

# Introduction to H<sub>2</sub>O

H<sub>2</sub>O.ai

Jakub Háva  
jakub@h2o.ai

Thanks for the slides to:  
Jo-fai (Joe) Chow  
Data Scientist  
joe@h2o.ai  
@matlabulous

The Amsterdam Pipeline Factory  
of Data Science, Amsterdam  
6<sup>th</sup> December, 2016

# About H<sub>2</sub>O.ai

What exactly is H<sub>2</sub>O?

# Company Overview

<b>Founded</b>	2011 Venture-backed, debuted in 2012
<b>Products</b>	<ul style="list-style-type: none"><li>• H2O Open Source In-Memory AI Prediction Engine</li><li>• Sparkling Water</li><li>• Steam</li></ul>
<b>Mission</b>	Operationalize Data Science, and provide a platform for users to build beautiful data products
<b>Team</b>	<p>70 employees</p> <ul style="list-style-type: none"><li>• Distributed Systems Engineers doing Machine Learning</li><li>• World-class visualization designers</li></ul>
<b>Headquarters</b>	Mountain View, CA



# Bring AI To Business Empower Transformation














H<sub>2</sub>O is an open source platform  
empowering business transformation

**Financial Services, Insurance and  
Healthcare as Our Vertical Focus**



**Community as Our Foundation**

# Users In Various Verticals Adore H<sub>2</sub>O

 	 	  	   	  
 Financial	 Telecom	 Insurance	 Healthcare	 Marketing



## Capital One



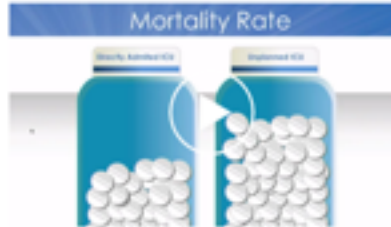
Capital One uses H2O open source machine learning for various use cases.

## MarketShare



H2O predictive analytics helps boost the impact and results of digital marketing.

## Kaiser



Kaiser uses H2O machine learning to save lives.

## Zurich Insurance



Zurich turned to H2O as a strategic differentiator for commercial insurance.

## Progressive



Progressive uses H2O predictive analytics for user-based insurance.

## Comcast



Comcast uses H2O to improve customer experience.

## Hospital Corporation of America



HCA uses H2O to predict patient outcomes in real-time.

## McKesson



McKesson discusses the adoption of artificial intelligence in healthcare.

## Macy's



Macy's uses H2O for personalized site recommendations.

## Transamerica



Transamerica turns to H2O to develop a product recommendation platform for insurance.

## Paypal



Paypal turned to H2O Deep Learning for fraud detection and customer churn.

## eBay



eBay chose H2O for open source machine learning.

# H<sub>2</sub>O.ai Makes A Difference as an AI Platform

## Open Source



- 100% open source

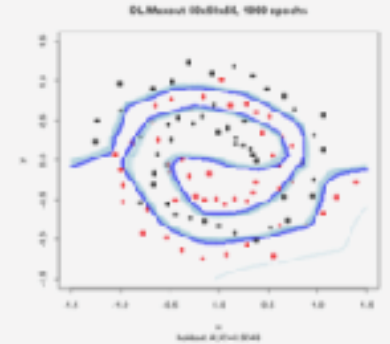
## Big Data Ecosystem



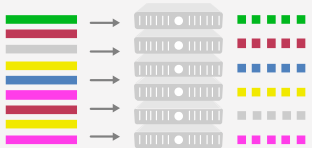
## Flexible Interface



## Smart and Fast Algorithms



## Scalability and Performance



- Distributed In-Memory Computing Platform
- Distributed Algorithms
- Fine-Grain MapReduce

## Rapid Model Deployment

- Highly portable models deployed in Java (POJO) and Model Object Optimized (MOJO)
- Automated and streamlined scoring service deployment with Rest API



## GPU Enablement

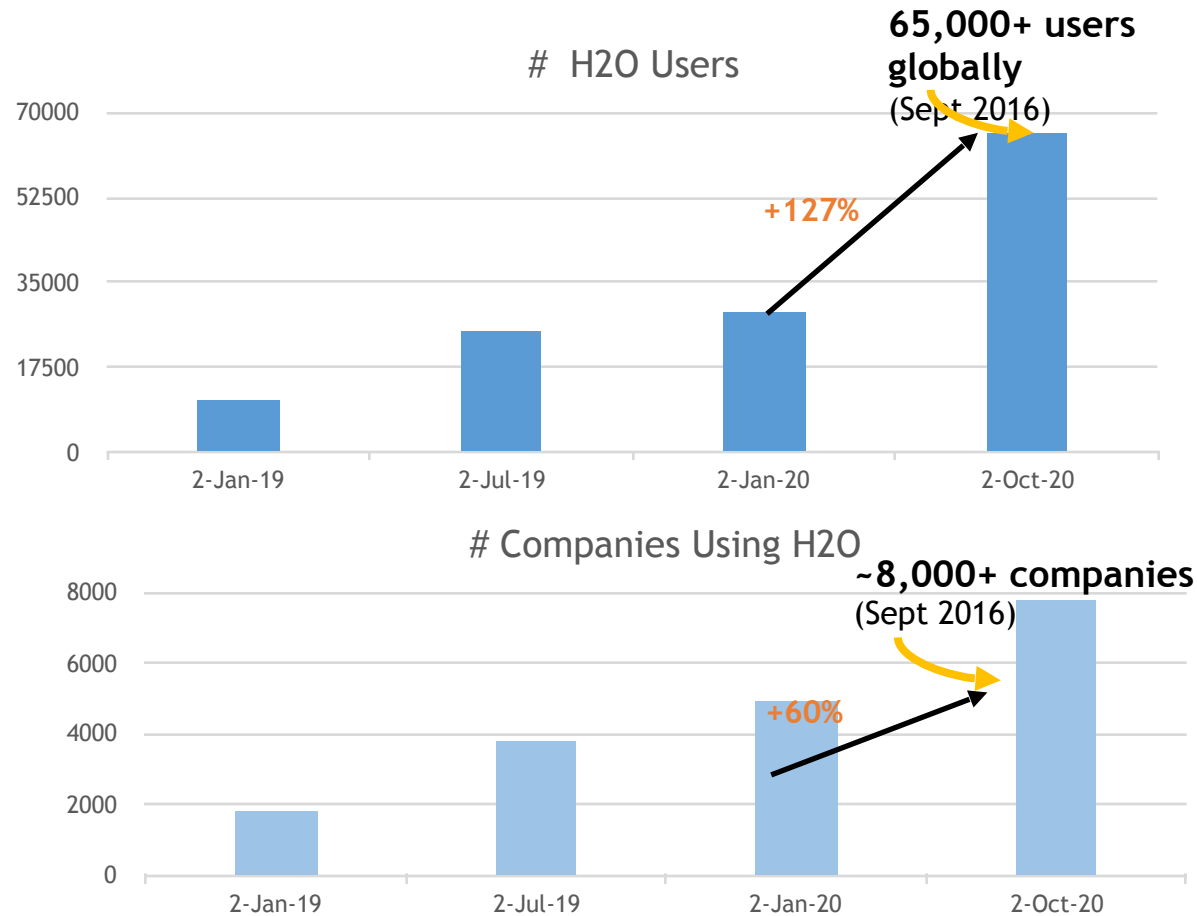


## Cloud Integration



# H<sub>2</sub>O Community Growth

## Tremendous Momentum Globally



### Large User Circle

- 65,000+ users from ~8,000 companies in 140 countries. Top 5 from:

