# H2O Driverless AI

Webinar

H2O.ai, Inc, Jul 27 2017

# Why Driverless?

- Recipes for Problem Solving
- Automatic Feature Engineering
- Fast GPU Editions of ML and DL -- Allows for automation.
- Model Interpretation
- Automatic Visualization
- Who needs Driverless AI?

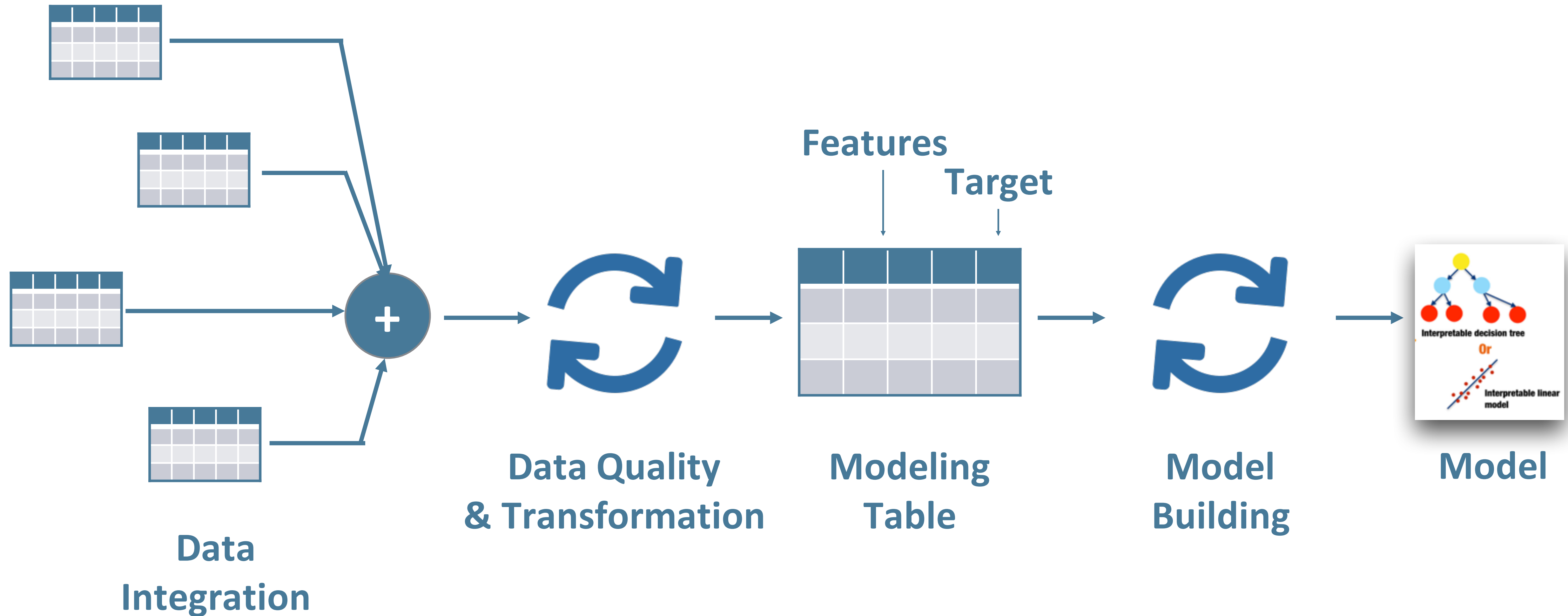# Shortage of Data Scientists

"The United States alone faces a shortage of 140,000 to 190,000 people with analytical expertise and 1.5 million managers and analysts"

–McKinsey Prediction for 2018

Competitions    Kernels    Discussion                                    Learn more about rankings ›

