# VO-CLOUD upgrade
# Integration of Spark, Jupyter and HDFS in a UWS-driven cloud service

## Petr Škoda, Jakub Koza

Astronomical Institute Academy of Sciences
Ondřejov
Czech Republic

IVOA Interoperability meeting , GWS Session 1
Shanghai, China, 16th May 2017

# Concept of scientific „CLOUD"

ITERATIVE REPEATING  of  SAME computation (workflow)

Machine Learning (of emission line profiles of LAMOST)

LARGE stable INPUT data + small changing PARAMS

Many runs on SAME data (tuning required)

Graphics visualization from postprocessed output (text) files

Using WWW browser - supercomputing in PDA/mobil

CZECH TECHNICAL UNIVERSITY IN PRAGUE

FACULTY OF INFORMATION TECHNOLOGY

DEPARTMENT OF SOFTWARE ENGINEERING

Bachelor's thesis

# VO-KOREL, server for astronomical cloud computing

*Lumír Mrkva*

Supervisor: RNDr. Petr Škoda, CSc.

18th May 2012

---

CZECH TECHNICAL UNIVERSITY IN PRAGUE

FACULTY OF INFORMATION TECHNOLOGY

DEPARTMENT OF SOFTWARE ENGINEERING

Bachelor's thesis

# Design and implementation of a distributed platform for data mining of big astronomical spectra archives

*Jakub Koza*

Supervisor: RNDr. Petr Škoda, CSc.

12th May 2015

Czech Technical University in Prague

Faculty of Information Technology

Department of Software Engineering

Master's thesis

# Interactive Cloud-Based Platform for Parallelized Machine Learning of Astronomical Big Data

*Bc. Jakub Koza*

# VO-CLOUD Architecture

Distributed engine

**MASTER** (frontend)

Database of users and their experiments

Visualization

Scheduling

Load balancing

SHARED DATA STORAGE  - controlled access (Big Data)

**WORKERS**  (backend)

Computation [+ output for visualization]

# Sources of Spectra

Getting spectra + store

(restricted access – big files)

Files

      UPLOAD   from given local directory (recursive)

      DOWNLOAD  by  http + index, FTP (recursive)

VOTable

      UPLOAD VOTable   (e.g. prepared in TOPCAT - meta)

      REMOTE VOTable

             SSAP query  + Accref

                       + DataLink  + SODA

SAMP  control  - send to SPLAT

# Machine Learning of BIG Archive

# VO-CLOUD deployment

# SOM Worker example

# VO-CLOUD spectra visualisation

- Big Data visualization  (thousands spectra)

- Implemented in Python using Matplotlib

- Can visualise multiple selected spectra files

- Figure generated on server-side and then transfered to client

- User can use panning, zooming and export to different formats

- Uses WebSocket for server-client communication

# VO-CLOUD spectra visualisation

# VO-CLOUD spectra visualisation

# VO-CLOUD spectra visualisation

# Hadoop cluster infrastructure

- HDFS – Distributed filesystem utilized as a storage in Hadoop/Spark jobs
  - NameNode – One process per cluster. Contains information about all files saved inside HDFS
  - DataNode – One process per each node in cluster. Stores individual file data blocks.
- YARN – Resource manager and scheduler for Hadoop/Spark jobs.
  - ResourceManager – Main process managing resources and scheduling jobs. One process per cluster.
  - NodeManager – Process executing assigned work on each node. One process per each cluster node.
  - Problem of millions of small files (FITS) – Apache AVRO (Sequence files)

# Hadoop cluster deployment

# Spark Worker in VO-CLOUD

- Accepts JSON configuration

- Downloads requested files from the Master server to HDFS

- Executes spark-submit script using implicit parameters or parameters present in JSON configuration

