

Jakub Háva
jakub@h2o.ai

Productionizing H2O Models using Sparkling Water

DataFest, Glasgow
March 21, 2018

Spark[★] + H₂O

SPARKLING
WATER

Bio

- Sr. Software Engineer at H2O.ai, Sparkling Water Core Team
- Finished Master's at Charles University in Prague
- Implemented high-performance cluster monitoring tool for JVM based languages (JNI, JVMTI, instrumentation)
- Enjoy climbing & tea

Distributed Sparkling Team

- Michal - Mountain View, CA
- Kuba - Prague, CZ

H₂O+Spark =
Sparkling
Water

Sparkling Water

- Transparent integration of H2O with Spark ecosystem - MLlib/ML and H2O side-by-side
- Transparent use of H2O data structures and algorithms with Spark API
- Platform for building Smarter Applications
- Excels in existing Spark workflows requiring advanced Machine Learning algorithms

Benefits



- Additional algorithms
 - NLP
- Powerful data munging
- ML Pipelines
- Advanced algorithms
 - speed & accuracy
 - advanced parameters
- Fully distributed and parallelized
- Graphical environment
- R/Python interface

**How to use
Sparkling
Water**

Start Spark with Sparkling Water

start.sh

```
1 $SPARK_HOME/bin/spark-submit \
2   --class water.SparklingWaterDriver \
3   --packages ai.h2o:sparkling-water-examples_2.10:1.6.3 \
4   --executor-memory=6g \
5   --driver-class-path scalastyle.jar /dev/null
```