1.	 United States
2.	 India
3.	 Japan
4.	 Germany
5.	 United Kingdom

\* DATA FROM GOOGLE ANALYTICS EMBEDDED IN THE END USER PRODUCT

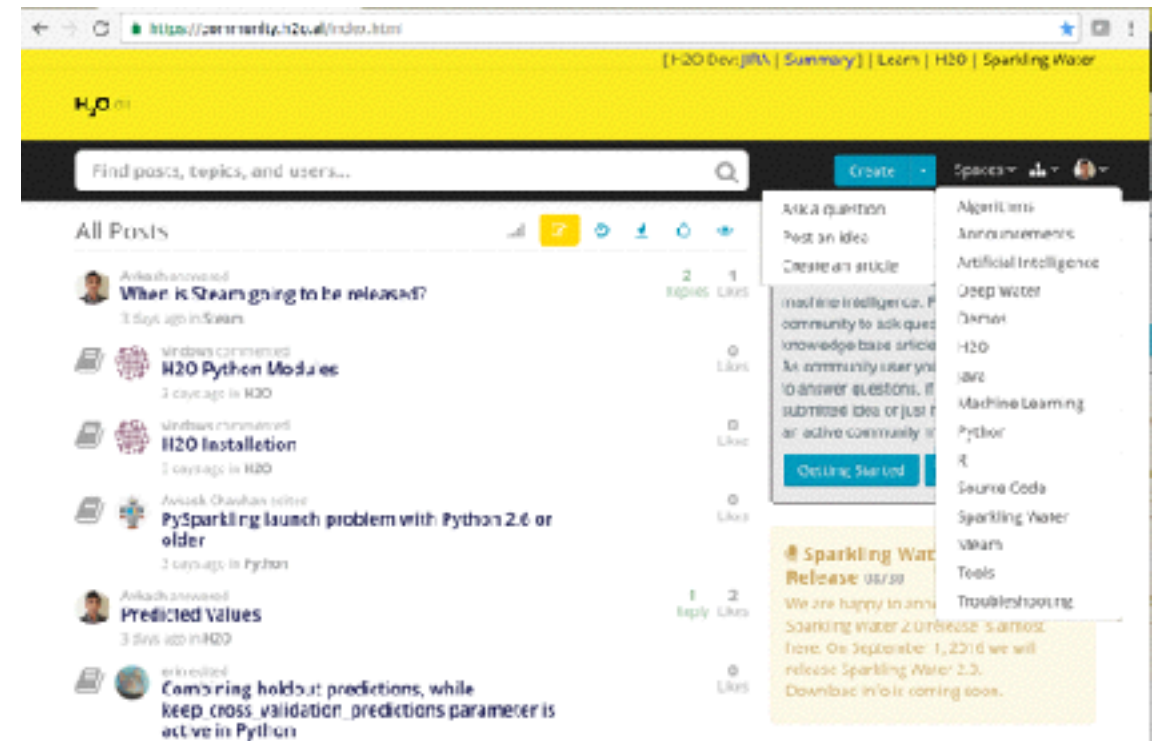
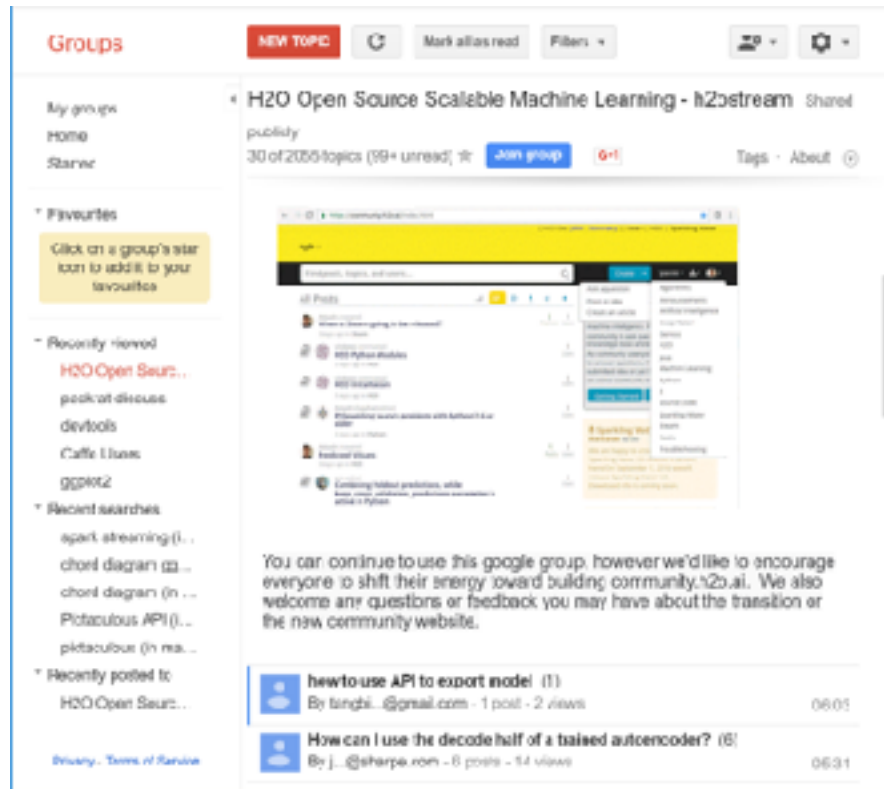


# H<sub>2</sub>O Community Support

Please try

Google forum - h2ostream

community.h2o.ai



H<sub>2</sub>O.ai

# H<sub>2</sub>O for Kaggle Competitions

**CIFAR-10 Competition  
Winners: Interviews with Dr.  
Ben Graham, Phil Culliton, &  
Zygmunt Zajac**

Triskellon | 01.02.2015

[READ MORE](#)

“I did really like H2O’s deep learning implementation in R, though - the interface was great, the back end extremely easy to understand, and it was scalable and flexible. Definitely a tool I’ll be going back to.”

**Kaggle challenge  
2nd place winner  
Colin Priest**

for creating this corpus. .  
do not contain Spanish sent-  
is a widespread major langu-  
reason was to create a corp-  
tasks. These tasks are com

Completed • Knowledge • 161 teams

**Denoising Dirty Documents**

Mon 1 Jun 2015 – Mon 5 Oct 2015 (3 months ago)

[READ MORE](#)

“For my final competition submission I used an ensemble of models, including 3 deep learning models built with R and h2o.”