92 Grandmasters    868 Masters    2,489 Experts    45,517 Contributors    13,156 Novices

H₂O.ai

# Typical Enterprise Machine Learning Workflow



**Features**

**Target**

**Data Quality & Transformation**

**Modeling Table**

**Model Building**

**Model**

**Data Integration**

# Typical Enterprise Machine Learning Workflow



**Data Integration**

**Data Quality & Transformation**

Features

Target

**Modeling Table**

**Model Building**

**Model**

Interpretable decision tree

Or

Interpretable linear model

**Driverless AI**

H2O.ai

# Driverless AI

**Data**

**VISUAL MODEL INTERPRETATION**

English Explanations
Reason Codes
K-Lime, LOCO, Partial Dependence

Recipes
**AUTO DL**

Kaggle Grandmaster in a Box
Automatic Feature Engineering
Pipeline Export

**H₂O**

Auto ML - Tuning + Ensembles
Deep Learning
Algorithms
data.table Munging

**Distributed Multi-CPU Multi-GPU**

**Model Repository**

**Deploy**

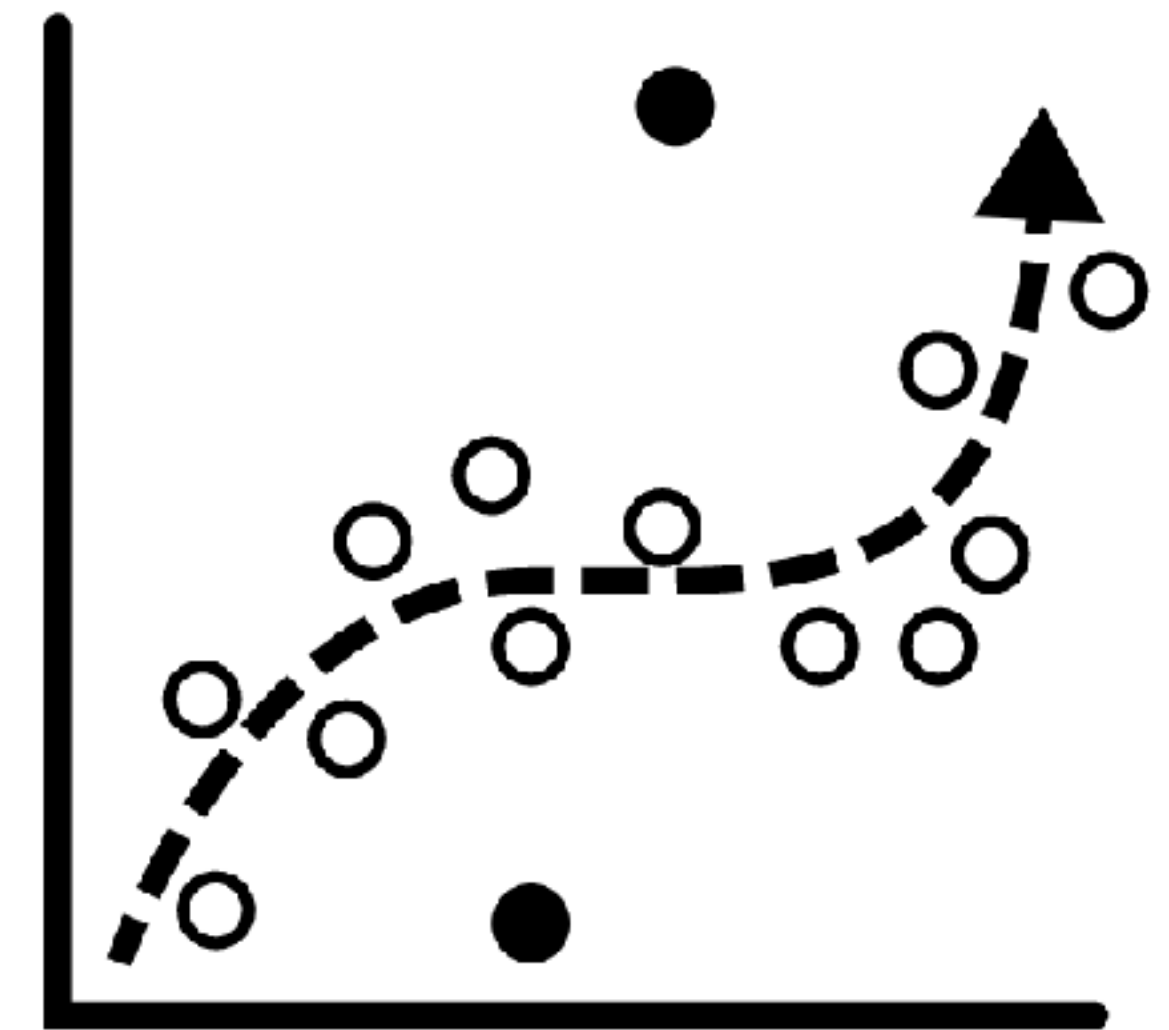**MODEL FITNESS**

H₂O.ai
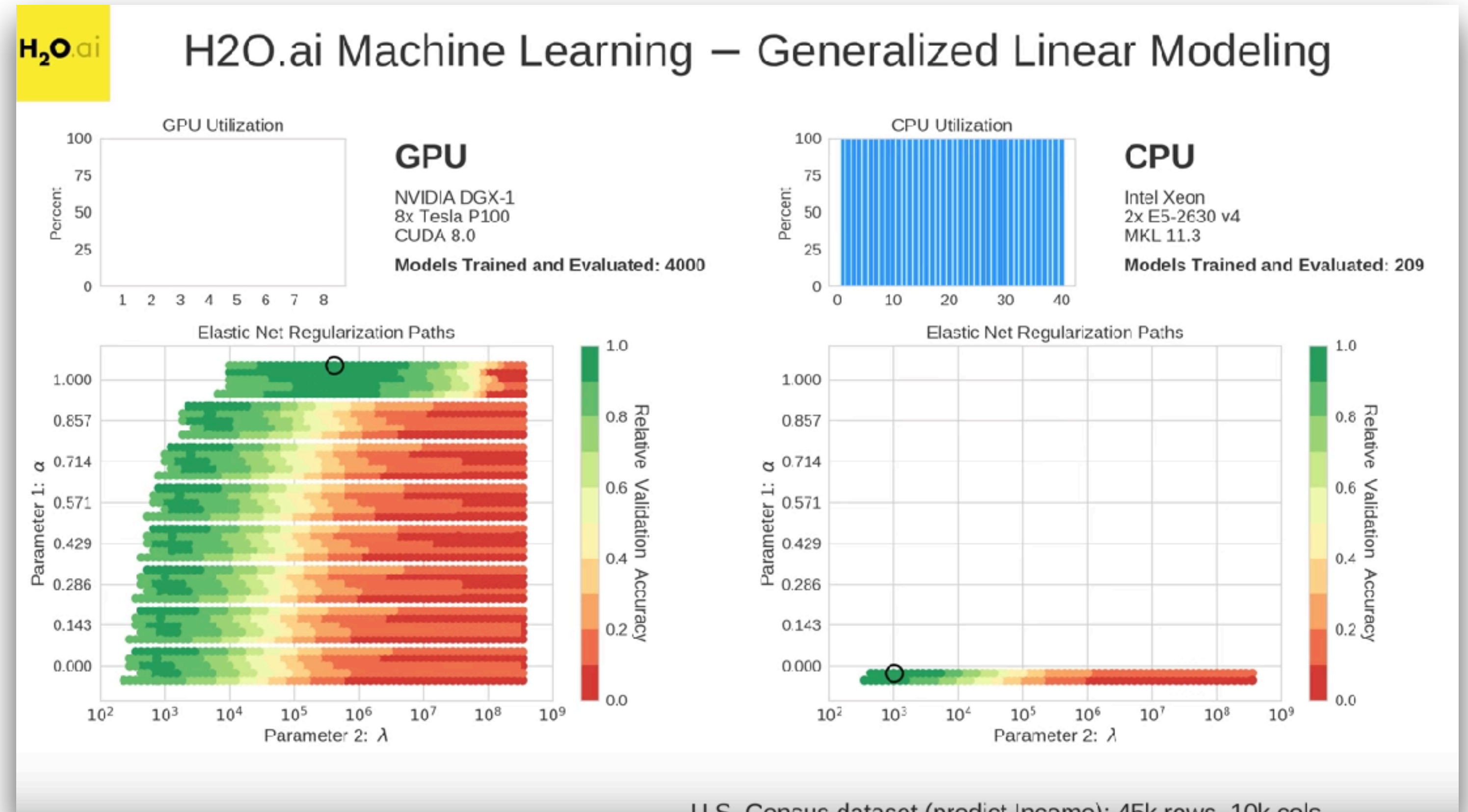
# 3 Pillars



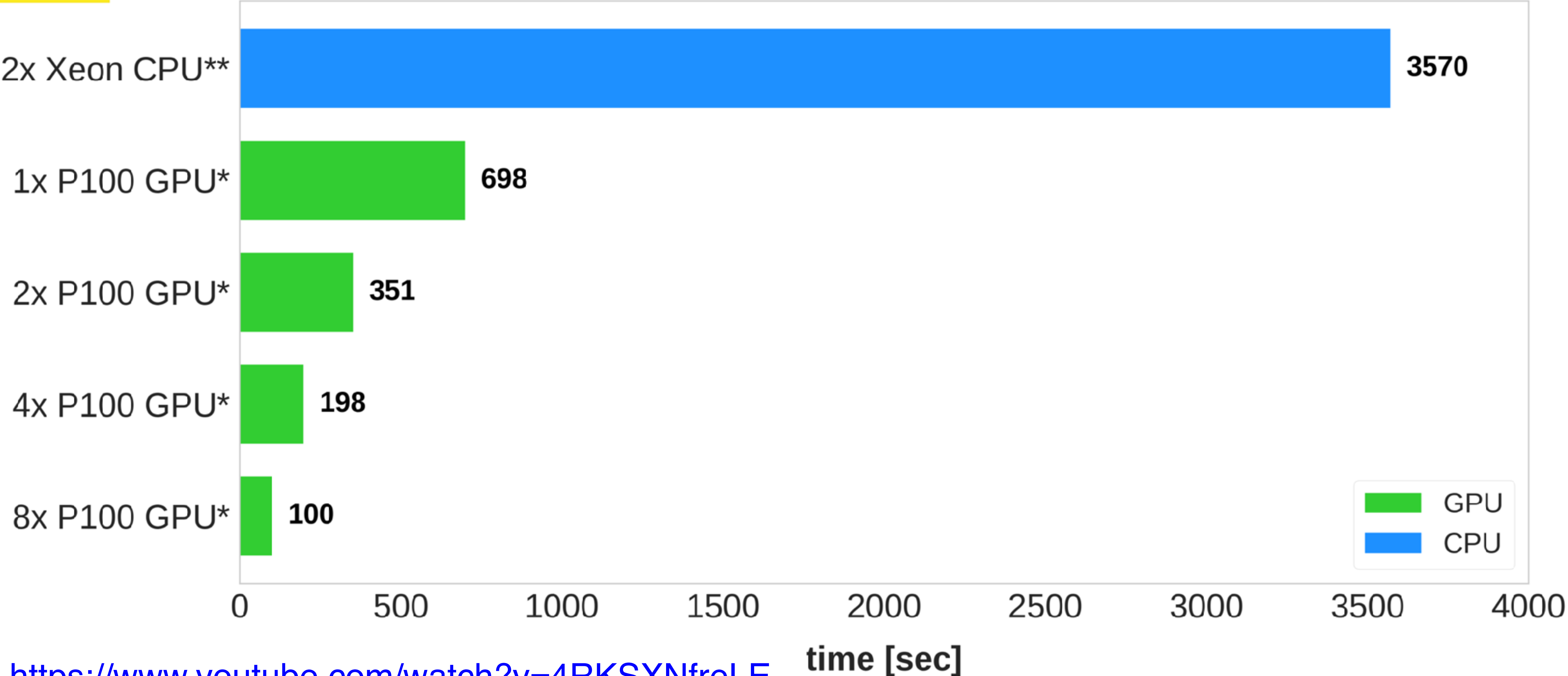**Speed**          **Accuracy**          **Interpretability**

# Speed

- GPU acceleration to achieve up to 40x speedups vs CPU

- Multi GPU - XGBoost, GLM, K-Means and more

- Achieve best performance in shortest time



**H2O.ai Machine Learning – Generalized Linear Modeling**

GPU
NVIDIA DGX-1
8x Tesla P100
CUDA 8.0
**Models Trained and Evaluated: 4000**

CPU
Intel Xeon
2x E5-2630 v4
MKL 11.3
**Models Trained and Evaluated: 209**

U.S. Census dataset (predict Income): 45k rows, 10k cols

# H2O.ai Machine Learning − Generalized Linear Modeling
## Time to Train and Evaluate 4000 Models

| Configuration | time [sec] |
|---|---|
| 2x Xeon CPU** | 3570 |
| 1x P100 GPU* | 698 |
| 2x P100 GPU* | 351 |
| 4x P100 GPU* | 198 |
| 8x P100 GPU* | 100 |

