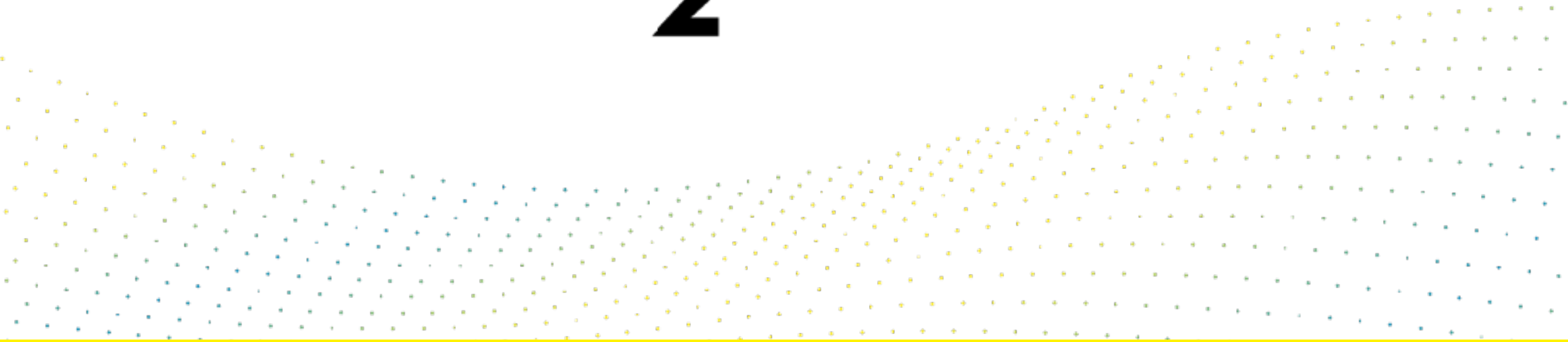


Introduction to Scalable & Automatic Machine Learning with H₂O

H₂O.ai



Agenda

- Talk 1: Introduction to H₂O
 - Company and People
 - H₂O Open Source ML Platform
 - Live Demos
- Talk 2: The Next Generation of Machine Learning on Apache Spark
 - Introduction to Sparkling Water
 - Live Demos
- Talk 3: Introduction to Driverless AI
 - Introduction to Driverless AI



About Me

- Graduated from Charles University in Prague (June 2017)
- Bachelor degree from Comenius University in Bratislava, Faculty of Mathematics, Physics and Informatics

- **Software Engineer**

- May 2015 - June 2017
 - Ruby on Rails and Android developer in min60 s.r.o.
- June 2014 - February 2015
 - Ruby on Rails developer in Staffino s.r.o.
- June 2017 - Present
 - H₂O.ai (working remotely from Bratislava)

Company Overview

Founded	2012, Series C in Nov, 2017
Products	<ul style="list-style-type: none">• Driverless AI – Automated Machine Learning• H₂O - Open Source Machine Learning Platform• H2O4GPU - Lightning Fast Machine Learning on GPUs• Sparkling Water - Integration of H₂O and Apache Spark
Mission	Democratize AI. Do Good.
Team	<p>~100 employees</p> <ul style="list-style-type: none">• Distributed Systems Engineers doing Machine Learning• World-class visualization designers
Offices	Mountain View, London, Prague



Our Mission:

Make Machine Learning Accessible to Everyone



Complexity is your enemy. Any fool
can make something complicated. It
is hard to keep things simple.

— *Richard Branson* —

AZ QUOTES

Scientific Advisory Council



Dr. Trevor Hastie

- John A. Overdeck Professor of Mathematics, Stanford University
- PhD in Statistics, Stanford University
- Co-author, *The Elements of Statistical Learning: Prediction, Inference and Data Mining*
- Co-author with John Chambers, *Statistical Models in S*
- Co-author, *Generalized Additive Models*



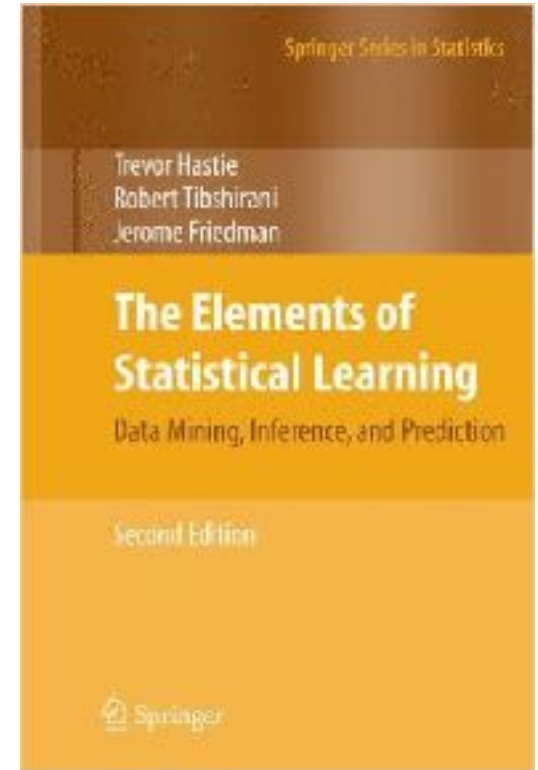
Dr. Robert Tibshirani

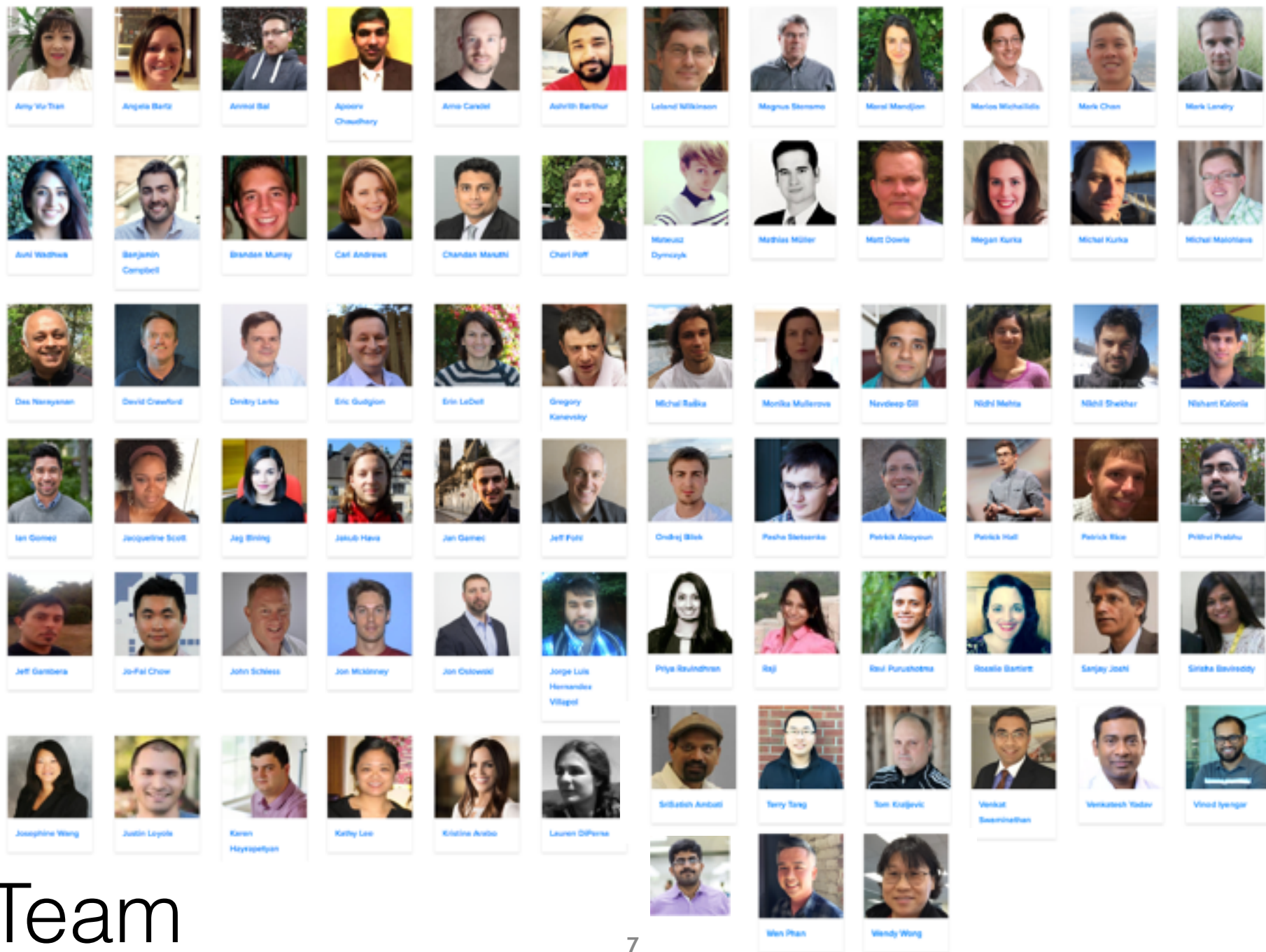
- Professor of Statistics and Health Research and Policy, Stanford University
- PhD in Statistics, Stanford University
- Co-author, *The Elements of Statistical Learning: Prediction, Inference and Data Mining*
- Author, *Regression Shrinkage and Selection via the Lasso*
- Co-author, *An Introduction to the Bootstrap*