**H<sub>2</sub>O.ai**

# H<sub>2</sub>O for Academic Research



<http://www.sciencedirect.com/science/article/pii/S0377221716308657>



<https://arxiv.org/abs/1509.01199>

H<sub>2</sub>O

democratizes

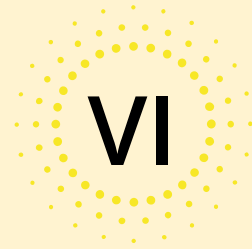
artificial intelligence & big data science

# Our Open Source Products

100% Open Source. Big Data Science for Everyone!

# H<sub>2</sub>O.ai Offers AI Open Source Platform

Product Suite to Operationalize Data Science with Visual Intelligence



Visual Intelligence and UX Framework For Data Interpretation and Story Telling on top of Beautiful Data Products

100% Open Source



---

In-Memory, Distributed  
Machine Learning  
Algorithms with Speed and  
Accuracy

Deep  
Water

---

State-of-the-art  
Deep Learning on GPUs  
with TensorFlow, MXNet or  
Caffe with the ease of use  
of H2O



---

H2O Integration with  
Spark. Best Machine  
Learning on Spark.

---

Operationalize and  
Streamline Model Building,  
Training and Deployment  
Automatically and  
Elastically



# H<sub>2</sub>O.ai Offers AI Open Source Platform

Product Suite to Operationalize Data Science with Visual Intelligence



Visual Intelligence and UX Framework For Data Interpretation and Story Telling on top of Beautiful Data Products

100% Open Source



---

In-Memory, Distributed  
Machine Learning  
Algorithms with Speed and  
Accuracy

Deep  
Water

---

State-of-the-art  
Deep Learning on GPUs  
with TensorFlow, MXNet or  
Caffe with the ease of use  
of H2O



---

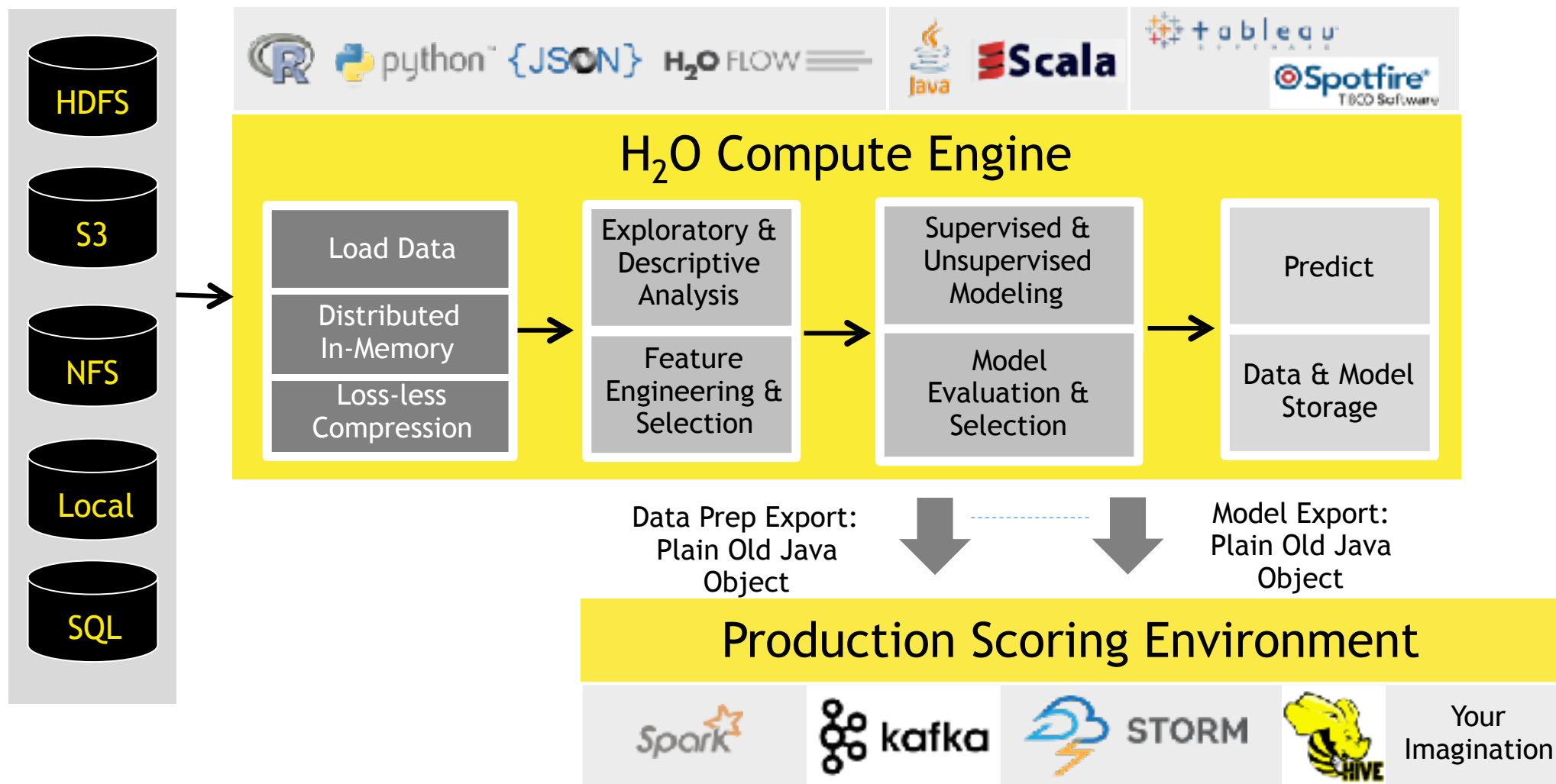
H2O Integration with  
Spark. Best Machine  
Learning on Spark.

Steam

---

Operationalize and  
Streamline Model Building,  
Training and Deployment  
Automatically and  
Elastically

# High Level Architecture



# 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

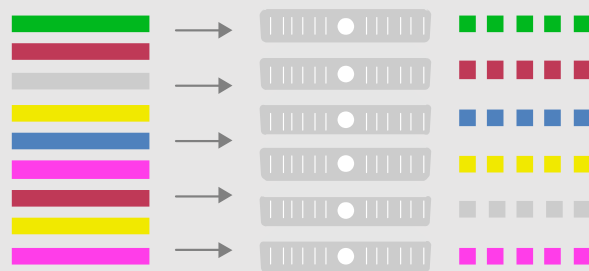
- **Principal Component Analysis:** Linearly transforms correlated variables to independent components
- **Generalized Low Rank Models:** extend the idea of PCA to handle arbitrary data consisting of numerical, Boolean, categorical, and missing data

### Anomaly Detection

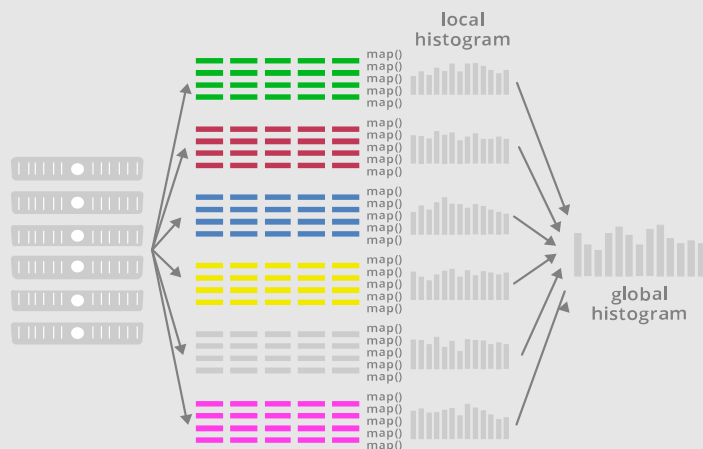
- **Autoencoders:** Find outliers using a nonlinear dimensionality reduction using deep learning

