

Jakub Háva  
[jakub@h2o.ai](mailto:jakub@h2o.ai)

# The Next Generation of Machine Learning on Apache Spark

PyData Bratislava  
Bratislava, May 29, 2018

 + H<sub>2</sub>O

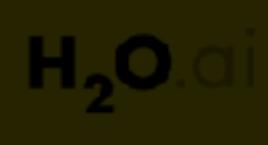
---

SPARKLING  
**WATER**

# **Who am I**

- Senior Software engineer at H2O.ai - Core Sparkling Water
- Finished Master's at Charles Uni in Prague
- Implemented high-performance cluster monitoring tool for JVM based languages ( JNI, JVMTI, instrumentation )
- Climbing & Tea lover, ( doesn't mean I don't like beer!)

# H<sub>2</sub>O Products



In-Memory, Distributed  
Machine Learning Algorithms  
with H2O Flow GUI



H2O AI Open Source Engine  
Integration with Spark



Lightning Fast machine  
learning on GPUs

DRIVERLESSAI

Automatic feature  
engineering, machine learning  
and interpretability

## Steam

Secure multi-tenant H2O clusters

# **Distributed Sparkling Team**

- Michal - Mt. View, CA
- Kuba - Prague, CZ

H<sub>2</sub>O+Spark =  
Sparkling  
Water

# Sparkling Water

- Transparent integration of H2O with Spark ecosystem - MLlib 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

Functionality missing in H2O can be replaced by Spark and vice versa

# Benefits



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

# How to use Sparkling Water?

# Start spark with Sparkling Water

start.sh

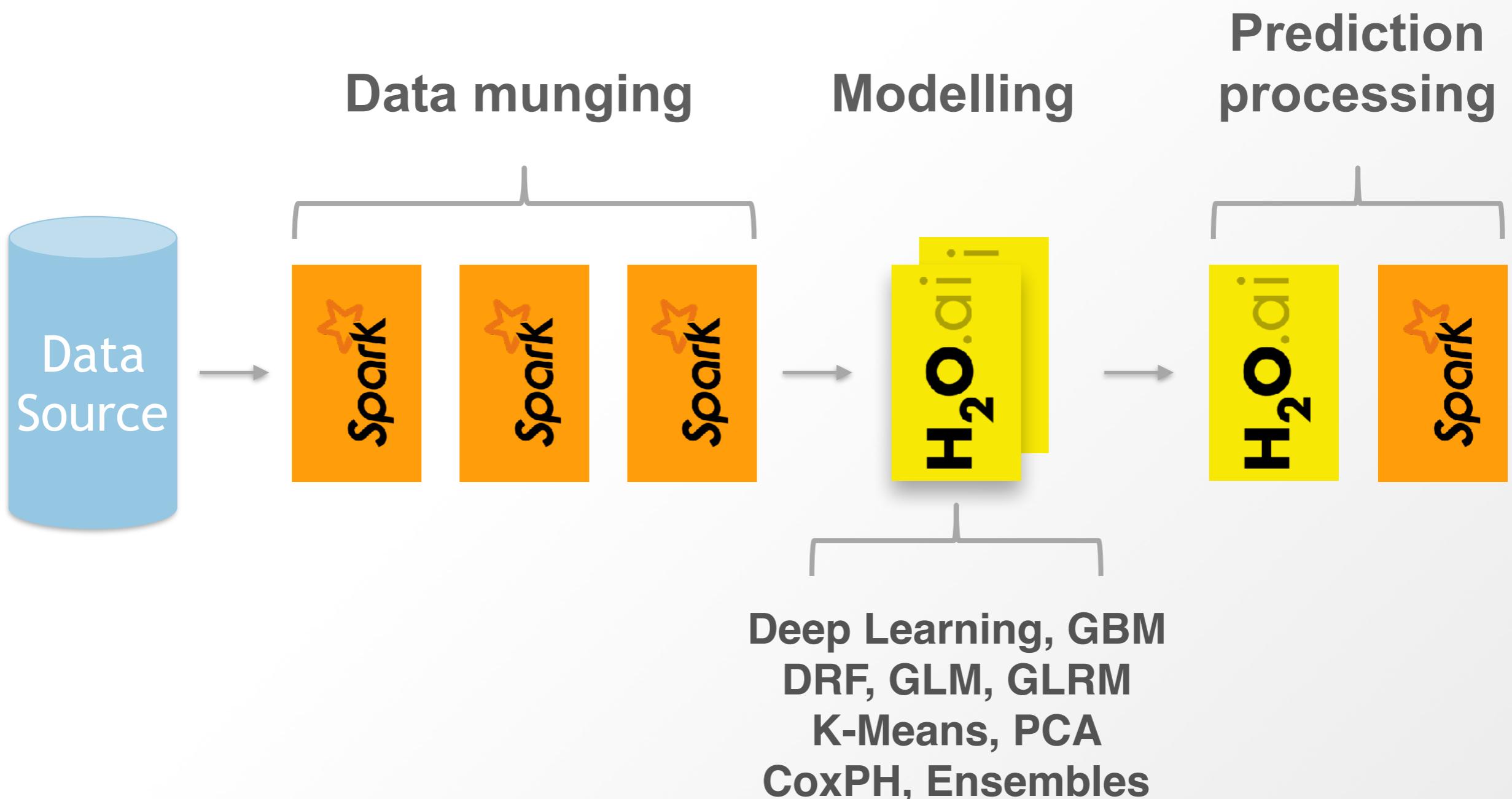
```
1 $SPARK_HOME/bin/spark-submit \
2   --class water.SparklingWaterDriver \
3   --packages ai.h2o:sparkling-water-examples_2.10:1.6.3 \
4   --executor-memory=6g \
5   --driver-class-path scalastyle.jar /dev/null
```

