

Hadoop

Jakub Podeszwik

infoShare Academy

05.09.2016

- 2 dni
- 20% prezentacja
- 40% live coding
- 40% zadań głównie programistycznych

① Wykład

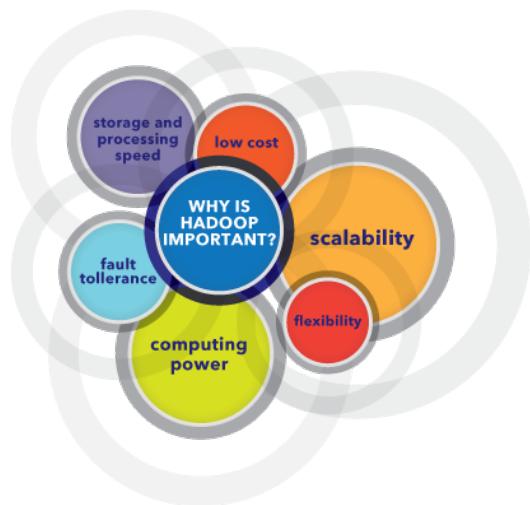
- ① Czym jest hadoop?
- ② Komponenty
- ③ Dystrybucje

② Zajęcia praktyczne

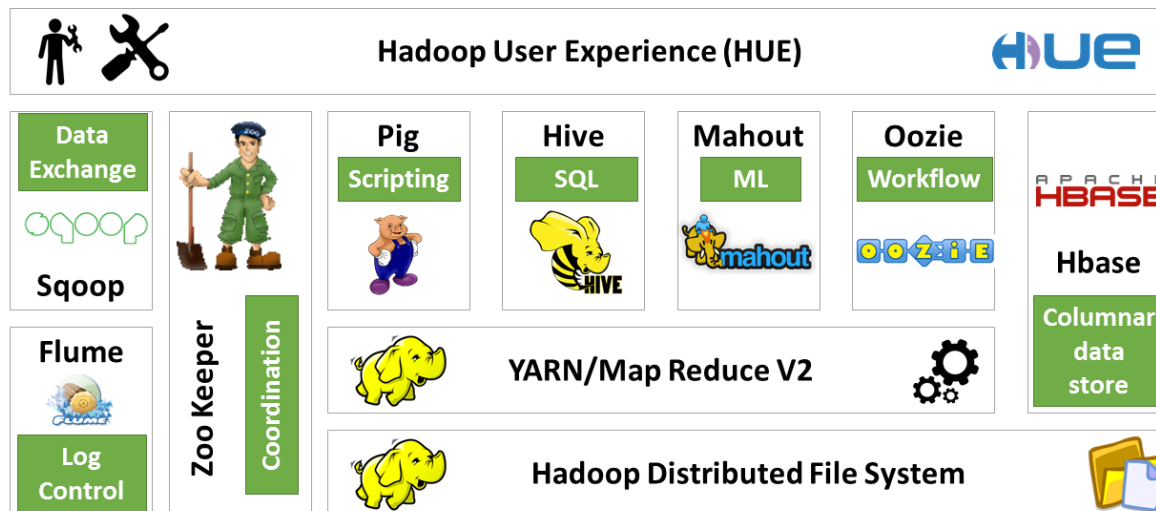
- ① Cloudera manager
- ② HDFS
- ③ Flume
- ④ Sqoop
- ⑤ Mapreduce
- ⑥ Camus*