# 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 H<sub>2</sub>O: 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**

# H<sub>2</sub>O Deep Learning in Action

116M rows, 6GB CSV file  
800+ predictors (numeric + categorical)

airlines\_all\_selected\_cols.hex

Actions: View Data Split Build Model Predict Download Export

Rows	Columns	Compressed Size
116495259	12	2GB

## Job

Run Time 00:00:36.712

Remaining Time 00:00:17.188

Type Model

Key Q deeplearning-dd2f42f7-81f7-42e8-9d98-e34437309828

Description DeepLearning

Status RUNNING

Progress 69%

Iterations: 12. Epochs: 0.628821. Speed: 2,243,735 samples/sec. Estimated time left: 21.849 sec

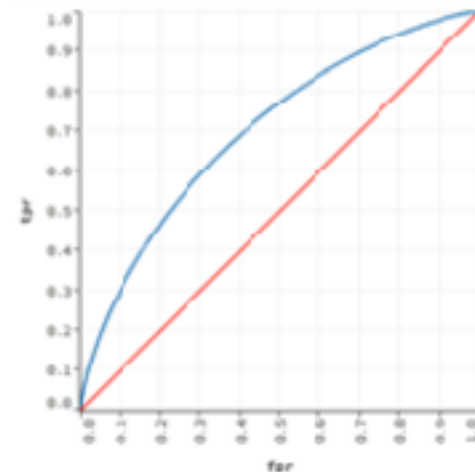
Actions View Cancel Job

model trained in <1 min:  
2M+ samples/second

OUTPUT - STATUS OF NEURON LAYERS (PREDICTING ISDEPDELAYED, 2-CLASS CLASSIFICATION, BERNOULLI DISTRIBUTION, CROSSENTROPY LOSS, 17.462 WEIGHTS/BIASES, 221.3 KB, 104,545,345 TRAINING SAMPLES, MINI-BATCH SIZE 1)

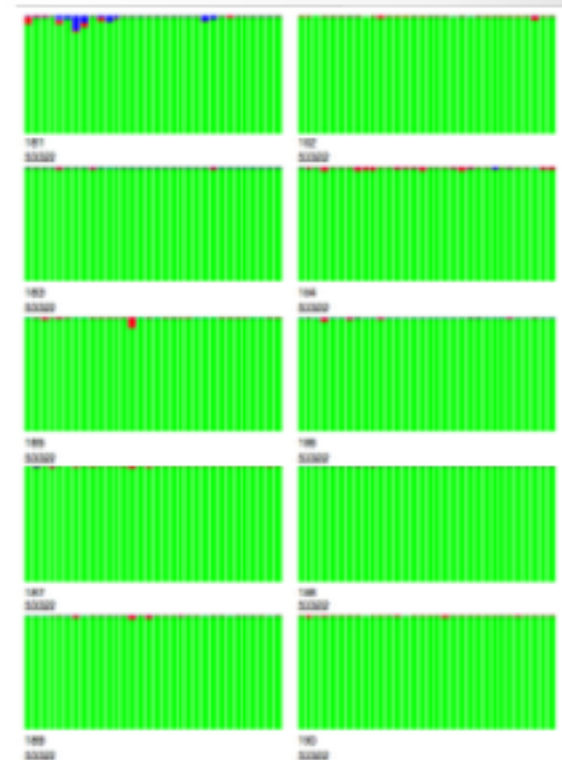
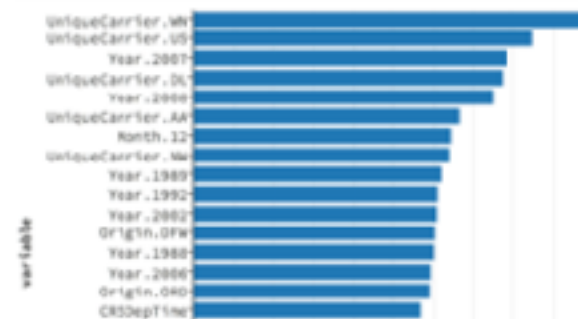
layer	units	type	dropout	l1	l2	mean_rate	rate_RMS	momentum	mean_weight	weight_RMS	mean_bias	bias_RMS
1	807	Input	0	0	0	0.0493	0.2029	0	-0.0021	0.2111	-0.9139	1.0036
2	24	Rectifier	0	0	0	0.9197	0.9227	0	-0.1033	0.9362	-1.3990	1.0299
3	24	Rectifier	0	0	0	0.9517	0.9440	0	-0.1575	0.9368	-0.0840	0.6040
4	24	Rectifier	0	0	0	0.9761	0.9844	0	-0.9374	0.2275	-0.2647	0.2481
5	2	Softmax	0	0	0	0.9161	0.9083	0	0.0741	0.7260	0.4269	0.2056

ROC CURVE - VALIDATION METRICS, AUC = 0.702560



Thresholds: Choose... Criteria: Choose...

VARIABLE IMPORTANCES



## Legend

Each bar represents one CPU.

Blue: idle time

Green: user time

Red: system time

White: other time (e.g. I/O)

10 nodes: all  
320 cores busy

H<sub>2</sub>O.ai Deep Learning Model

real-time, interactive  
model inspection in Flow



# H<sub>2</sub>O + R

```
0 h2o_iris_demo.R
Source on Save
Run
Source

1- # .....
2 # Build a simple classification model using iris dataset
3- # .....
4
5 # Start and connect to a local H2O cluster
6 library(h2o)
7 h2o.init(nthreads = -1)
8
9 # Import data from a R data frame
10 data(iris)
11 d_iris <- as.h2o(iris)
12
13 # Define Targets and Features
14 target <- "Species"
15 features <- setdiff(colnames(d_iris), c("Species"))
16
17- # .....
18 # Train a H2O Model
19- # .....
20
21 # Train three basic H2O models
22 model_drf <- h2o.randomForest(x = features,
23 | | | | | y = target,
24 | | | | | model_id = "iris_random_forest",
25 | | | | | training_frame = d_iris)
26
27 model_gbm <- h2o.gbm(x = features,
28 | | | | | y = target,
29 | | | | | model_id = "iris_gbm",
30 | | | | | training_frame = d_iris)
31
32 model_dnn <- h2o.deeplearning(x = features,
33 | | | | | y = target,
34 | | | | | model_id = "iris_deep_learning",
35 | | | | | training_frame = d_iris)
36
```