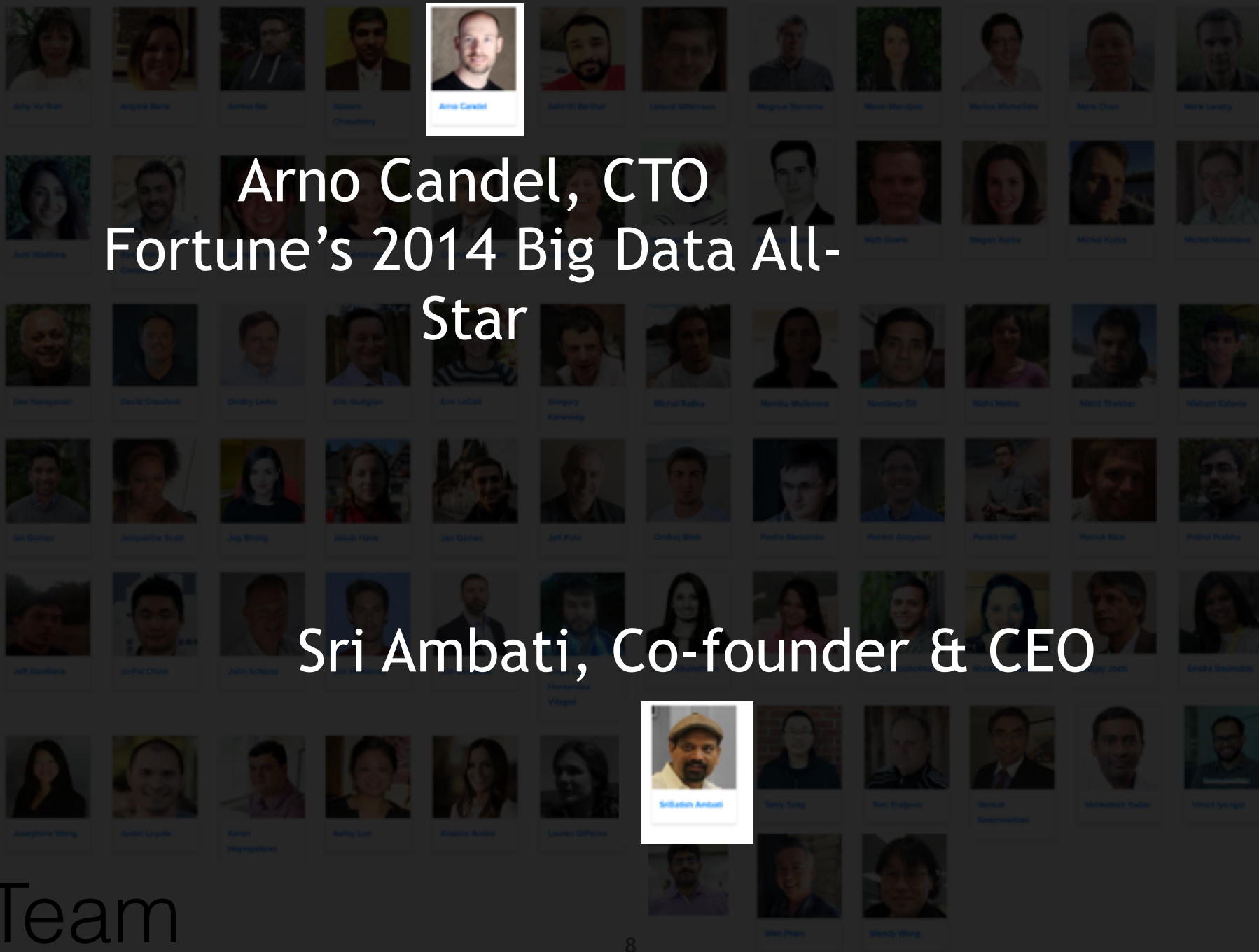


Dr. Steven Boyd

- Professor of Electrical Engineering and Computer Science, Stanford University
- PhD in Electrical Engineering and Computer Science, UC Berkeley
- Co-author, *Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers*
- Co-author, *Linear Matrix Inequalities in System and Control Theory*
- Co-author, *Convex Optimization*



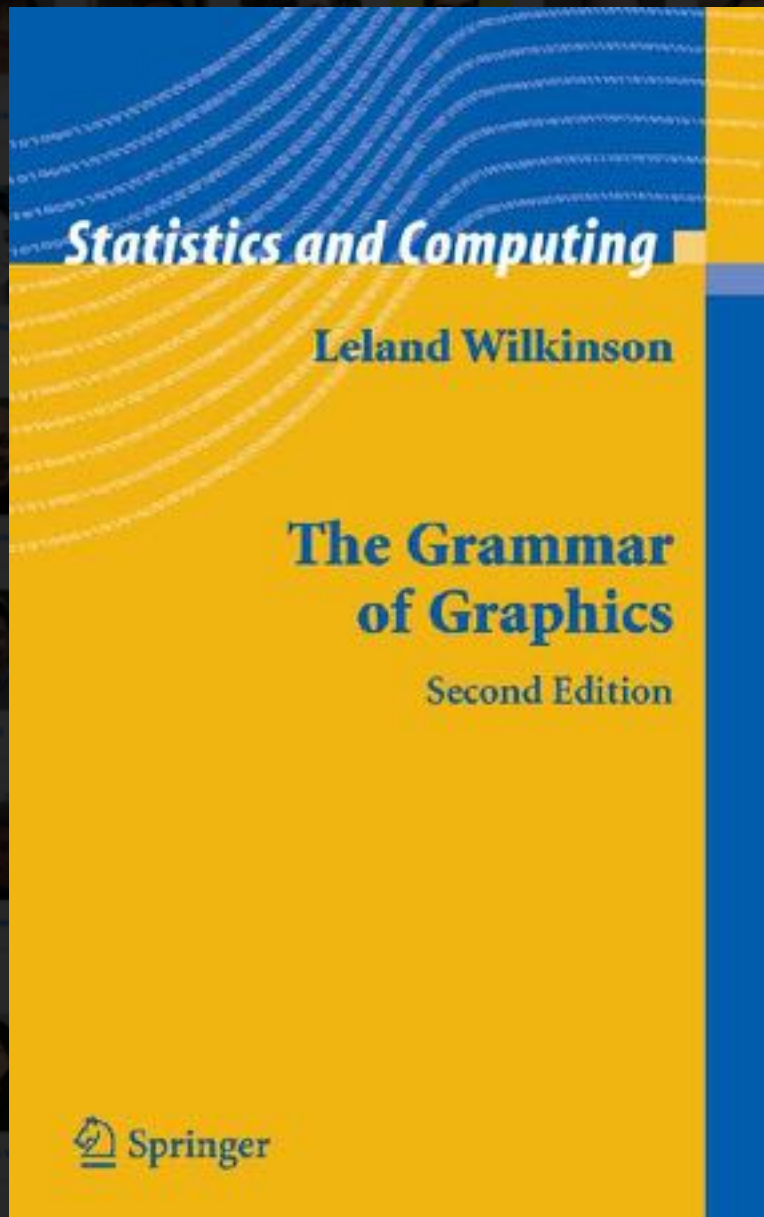




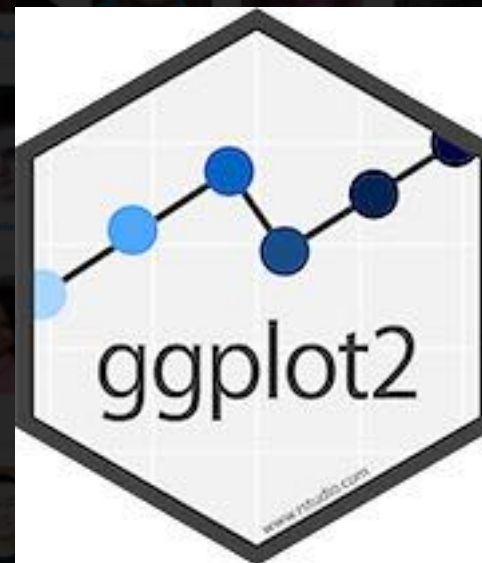
Arno Candel, CTO
Fortune's 2014 Big Data All-
Star

Sri Ambati, Co-founder & CEO

H₂O Team



Origin of R Package 'ggplot2'