Wykład

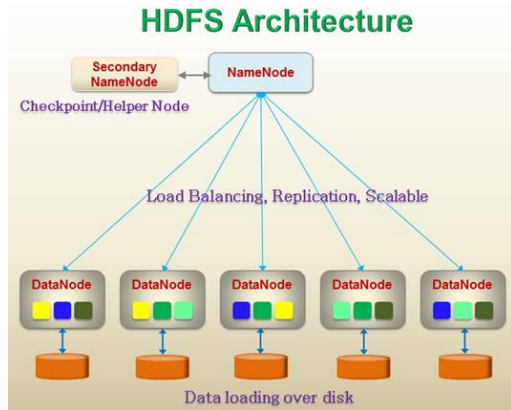
Czym jest hadoop



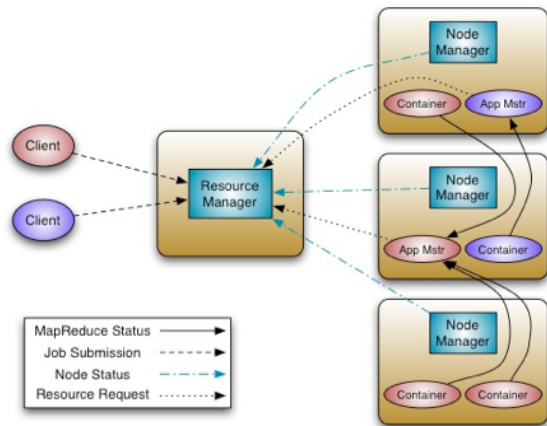
The Apache Hadoop Stack



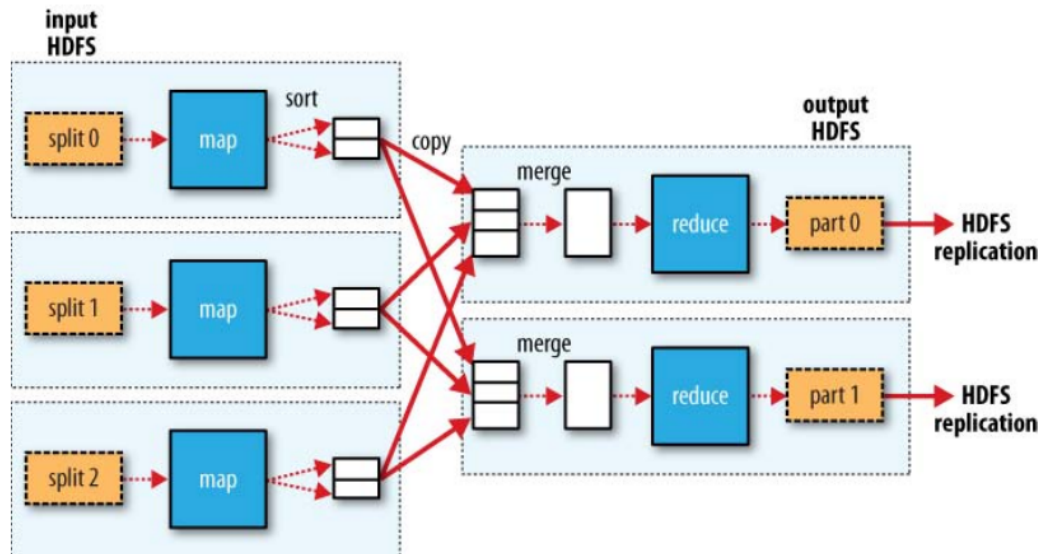
Hdfs



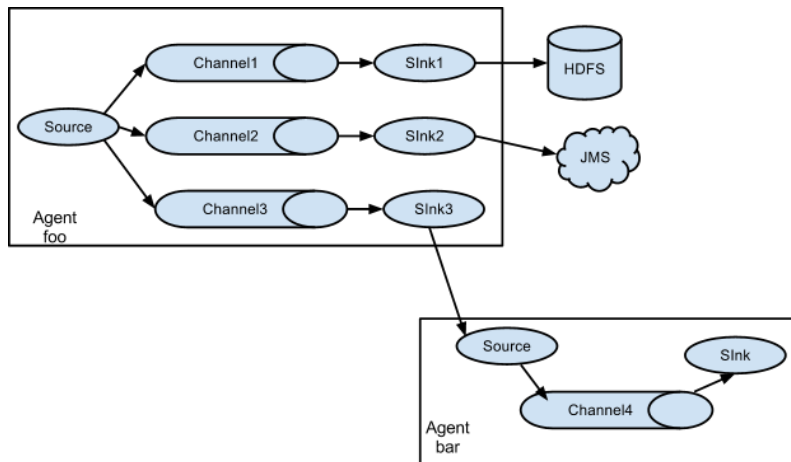
YARN



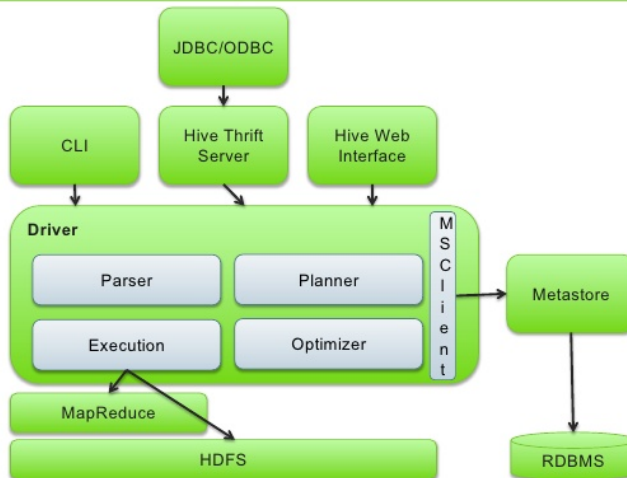
Mapreduce



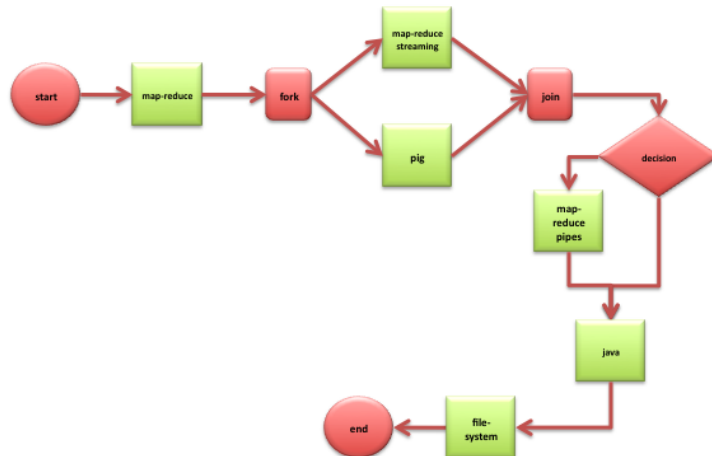
Flume

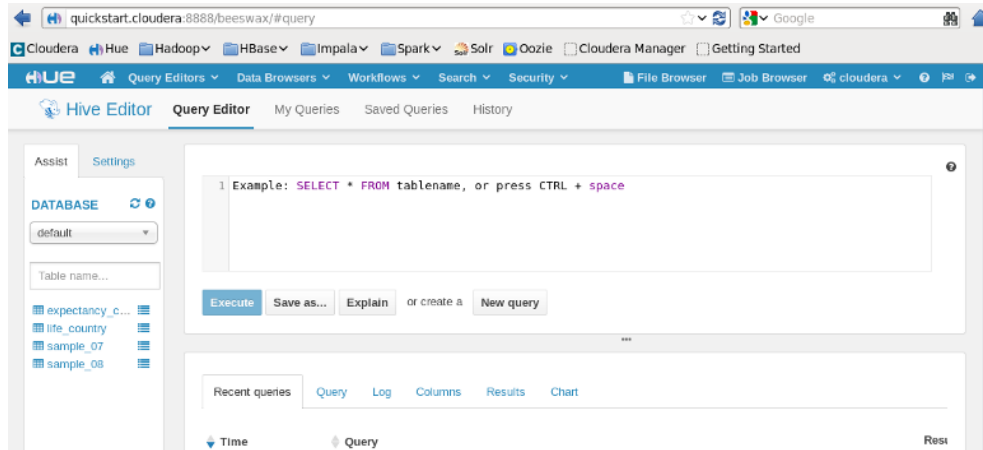


Apache Hive Architecture



Oozie





cloudera

MAPRTM
TECHNOLOGIES



Hortonworks

Zajęcia praktyczne

- <https://github.com/jpodeszwik/Fake-Apache-Log-Generator.git>
- <https://github.com/jpodeszwik/hadoop-workshop-day1.git>

Cloudera manager

- vagrant:vagrant /user/vagrant
- vagrant:vagrant /user/vagrant/inputs,outputs
- flume:flume /user/flume
- kursant:kursant /user/kursant
- wrzucić pliki loremipsum i apache_logs do /user/vagrant/inputs

użyj zmiennej HADOOP_USER_NAME

hdfs - root

- zmień prawa do odczytu do katalogu `/user/kursant` tak, żeby tylko kursant mógł go odczytać
- spróbuj odczytać katalog jako user `'vagrant'`
- spróbuj odczytać katalog jako user `'hdfs'`

- utwórz pusty plik w katalogu `/user/kursant`
- przeczytaj plik `/user/vagrant/apache_logs`

użyj zmiennej `HADOOP_CONF_DIR`

- wylistuj katalog
- utwórz plik, zapisz do niego jakieś dane
- usuń przed chwilą utworzony plik

Flume

Mapreduce

- Zainstaluj na vm-cluster-node1 server mysql (sudo apt-get install mysql-server)
- zaimportuj do niego dane z pliku dump.sql (mysql -u root < dump.sql)
- zaimportuj dane z tej bazy na hdfs

Mapreduce

Camus*

Pytania?