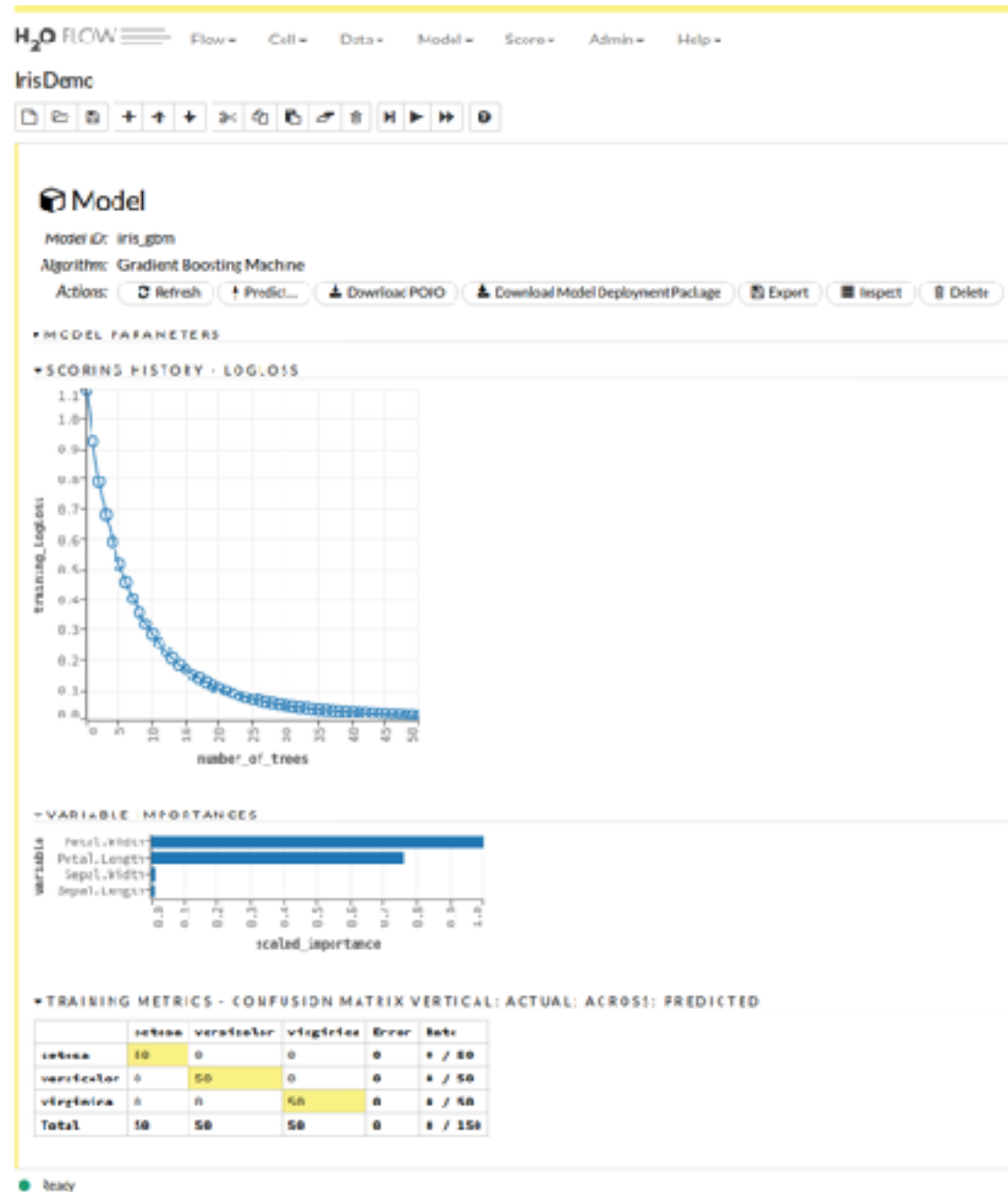
Please try



# H<sub>2</sub>O Flow (Web)

The screenshot displays the H2O Flow web interface in a browser window. The address bar shows the URL `localhost:54321/demos/index.html`. The top navigation bar includes the H2O FLOW logo and menu items: Flow, Cell, Data, Model, Score, Admin, and Help. The main content area is titled "Iris Demo" and features a toolbar with icons for file operations, execution, and debugging. Below the toolbar, the "Expression" field contains the code `getModels`. A sidebar on the left, labeled "Models", lists three models: `iris_deep_learning` (Deep Learning), `iris_gbm` (Gradient Boosting Machine), and `iris_random_forest` (Distributed Random Forest). Each model has an "Inspect" button. At the bottom right, there are "Predict..." and "Inspect" buttons for each model. The status bar at the bottom indicates "Ready" and shows a "Connections" count of 0.

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



# Key Learning Resources

- Help Documentations

- [docs.h2o.ai](https://docs.h2o.ai)