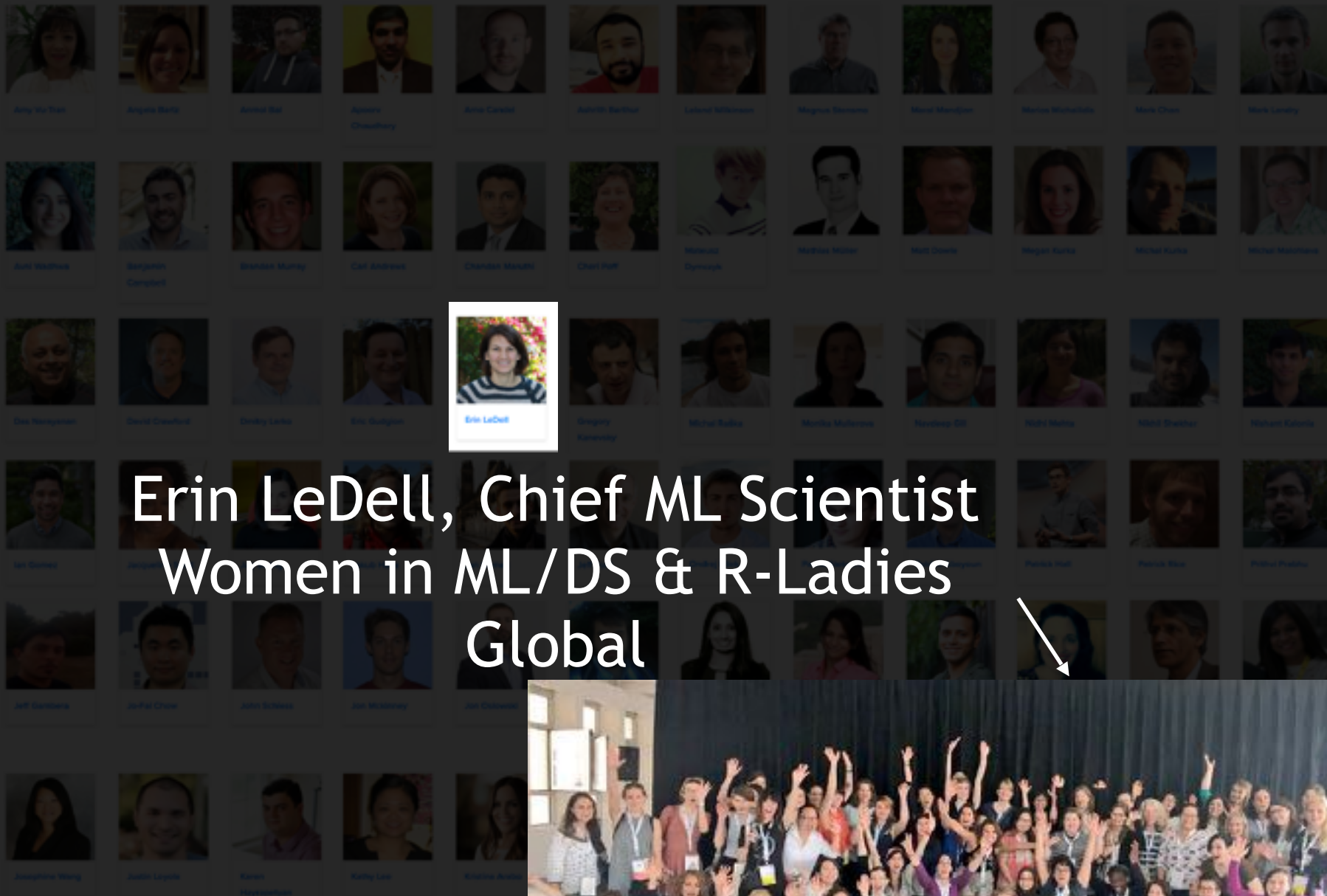




Matt Dowle



H₂O Team



Erin LeDell, Chief ML Scientist Women in ML/DS & R-Ladies Global

H₂O Team



H₂O.ai



48th

1st

33rd

4th

25th

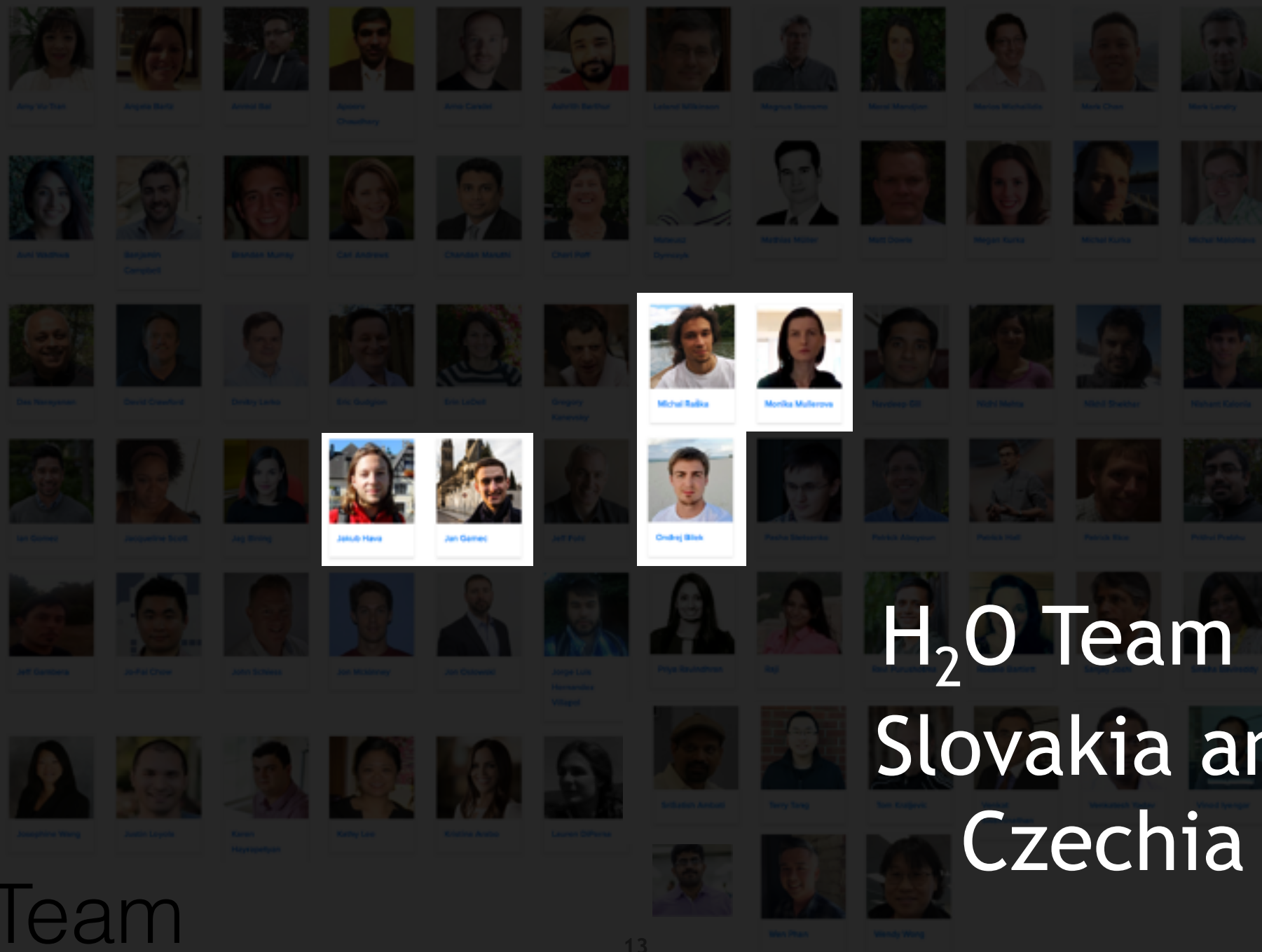
Kaggle Grand Masters

Their Highest Rank in
Kaggle
(about 80,000
competitors)

