# Force multiplier for data science: Introduction to H2O

**H₂O**.ai

Hank Roark
@hankroark
hank@h2o.ai

# DEMONSTRATION

# WHO AMI I

Lead, Customer Data Science @ H2O.ai

John Deere: Research, Software Product Development, High Tech Ventures

Lots of time dealing with data off of machines, equipment, satellites, weather, radar, hand sampled, and on.

Geospatial, temporal / time series data almost all from sensors.

Previously at startups and consulting (Red Sky Interactive, Nuforia, NetExplorer, Perot Systems, a few of my own)

Engineering & Management MIT
Physics Georgia Tech

hank@h2oai.com
@hankroark
https://www.linkedin.com/in/hankroark

# WHAT IS H2O?

**Math Platform** — Open source in-memory prediction engine

- Parallelized and distributed algorithms making the most use out of multithreaded systems
- GLM, Random Forest, GBM, Deep Learning, etc.

**API** — Easy to use and adopt

- Written in Java – perfect for Java Programmers
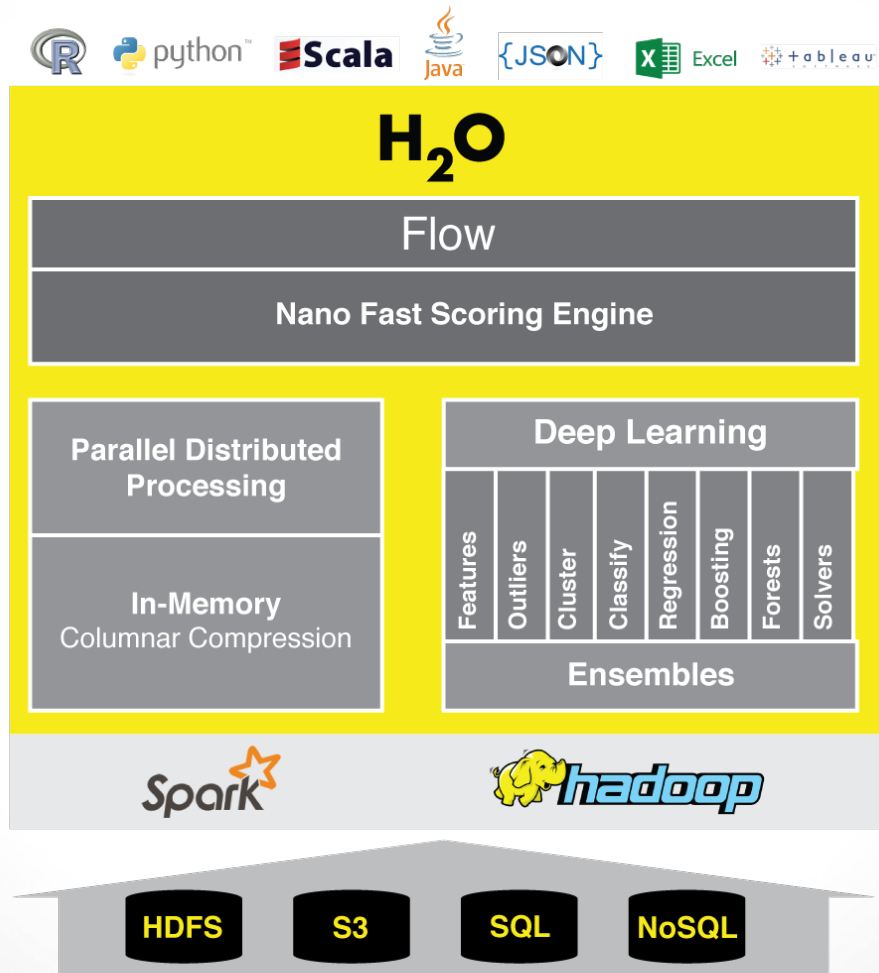- REST API (JSON) – drives H2O from Browser UI, R, Python, Tableau
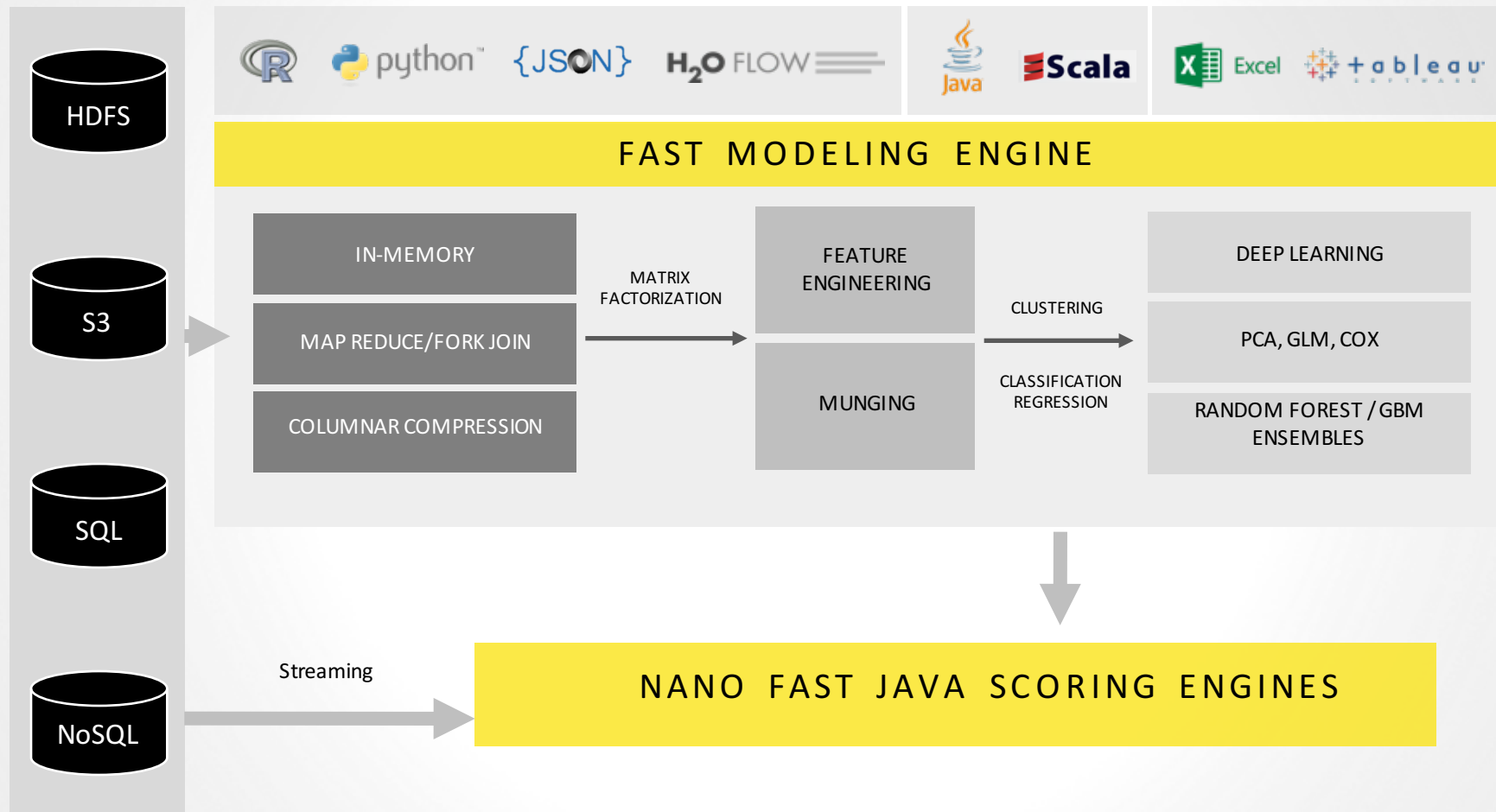
**Big Data** — More data? Or better models? BOTH

