

Administracja

Agenda

1. Topologia klastrów obliczeniowych w Big Data
2. Metody instalacji narzędzi
3. Zarządzanie i Monitoring

Topologia klastrów

- “Active-Standby”
- “Gold-silver”
- Active - Active ?

Gdzie można uruchamiać?

- Własna serwerownia
- Google Cloud
- Amazon AWS
- Microsoft Azure

Zarządzanie infrastrukturą

- Terraform
- AWS CloudFormation

Terraform

```
provider "aws" {  
    region = "eu-west-1"  
}
```

```
module "ec2" {  
    instance_count = 2  
    name           = "example-normal"  
    ami            = "<id_obrazu_do_uruchomienia>"  
    instance_type  = "m4.large"  
    subnet_id      = "<subnet_id>"  
    vpc_security_group_ids = ["<security_group_id>"]  
    associate_public_ip_address = true  
}
```

Zarządzanie konfiguracją

- Puppet
- Chef
- Ansible

Puppet

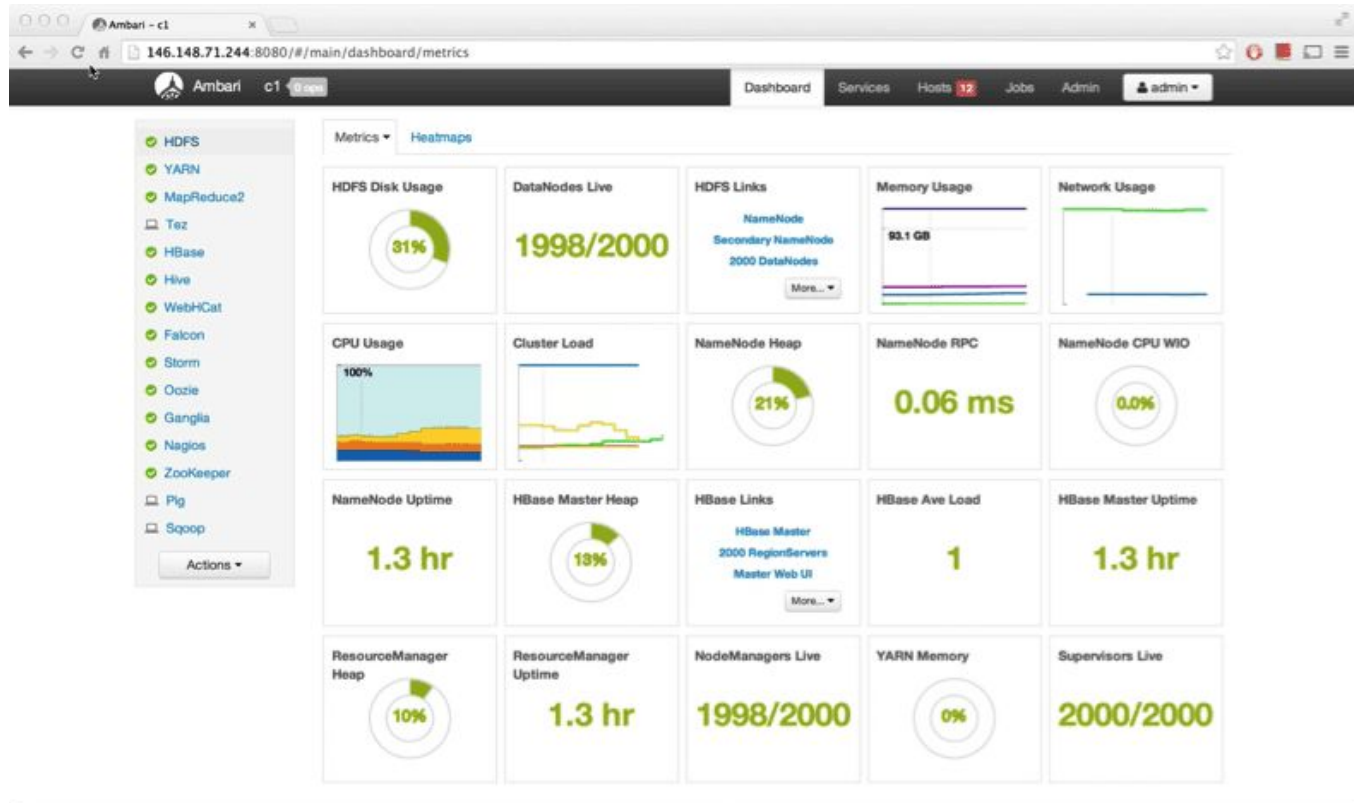
```
package { 'mysql-server':  
  ensure => installed,  
}
```

```
service { 'hadoop-yarn-resourcemanager':  
  ensure => running,  
}
```

```
user { 'janusz':  
  ensure    => present,  
  uid       => '1521',  
  gid       => '1345',  
  shell     => '/bin/bash',  
  home      => '/home/janusz'  
}
```