Legend: GPU (green), CPU (blue)

https://www.youtube.com/watch?v=4RKSXNfreLE
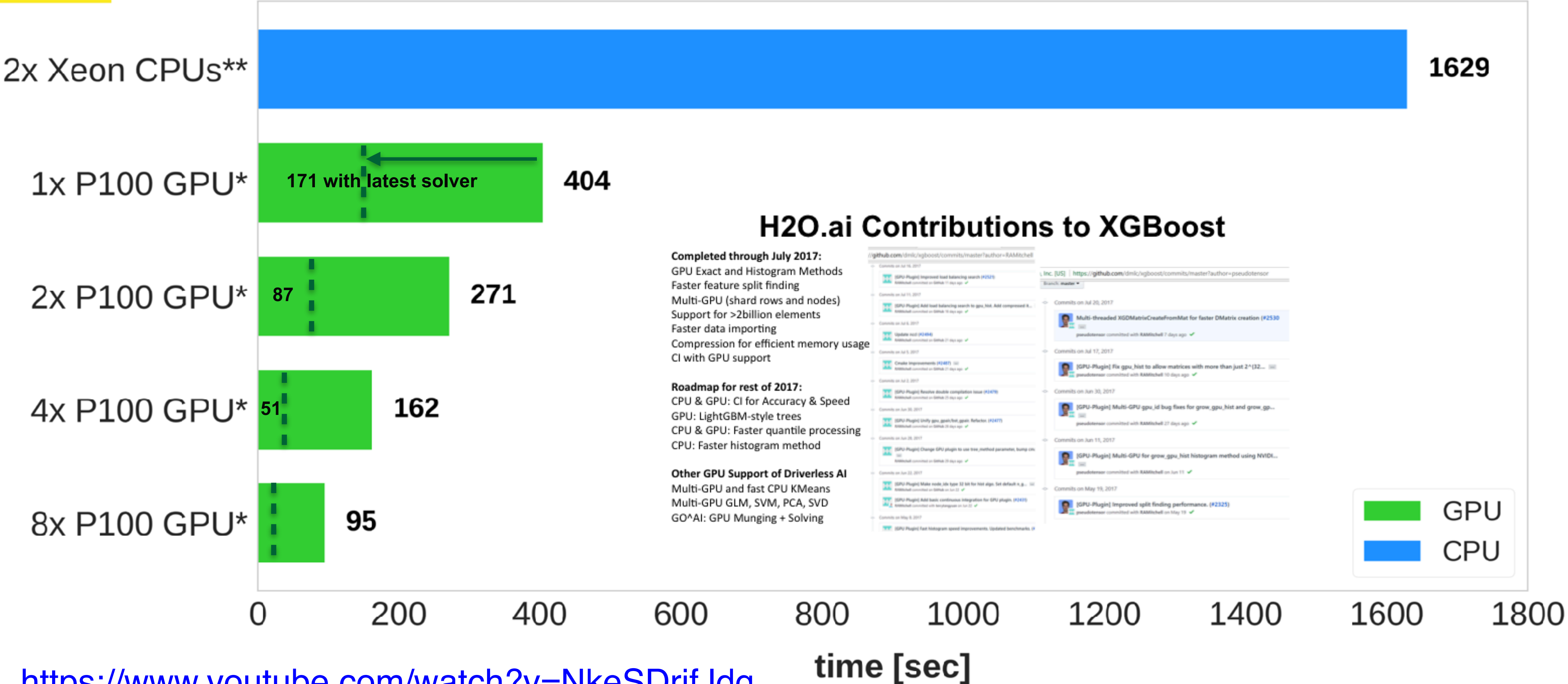
http://github.com/h2oai/perf/

*NVIDIA DGX-1, **Dual Intel Xeon E5-2630 v4
U.S. Census dataset (predict Income): 45k rows, 10k cols
Elastic Net Model Parameters: 5-fold cross-validation, $\alpha = \{\frac{i}{7}, i = 0...7\}$, full $\lambda$-search

