

# Deep Water



Jo-fai (Joe) Chow  
Data Scientist  
[joe@h2o.ai](mailto:joe@h2o.ai)  
@matlabulous

Data Science Milan  
Politecnico di Milano  
10<sup>th</sup> October, 2016

# Agenda

- First Talk (25 mins)
  - About H2O.ai
  - Demo
    - A Simple Classification Task
    - H2O's Web Interface
  - Why H2O?
    - Our Community
    - Our Customers
  - What's Next?
    - New H2O Features
- Second Talk (25 mins)
  - H2O for IoT
    - Predictive Maintenance
    - Anomaly Detection
    - H2O's R Interface
- Third Talk (25 mins)
  - Deep Water
  - Demo
    - H2O + mxnet on GPU
    - H2O's Python Interface

# Deep Learning in H2O



# H2O Overview

Computer Science (CS)

Artificial Intelligence (A.I.)

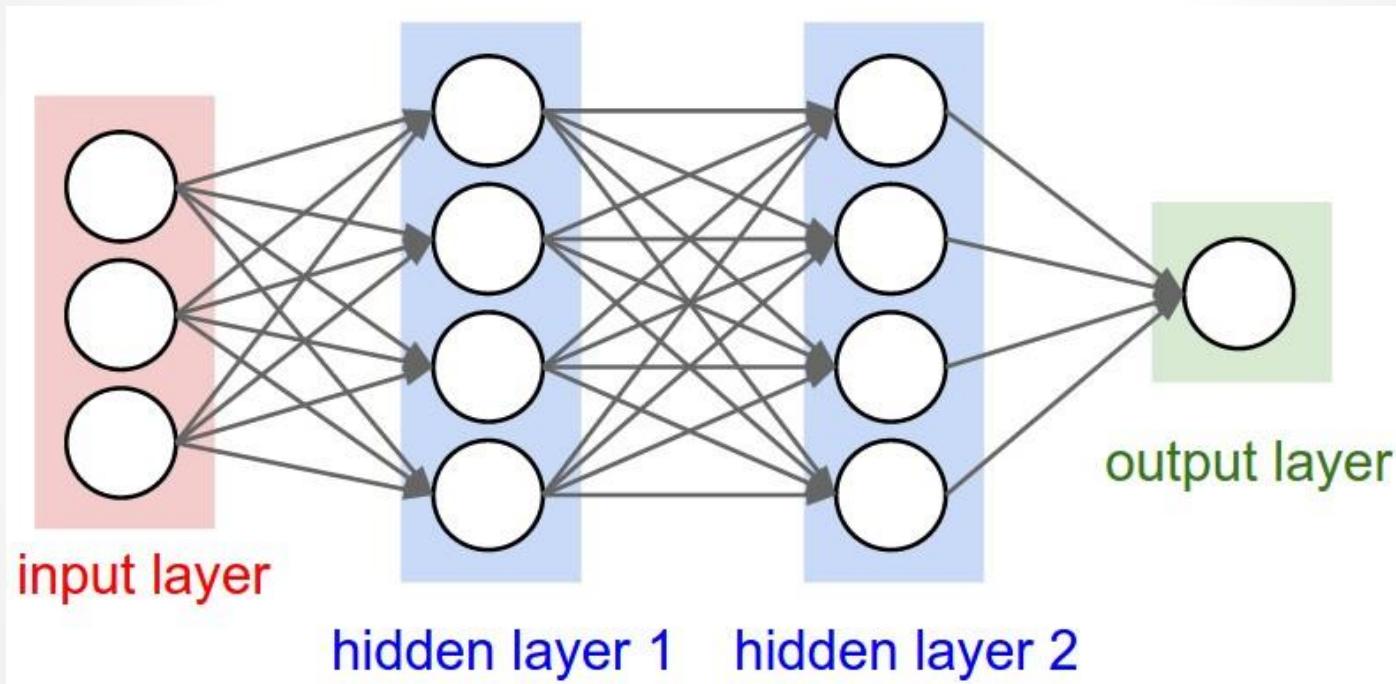
Machine Learning (ML)

Deep Learning (DL)

hot hot hot hot hot



# A Simple Neural Network



# H2O 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
116695259	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

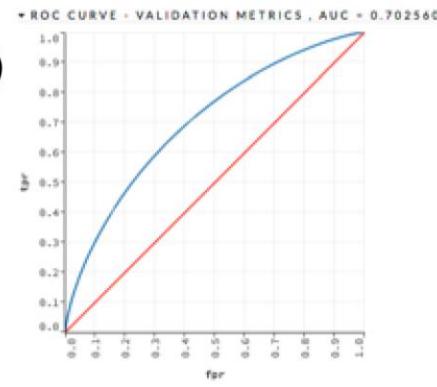
Actions View Cancel Job

\* OUTPUT - STATUS OF NEURON LAYERS (PREDICTING ISDEPDELAYED, 2-CLASS CLASSIFICATION, BERNoulli DISTRIBUTION, CROSSENTROPY LOSS, 17,462 WEIGHTS/BIASES, 221.3 KB, 106,585,365 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	897	Input	0									
2	20	Rectifier	0	0	0	0.0463	0.2020	0	-0.0021	0.2111	-0.9139	1.0036
3	20	Rectifier	0	0	0	0.0197	0.2027	0	-0.1053	0.5362	-1.3908	1.5259
4	20	Rectifier	0	0	0	0.0517	0.0446	0	-0.1575	0.3068	-0.0846	0.0046
5	20	Rectifier	0	0	0	0.0761	0.0844	0	-0.0374	0.2275	-0.2647	0.2481
6	2	Softmax	0	0	0	0.0161	0.0083	0	0.0741	0.7260	0.4269	0.2056