13th

H₂O Team

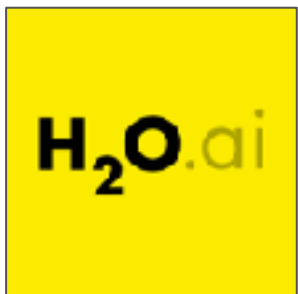
H₂O.ai



H₂O Team in Slovakia and Czechia

H₂O Team

H₂O 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

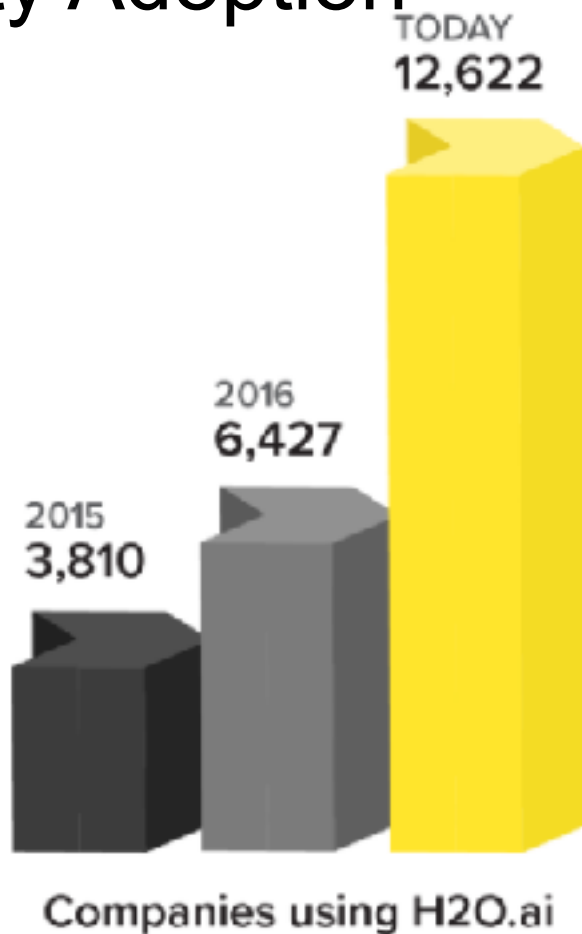


Automatic feature
engineering, machine learning
and interpretability

Steam

Secure multi-tenant H2O clusters

Worldwide Community Adoption

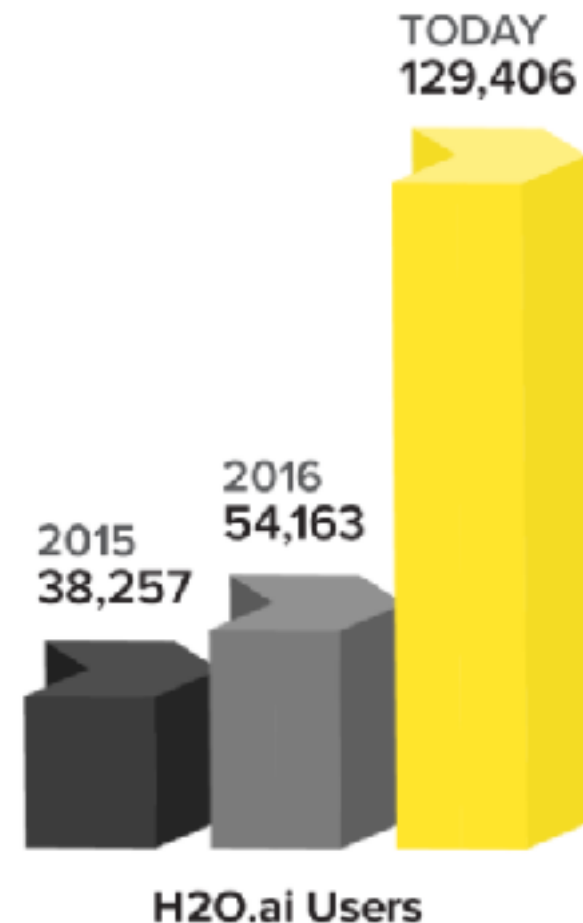


222 OF **THE 500** FORTUNE
❤️ **H₂O**

8 OF TOP 10
BANKS

7 OF TOP 10
INSURANCE COMPANIES

4 OF TOP 10
HEALTHCARE COMPANIES



* DATA FROM GOOGLE ANALYTICS EMBEDDED IN THE END USER PRODUCT

H2O.ai Solution Leadership Across Verticals



Gartner names H2O as **Leader** with the **most completeness of vision**

- H2O.ai recognized as a **technology leader with most completeness of vision**
- H2O.ai was recognized for the mindshare, partner network and status as a **quasi-industry standard** for machine learning and AI.
- **H2O customers gave the highest overall score** among all the vendors for sales relationship and account management, customer support (onboarding, troubleshooting, etc.) and overall service and support.

Figure 1. Magic Quadrant for Data Science and Machine-Learning Platforms



Source: Gartner (February 2018)

Platforms with H₂O integration



srisatish
@srisatish

Following

Replying to @BobMuenchen @knime @h2oai

@KNIME gained the ability to run @H2O.ai algorithms, so these two may be viewed as complementary, not competitors

#Ecosystem #OpenSource

3:32 PM - 2 Mar 2018

H₂O + KNIME Talk
at KNIME Summit
Mar 2017



1:54 PM - 7 Mar 2018 from Lionel Eedins

Figure 1. Magic Quadrant for Data Science and Machine-Learning Platforms



Source: Gartner (February 2018)

H₂O.ai

Community Expansion



meetup

88,286
members

43
interested

52
Meetups

48
cities

18
countries

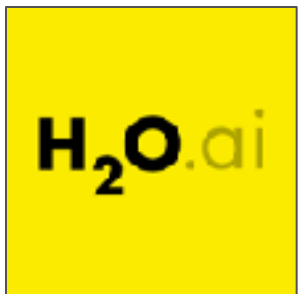
Find out more: www.h2o.ai/community/

H2O

H₂O Machine Learning Platform



H₂O 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



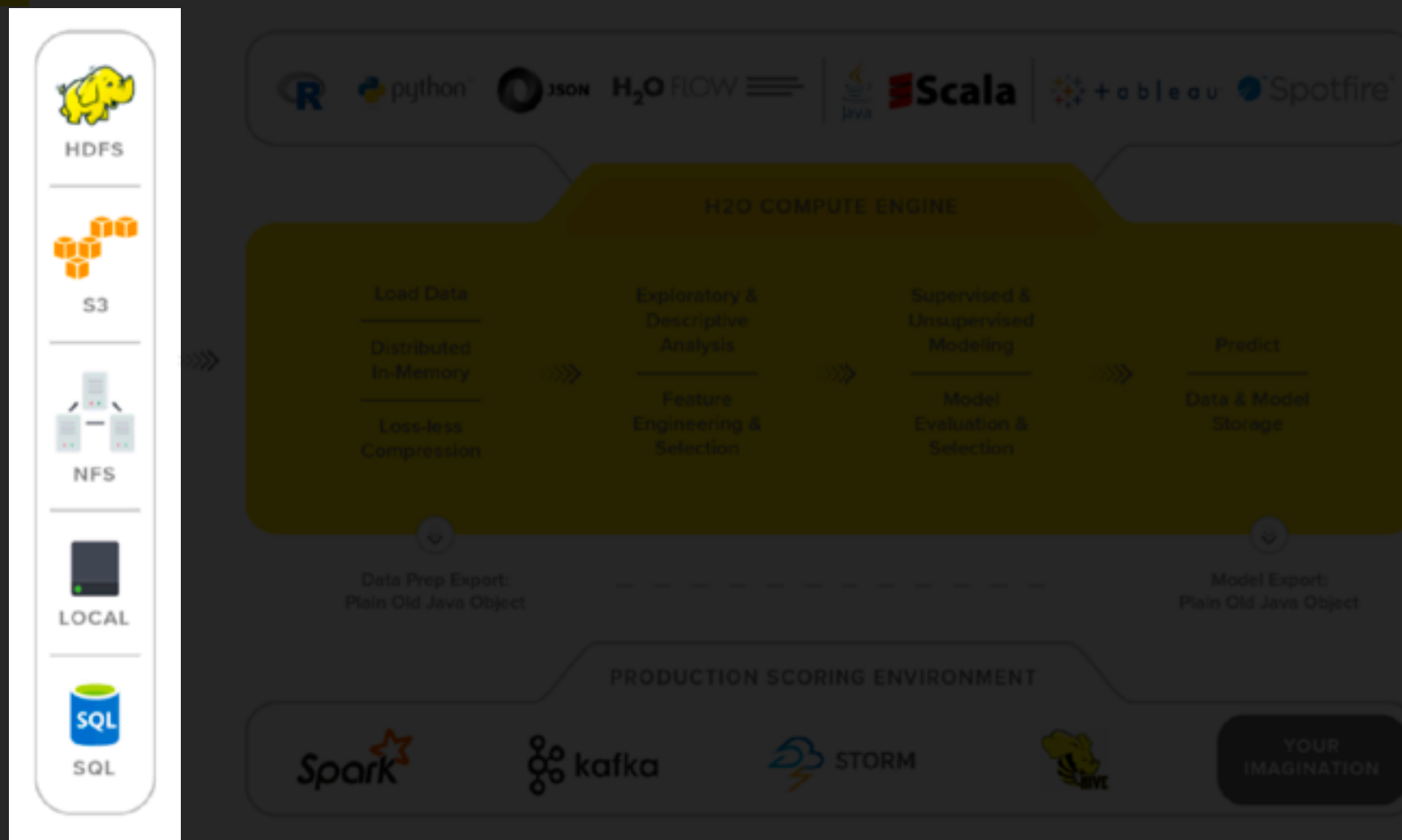
Automatic feature
engineering, machine learning
and interpretability