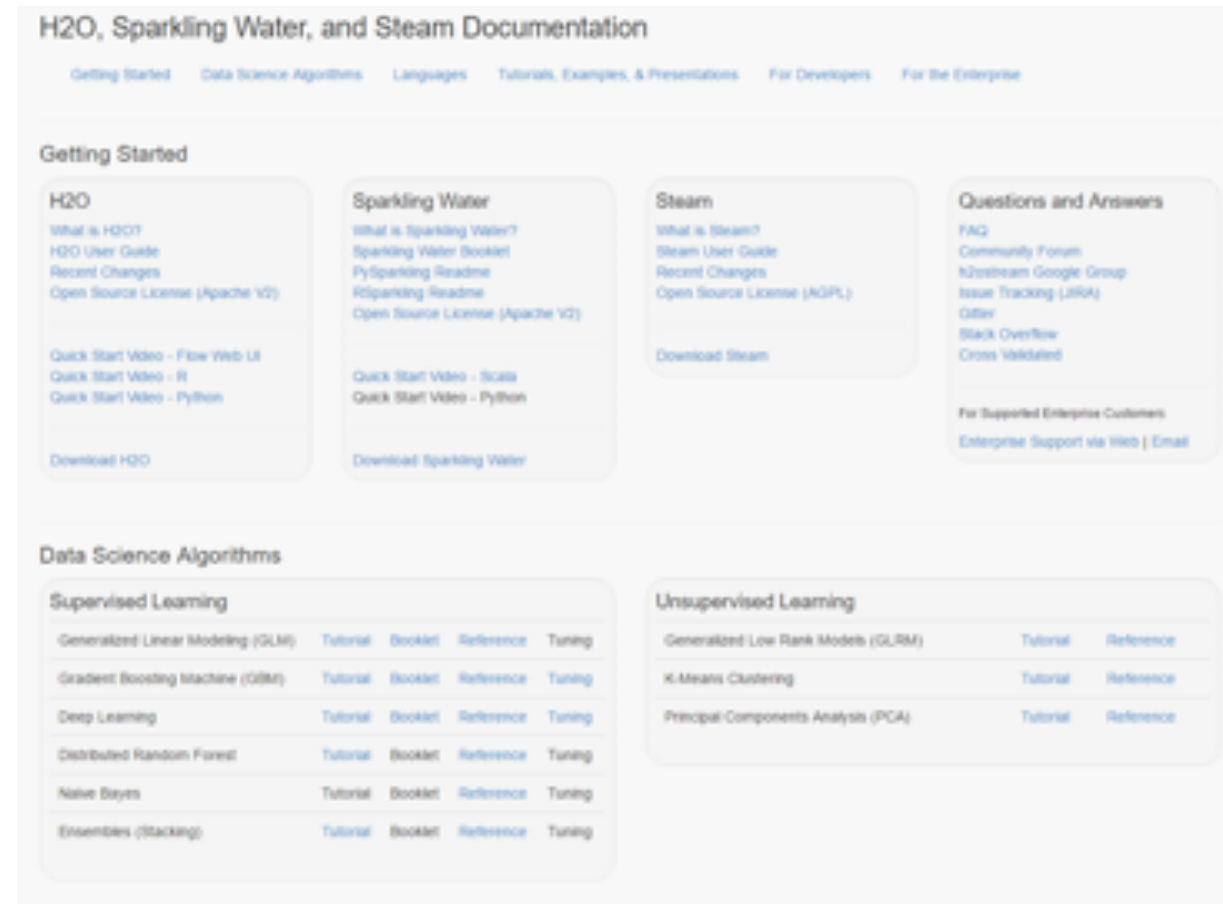


- Meetups

- [bit.ly/h2o\\_meetups](https://bit.ly/h2o_meetups)

- YouTube Channel

- [bit.ly/h2o\\_youtube](https://bit.ly/h2o_youtube)



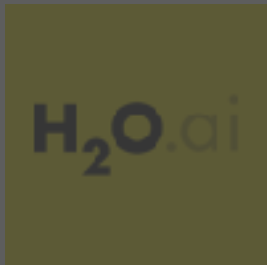
# H<sub>2</sub>O.ai Offers AI Open Source Platform

Product Suite to Operationalize Data Science with Visual Intelligence



Visual Intelligence and UX Framework For Data Interpretation and Story Telling on top of Beautiful Data Products

100% Open Source



---

In-Memory, Distributed  
Machine Learning  
Algorithms with Speed and  
Accuracy

**Deep  
Water**

---

State-of-the-art  
Deep Learning on GPUs  
with TensorFlow, MXNet or  
Caffe with the ease of use  
of H2O



---

H2O Integration with  
Spark. Best Machine  
Learning on Spark.

**Steam**

---

Operationalize and  
Streamline Model Building,  
Training and Deployment  
Automatically and  
Elastically

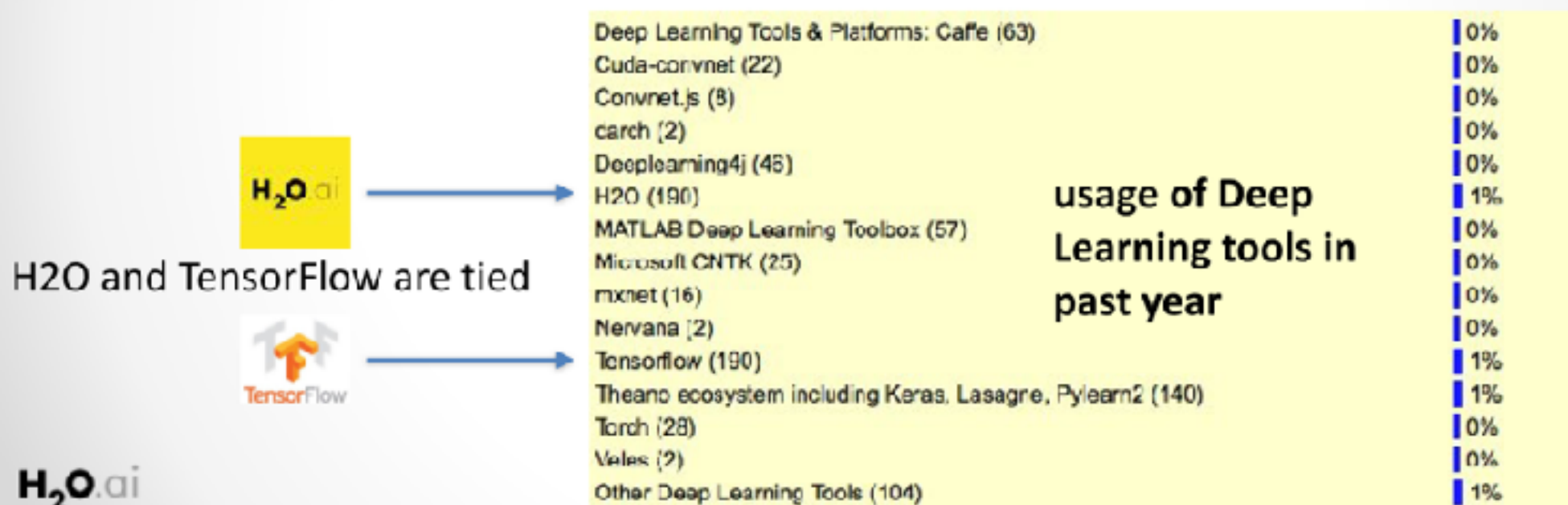
# Both TensorFlow and H<sub>2</sub>O are widely used

The usage of Hadoop/Big Data tools grew to 39%, up from 29% in 2015 (and 17% in 2014), driven by Apache Spark, MLlib (Spark Machine Learning Library) and H2O.

See also

- KDnuggets interview with Spark Creator Matei Zaharia
- KDnuggets interview with Arno Candel, H2O.ai on How to Quick Start Deep Learning with H2O

<http://www.kdnuggets.com>



TensorFlow democratizes the power of deep learning.

H<sub>2</sub>O democratizes artificial intelligence & big data science.

There are other open source libraries like MXNet and Caffe too.  
Let's have a party, this will be fun!

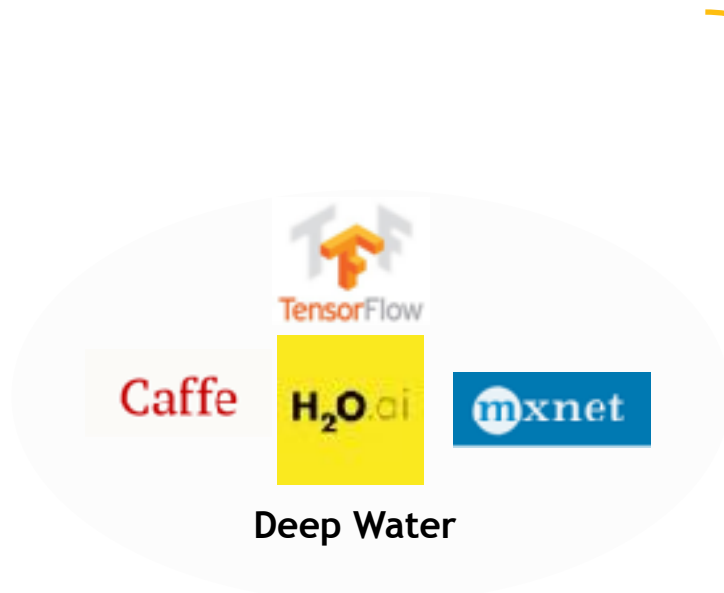


# Deep Water

## Next-Gen Distributed Deep Learning with H<sub>2</sub>O

**One Interface - GPU Enabled - Significant Performance Gains**

Inherits All H<sub>2</sub>O Properties in Scalability, Ease of Use and Deployment