Raw

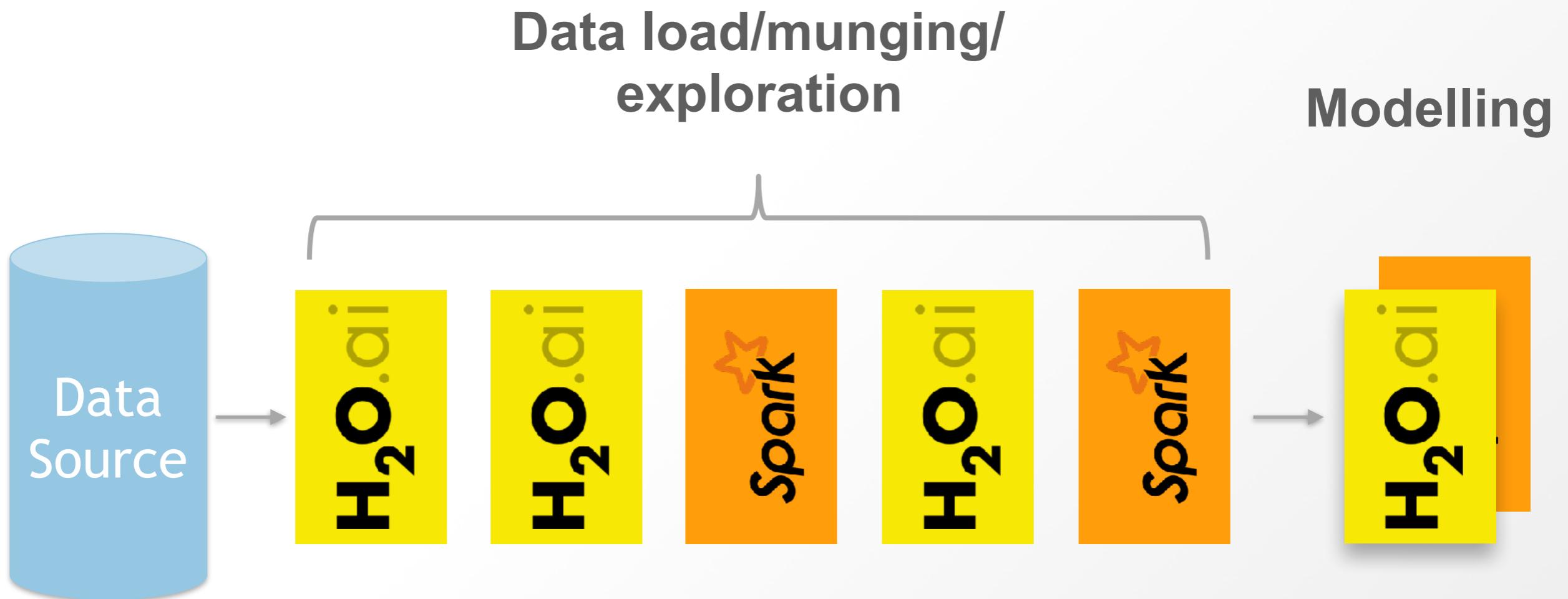


The screenshot shows the H2O Flow web application. At the top, there's a navigation bar with tabs for 'Flow', 'Cell', 'Data', 'Model', 'Score', 'Admin', and 'Help'. Below the navigation is a toolbar with various icons. The main area is titled 'Untitled Flow' and contains a search bar and a list of flow nodes. To the right of the flow area is a sidebar with sections for 'Assistance' (listing methods like importFrame, getFrames, splitFrame, etc.) and 'Help' (with links to quickstart videos and examples). The bottom of the screen shows a footer with links for 'Flow Web UI', 'Importing Data', 'Building Models', 'Making Predictions', 'Using Flows', and 'Troubleshooting Flow'.