Steam

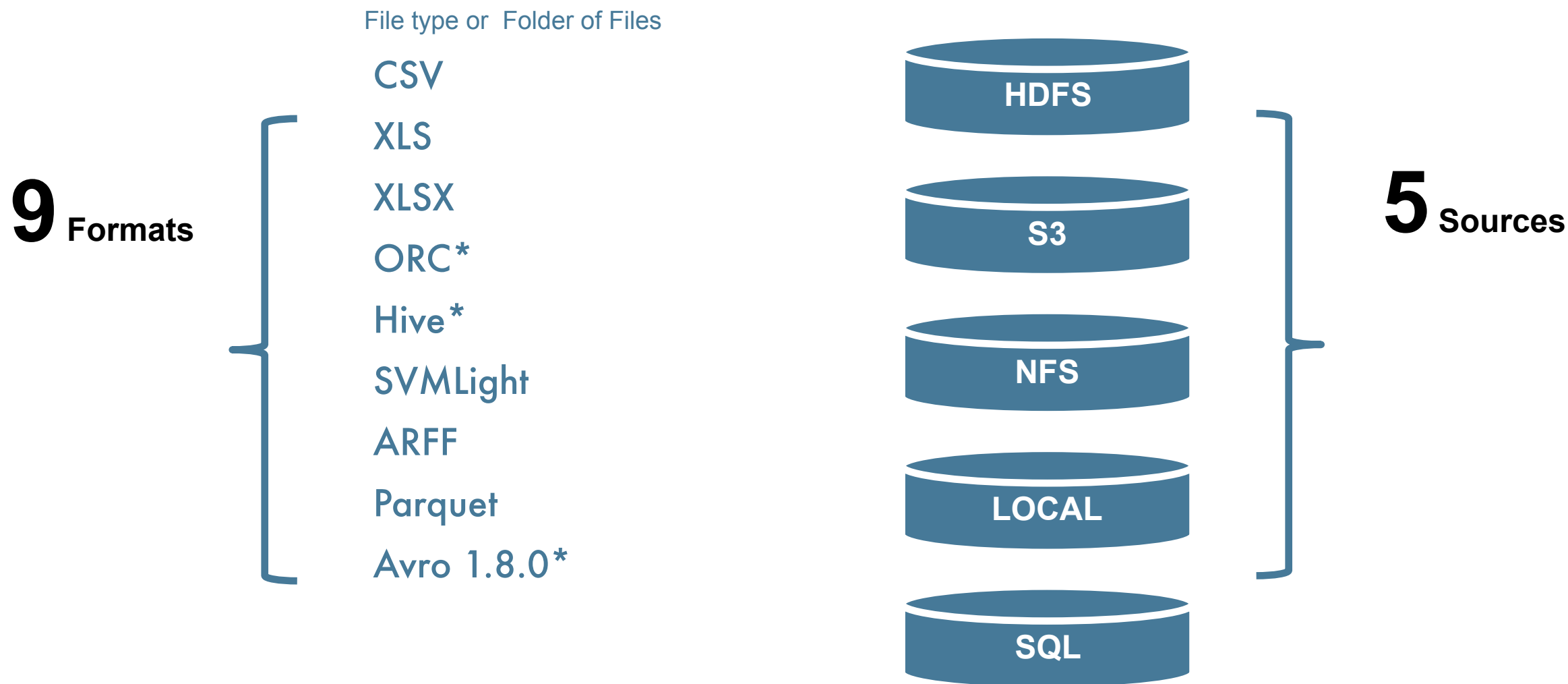
Secure multi-tenant H2O clusters

High Level Architecture





Supported Formats & Data Sources

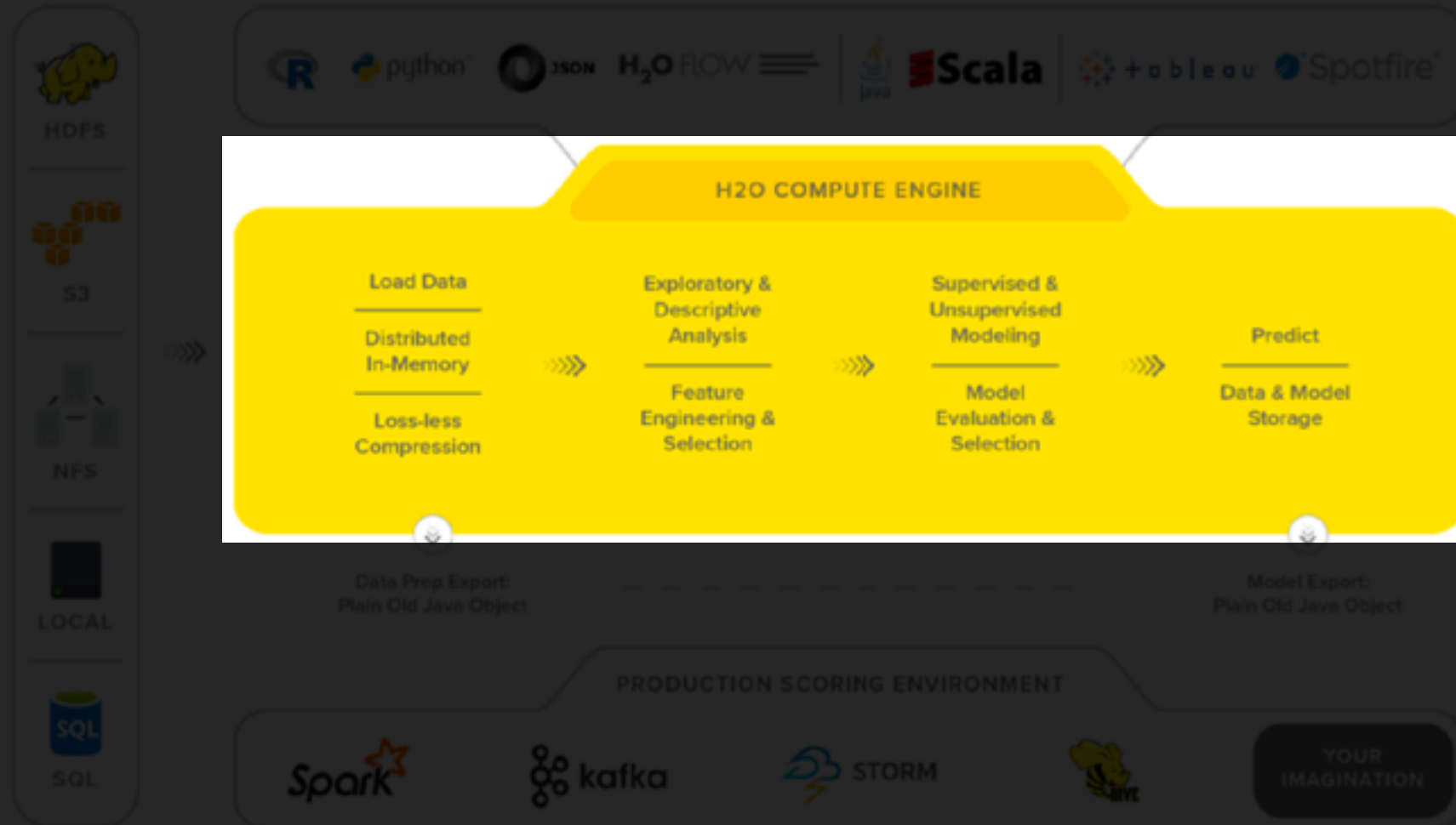


* 1. only if H2O is running as a Hadoop job

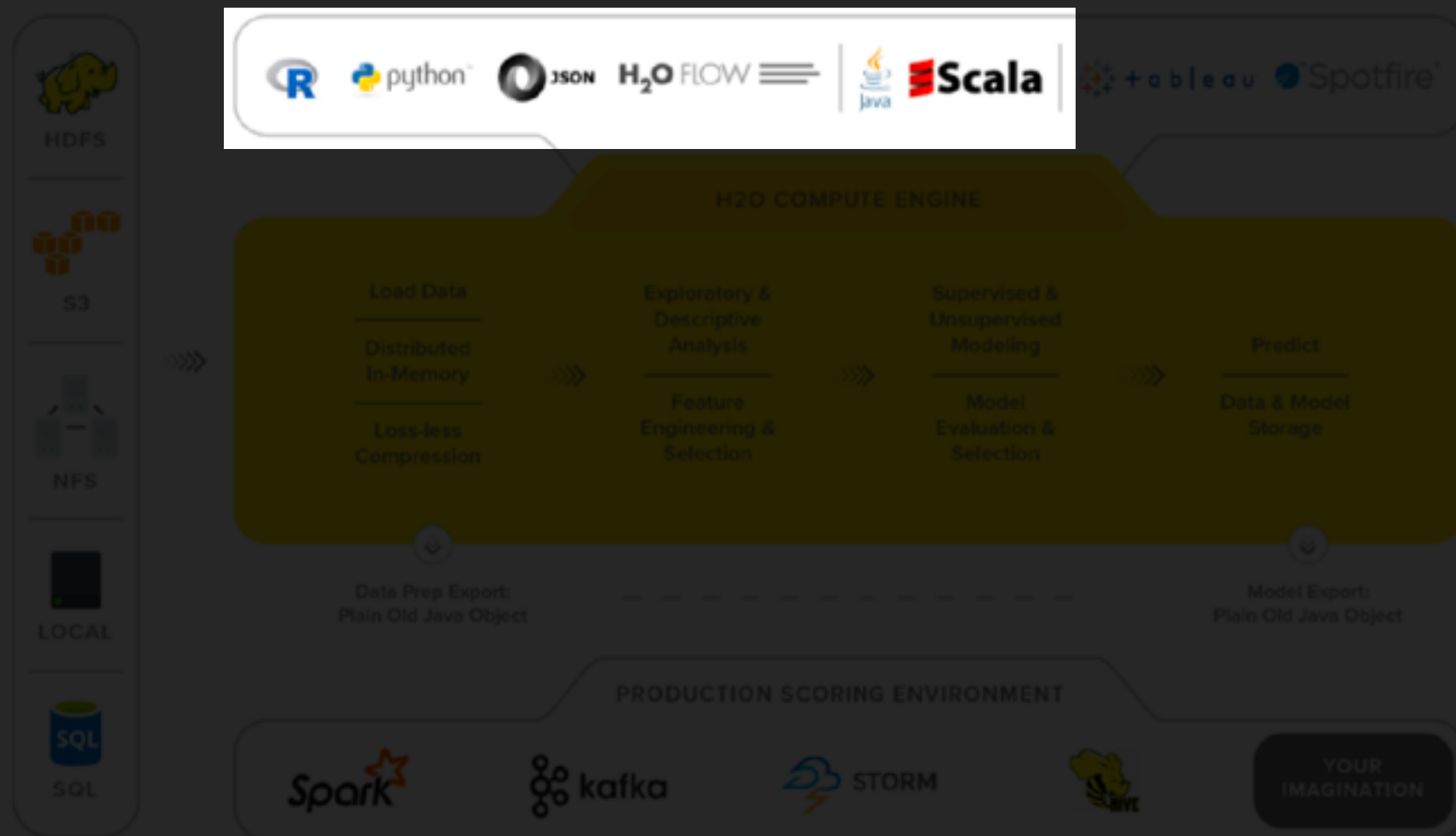
* 2. Hive files that are saved in ORC format

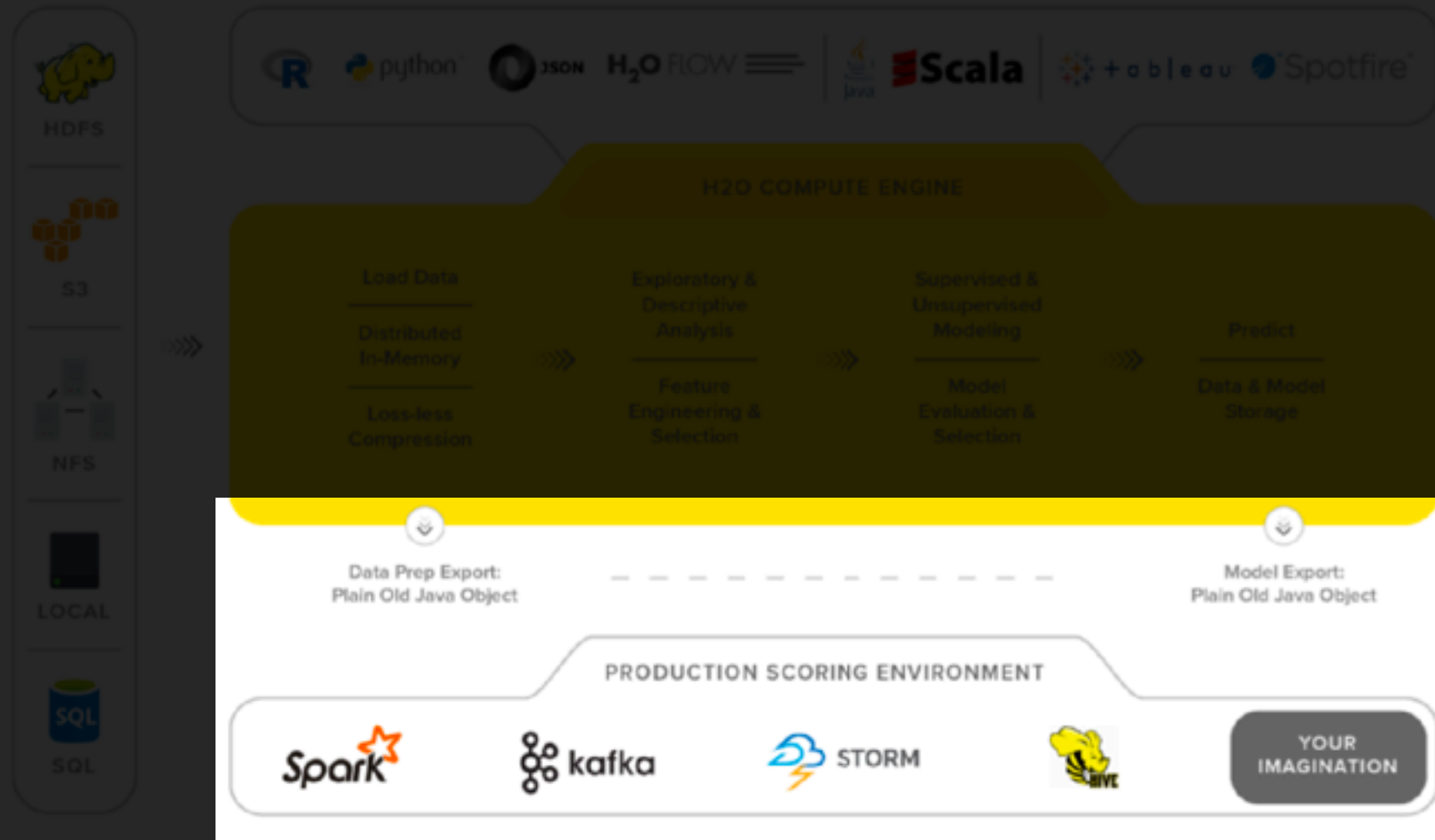
* 3. without multi-file parsing or column type modification

Fast, Scalable &
Distributed Compute
Engine Written in Java



High Level Architecture

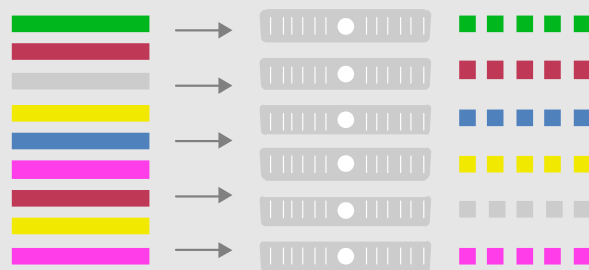




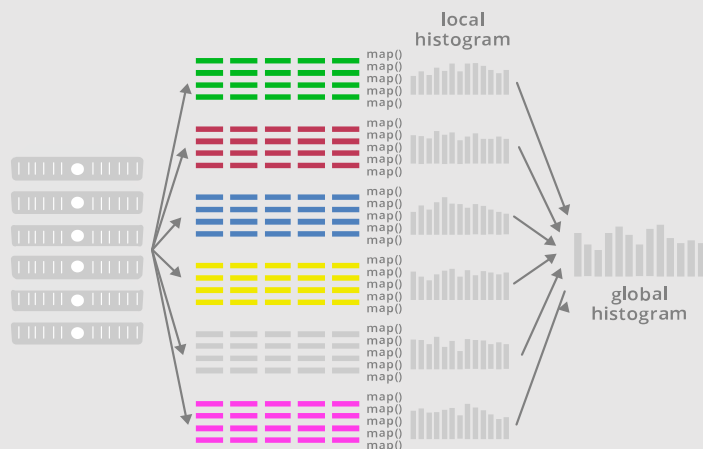