- Use all of your data – model without down sampling
- Run a simple GLM or a more complex GBM to find the best fit for the data
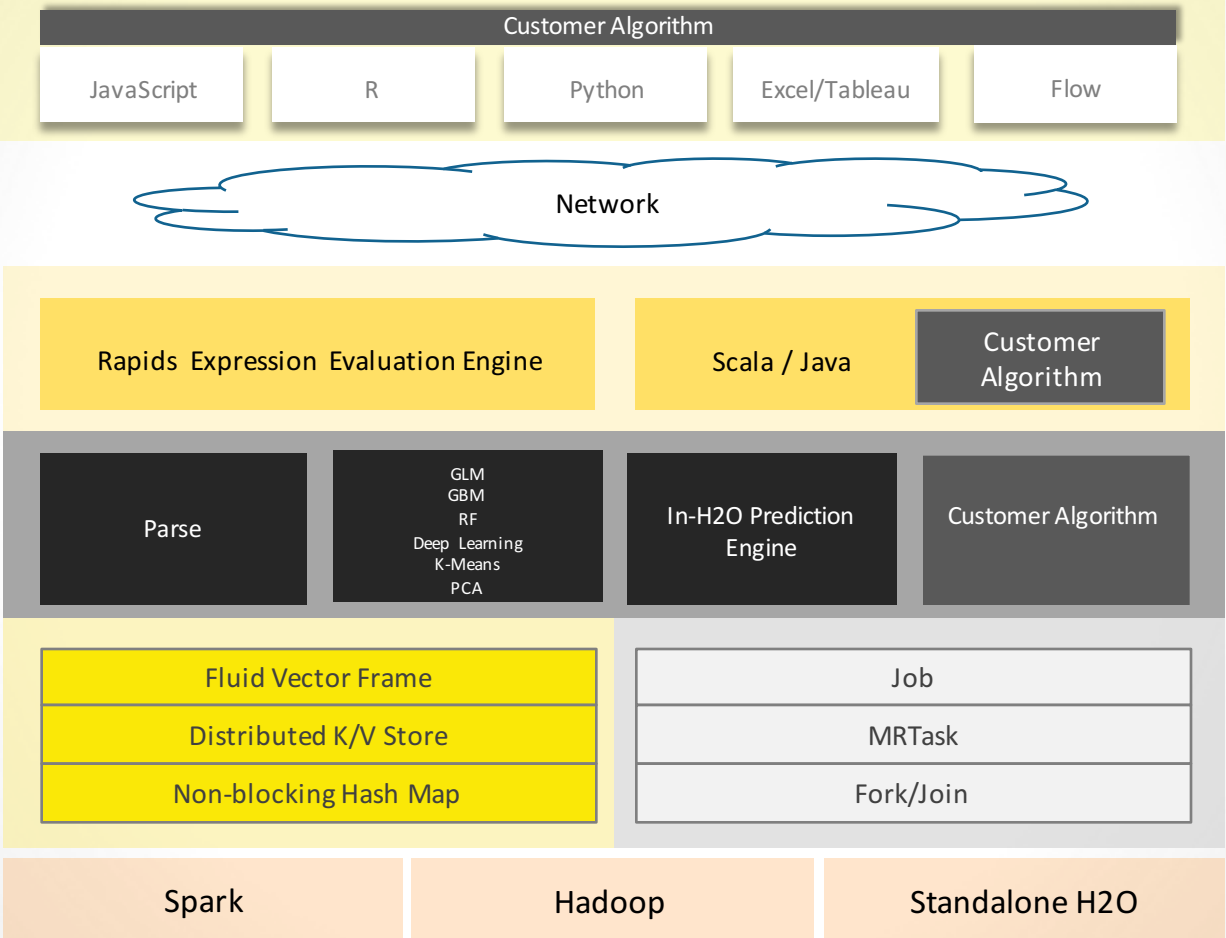- More Data + Better Models = Better Predictions