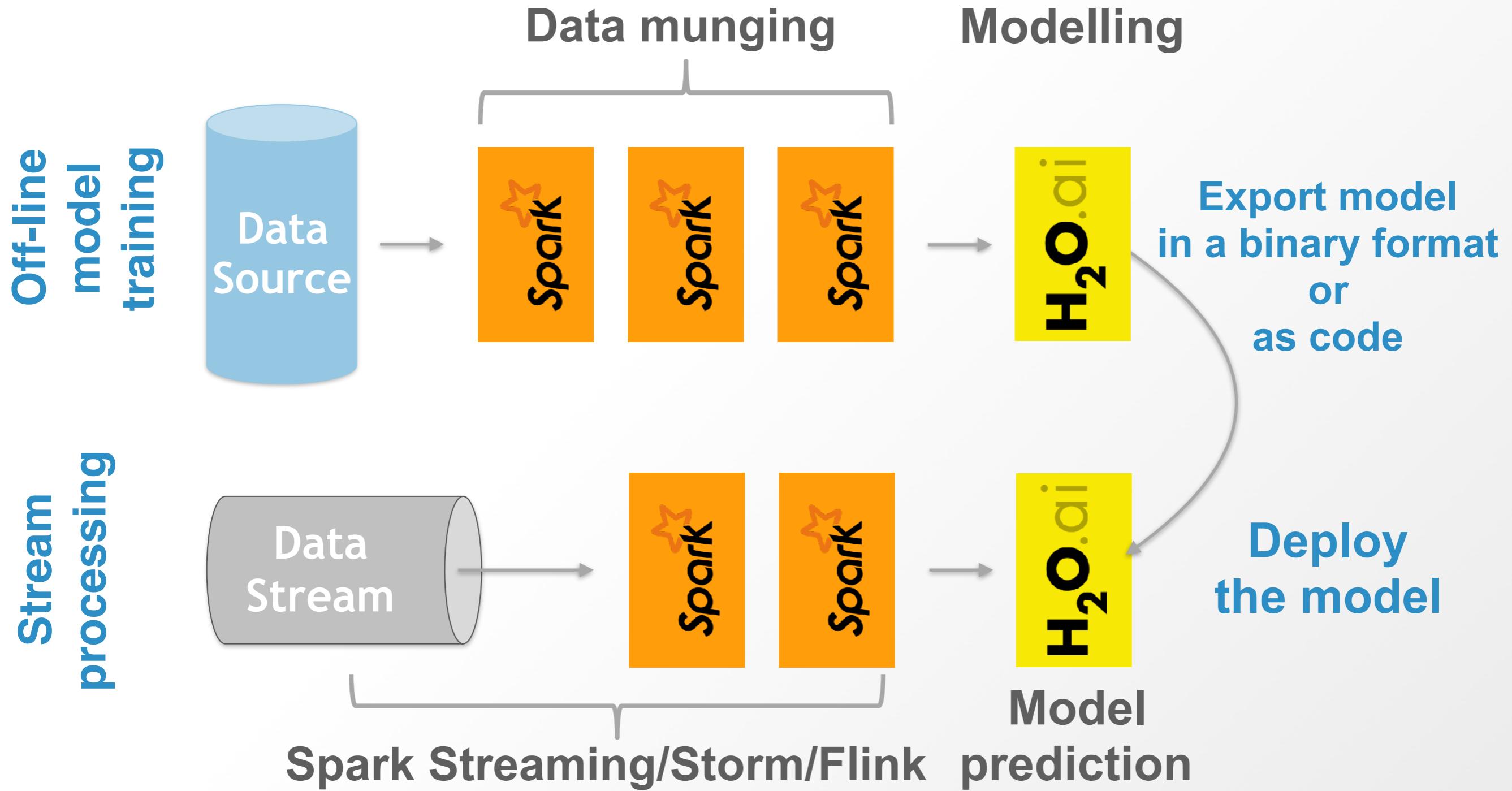
# Model Building



# Data Munging



# Stream Processing



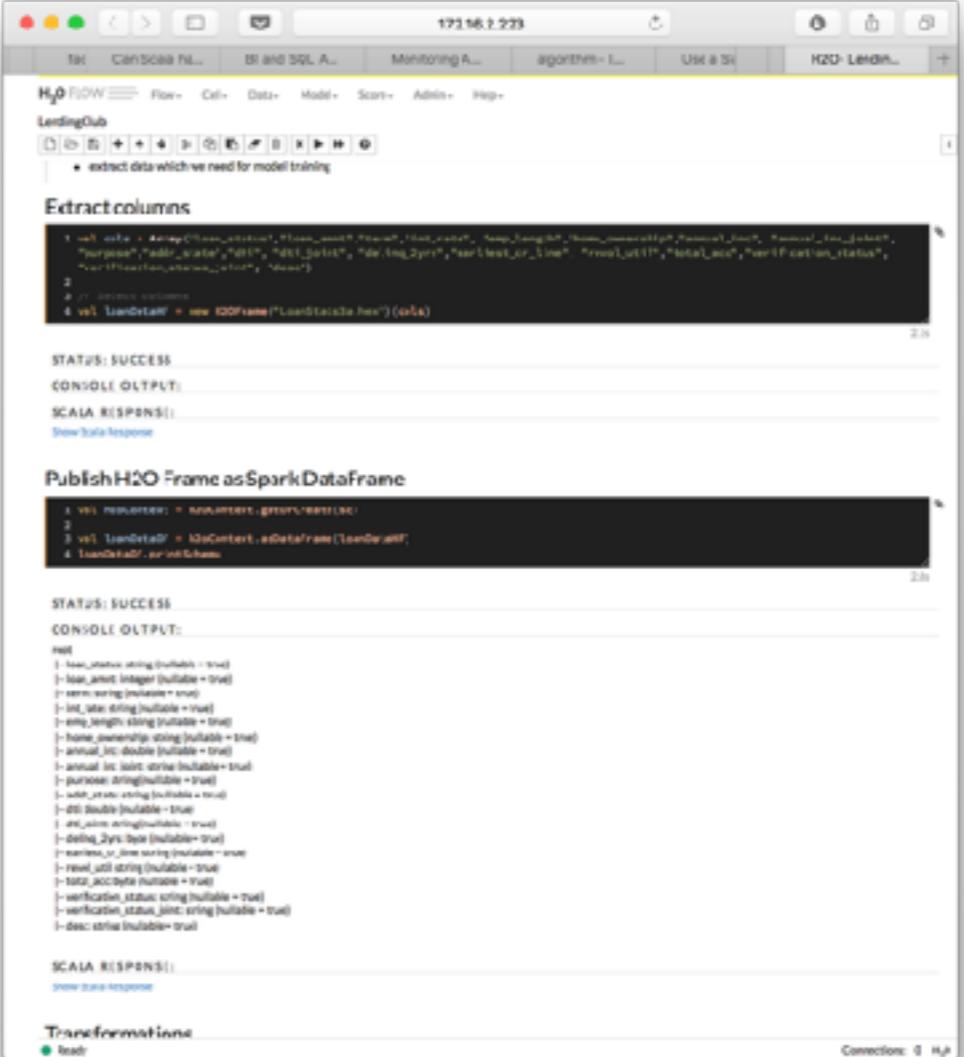
# Scoring

- POJO
  - Plain Old Java Object
- MOJO
  - Model Object Optimised
- 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



The screenshot shows the H2O Flow interface with a Scala code cell. The code is:

```
1 val df = spark.read.option("header","true").option("inferSchema","true").option("sep",",").load("hdfs://localhost:9000/SparkData/loan_train.csv")
2 df.printSchema()
3 df.show(5)
4 df.createOrReplaceTempView("loan")
5 val trainData = new H2OFrame("LoadData(hdfs://localhost:9000/SparkData/loan_train.csv) (data)")
```

Below the code, the status is shown as "STATUS: SUCCESS". The "SCALA RESPONSE" section contains:

```
1 val responder = Respondent.getOrCreate()
2 responder.setContext(df)
3 responder.setH2OFrame(trainData)
4 responder.publish()
```

The "CONSOLE OUTPUT" section shows the output of the Scala code, including the schema and first five rows of the DataFrame.

# **H2O Frame as Spark's Datasource**

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

# Machine learning pipelines

- Wrap our algorithms as Transformers and Estimators
- Support for embedding them into Spark ML Pipelines
- Can serialise fitted/unfitted pipelines
- Unified API => Arguments are set in the same way for Spark and H2O Models
- Integration with Mojo pipelines

# **PySparkling made easy**

- PySparkling is on PyPi
- Contains all H2O and Sparkling Water dependencies, no need to worry about them
- Just add it to your Python path and that's it
- Correctly specifies dependency on PySpark

# 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..

# **Coming features**

- Integration with Steam
- More advanced pipelines with MOJOs
- ...