H<sub>2</sub>O integrates with existing **GPU** backends for **significant performance gains**



Convolutional Neural Networks enabling **Image, video, speech** recognition

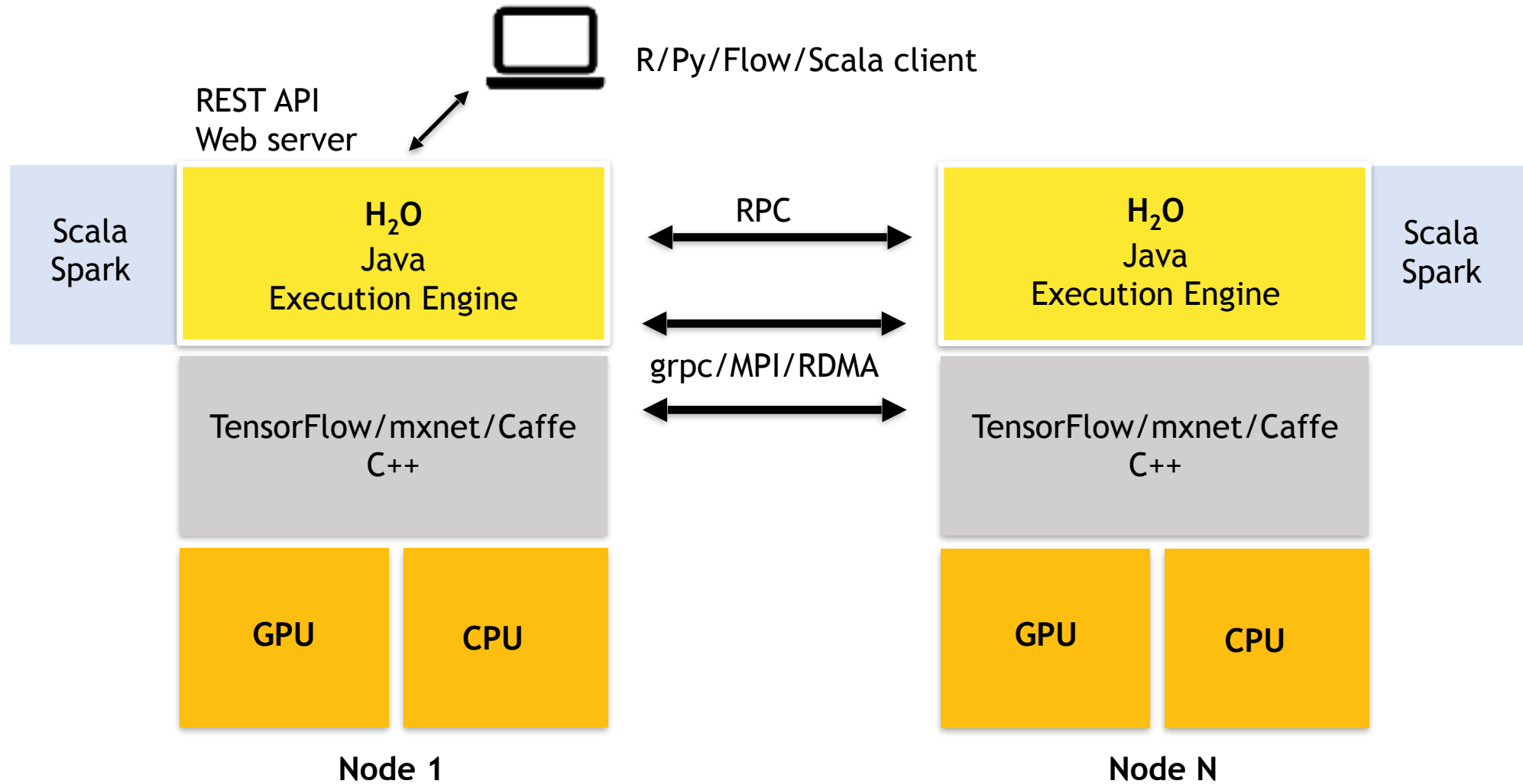


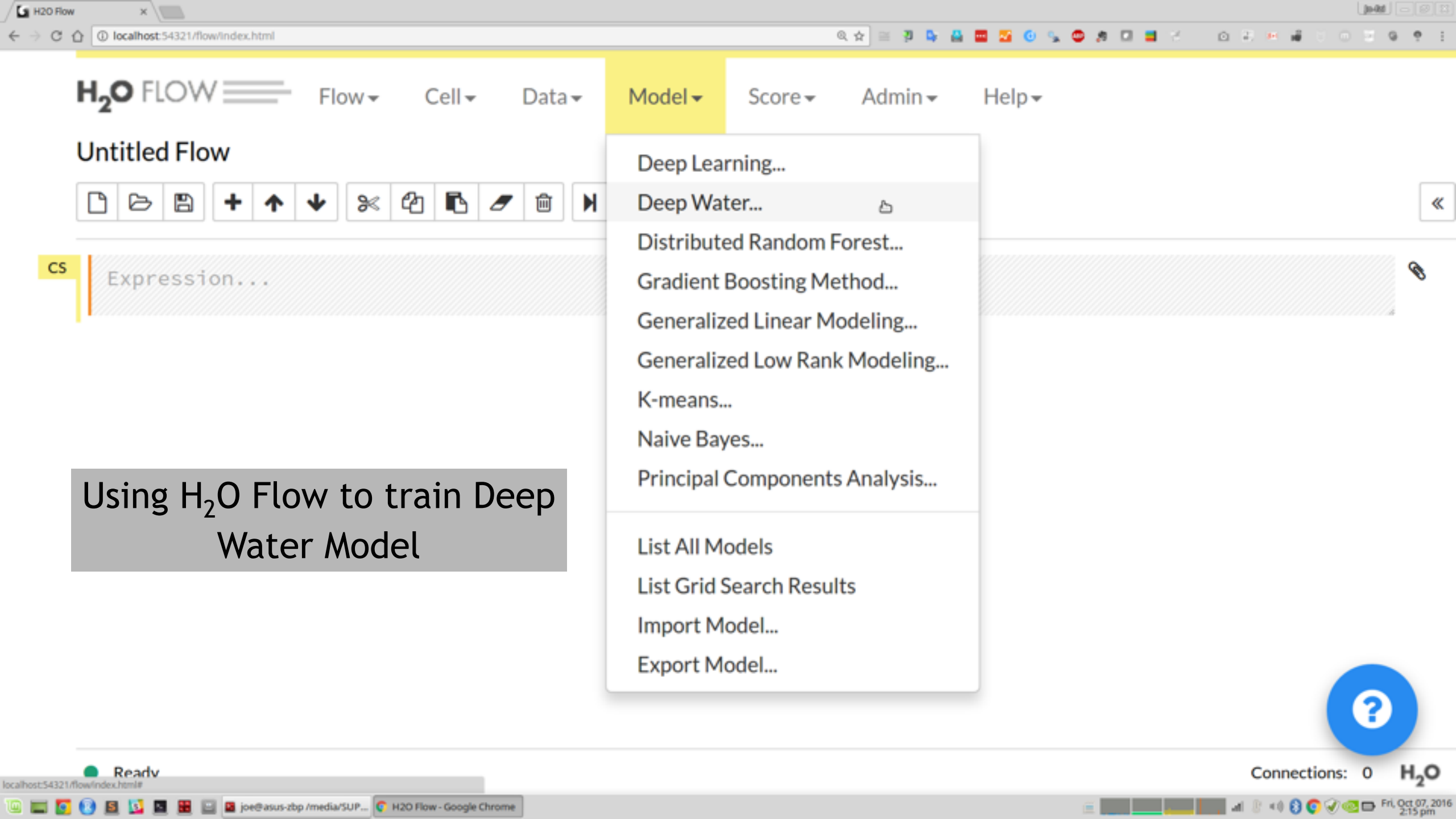
Recurrent Neural Networks enabling **natural language processing, sequences, time series, and more**



