# H₂O.ai

Scalable Machine Learning
For Smarter Applications

Hank Roark
Data Scientist / Hacker
hank@h2o.ai
@hankroark
https://www.linkedin.com/in/hankroark

# H2O.ai Overview

- Founded: 2011 venture-backed, debuted in 2012
- Product: H2O open source in-memory prediction engine
- Team: 34
- HQ: Mountain View, CA
- SriSatish Ambati – CEO & Co-founder (Founder Platfora, DataStax; Azul)
- Cliff Click – CTO & Co-founder (Creator Hotspot, Azul, Sun, Motorola, HP)
- Tom Kraljevic – VP of Engineering (CTO & Founder Luminix, Azul, Chromatic)

H2O.ai
Machine Intelligence

**Distributed
Systems
Engineers
Making
ML Scale!**

H₂O.ai
Machine Intelligence

# Scientific Advisory Council

**Stephen Boyd**
Professor of EE Engineering
Stanford University

**Rob Tibshirani**
Professor of Health Research
and Policy, and Statistics
Stanford University

**Trevor Hastie**
Professor of Statistics
Stanford University

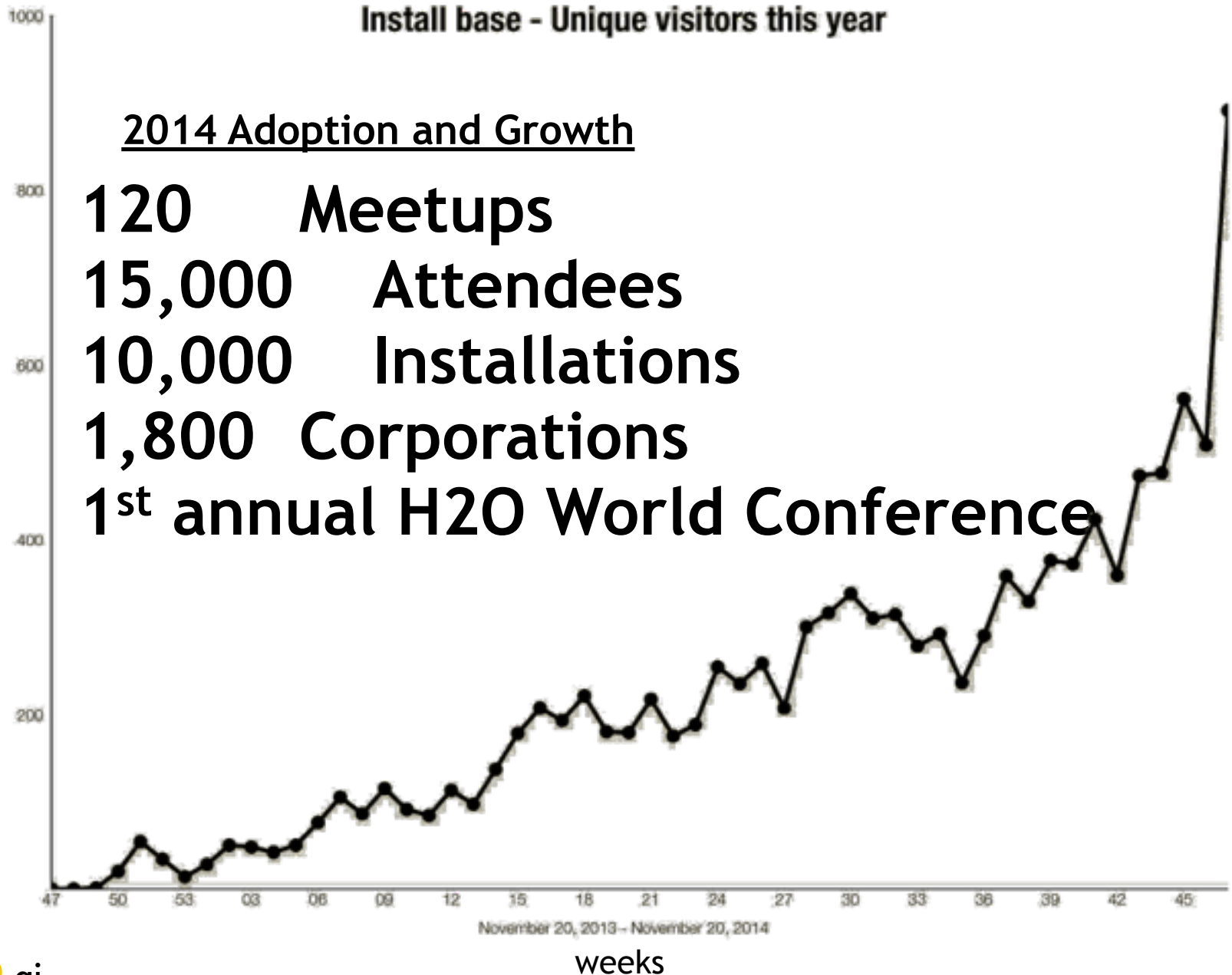# Install base - Unique visitors this year