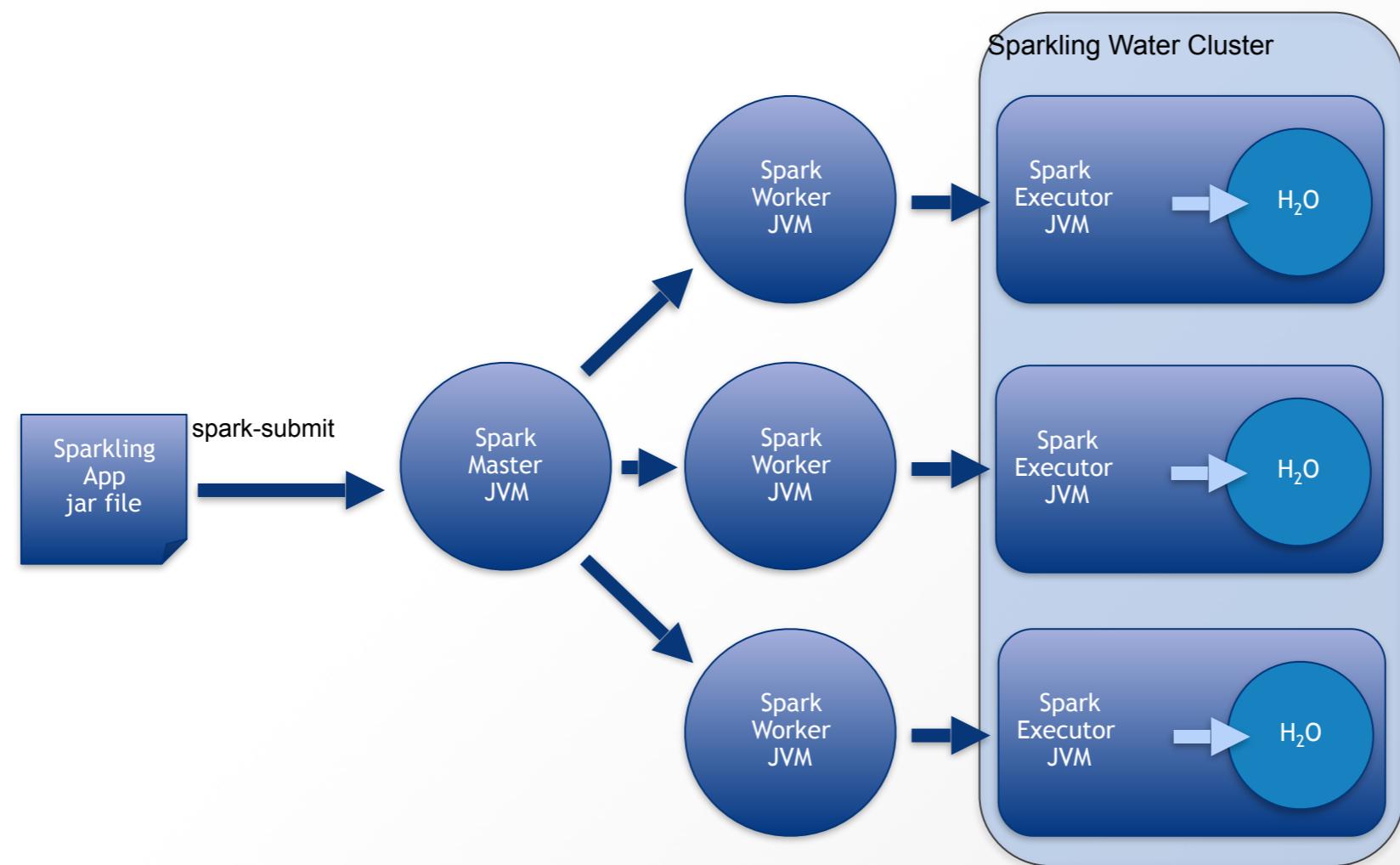
**What is  
inside?**

# Two Backends

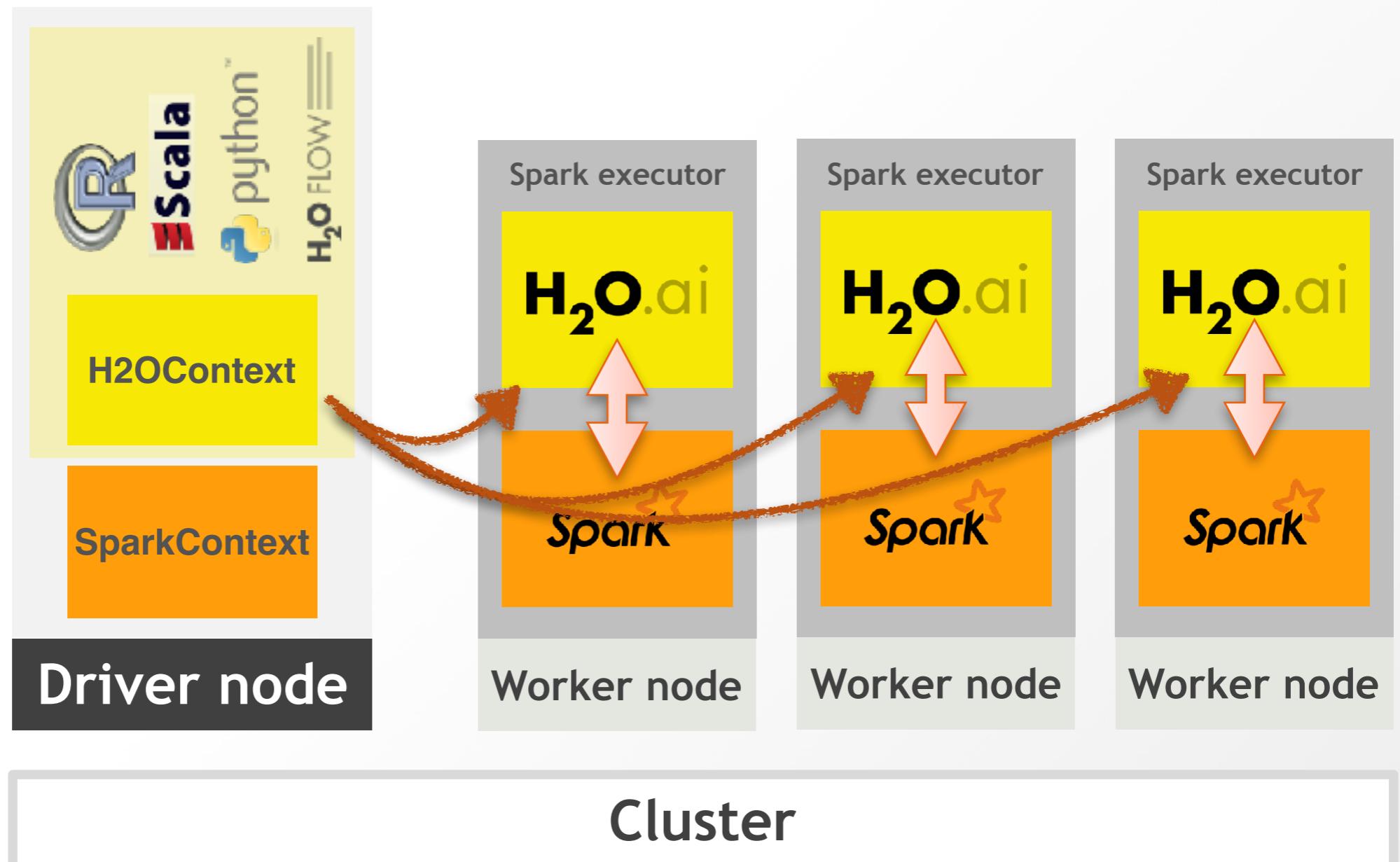
- Sparkling Water has two backends
  - External
  - Internal
- The backend determines where H2O cluster is located
- Each backend is good for different use-cases