# ACCURACY WITH SPEED AND SCALE

# ACCURACY WITH SPEED AND SCALE

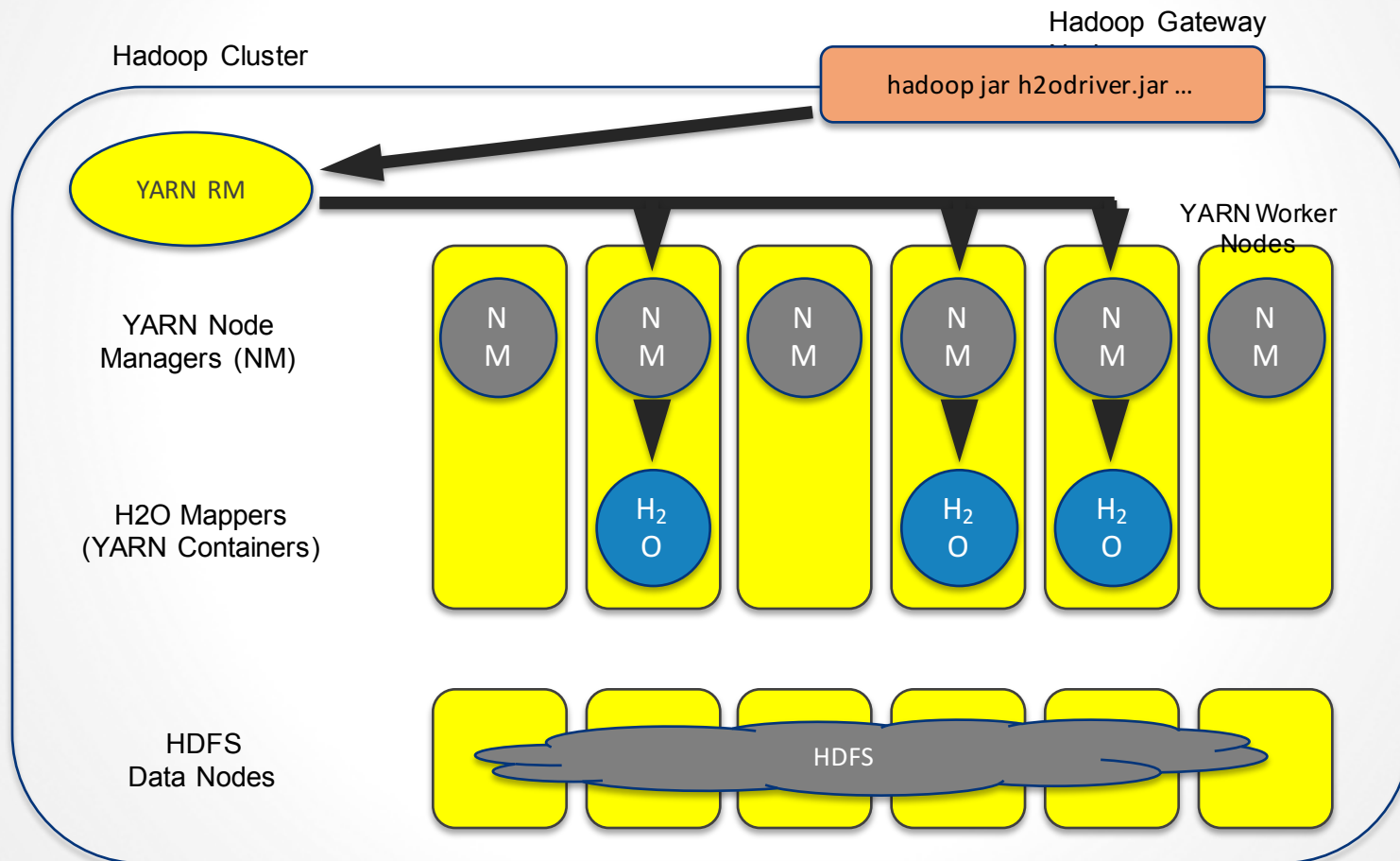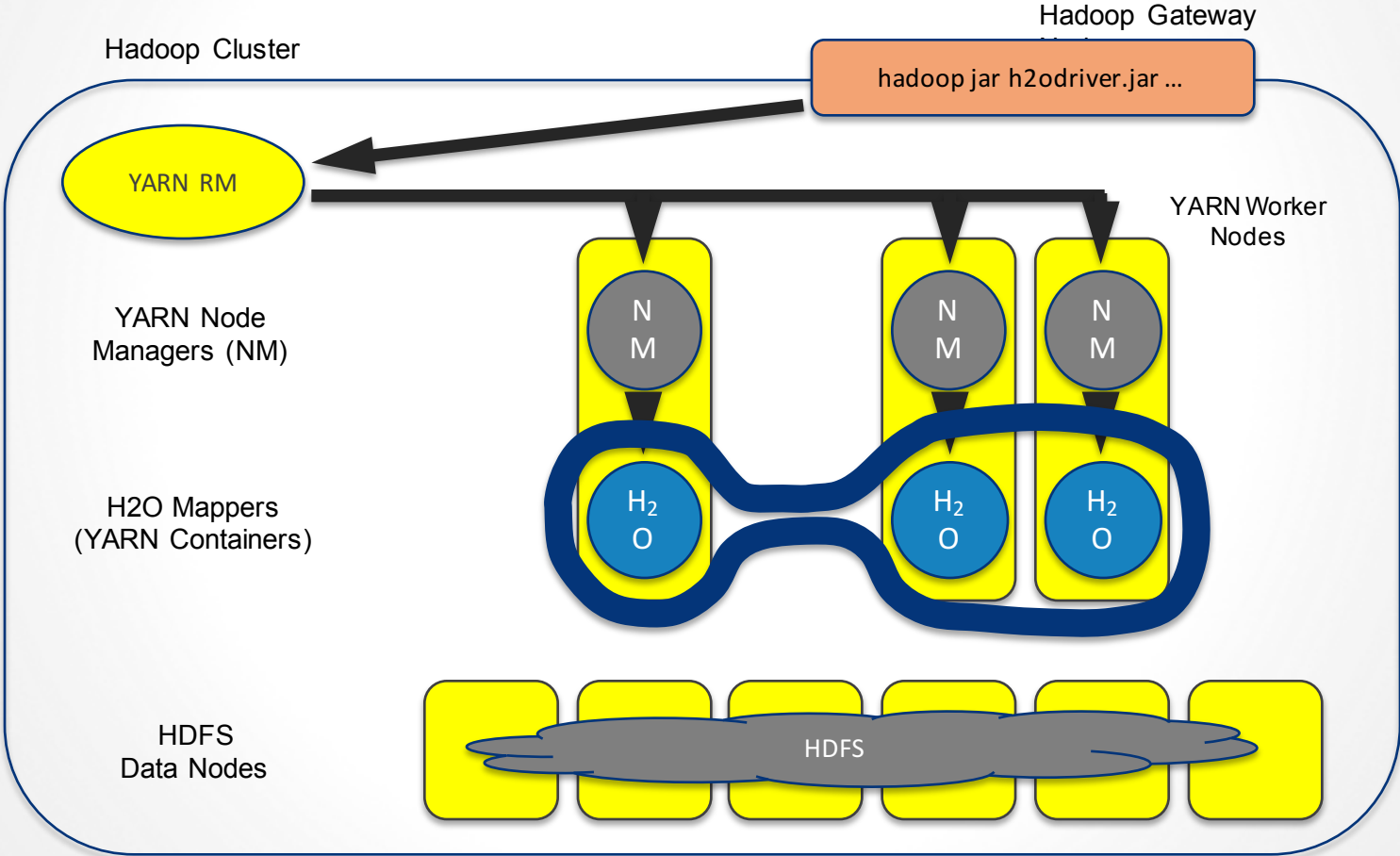| HDFS | R  python™  {JSON}  H₂O FLOW | Java  Scala | Excel  tableau |
|------|------|------|------|

