

IHEP Jupyter Service for physics analysis in the future

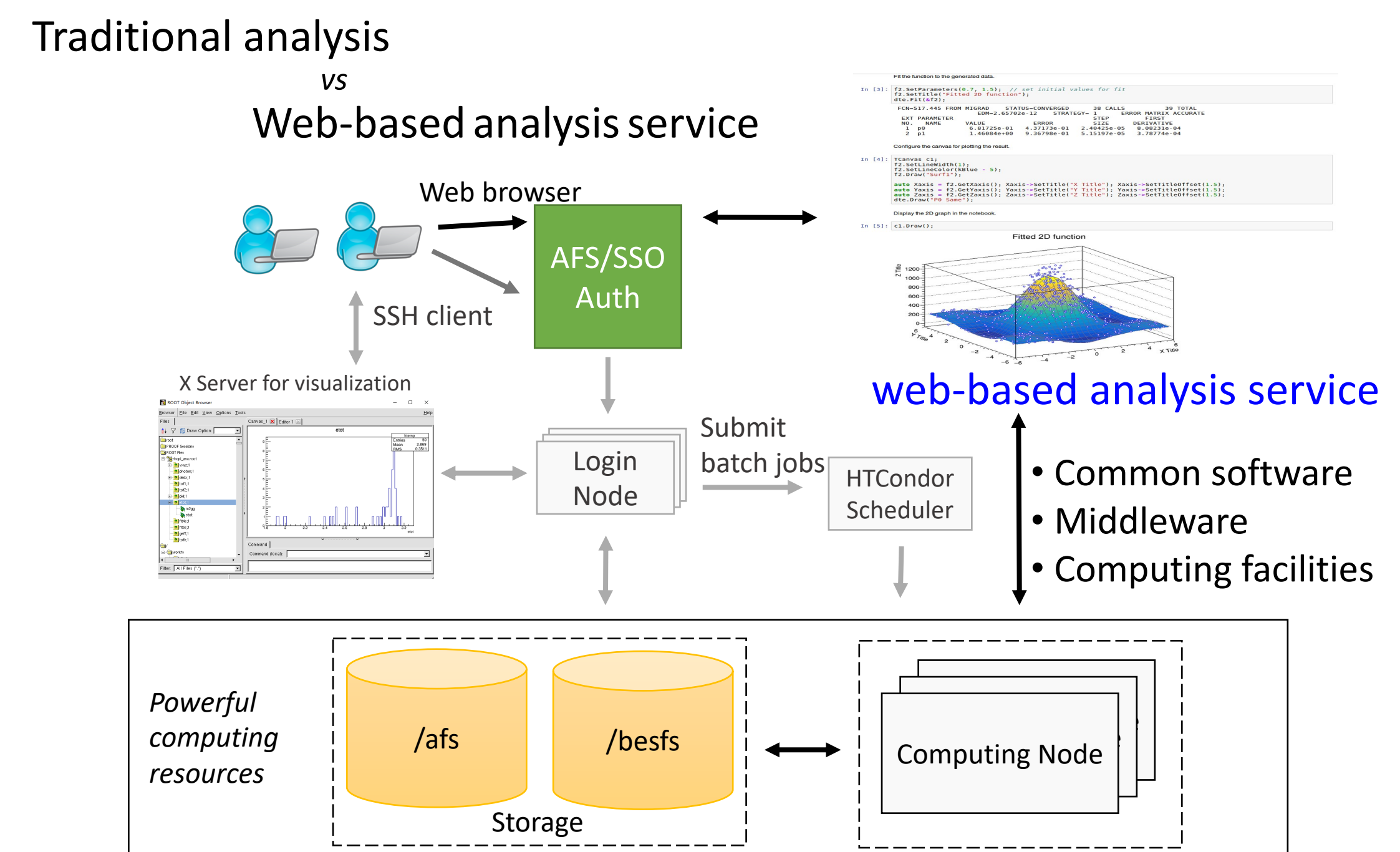
Tao Lin, Mengyao Qi, Yan Liu, Qiulan Huang, Wei Zheng, Weidong Li
Computing Center, Institute of High Energy Physics, CAS

• Motivation

- Physics analysis tasks is challenge due to “big data” (PB scale).
- A physics analysis task is split into many jobs manually, while each job will process a part of dataset.
- For the physics analysis, the same dataset will be loaded many times due to iteration.
 - Write code, edit job option and update selection criteria.
 - Submit batch jobs.
 - After jobs finished, plot data, analyze result.
 - Go to step 1.
- Submitting jobs and plotting are the most time consuming part.
- Loading dataset many times also cause I/O performance.