Zarządzanie klastrami Hadoop

- Apache Ambari
- Cloudera Manager



Home

30 minutes preceding November 3, 2015, 1:44 PM PST

Status

All Health Issues

Configuration 16 ▾

All Recent Commands

Add Cluster

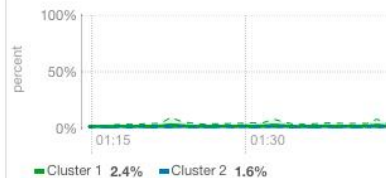
Cluster 1 (CDH 5.5.0, Packages) ▾

	Hosts	4	
	FLUME-1		▾
	HBASE-1		▾
	HDFS-1	1	▾
	HIVE-1		▾
	HUE-1	1	▾
	IMPALA-1		▾
	KAFKA-1		▾
	KS_INDEXER-1		▾
	KUDU-1		▾
	MAPREDUCE-1		▾
	OOZIE-1		▾
	SOLR-1		▾
	SPARK_ON_YARN-1		▾
	SQOOP-1		▾
	SQOOP_CLIENT-1		▾
	YARN-1		▾
	ZOOKEEPER-1	1	▾

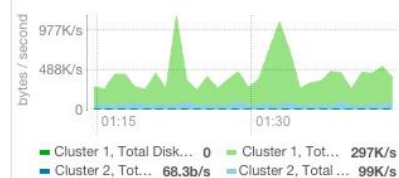
Charts

30m 1h 2h 6h 12h 1d 7d 30d ▾

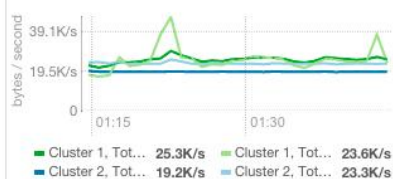
Cluster CPU



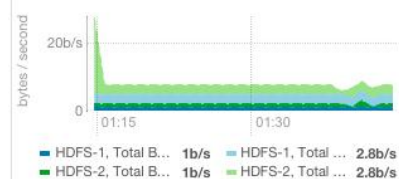
Cluster Disk IO



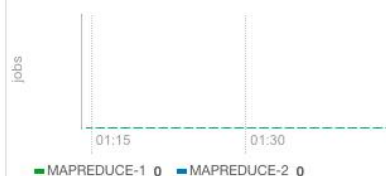
Cluster Network IO



HDFS IO



Running MapReduce Jobs



Completed Impala Queries



Zabezpieczenie danych

- Non secure mode
- Kerberos
- HiveServer

Monitorowanie

- Działanie procesów (czy np. odpowiada na porcie?)
- Statystyki systemu operacyjnego
- Utylizacja zasobów OS
- Monitoring JMX (procesy Javy)
- Testy funkcjonalne

Jakie parametry monitorować?

- Ilość bloków
- Ilość plików
- Utylizacja pamięci / YARN
- Ilość aplikacji w różnych stanach
- Ilość błędnych aplikacji oraz kontenerów
- Użycie pamięci + GC dla procesów HDFS+YARN

Narzędzia do metryk

- “Time Series DB”
- Graphite
- InfluxDB
- OpenTSDB
- Grafana



Big Dashboard



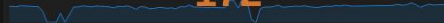
Zoom Out

Oct 4, 2015 11:29:09 to Oct 4, 2015 14:13:23 UTC



Logins

172



Sign ups

263



Sign outs

268

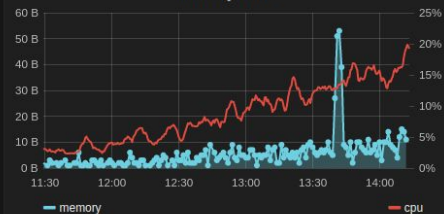


Support calls

80



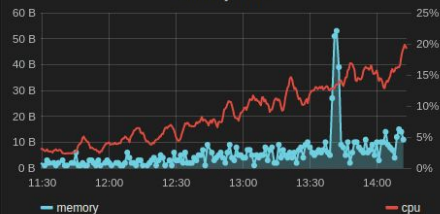
Memory / CPU



logins



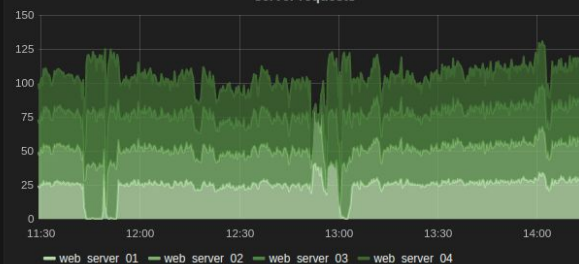
Memory / CPU



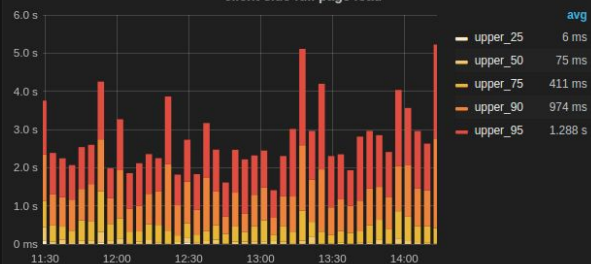
logins



server requests



client side full page load



server requests