Distributed Algorithms

Foundation for Distributed Algorithms



Parallel Parse into Distributed Rows



Fine Grain Map Reduce Illustration:
Scalable Distributed Histogram Calculation
for GBM

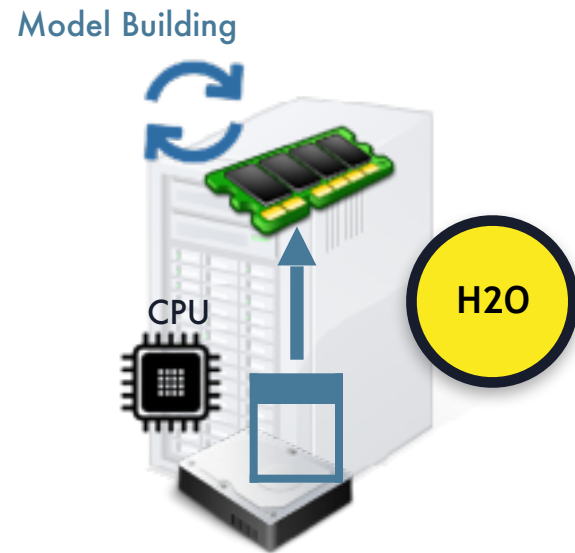
Advantageous Foundation

- Foundation for In-Memory Distributed Algorithm Calculation - **Distributed Data Frames** and **columnar compression**
- All algorithms are distributed in H₂O: GBM, GLM, DRF, Deep Learning and more. Fine-grained map-reduce iterations.
- Only enterprise-grade, open-source distributed algorithms in the market