Raw

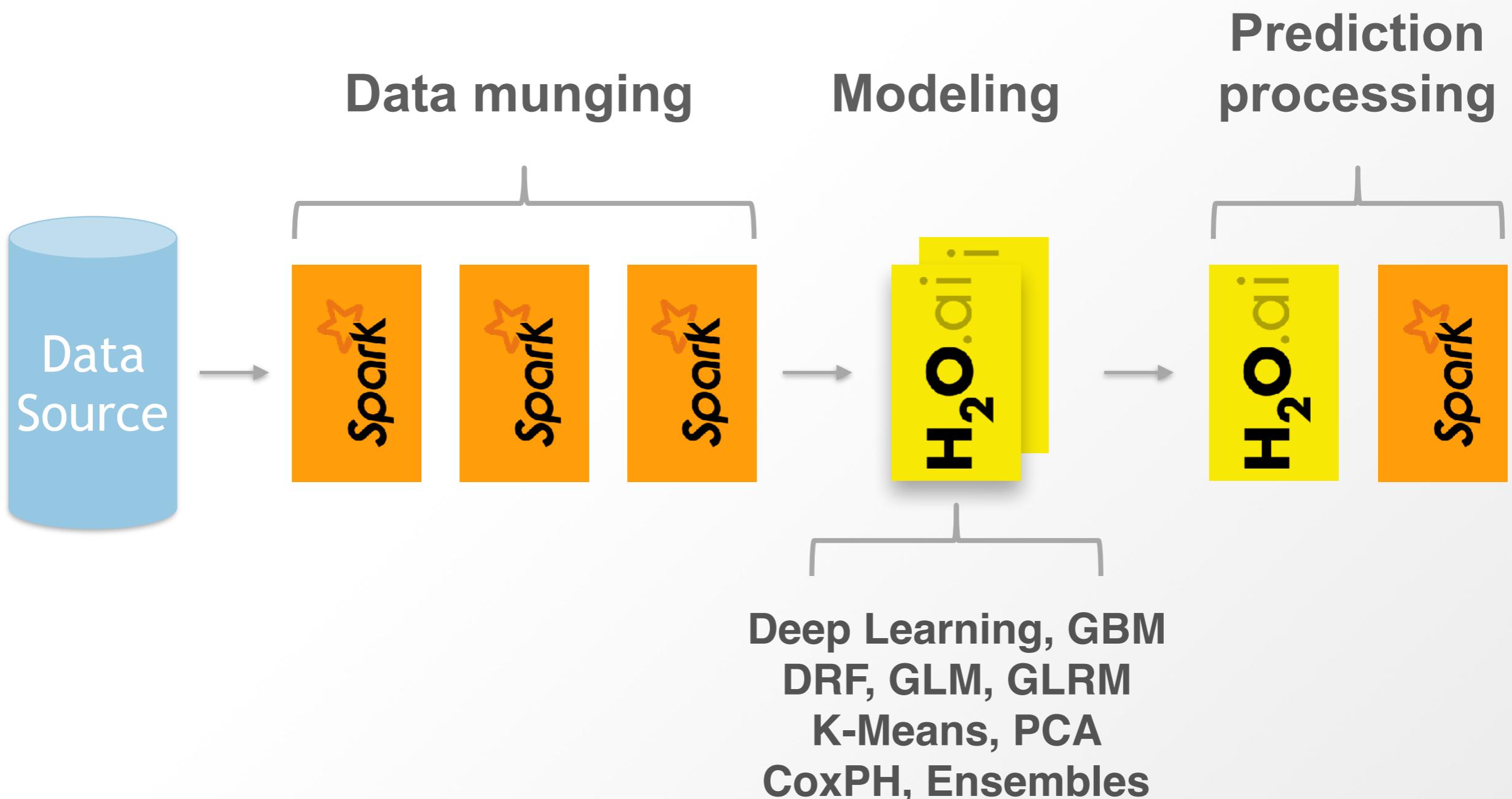


The screenshot shows the H2O Flow web application. At the top, there's a navigation bar with tabs for 'Flow', 'Cell', 'Data', 'Model', 'Score', 'Admin', and 'Help'. Below the navigation is a toolbar with various icons. The main area is titled 'Untitled Flow' and contains a list of H2O API routes under the heading 'Assistance'. The routes listed are:

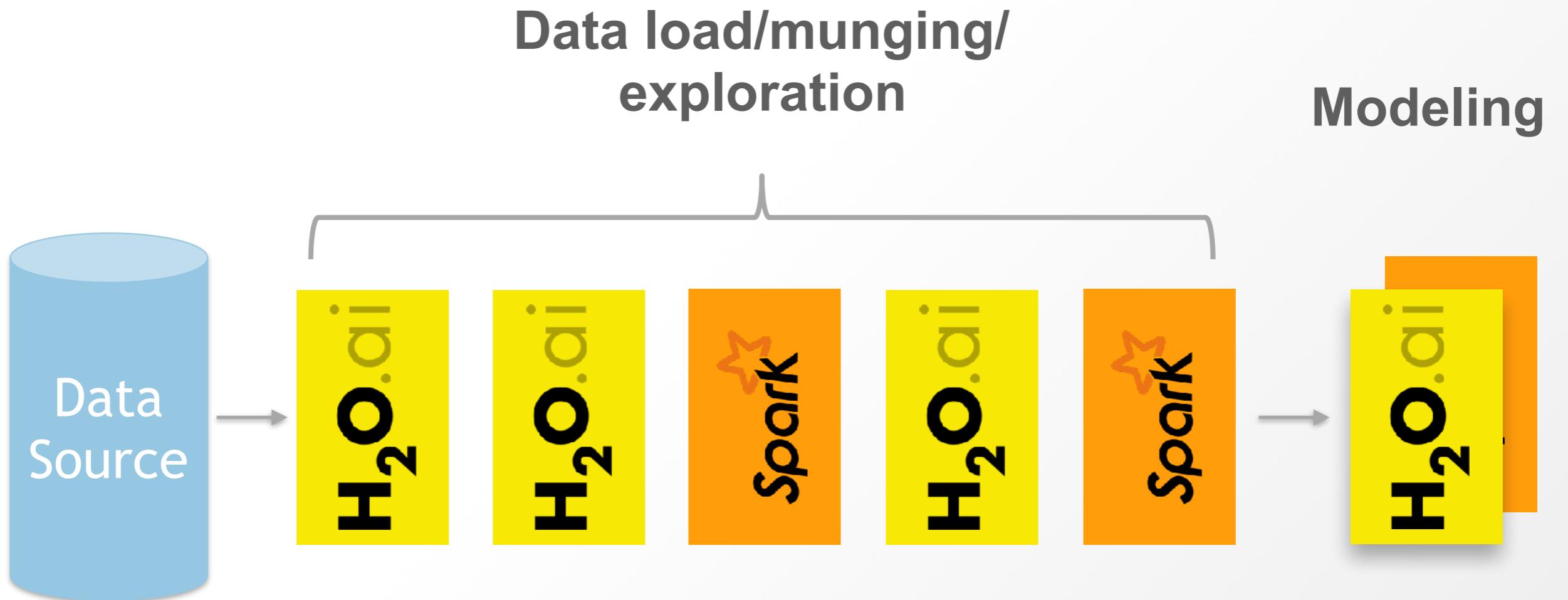
- Import File
- getFrames
- splitFrame
- getModels
- getGrids
- getPredictions
- getJobs
- buildModel
- Input Model
- predict
- getRDDs
- getDataframes

To the right of the assistance panel is a 'Help' section with links to 'Quickstart Videos' and 'View example Flow'. It also includes sections for 'GENERAL' and 'EXAMPLE' flows.

