



Gaia Spark analysis platform Technology choices





Original GDAF system
Zeppelin & Spark deployed on physical hardware
Based in part on Cloudera deployment

Initial goal for our project

- replicate GDAF in the cloud

Subsequent goals

- additional analysis tools
- scalable deployment
- more users
- more data
 - DR3
 - DR4





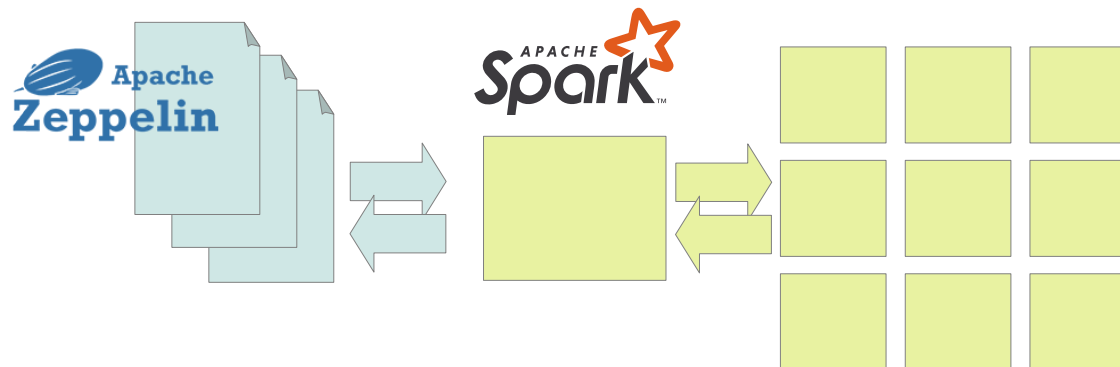
Technology choice #1 Hadoop/Yarn scheduler



“known unknowns”

Hadoop/Yarn

- GDAF system – working example
- Cloudera deployment - lots of documentation
- Standard deployment - lots of blogs and howtos
- Virtual machine based
- Spark cluster deployed on a static set of resources
- Zeppelin notebooks all interact with the same Spark cluster





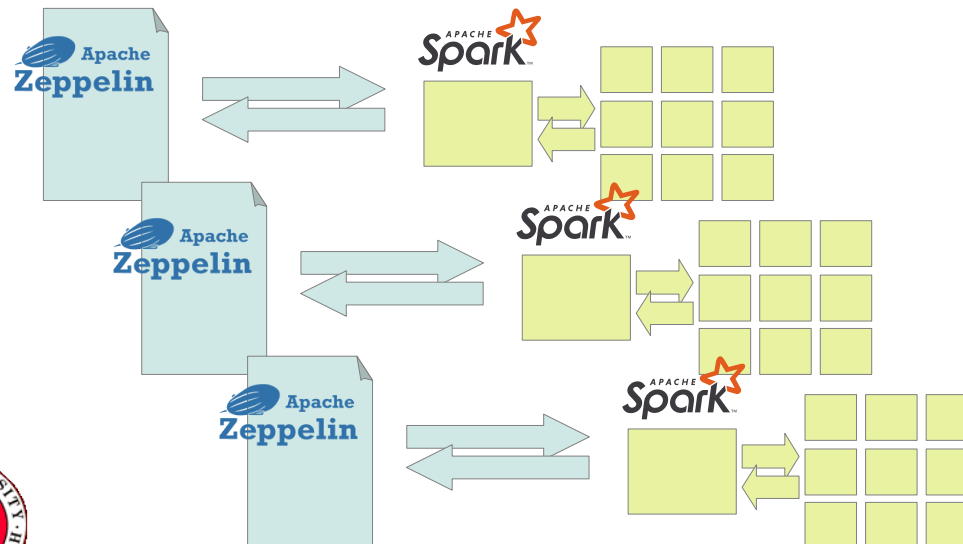
Technology choice #2 Kubernetes scheduler



kubernetes

“unknown unknowns”

- Experimental in 2020
- Zeppelin pre-release only
- Not recommended for production
- Spark cluster on demand
- Notebooks launch their own Spark cluster
- Kubernetes container based
- Standard technology in 2021/2022



Technology choices Spark deployments

Developing both systems in parallel



Hadoop/Yarn

- Up and running early
- Spark and Zeppelin by Nov 2019
- Gaia DR2 parquet data by Dec 2019
- Automated with Ansible



Working system to develop science cases

- Python libraries
- AXS extension
- Parquet partitions



kubernetes

- Development platform for Spark and Zeppelin on Kubernetes
- Automated with Helm



Technology experiments

- Openstack Magnum & Manila
- Kubernetes Helm
- Terraform
- OAuth login
- Drupal CMS



Technology choices deployment tools



- Good documentation
- Simple client interface
- Machine readable results

Working towards full automation



- Well documented
- Well supported on StackOverflow
- Lots of plugins



- De facto standard for Kubernetes deployments
- Good version control for dependencies



- Didn't live up to the hype
- Relies on local state
- Clear text secrets