# H2O.ai Machine Learning − Gradient Boosting Machine
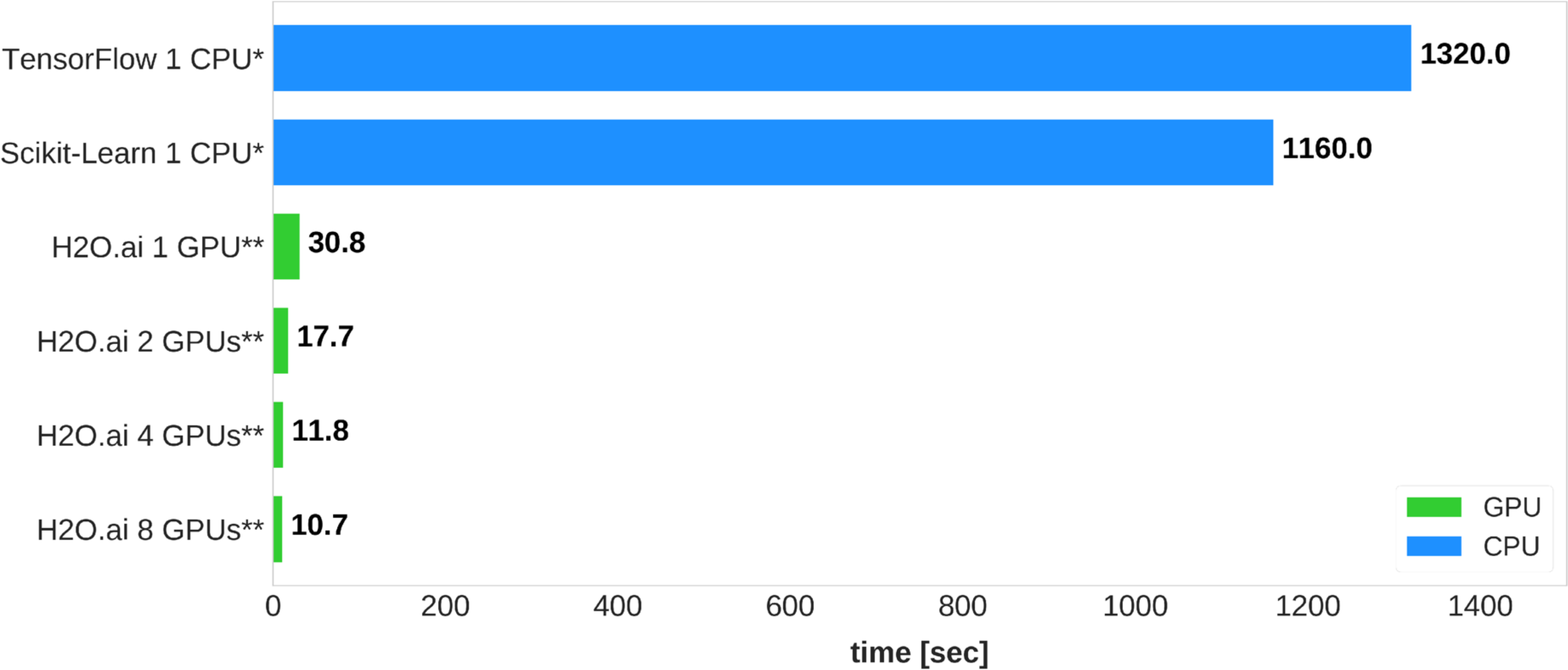## Time to Train 16 H2O XGBoost Models



| Configuration | Time [sec] |
|---|---|
| 2x Xeon CPUs** | 1629 |
| 1x P100 GPU* | 404 (171 with latest solver) |
| 2x P100 GPU* | 271 (87) |
| 4x P100 GPU* | 162 (51) |
| 8x P100 GPU* | 95 |

Legend: GPU (green), CPU (blue)

### H2O.ai Contributions to XGBoost

**Completed through July 2017:**
GPU Exact and Histogram Methods
Faster feature split finding
Multi-GPU (shard rows and nodes)
Support for >2billion elements
Faster data importing
Compression for efficient memory usage
CI with GPU support