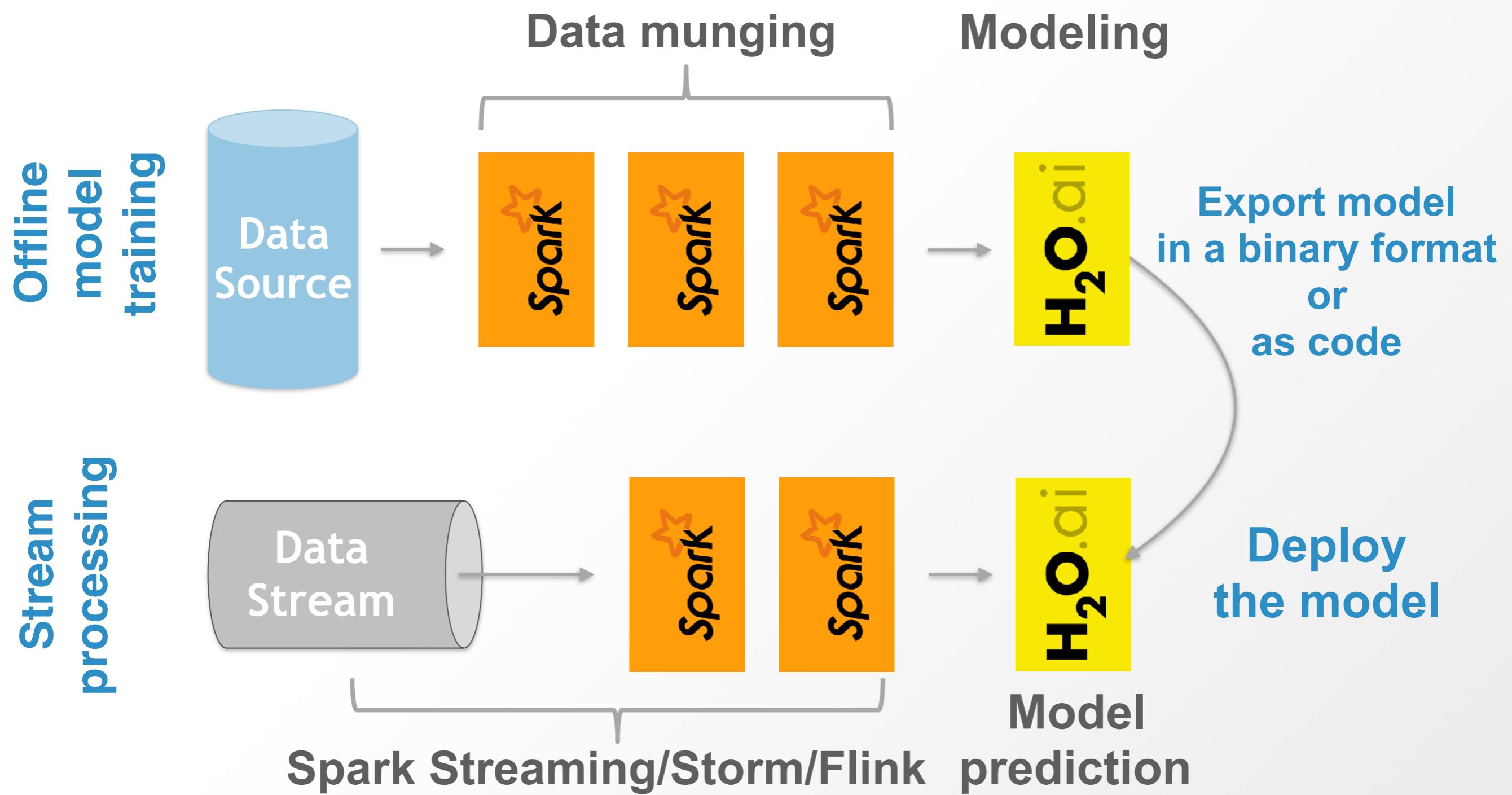
Model Building



Data Munging



Stream Processing



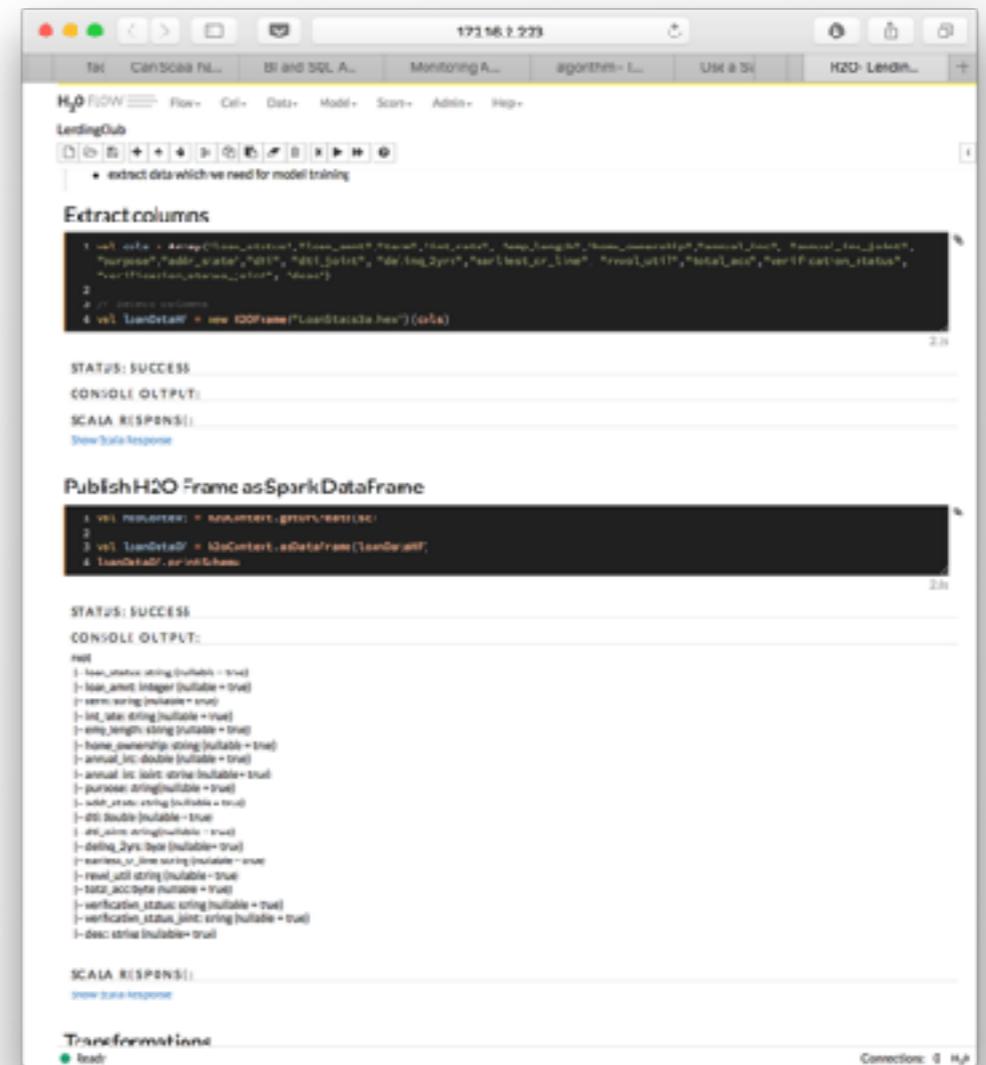
Scoring

- POJO
 - Plain Old Java Object
- MOJO
 - Model Object Optimized
- No runtime dependency on H2O framework

Features Overview

Scala Code in H2O Flow

- New type of cell
 - Access Spark from Flow UI
 - Experimenting made easy



H2O Frame as Spark's Datasource

- Use native Spark API to load and save data
- Spark can optimize the queries when loading data from H2O Frame
- Use of Spark query optimizer

Machine Learning Pipelines

- Wrap our algorithms as Transformers and Estimators
- Support for embedding them into Spark ML Pipelines
- Can serialize fitted/unfitted pipelines
- Unified API => Arguments are set in the same way for Spark and H2O Models
- Using MOJO when possible

MLlib Algorithms in Flow UI

- Can examine them in H2O Flow
- Can generate POJO/MOJO out of them

PySparkling Made Easy

- PySparkling is on PyPi
- Contains all H2O and Sparkling Water dependencies, no need to worry about them
- Just add in on your Python path and that's it
- Python independent PySparkling distribution - implemented by [SW-341]

RSparkling

- Sparkling Water for R
- Based on Sparklyr package
- Independent on Spark and Sparkling Water version

And others!