User Benefits

- “Out-of-box” functionalities for all algorithms (**NO MORE SCRIPTING**) and uniform interface across all languages: R, Python, Java
- Designed for all sizes of data sets, especially **large data**
- Highly optimized Java code for model exports
- **In-house expertise for all algorithms**

H₂O Core

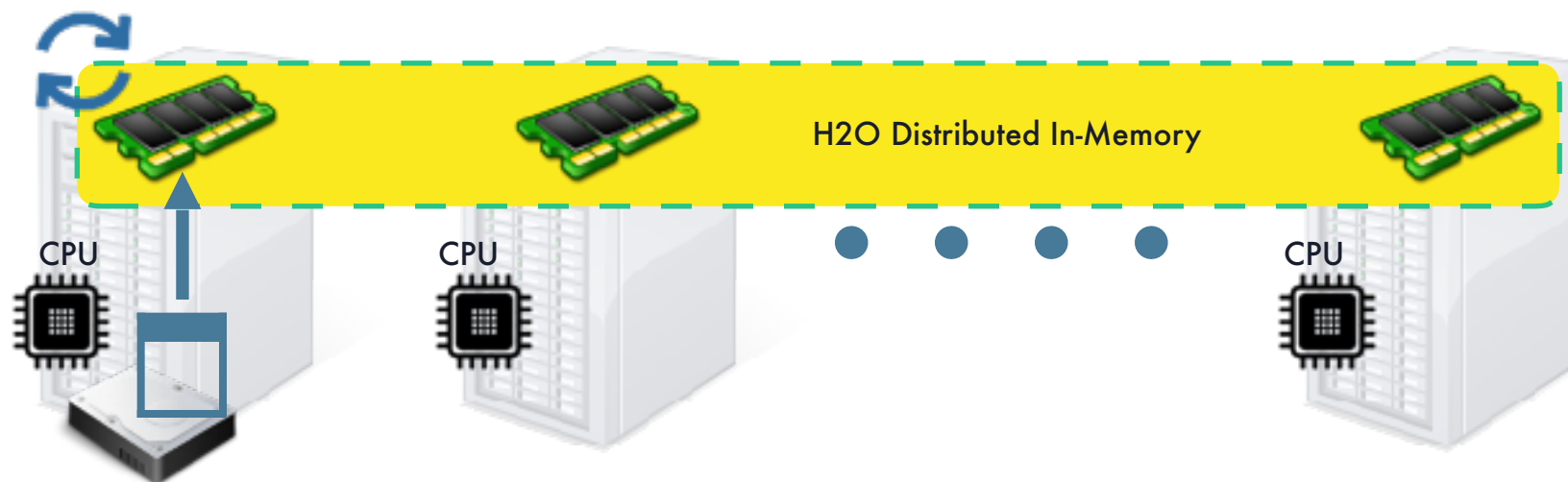


H₂O Core



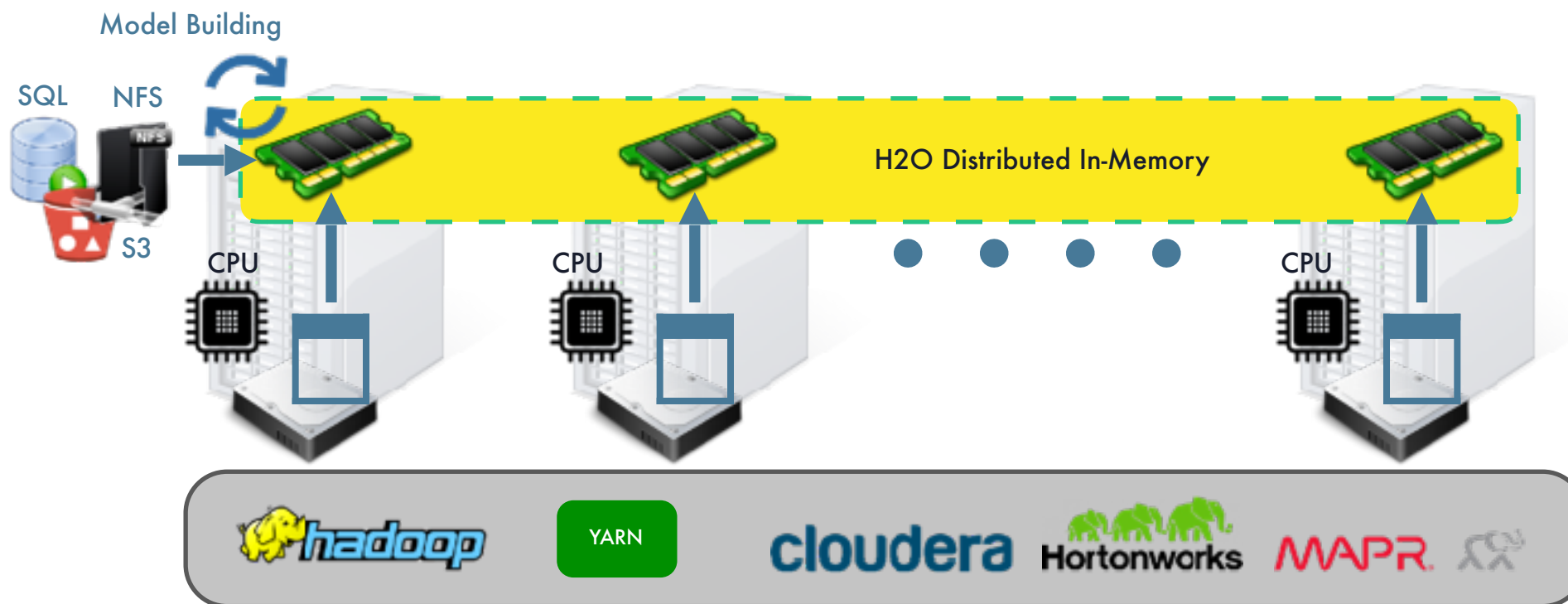
H₂O Core

Model Building



H₂O Core

Firewall or Cloud



H₂O-3 Algorithms Overview

Supervised Learning

Statistical Analysis

- **Generalized Linear Models:** Binomial, Gaussian, Gamma, Poisson and Tweedie
- **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

Unsupervised Learning

Clustering

- **K-means:** Partitions observations into k clusters/groups of the same spatial size. Automatically detect optimal k

Dimensionality Reduction

- **Principal Component Analysis:** Linearly transforms correlated variables to independent components
- **Generalized Low Rank Models:** extend the idea of PCA to handle arbitrary data consisting of numerical, Boolean, categorical, and missing data