## FAST MODELING ENGINE

| IN-MEMORY | | FEATURE ENGINEERING | | DEEP LEARNING |
|-----------|--|---------------------|--|---------------|
| MAP REDUCE/FORK JOIN | MATRIX FACTORIZATION → | | CLUSTERING → | PCA, GLM, COX |
| COLUMNAR COMPRESSION | | MUNGING | CLASSIFICATION REGRESSION | RANDOM FOREST / GBM ENSEMBLES |

S3

SQL

Streaming

## NANO FAST JAVA SCORING ENGINES

NoSQL

# H2O SOFTWARE STACK

Customer Algorithm

| JavaScript | R | Python | Excel/Tableau | Flow |
|---|---|---|---|---|

Network

| Rapids Expression Evaluation Engine | Scala / Java | Customer Algorithm |
|---|---|---|

| Parse | GLM<br>GBM<br>RF<br>Deep Learning<br>K-Means<br>PCA | In-H2O Prediction Engine | Customer Algorithm |
|---|---|---|---|

| Fluid Vector Frame | Job |
|---|---|
| Distributed K/V Store | MRTask |
| Non-blocking Hash Map | Fork/Join |

| Spark | Hadoop | Standalone H2O |
|---|---|---|

# Hadoop (and YARN)

- You can launch H2O directly on Hadoop:
  *$ hadoop jar h2odriver.jar … –nodes 3 –mapperXmx 50g*

- H2O uses Hadoop MapReduce to get CPU and Memory on the cluster, not to manage work
  - o  H2O mappers stay at 0% progress forever
    - Until you shut down the H2O job yourself
  - o  All mappers (3 in this case) must be running at the same time
  - o  The mappers communicate with each other
    - Form an H2O cluster on-the-spot within your Hadoop environment
  - o  No Hadoop reducers(!)