- Awaits spark-submit script completion

- Downloads requested files from the HDFS

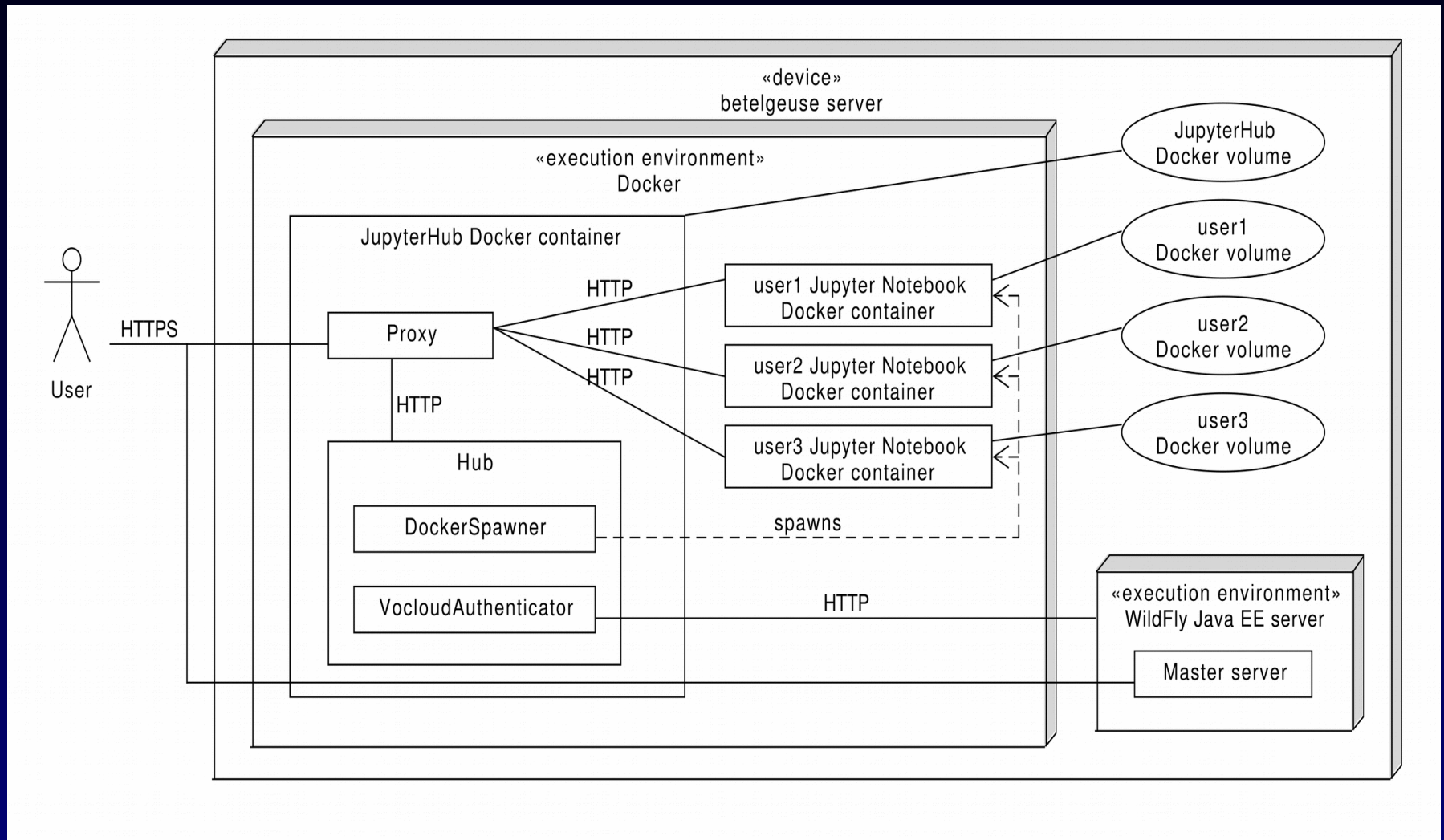- Master FILESYSTEM – NFS emulation for HDFS

# JupyterHub

- Consists of Proxy, Hub and individual Jupyter Notebook server instances

- One Jupyter Notebook server instance per authenticated user

- JupyterHub runs as Docker container

- JupyterHub spawns additional Docker containers – one for each Jupyter Notebook server instance

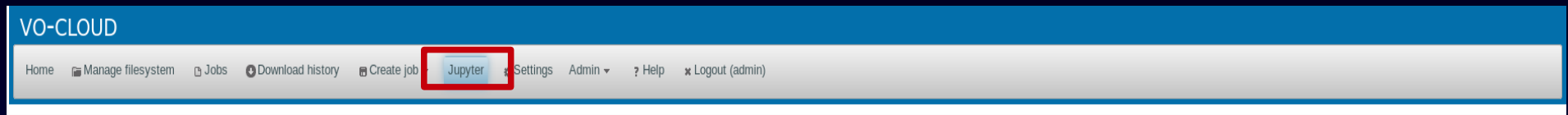- Users are isolated from each other and from hosting system itself