<u>2014 Adoption and Growth</u>

**120        Meetups**
**15,000     Attendees**
**10,000     Installations**
**1,800   Corporations**
**1<sup>st</sup> annual H2O World Conference**

November 20, 2013 – November 20, 2014

weeks

H2O.ai
Machine Intelligence

# 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, PCA, etc.

**API**

- Easy to use and adopt

- Written in Java – perfect for Java Programmers

- REST API (JSON) – drives H2O from R, Python, Excel, 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

$H_2O$.ai
Machine Intelligence

Python JSON R Scala Java Tableau Excel

# H$_2$O Prediction Engine

## SDK / API

| Rapids Query R-engine | Nano Fast Scoring Engine |

**In-Mem Map Reduce**
Distributed fork/join

**Memory Manager**
Columnar Compression

### Deep Learning

| Cluster | Classify | Regression | Trees | Boosting | Forests | Solvers | Gradients |

**Ensembles**

On Premise
On Hadoop & Spark
On EC2

Per Node
2M      Row ingest/sec
50M     Row Regression/sec
750M    Row Aggregates / sec

HDFS    S3    SQL    NoSQL

H$_2$O.ai
Machine Intelligence

# Algorithms on H2O

*Supervised Learning*

**Statistical Analysis**

- **Generalized Linear Models**: Binomial, Gaussian, Gamma, Poisson and Tweedie
- **Cox Proportional Hazards Models**
- **Naïve Bayes**

**Ensembles**

- **Distributed Random Forest**: Classification or regression models
- **Gradient Boosting Machine**: Produces an ensemble of decision trees with increasing refined approximations

**Deep Neural Networks**

- **Deep learning**: Create multi-layer feed forward neural networks starting with an input layer followed by multiple layers of nonlinear transformations

# Algorithms on H$_2$O

*Unsupervised Learning*

| Clustering |
|---|

- **K-means**: Partitions observations into k clusters/groups of the same spatial size
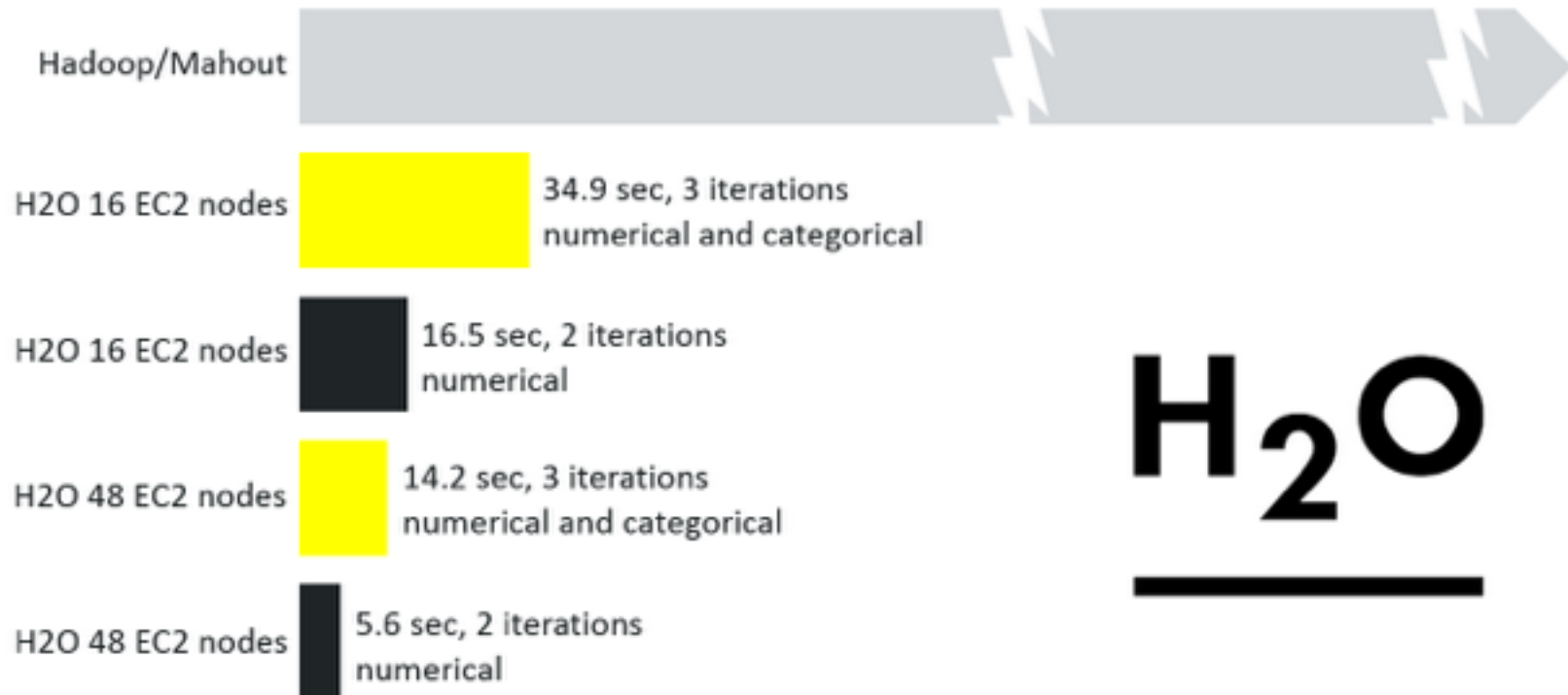
| Dimensionality Reduction |
|---|