- Special YARN memory settings for large mappers
  - o  yarn.nodemanager.resource.memory-mb
  - o  yarn.scheduler.maximum-allocation-mb

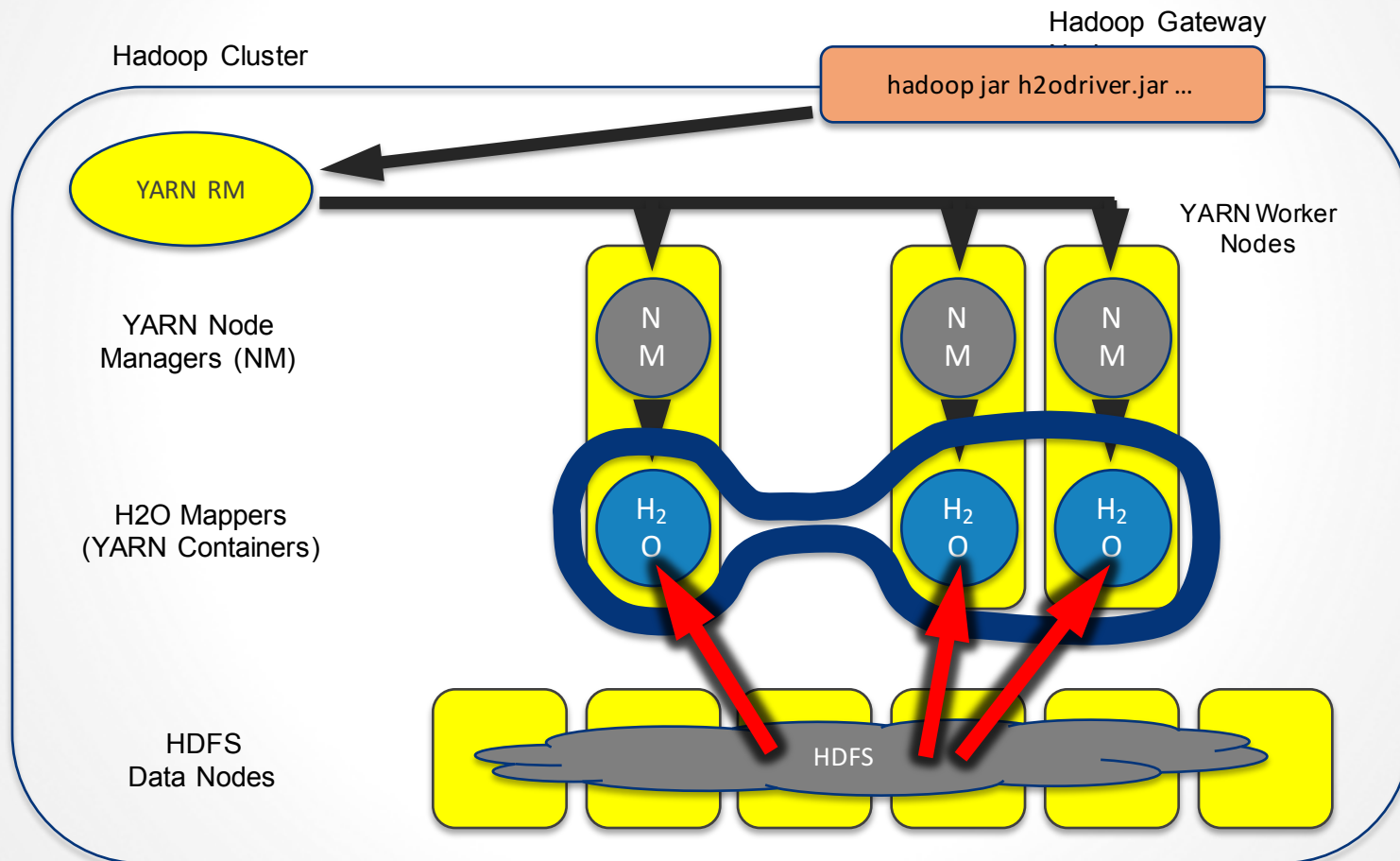- CPU resources controlled via –nthreads h2o command line argument