- Support for Datasets
- Zeppelin notebook support
- XGBoost Support (local mode so far)
- Support for high-cardinality fast joins
- Secure Communication - SSL
- Support for Sparse Data conversions
- Bug fixes

Upcoming Features

- Integration with Steam
- More advanced pipelines with MOJOs
- Integration with Driverless AI
- Advanced Telemetry

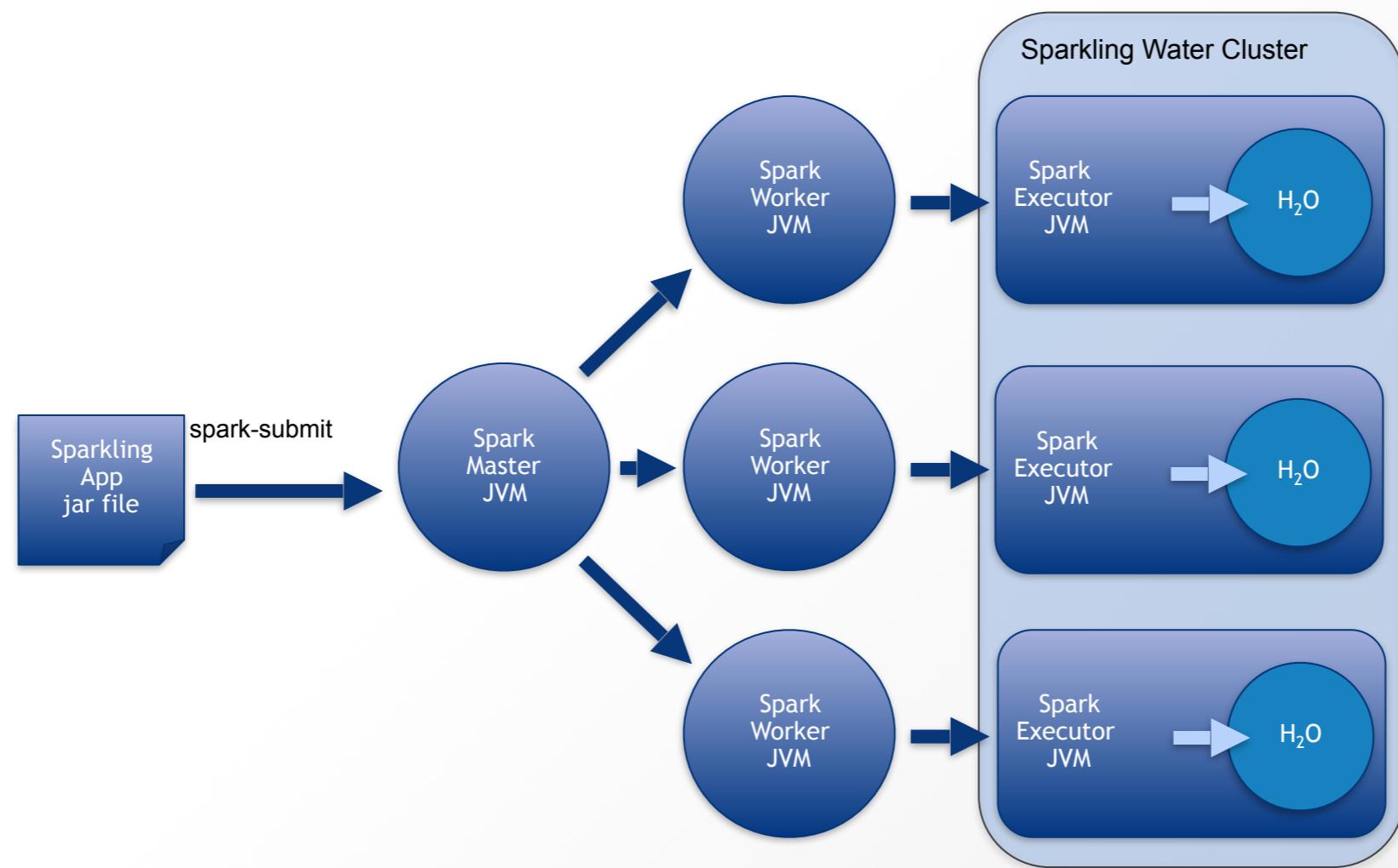
**What's
inside?**

Two Backends

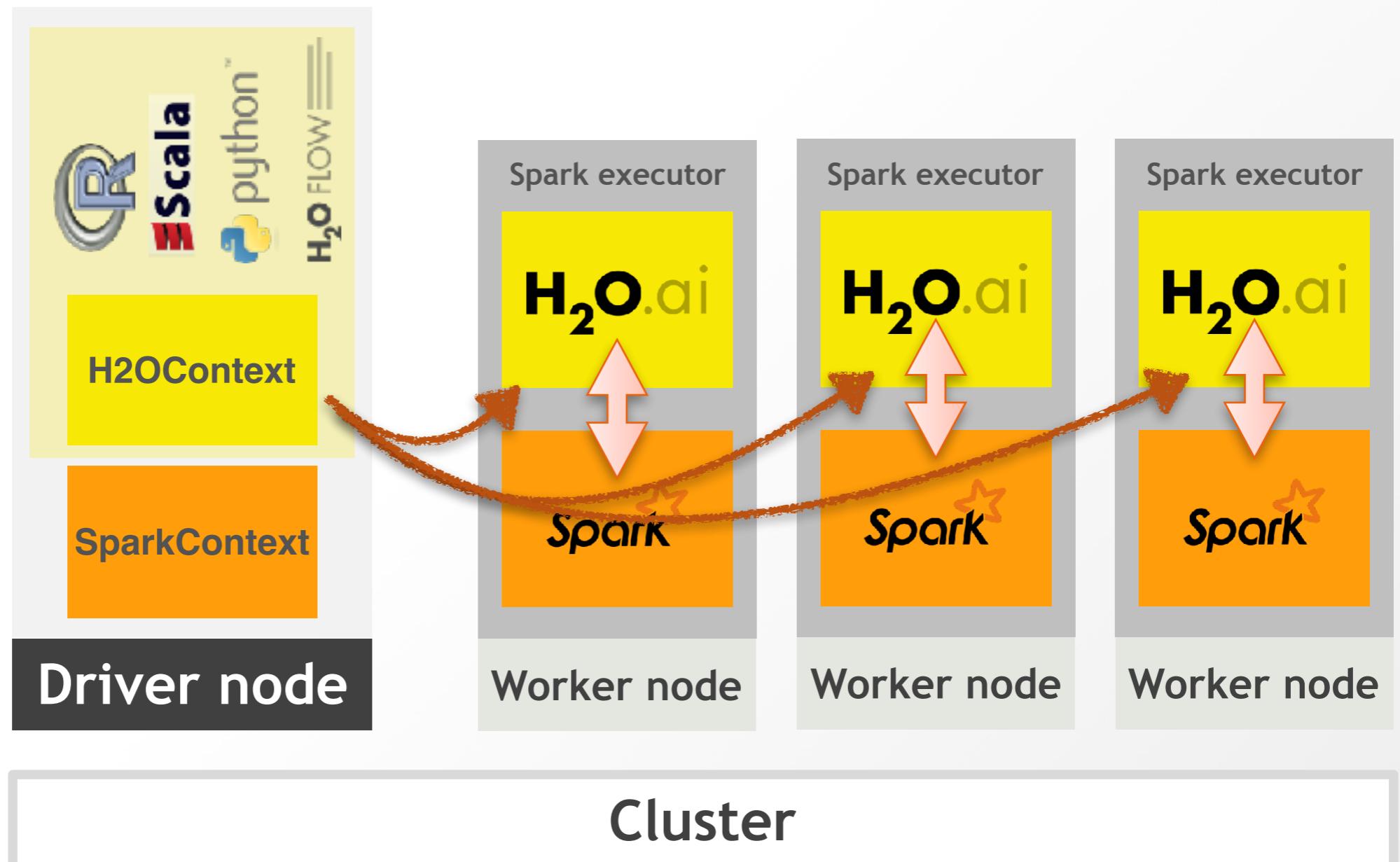
- Sparkling Water has two backends
 - External
 - Internal
- The backend determines where the H2O cluster is located
- Each backend is good for different use cases

Internal Backend

Sparkling Water Internal Backend



Internal Backend



Pros & Cons

- Advantages
 - Easy to configure
 - Faster (no need to send data to different cluster)
- Disadvantages
 - Spark kills or joins new executor => H2O goes down
 - No way how to discover all executors