# Internal Backend

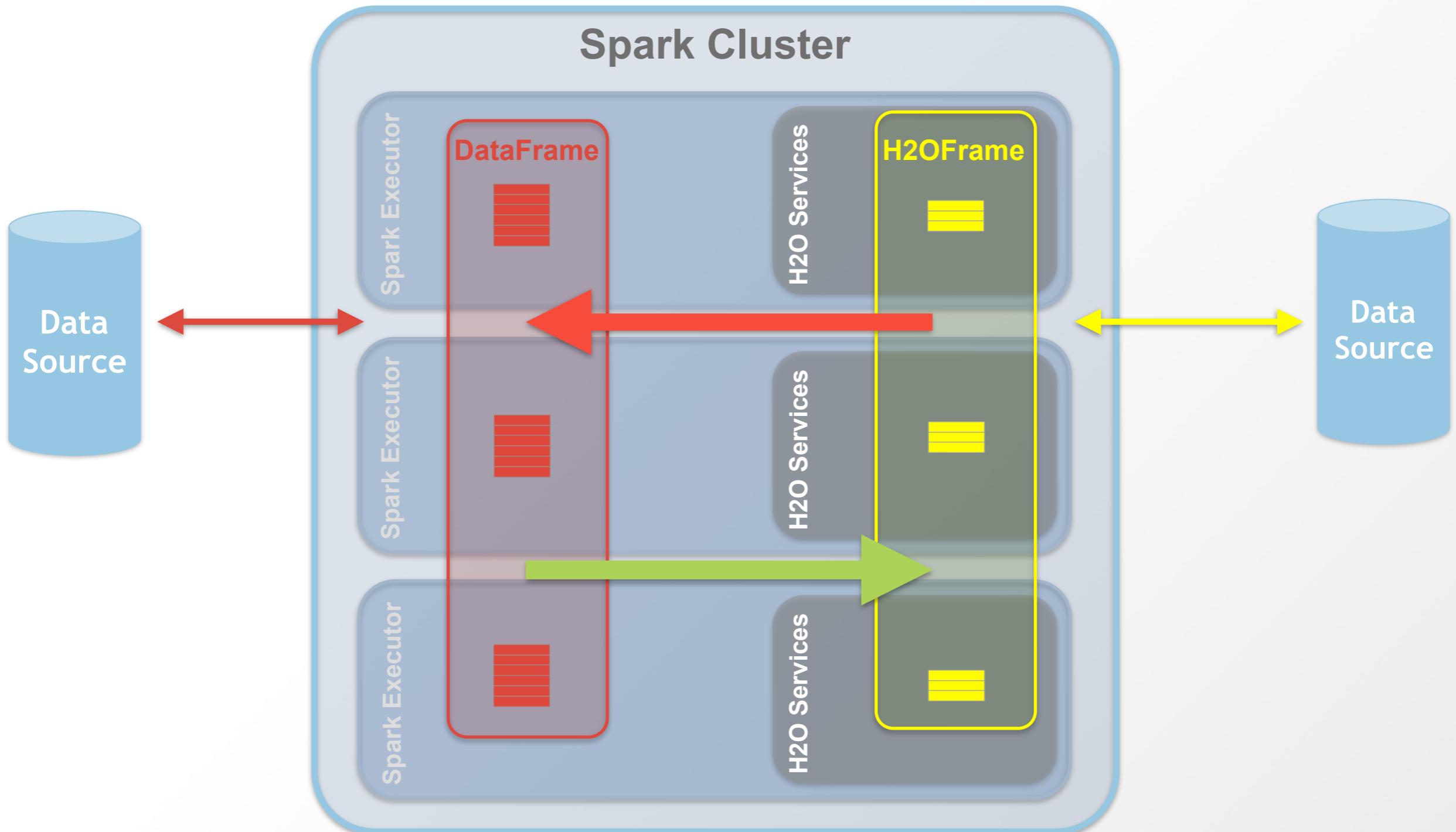
# Sparkling Water Internal Backend



# Internal Backend



# Data Transfers



## **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

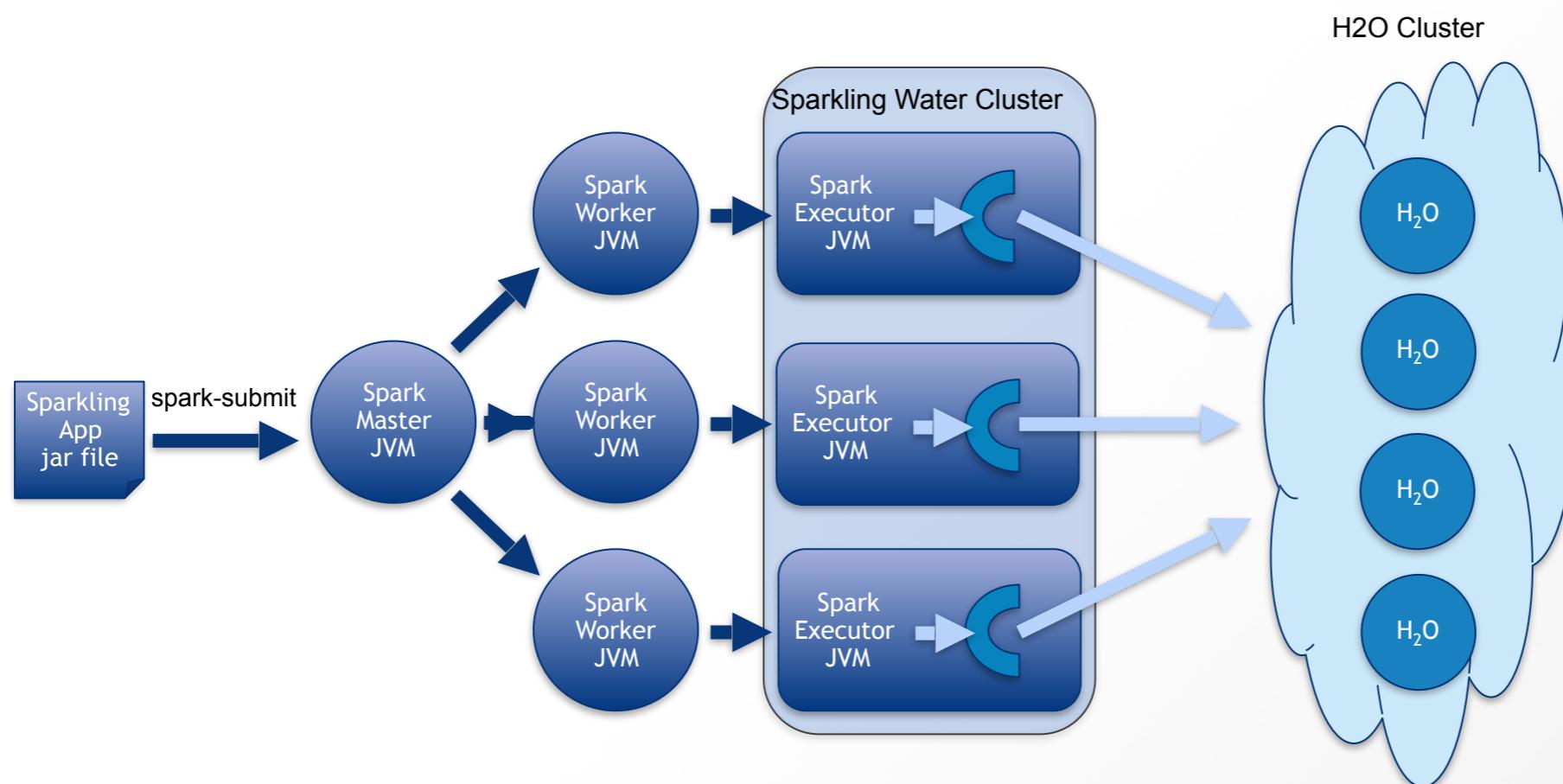
# Overview

- 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

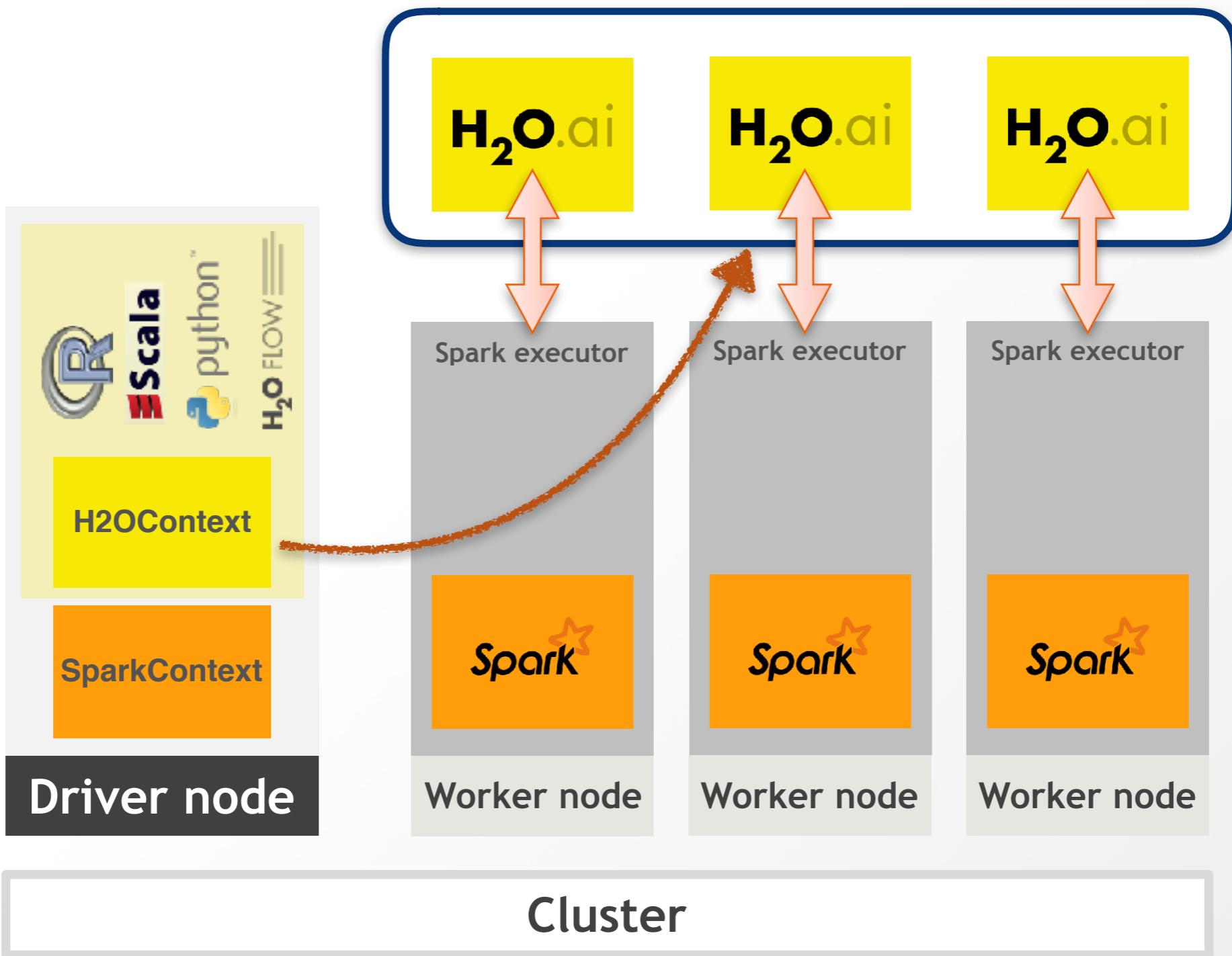
# Separation Approach

- Separating Spark and H2O
  - But preserving same API:  
`val h2oContext = H2OContext.getOrCreate("http://h2o:54321/")`
- Spark and H2O can be submitted as 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

# Demo Time

# Upcoming Events

- Productionizing H2O Models with Apache Spark



# More Info

- Documentation: <http://docs.h2o.ai>
- Tutorials: <https://github.com/h2oai/h2o-tutorials>
- Slidedecks: <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>

# Thank you!

Sparkling Water is  
open-source  
ML application platform  
combining  
power of Spark and H2O

Learn more at [h2o.ai](https://h2o.ai)

Follow us at [@h2oai](https://twitter.com/h2oai)

PS: We are hiring!