Hybrid Neural Network Architectures enabling **speech to text translation, image captioning, scene parsing** and more

# Deep Water Architecture





Using H<sub>2</sub>O Flow to train Deep Water Model

# Same H2O R/Python Interface

To build a LeNet image classification model in H2O, simply specify network = "lenet":

```
model <- h2o.deeppwater(x=path, y=response,  
                        training_frame=df, epochs=50,  
                        learning_rate=1e-3, network = "lenet")  
model
```

```
|-----| 100%
```

Model Details:

-----

H2OMultinomialModel: deeppwater

Model ID: DeepWater\_model\_R\_1477378962436\_2

Status of Deep Learning Model: lenet, 1.6 MB, predicting C2, 3-class classification, 14,335 training samples, mini-batch size 32

	input_neurons	rate	momentum
1	2352	0.000900	0.990000

H2OMultinomialMetrics: deeppwater

\*\* Reported on training data. \*\*

\*\* Metrics reported on full training frame \*\*

Training Set Metrics:

-----

Extract training frame with `h2o.getFrame("cat\_dog\_mouse.hex\_sid\_95f8\_1")`

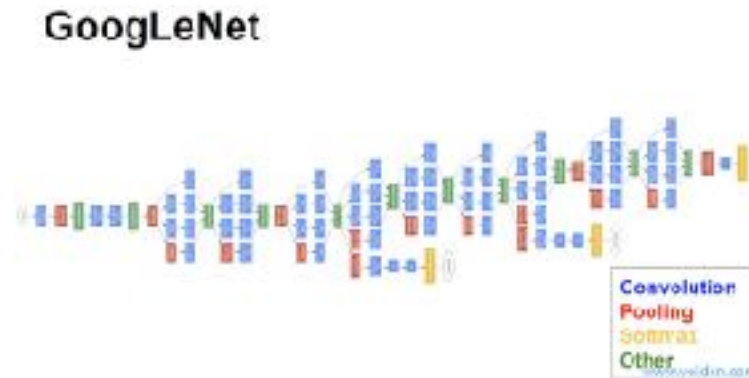
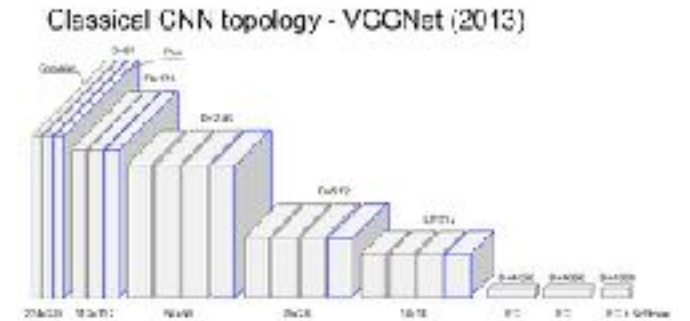
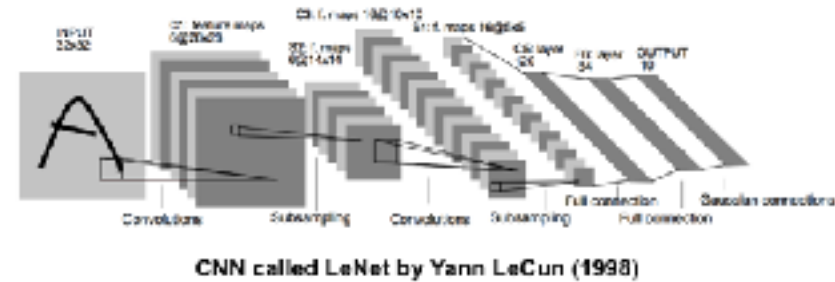
MSE: (Extract with `h2o.mse`) 0.131972

RMSE: (Extract with `h2o.rmse`) 0.3628386

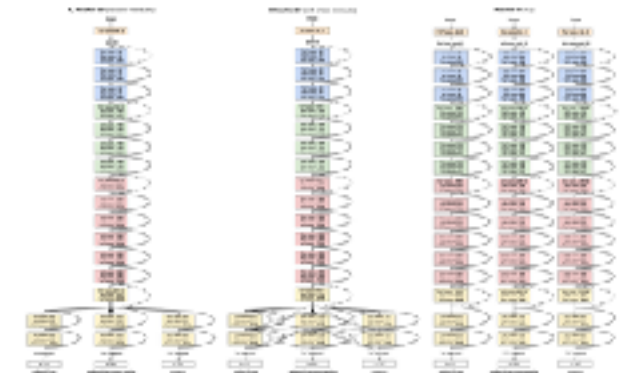
Logloss: (Extract with `h2o.logloss`) 0.4176429

# Available Networks in Deep Water

- LeNet (This Demo)
- AlexNet
- VGGNet
- Inception  
(GoogLeNet)
- ResNet (Deep  
Residual Learning)
- Build Your Own



# ResNet



# Want to try Deep Water?

- Build it
  - [github.com/h2oai/deepwater](https://github.com/h2oai/deepwater)
  - Ubuntu 16.04
  - CUDA 8
  - cuDNN 5
  - ...
- Pre-built Amazon Machine Images (AMIs)
  - Info to be confirmed

## Python/R Jupyter Notebooks

Check out a sample of cool Deep Learning Jupyter notebooks!

## PreRelease Downloads

For the following system dependencies, we provide recent builds for your convenience.

- Ubuntu 16.04 LTS
- Latest NVIDIA Display driver
- CUDA 8 (latest available) in /usr/local/cuda
- cuDNN 5 (inside of lib and include directories in /usr/local/cuda/)

In the future, we'll have more pre-built jars for more OS/CUDA combinations.

- Required to run Jupyter notebook: [H2O Deep Water enabled Python module](#) -- install via `pip install <file>`
- To build custom network: [Matching MXNet Python egg](#) -- install via `easy_install <file>`
- To run from Flow only: [H2O Standalone h2o.jar](#) -- launch via `java -jar h2o.jar`

If you are interested in running H2O Deep Water on a different infrastructure, see the [DIY build instructions](#) below



# H<sub>2</sub>O's Mission



## Making Machine Learning Accessible to Everyone

*Photo credit: Virgin Media*