# H2O on YARN Deployment
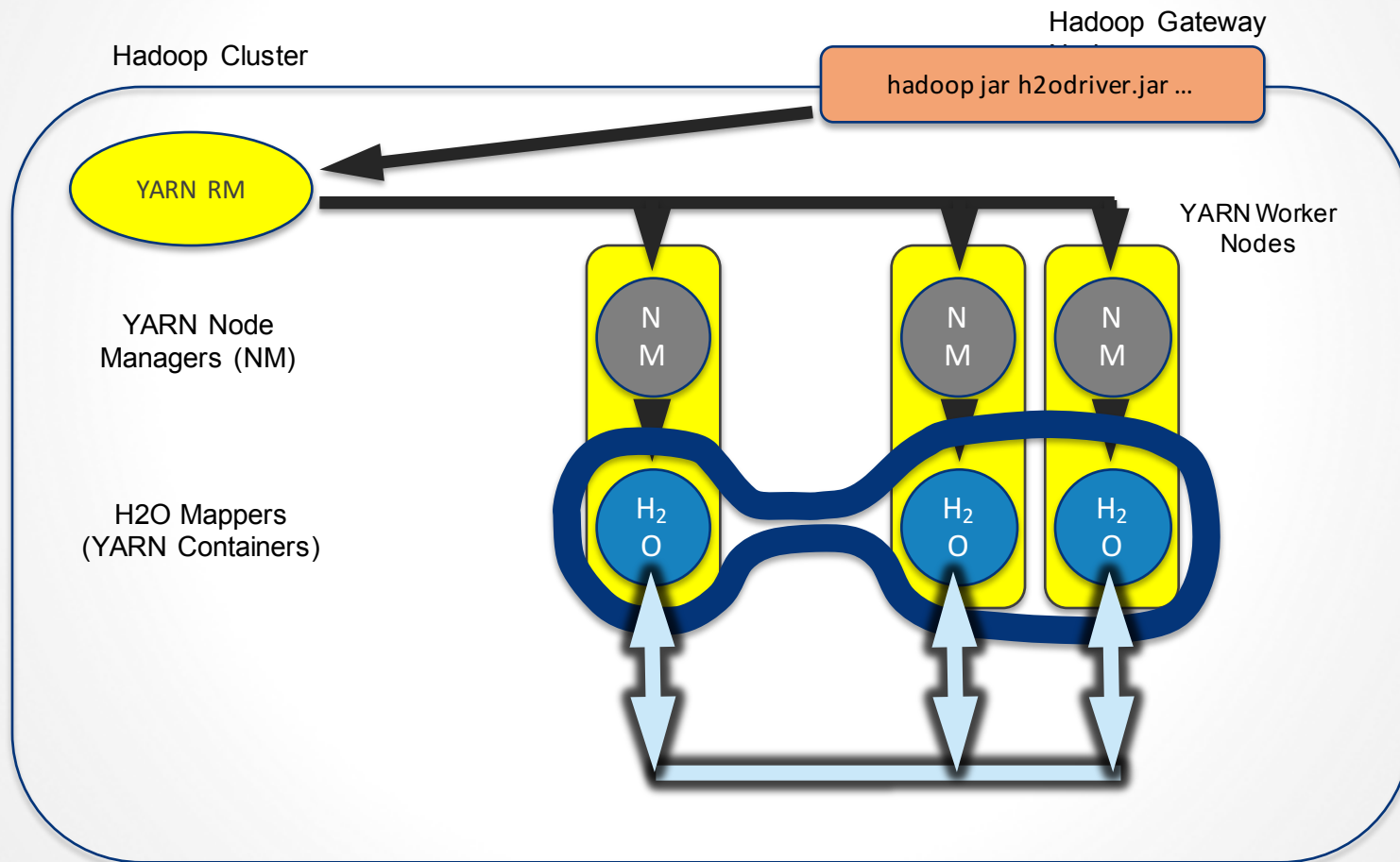
# Now You Have an H2O Cluster

# Read Data from HDFS *Once*

# PYTHON AND R OBJECTS ARE PROXIES FOR BIG DATA

## STEP 1



h2o_df = h2o.import_file("hdfs://path/to/data.csv")

Python user

# PYTHON AND R OBJECTS ARE PROXIES FOR BIG DATA

**STEP 2**

**2.3**
Initiate distributed ingest

**H2O Cluster**

**2.2**

Python

HTTP REST API request to $H_2O$ has HDFS path

h2o.import_file()

$H_2O$

$H_2O$

$H_2O$

**2.1**

Python function call

**2.4**

HDFS

data.csv

Request data from HDFS

# PYTHON AND R OBJECTS ARE PROXIES FOR BIG DATA

**STEP 3**

**3.2**
Distributed $H_2O$ Frame in DKV

**H2O Cluster**

Python

**3.3**

h2o_df

Return pointer to data in REST API JSON Response

Cluster IP
Cluster Port
Pointer to Data

**3.4**

h2o_df object created in Python

$H_2O$
$H_2O$
$H_2O$
Frame
$H_2O$

**3.1**

HDFS

data.csv

HDFS provides data

# PYTHON SCRIPT STARTING H2O GLM



| | |
|---|---|
| TCP/IP | TCP/IP |
| HTTP | HTTP |
| REST/JSON | REST/JSON |
| GET /3/Models/model_id | /3/Models endpoint |
| glm_estimator.show() | |
| Python script | Fork/Join framework · K/V store framework |
| Standard Python process | H2O process |

**Legend**

| |
|---|
| Network layer |
| REST layer |
| H2O - algos |
| H2O - core |
| User process |
| H2O process |

# H$_2$O on Storm

# H$_2$O – The Killer-App for Spark

## H2O Sparkling Water

| | | |
|---|---|---|
| MLlib | H$_2$O | SQL |

**H$_2$ORDD**

**HDFS=DATA**

| | |
|---|---|
| In-Memory | Big Data, Columnar |
| ML | 100x faster Algos |
| R | CRAN, API, fast engine |
| API | Spark API, Java MM |
| Community | Devs, Data Science |