- **Principal Component Analysis**: Linearly transforms correlated variables to independent components

| Anomaly Detection |
|---|

- **Autoencoders**: Find outliers using a nonlinear dimensionality reduction using deep learning

# H2O Billion Row Machine Learning Benchmark
## GLM Logistic Regression



| | |
|---|---|
| Hadoop/Mahout | |
| H2O 16 EC2 nodes | 34.9 sec, 3 iterations, numerical and categorical |
| H2O 16 EC2 nodes | 16.5 sec, 2 iterations, numerical |
| H2O 48 EC2 nodes | 14.2 sec, 3 iterations, numerical and categorical |
| H2O 48 EC2 nodes | 5.6 sec, 2 iterations, numerical |

H2O

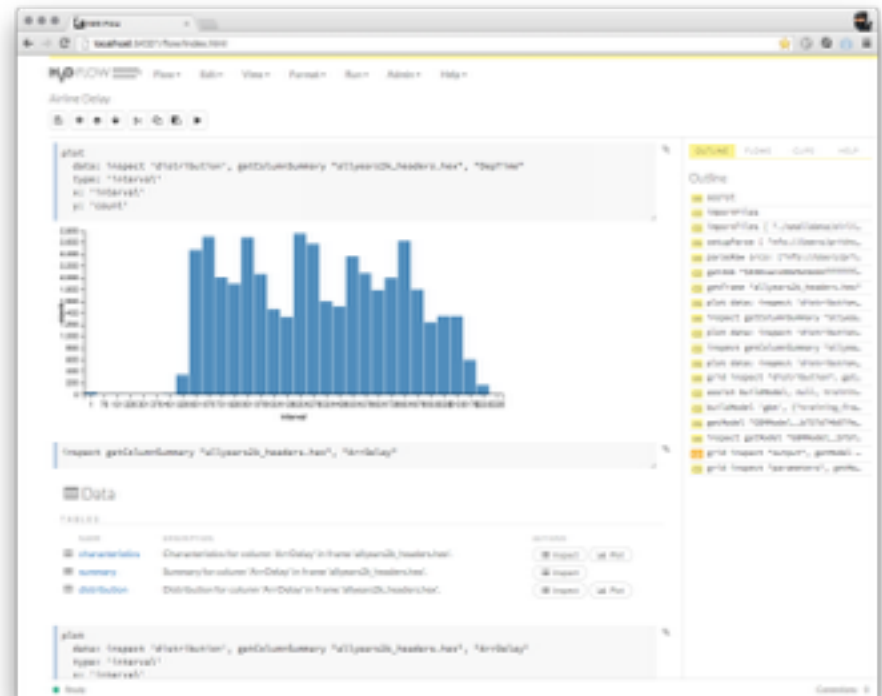Compute Hardware: AWS EC2 c3.2xlarge - 8 cores and 15 GB per node, 1 GbE interconnect