Anomaly Detection

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

Downloading H₂O

H₂O

H2O works with R, Python, Scala on Hadoop/Yarn, Spark or your laptop

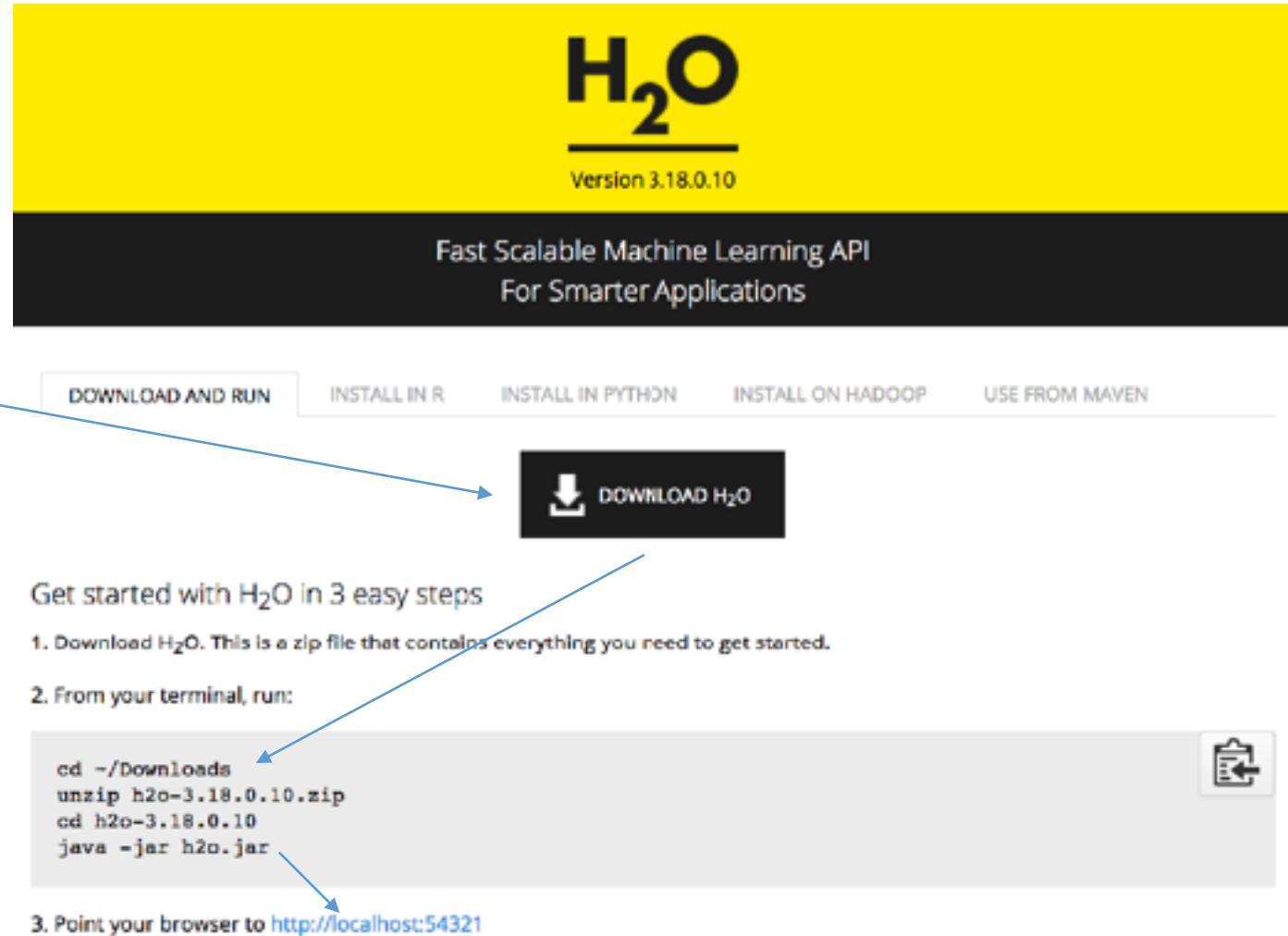
Latest Stable Release

Nightly Bleeding Edge

H₂O in the Cloud



H₂O is licensed under the [Apache License, Version 2.0](#)



The screenshot shows the H2O download page. At the top, the H2O logo is displayed with the version 3.18.0.10. Below the logo, the text 'Fast Scalable Machine Learning API For Smarter Applications' is shown. A navigation bar contains links: 'DOWNLOAD AND RUN', 'INSTALL IN R', 'INSTALL IN PYTHON', 'INSTALL ON HADOOP', and 'USE FROM MAVEN'. The 'DOWNLOAD AND RUN' link is highlighted, and a blue arrow points from the 'Latest Stable Release' button in the left sidebar to the 'DOWNLOAD H2O' button. Below the navigation bar, the text 'Get started with H₂O in 3 easy steps' is followed by three steps: 1. Download H₂O. This is a zip file that contains everything you need to get started. 2. From your terminal, run:

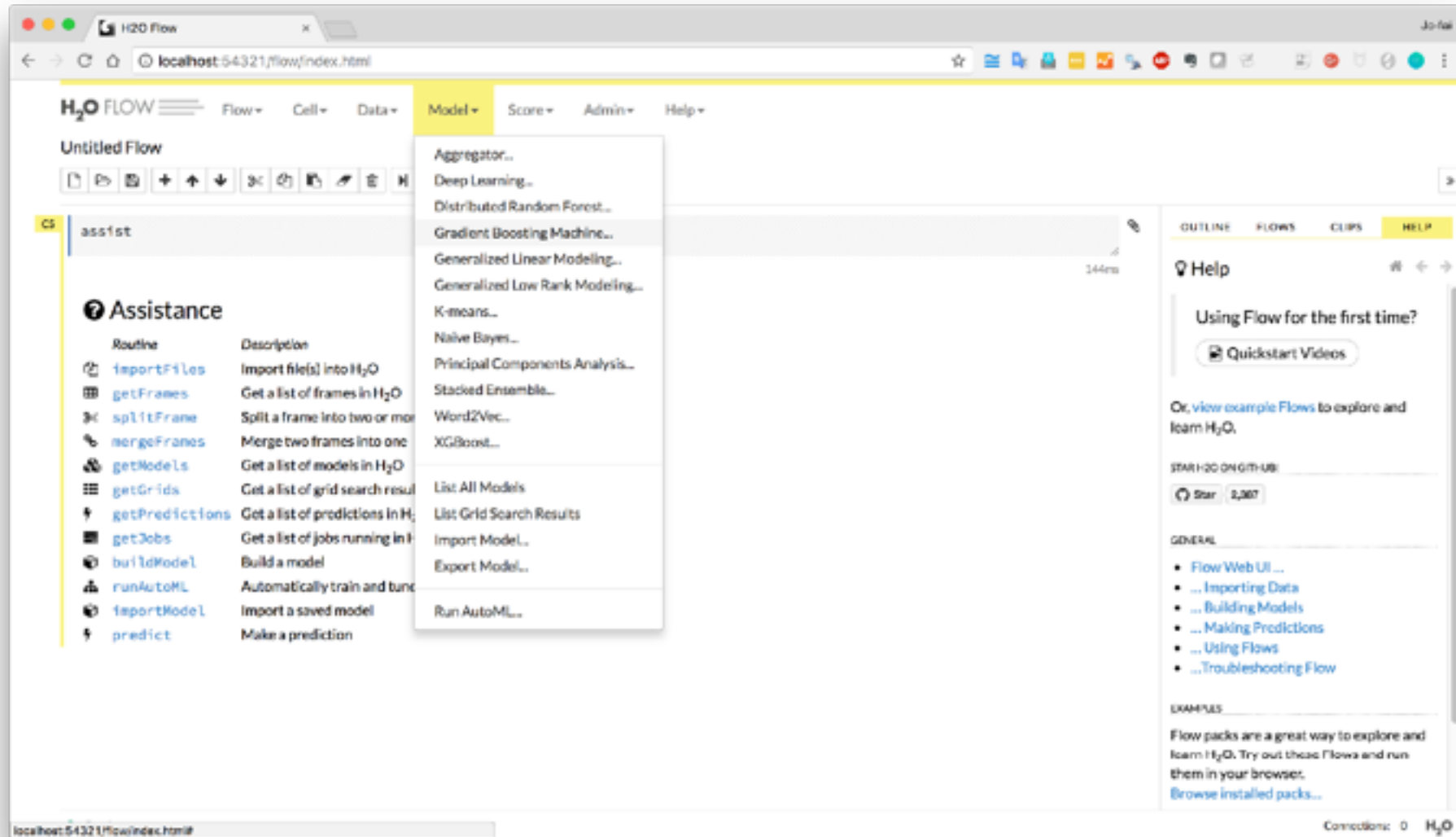
```
cd ~/Downloads
unzip h2o-3.18.0.10.zip
cd h2o-3.18.0.10
java -jar h2o.jar
```

 A blue arrow points from the 'DOWNLOAD H2O' button to the first line of the terminal command. 3. Point your browser to <http://localhost:54321>. A blue arrow points from the second line of the terminal command to the URL.

h2o.ai/download/

H₂O.ai

H₂O Flow and Python Client – First Demo



H₂O Documentation

[Getting Started & User Guides](#) | [Q & A](#) | [Algorithms](#) | [Languages](#) | [Tutorials, Examples, & Presentations](#) | [API & Developer Docs](#) | [For the Enterprise](#)

Getting Started & User Guides

 Open Source |  Commercial

H₂O

What is H₂O?

[H₂O User Guide](#) (Main docs)

[H₂O Book](#) (O'Reilly)

[Recent Changes](#)

[Open Source License](#) (Apache V2)

[Quick Start Video - Flow Web UI](#)

[Quick Start Video - R](#)

[Quick Start Video - Python](#)

[Download H₂O](#)

Sparkling Water

What is Sparkling Water?

[Sparkling Water User Guide](#) 2.3 2.2 2.1

[Sparkling Water Booklet](#)

[RSparkling Readme](#)

[PySparkling User Guide](#) 2.3 2.2 2.1

[Recent Changes](#) 2.3 2.2 2.1

[Open Source License](#) (Apache V2)

[Quick Start Video - Scala](#)

[Download Sparkling Water](#)

Driverless AI

What is Driverless AI?

[Driverless AI User Guide](#) [HTML](#) [PDF](#)

[Recent Changes](#)

[Driverless AI Booklet](#)

[MLI with Driverless AI Booklet](#)

[Quick Start Video - Downloading Driverless AI](#)

[Quick Start Video - Launching an Experiment](#)

[Driverless AI Webinars](#)

[Download Driverless AI](#)

H₂O4GPU (alpha)

[H₂O4GPU Readme](#)

[Open Source License](#) (Apache V2)

[Download H₂O4GPU](#)

Thanks!

- Organisers & Sponsors

- GapData Institute
- PyData Bratislava

- H₂O's Mission

- Democratize AI
- Make Machine Learning Accessible to Everyone

- Code, Slides & Documents

- bit.ly/h2o_meetups
- docs.h2o.ai

- Contact

- michalr@h2o.ai
- github.com/michal-raska

- Please search/ask questions on **Stack Overflow**

- Use the tag `h2o` (not h2 zero)