**H$_2$O.ai**
Machine Intelligence
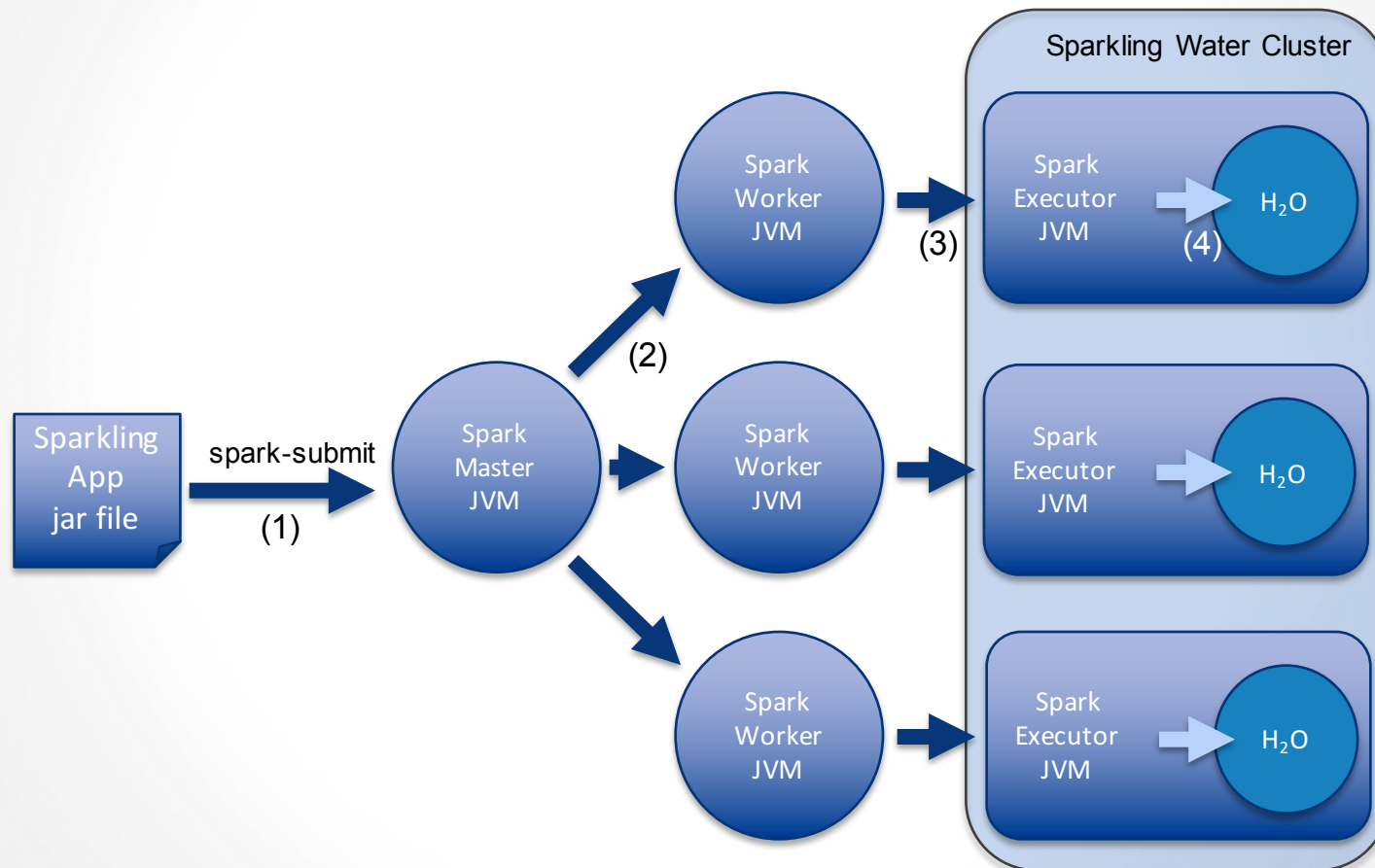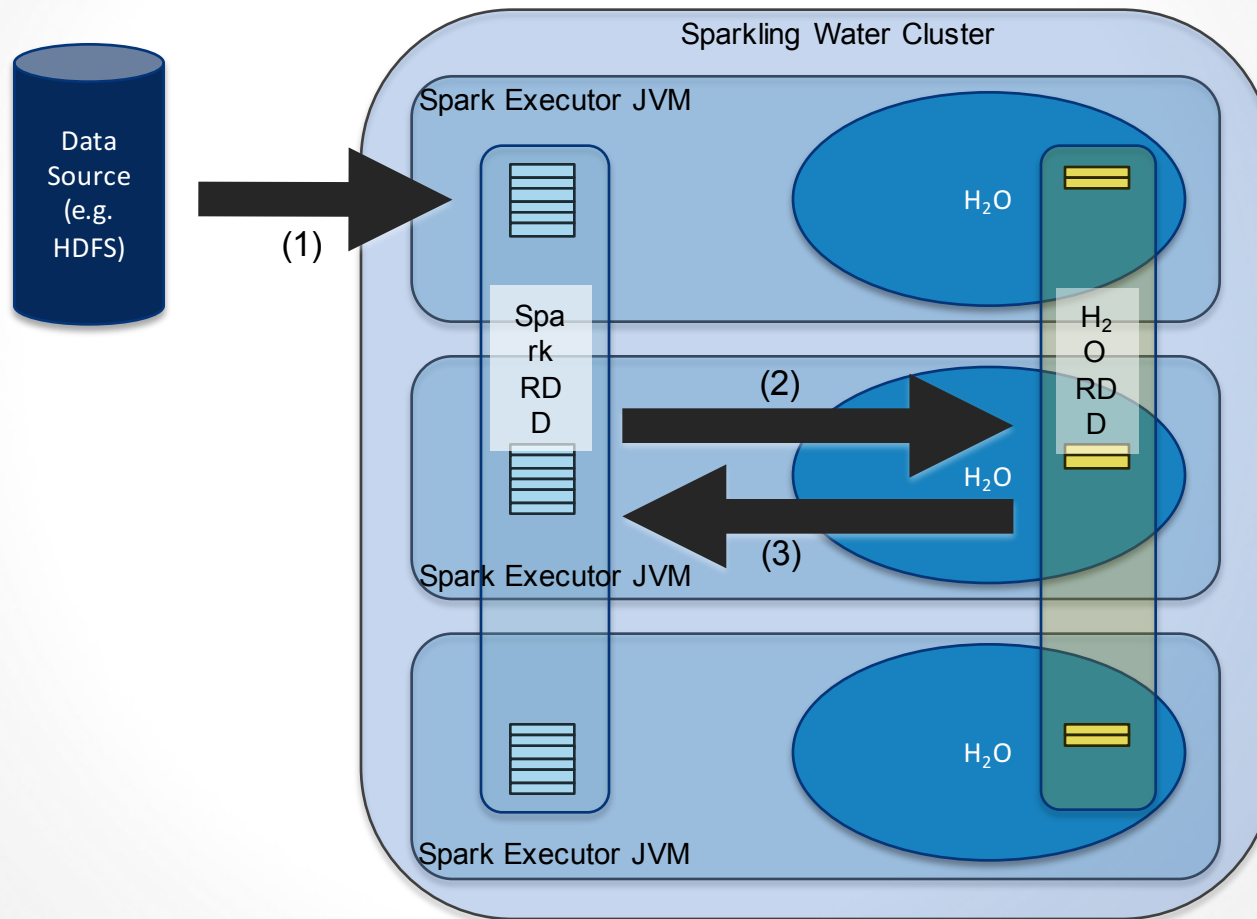
# Spark (and Sparkling Water)