Airline Dataset 1987-2013, 42 GB CSV, 1 billion rows, 12 input columns, 1 outcome column
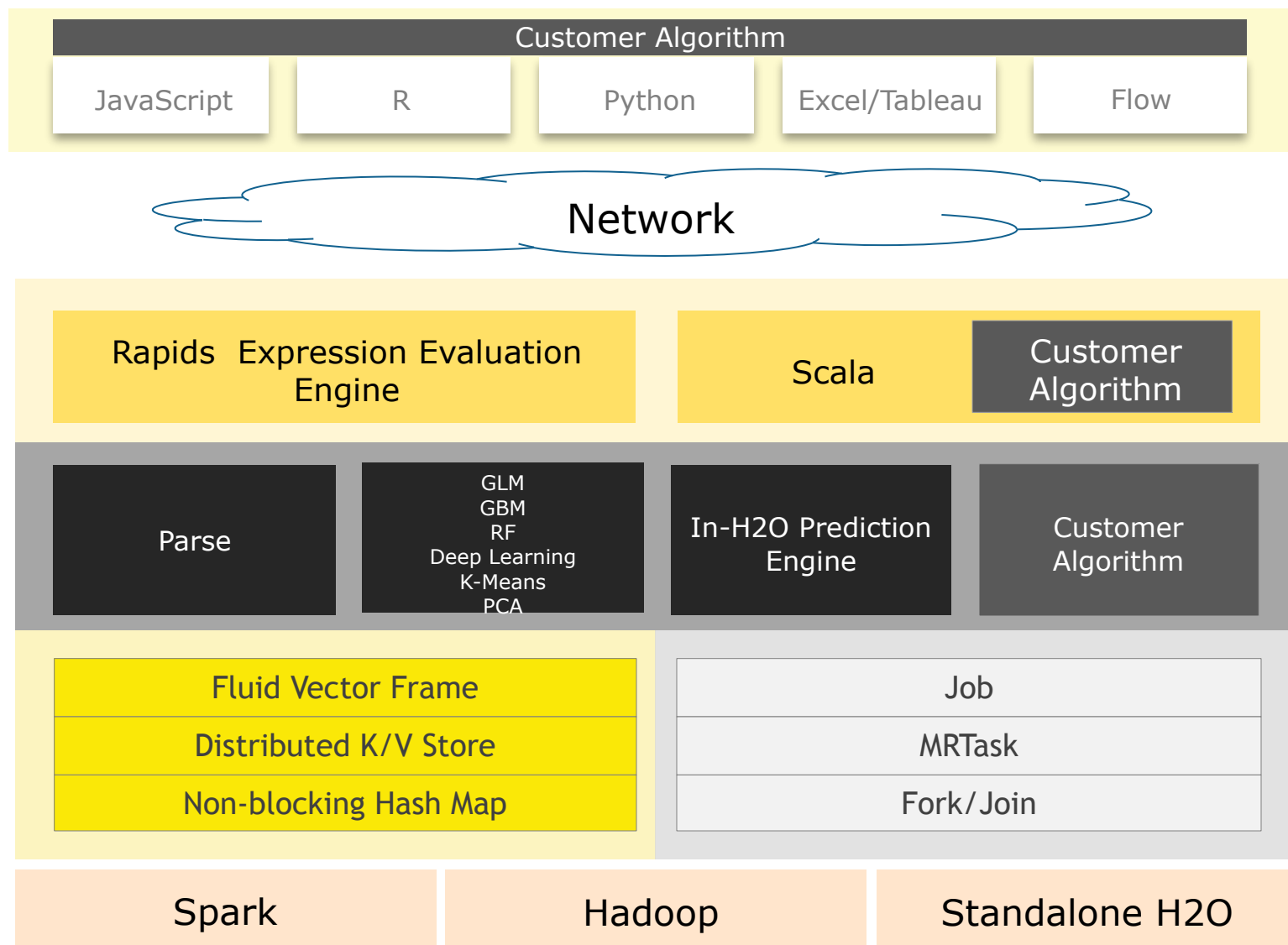9 numerical features, 3 categorical features with cardinalities 30, 376 and 380

H2O.ai
Machine Intelligence

# H2O Flow

- A Web-based interactive computing environment for Big Data Machine Learning

- New web interface of H2O

- Model comparisons

- Mixed environment for
  - Coffeescript
  - Text & Markdown
  - Charts & Visualization (more to come)
  - R/Spark/Python code (coming soon)
  - Mathematics Equations (coming soon)
  - Video & Rich Media

# H2O Software Stack

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

**Network**

| Rapids Expression Evaluation Engine | Scala | 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 |
|---|---|---|