External Backend

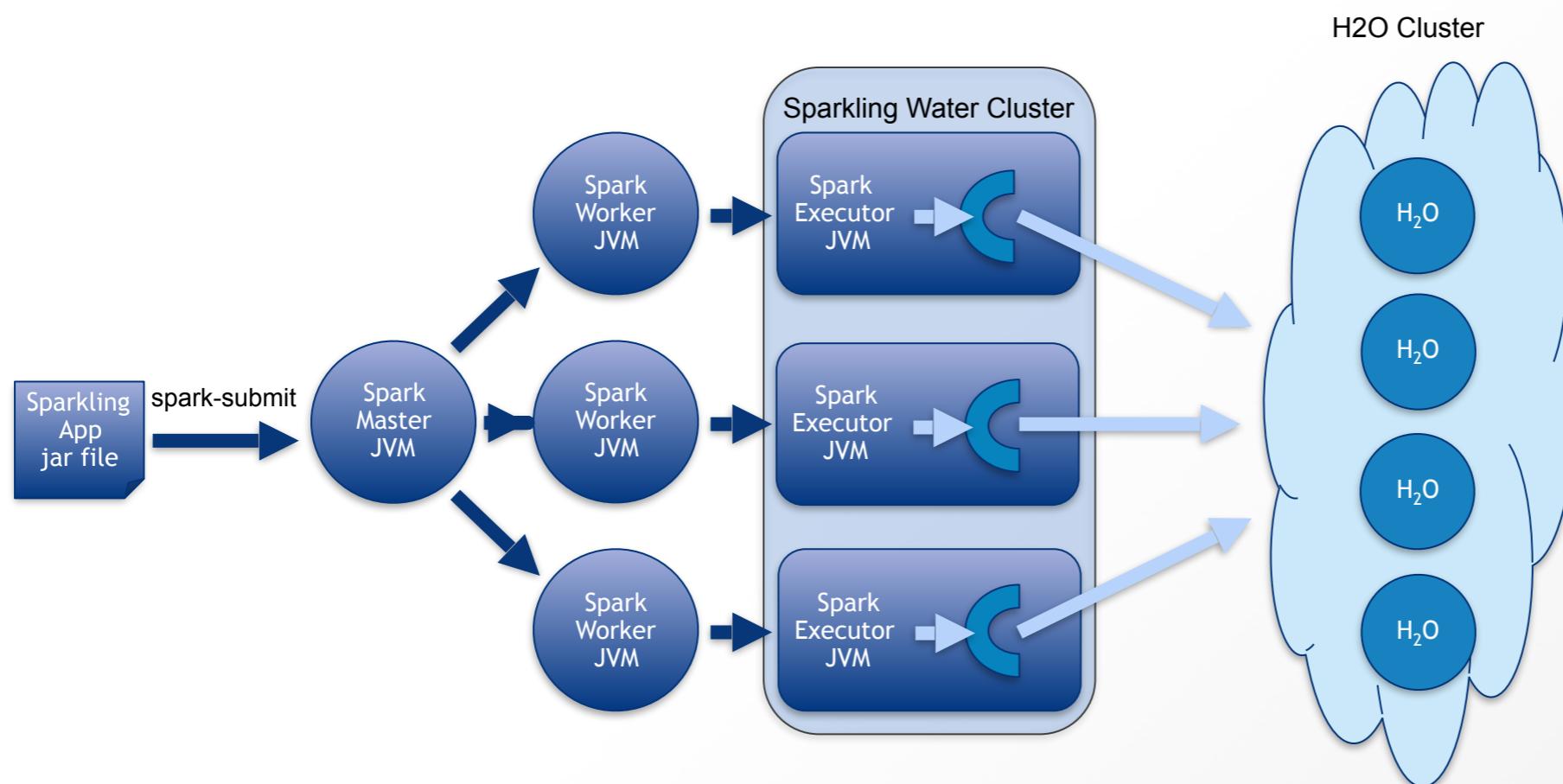
Overview

- Sparkling Water is using external H2O cluster instead of starting H2O in each executor
- Spark executors can come and go and H2O won't be affected
- Start H2O cluster on YARN automatically