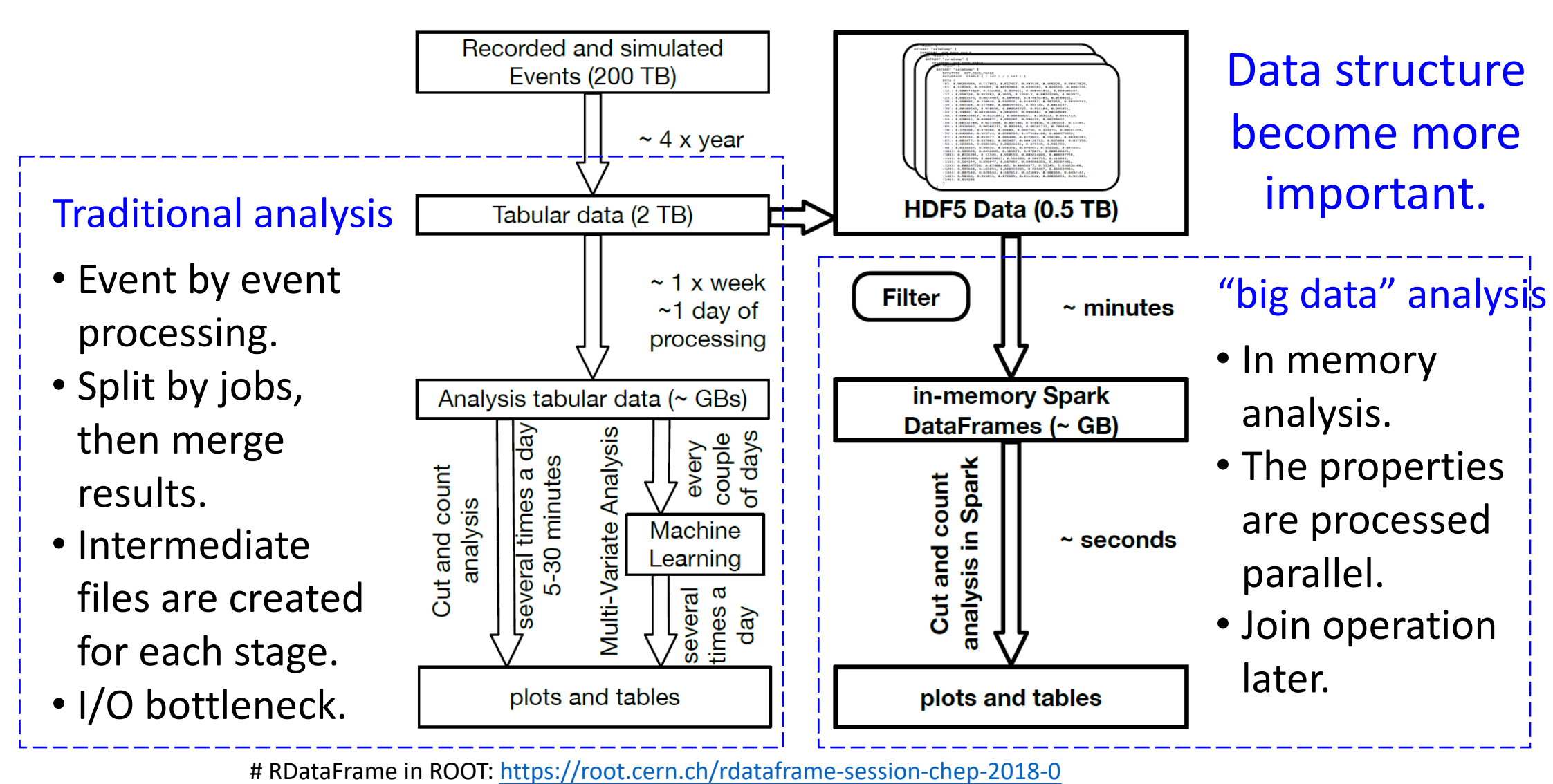
Q: How to reduce the iteration time? A: **“Analysis as a Service”**