H₂O.ai
Machine Intelligence

# Reading Data from HDFS into H2O with R

**STEP 1**



h2o_df = h2o.importFile("hdfs://path/to/data.cs

R user
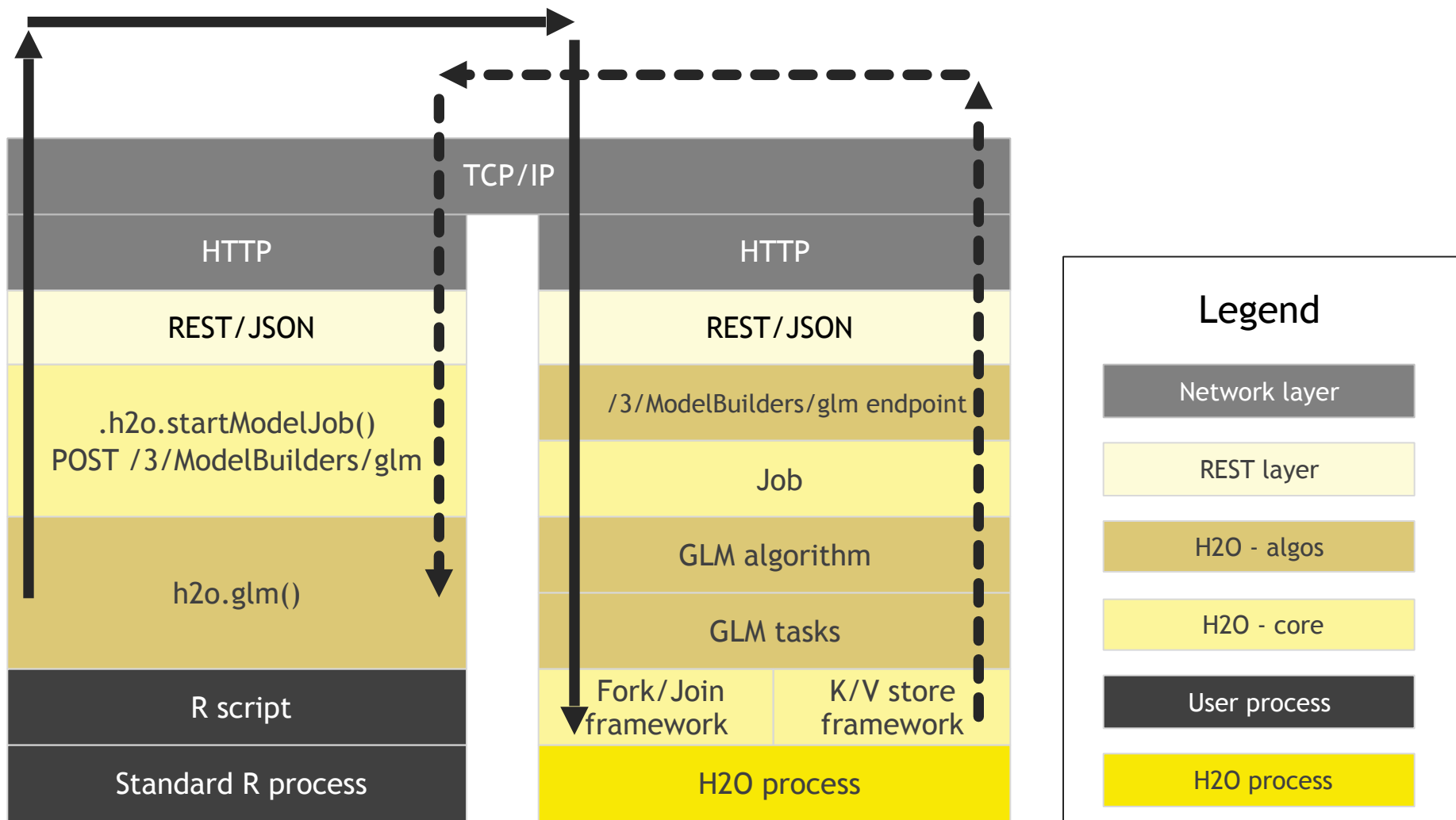
# Reading Data from HDFS into H2O with R



**STEP 2**

**R**

h2o.importFile()

**2.1**

R function call

**2.2**

HTTP REST API
request to $H_2O$
has HDFS path

**2.3**

Initiate distributed
ingest

**H2O Cluster**

$H_2O$

$H_2O$

$H_2O$

**2.4**

Request data
from HDFS

**HDFS**

data.csv

# Reading Data from HDFS into H2O with R

# R Script Starting H2O GLM

# R Script Retrieving H2O GLM Result

# H$_2$O – 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* |

# H₂O.ai

Scalable Machine Learning
For Smarter Applications