**Roadmap for rest of 2017:**
CPU & GPU: CI for Accuracy & Speed
GPU: LightGBM-style trees
CPU & GPU: Faster quantile processing
CPU: Faster histogram method

**Other GPU Support of Driverless AI**
Multi-GPU and fast CPU KMeans
Multi-GPU GLM, SVM, PCA, SVD
GO^AI: GPU Munging + Solving

https://www.youtube.com/watch?v=NkeSDrifJdg

http://github.com/h2oai/perf/

*NVIDIA DGX-1, **Dual Intel Xeon E5-2630 v4
Higgs dataset (binary classification): 1M rows, 29 cols; max_depth: {6,8,10,12}, sample_rate: {0.7,0.8,0.9,1.0}

# H2O.ai Machine Learning − k-Means Clustering

## Time to run 1000 Lloyds iterations for k=1000 clusters

| Method | time [sec] |
|---|---|
| TensorFlow 1 CPU* | 1320.0 |
| Scikit-Learn 1 CPU* | 1160.0 |
| H2O.ai 1 GPU** | 30.8 |
| H2O.ai 2 GPUs** | 17.7 |
| H2O.ai 4 GPUs** | 11.8 |
| H2O.ai 8 GPUs** | 10.7 |

Legend: GPU, CPU

Kaggle Homesite Home Insurance Claims Predictions Dataset (261k rows, 298 cols)
k-Means Clustering (Lloyds), random initialization, 1000 centroids, 1000 iterations
Hardware: *Intel i7 5820k (6-core), **NVIDIA Tesla P100 (DGX-1)

http://github.com/h2oai/perf/

# Accuracy

- Automatic feature engineering to increase accuracy - AlphaGo for AI

- Automatic Kaggle Grandmaster recipes in a box for solving wide variety of use-cases

- Automatic machine learning to find and tune the right ensemble of models

**Relative Error: Lower is Better**

Driverless AI



Preliminary results - untuned, single model

H₂O.ai

# What's the "Secret" to Data Science?

"Coming up with features is difficult, time-consuming, requires expert knowledge. "Applied machine learning" is basically feature engineering."

–Andrew Ng

"… some machine learning projects succeed and some fail. What makes the difference? Easily the most important factor is the features used."

–Pedro Domingos

"Good data preparation and feature engineering is integral to better predictions."

–Kazanova, H2O.ai Kaggle Grandmaster #2 (former #1)

H₂O.ai

# Why is it so difficult to be a good Data Scientist?

Data matters: need access to business-relevant data
and need domain knowledge about how features
interact with each other       *Data Science Recipes*

Powerful feature transformations (like target encoding) can lead to
overfit models if done wrong, need strong math & statistics skills

Need to run thousands of experiments to reach robust
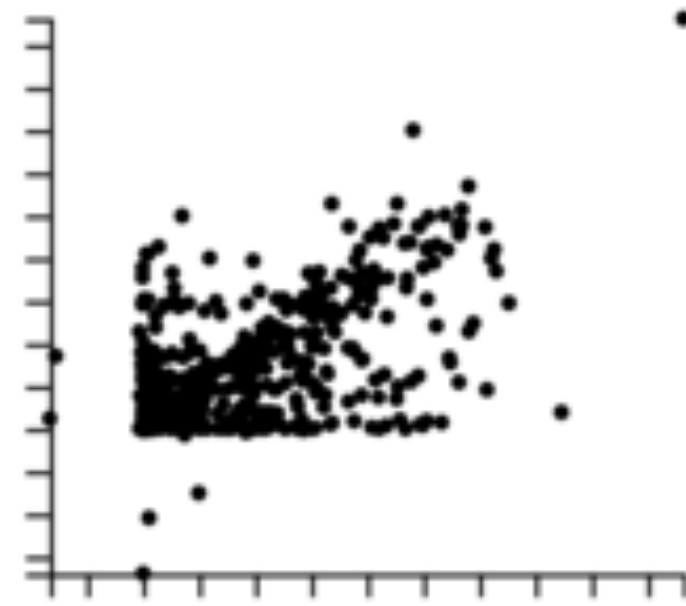conclusions, need computer science skills and access to
compute hardware       *GPUs, Automation*