- H2O runs as an application on a Spark cluster using spark-submit
  - Standard Spark 1.3+
  - Includes H2O on Spark on YARN
- H2O and Spark nodes share a JVM process
- H2ORDD facilitates easy data sharing between Spark (e.g. Spark SQL, MLlib) and H2O (e.g. Deep Learning)
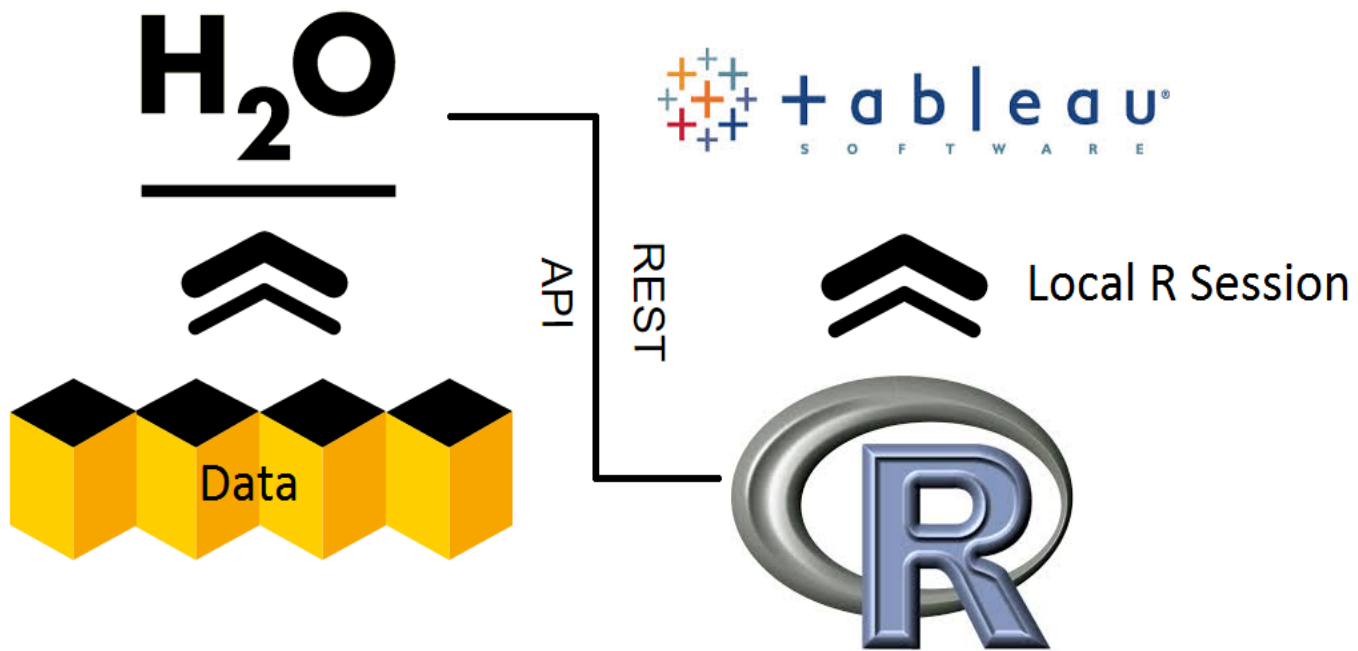- Scala & PySpark support

# Sparkling Water Data Distribution

# RESOURCES

- Download and go: http://www.h2o.ai/download
- Documentation: http://docs.h2o.ai/
- Booklets, Datasheet: http://www.h2o.ai/resources/
- Github: http://github.com/h2oai/
- Training: http://learn.h2o.ai/

# THANK YOU