# JupyterHub authentication

1) VO-CLOUD generates temporal token and passes it to user's web browser

2) Web browser uses username and token for JupyterHub authentication

3) JupyterHub uses vocloud-authenticator package for authentication

4) Authenticator asks VO-CLOUD whether token is valid for relevant username

5) JupyterHub spawns new Docker container with the new Jupyter Notebook server instance for the user
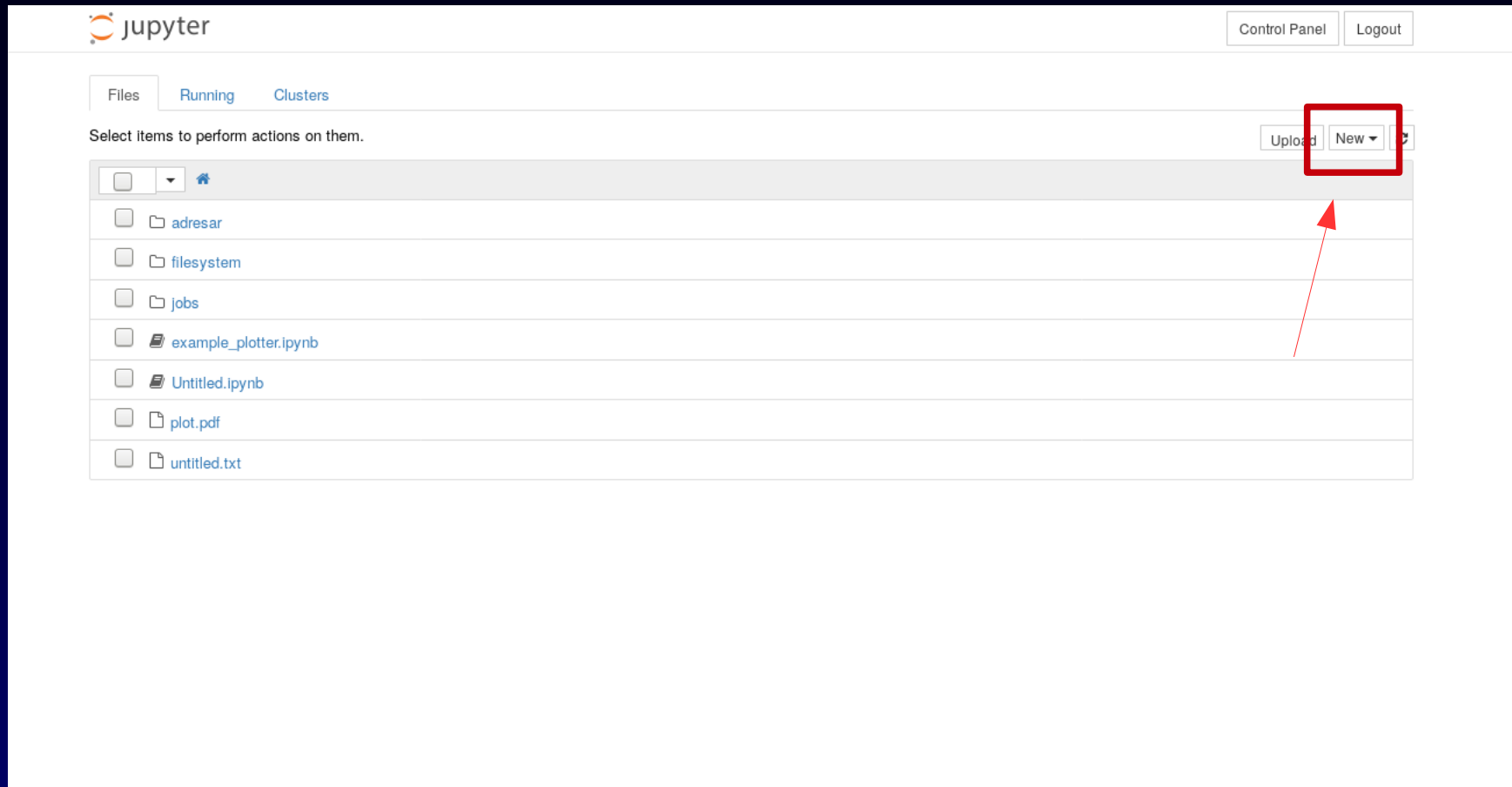
# JupyterHub deployment

# JupyterHub example



**VO-CLOUD**

Home  Manage filesystem  Jobs  Download history  Create job  **Jupyter**  Settings  Admin ▾  ? Help  Logout (admin)

Logging into JupyterHub. Please wait...
If nothing happens, click here to return to VOCLOUD.

# JupyterHub example

# JupyterHub example

# Conclusions

- VO-CLOUD is now very powerful machine learning environment capable of visualization of Big data

- It can spawn jobs on remote Spark cluster

- Provides sandbox for playing with big data in Jupyter notebook ON BIG SERVER (memory, CPU, GPU)

  but

- Still missing important capabiliites
  - AVRO not part of Spark-worker
  - Using Docker but so far not deployable as a docker (compose)
  - Combines Java EE + Python (+ Scala)

- Focused on Machine learning of 1D vectors (spectra, time series) employing VO  technology and protocols

# Source Code

https://github.com/vodev/vocloud