim)
  println("setProperty " + fullKey + " set as [" + value.trim + "]")
}
}
println()
```

View

Algorithm



Zeppelin provided Resources



Example of resource

Data

- Result of last execution
- JDBC connection (from JDBC Interpreter)*

Computing

- SparkContext (from SparkInterpreter)
- Flink environment (from FlinkInterpreter)*

Any java object

- Provided by user created Interpreter
- Provided by user created Helium application

데이터

- ex) get git commit log data

<https://github.com/Leemoonsoo/zeppelin-gitcommitdata>

컴퓨팅

- ex) run cpu usage monitoring code across spark cluster, using SparkContext

<https://github.com/Leemoonsoo/zeppelin-sparkmon>

시각화

- ex) display result data as a wordcloud

<https://github.com/Leemoonsoo/zeppelin-wordcloud>

동작 방식

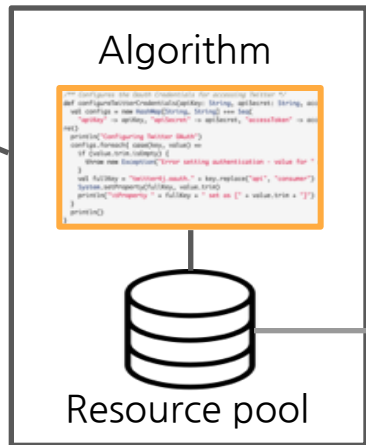
DEVIEW
2015



Web browser



Zeppelin Server



Interpreter Process

Resource
pools are
connected

“Algorithm runs where resource exists”

쉽다

```
class YourApplication extends org.apache.zepplin.helium.Application {  
  
    @Override  
    public void run(ApplicationArgument arg, InterpreterContext context) {  
        .....  
    }  
}
```

하나의 class 만 extend 하면 만들수 있다.

간단하다

```
{  
  mavenArtifact : "groupId:artifactId:version",  
  className : "your.helium.application.Class",  
  icon : "fa fa-cloud",  
  name : "My app name",  
  description : "some description",  
  consume : [  
    "org.apache.spark.SparkContext"  
  ]  
}
```

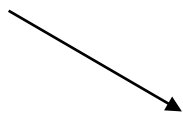
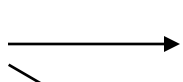
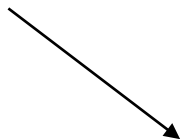
Spec 파일만 정의하면 바로 사용

Deploy

DEVIEW
2015

편리하다

공개



사내

maven 사용해서 배포되고
설치되므로 사내 라이브러
리를 따로 구축할 수 있다.

사용자가 클릭하면 바로 다
운로드되어 실행

maven 통해서 배포하여 편리하게 사용



Demo

Get involved

DEVIEW
2015

홈페이지

<http://zeppelin.incubator.apache.org/>

메일링리스트

users@zeppelin.incubator.apache.org

dev@zeppelin.incubator.apache.org

이슈트래커

<https://issues.apache.org/jira/browse/ZEPPELIN>

소스리파지토리

<https://github.com/apache/incubator-zeppelin>

Q&A

Thank You