**H₂O**.ai

# Interpretability

- Interpretability for debugging, not just for regulators

- Get reason codes and model interpretability in plain english

- K-Lime, LOCO, partial dependence and more
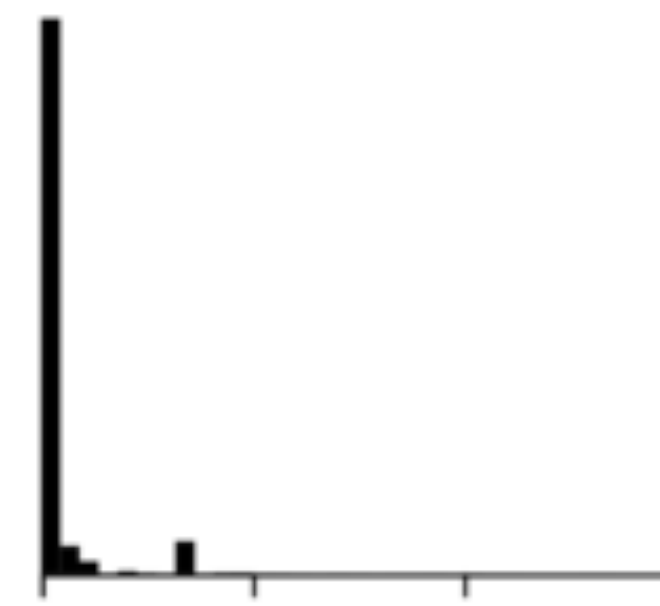
# Automatic Visualization

H2O.ai

# Next: Live Demo!
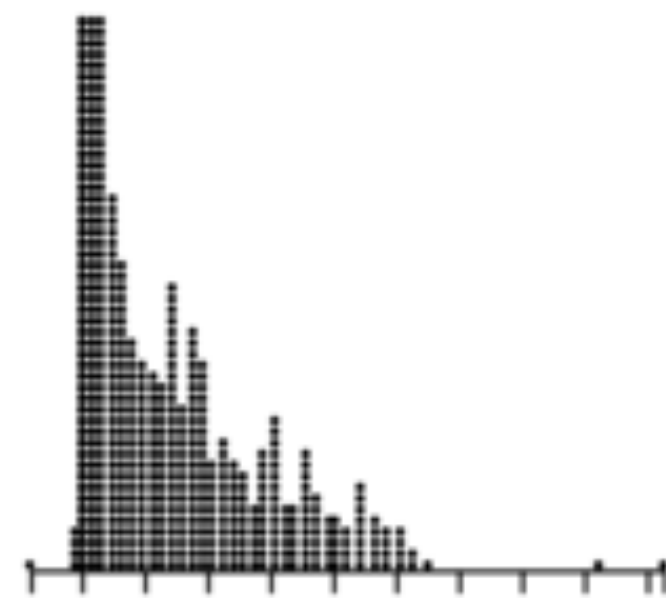
# Installation and Start

```
# Load the Driverless AI docker image
docker load < downloaded-h2oai-driverless-ai-image.tar
```

```
# Start the Driverless AI docker image
nvidia-docker run --rm -u `id -u`:`id -g` -p 12345:12345 -v `pwd`/data:/data -v `pwd`/
 ↪log:/log opsh2oai/h2oai-runtime
```

**H₂O**.ai

# The Team Behind Driverless AI



**AutoDL**
Marios Michailidis, Dmitry Larko, Branden Murray, Mark Landry

**Datatable**
Matt Dowle, Pasha Stetsenko

**MLI**
Patrick Hall, Mark Chan, Navdeep Gill

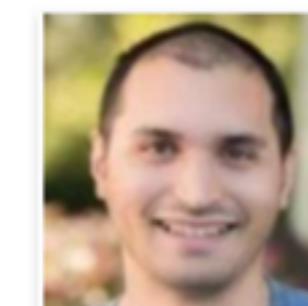**GPU Acceleration**
Arno Candel, Jon Mckinney, Rory Mitchell

**AutoViz**
Leland Wilkinson, Prithvi Prabhu, Justin Loyola

**H₂O**.ai

# Upcoming Webinars

- Machine Learning Interpretability - Patrick Hall - **Aug 17**

- AutoDL - Dmitry Larko - **Aug 24**

- Automatic Visualization - **Sep (TBD)**