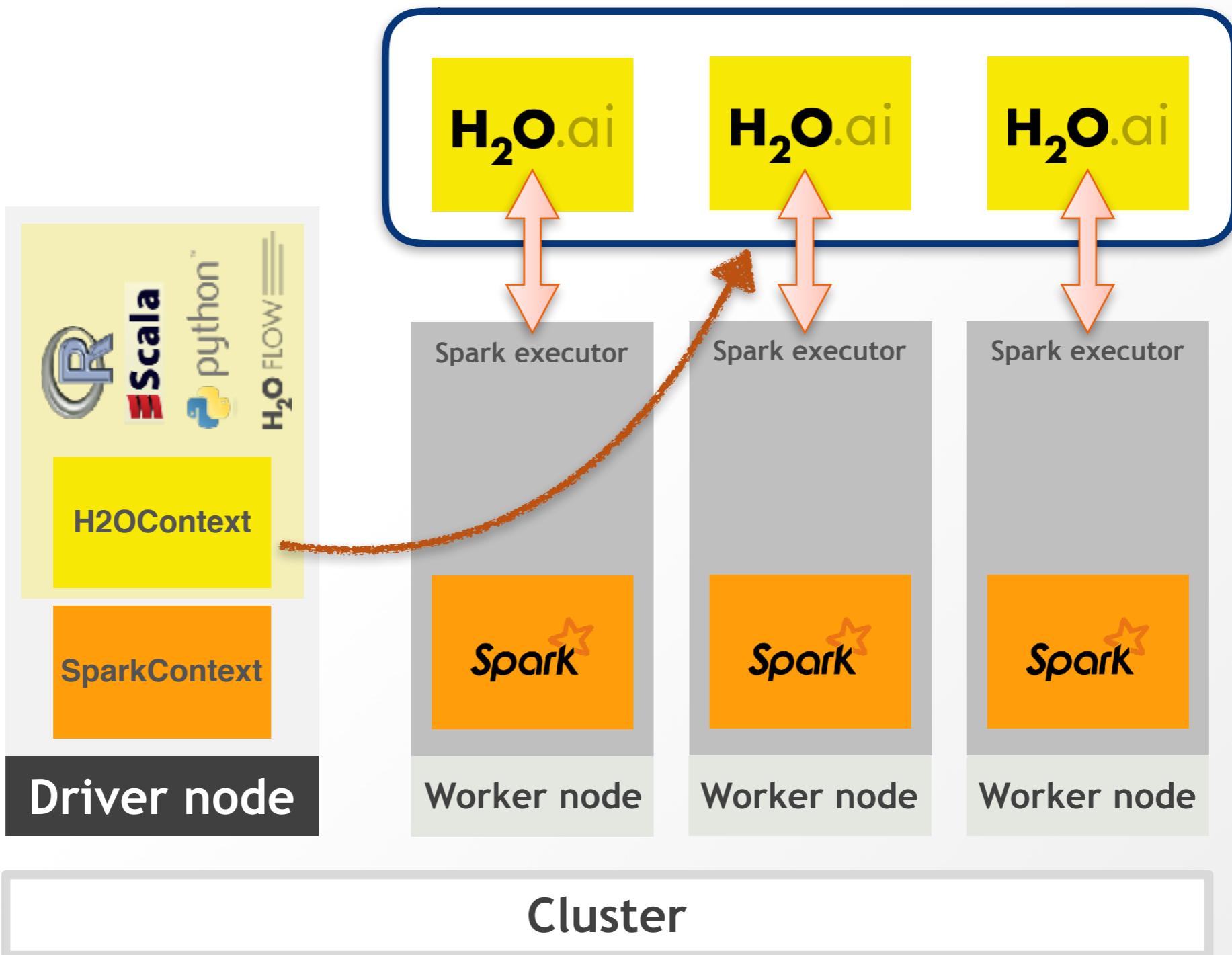
Separation Approach

- Separating Spark and H2O
 - But preserving same API:
`val h2oContext = H2OContext.getOrCreate(http://h2o:
54321/)`
- Spark and H2O can be submitted as a YARN job and controlled in separation
 - But H2O still needs non-elastic environment (H2O itself does not implement HA yet)

Sparkling Water External Backend



External Backend



Pros & Cons

- **Advantages**
 - H2O does not crash when Spark executor goes down
 - Better resource management since resources can be planned per tool
- **Disadvantages**
 - Transfer overhead between Spark and H2O processes
 - Under measurement with cooperation of a customer

Modes

- **Auto Start Mode**
 - Start h2o cluster automatically on YARN
- **Manual Start Mode**
 - User is responsible for starting the cluster manually

Thank you!

Learn more at h2o.ai
Follow us at [@h2oai](https://twitter.com/h2oai)

More Info

- Documentation: <http://docs.h2o.ai>
- Tutorials: <https://github.com/h2oai/h2o-tutorials>
- Slides: <https://github.com/h2oai/h2o-meetups>
- Videos: <https://www.youtube.com/user/0xdata>
- Events & Meetups: <http://h2o.ai/events>
- Stack Overflow: <https://stackoverflow.com/tags/sparkling-water>
- Google Group: <https://tinyurl.com/h2ostream>
- Gitter: <http://gitter.im/h2oai/sparkling-water>