• Analysis workflow

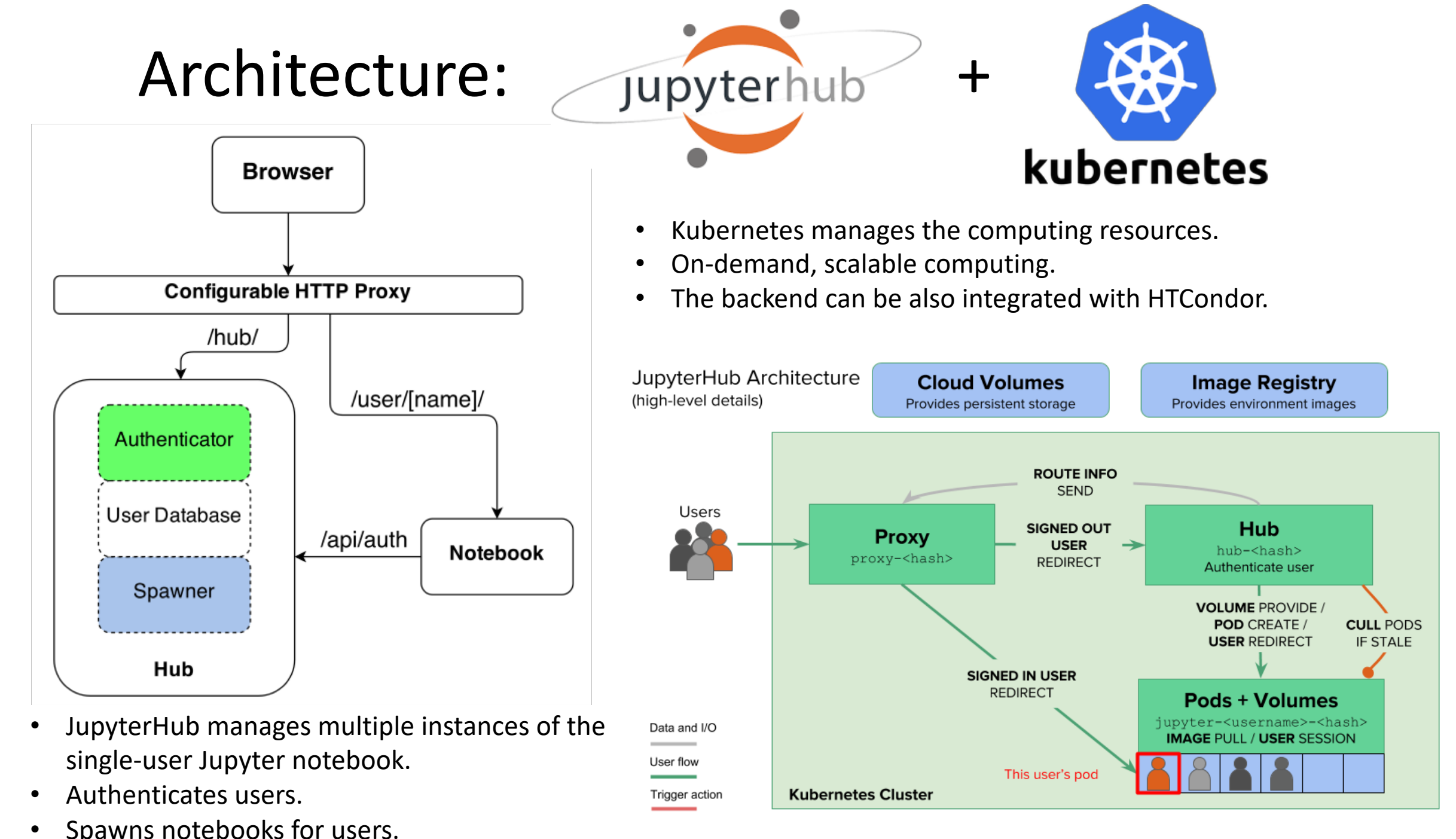


• Analysis with BIG data

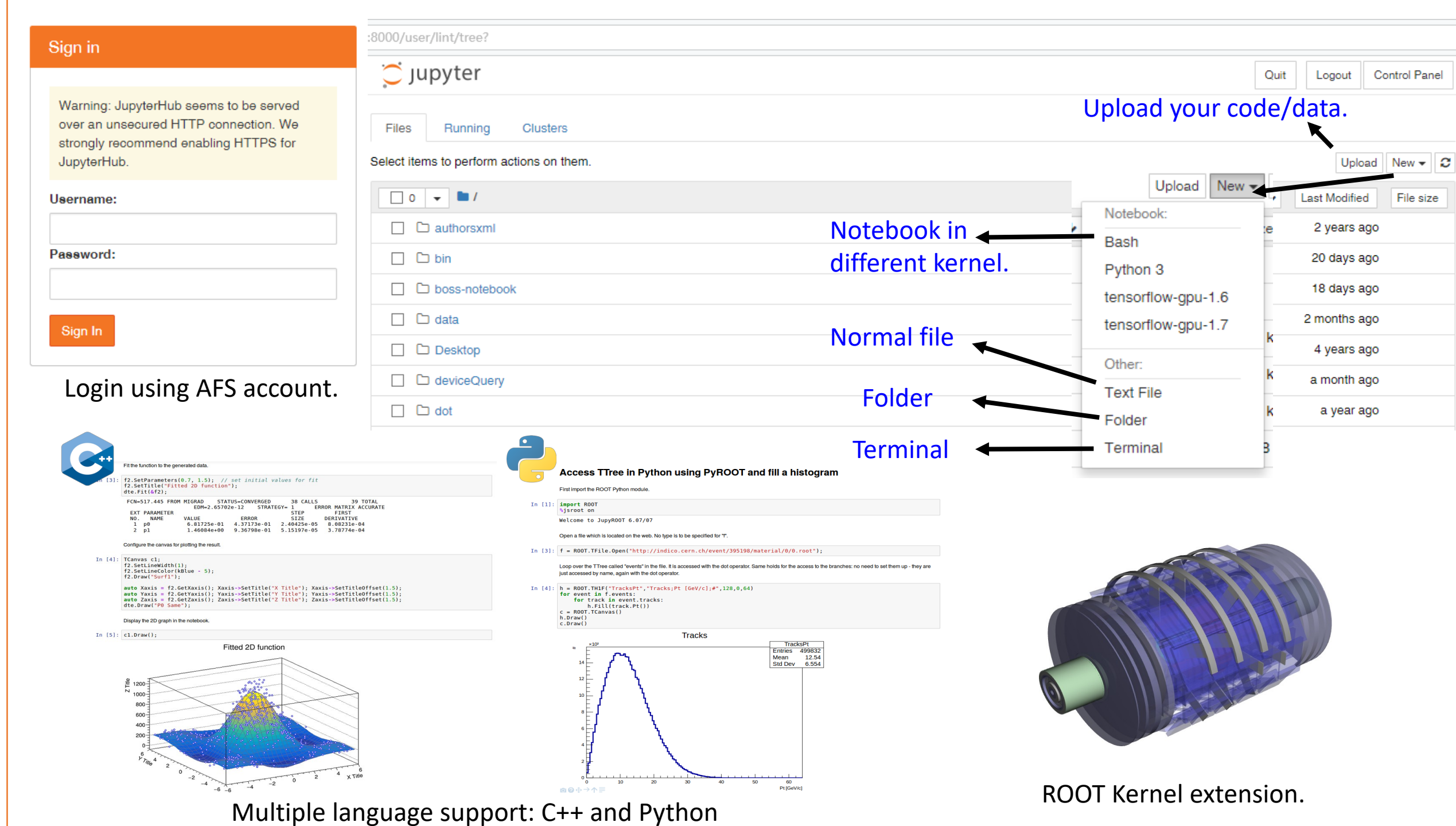
Using “big data” technology: in-memory analysis.
Read once, Analysis multiple times



• IHEP Jupyter Service



• Uses case



• Status and Challenges

Status

- Jupyter software stack is deployed at AFS.
- Users can start Jupyter and ROOT 6 in their own space.
- Setup JupyterHub in a virtual machine and setup Kubernetes in two blade servers.

Challenges

- Need to klog manually to get AFS token.
- Support multiple users and experiments.
- Unified data access between SSH and Web.
- How to benefit from “big data” technology.

Make the service available as soon as possible!
Get more feedbacks from physicists.

References:

- SWAN: <http://dx.doi.org/10.1016/j.future.2016.11.035>
- DIANA: <http://diana-hep.org>
- Jupyter: <http://jupyter.org/>
- JupyterHub: <https://jupyterhub.readthedocs.io/en/latest/index.html>
- Jupyter and Spark: <https://mapr.com/blog/configure-jupyter-spark-python/>
- Fermilab big data: <http://computing.fnal.gov/big-data/>
- NERSC: <http://www.nersc.gov/users/data-analytics/data-analytics-2/jupyter-and-rstudio/>
- Spark for HPC: <https://ieeexplore.ieee.org/document/7965154/>
- ROOT RDataFrame: <https://indico.cern.ch/event/743070/>

• Conclusions

- Jupyter service will be available at IHEP, providing users web based interactive analysis.
- A general solution to support different users and experiments.
- It can speedup the analysis such as in-memory big data analysis.
- A prototype based on JupyterHub is deployed in a virtual machine.
- User authentication works, however AFS token does not work properly.
- Kubernetes cluster is also setup using two blade servers.