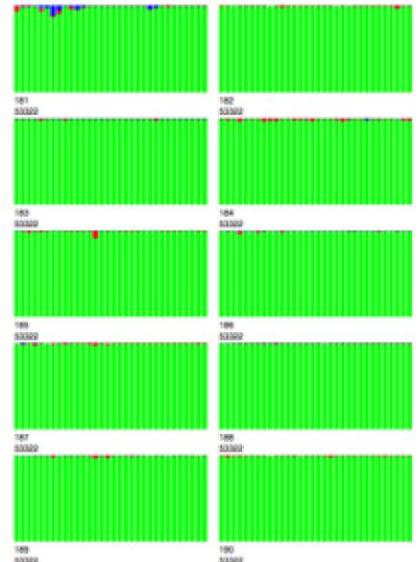
H<sub>2</sub>O.ai

Deep Learning Model



Threshold: Choose... Criterion: Choose...

## VARIABLE IMPORTANCES



## Legend

Each bar represents one CPU.

Blue: idle time

Green: user time

Red: system time

White: other time (e.g. Vo)

10 nodes: all 320 cores busy



real-time, interactive model inspection in Flow

# H2O Deep Learning Community Quotes

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

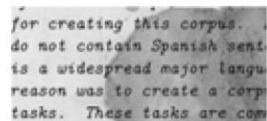
Triskelion | 01.02.2015

[READ MORE](#)

## Kaggle challenge 2nd place winner Colin Priest

[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.”



Completed • Knowledge • 161 teams

## Denoising Dirty Documents

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

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



# Why Deep Water?



# Deep Water: Best Open-Source Deep Learning

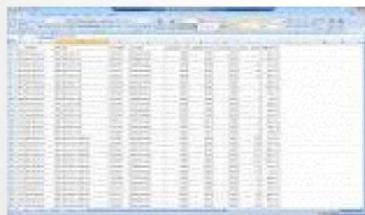
## Enterprise Deep Learning for Business Transformation

<b>Deep Water = THE Deep Learning Platform</b>	H2O integrates the top open-source DL tools	
<b>Native GPU support</b>	  is up to 100x faster than	
<b>Enterprise Ready</b>	Easy to train and deploy, interactive, scalable, etc. Flow, R, Python, Spark/Scala, Java, REST, POJO, <b>Steam</b>	
<b>New Big Data Use Cases (previously impossible or difficult in H2O)</b>	<b>Image</b> - social media, manufacturing, healthcare, ... <b>Video</b> - UX/UI, security, automotive, social media, ... <b>Sound</b> - automotive, security, call centers, healthcare, ... <b>Text</b> - NLP, sentiment, security, finance, fraud, ... <b>Time Series</b> - security, IoT, finance, e-commerce, ...	

# Deep Water opens the Floodgates for state-of-the-art Deep Learning

## H2O Deep Learning: simple multi-layer neural networks

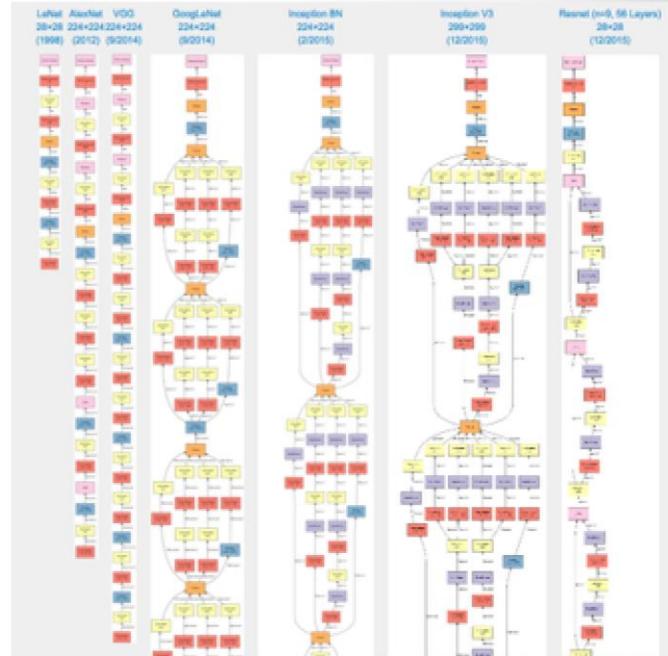
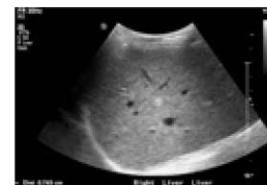
1-5 layers  
MBs/GBs of data



Limited to business analytics,  
statistical models (CSV data)

## Deep Water: deep complex networks

5-1000 layers  
GBs/TBs of data



Large networks for big data  
(e.g. image 1000x1000x3 -> 3m inputs)

## Current Contributors (more H2O.ai folks joining soon)



Fabrizio Milo



Cyprien Noel



Qiang Kou



Arno Candel



Caffe



This repository

mxnet



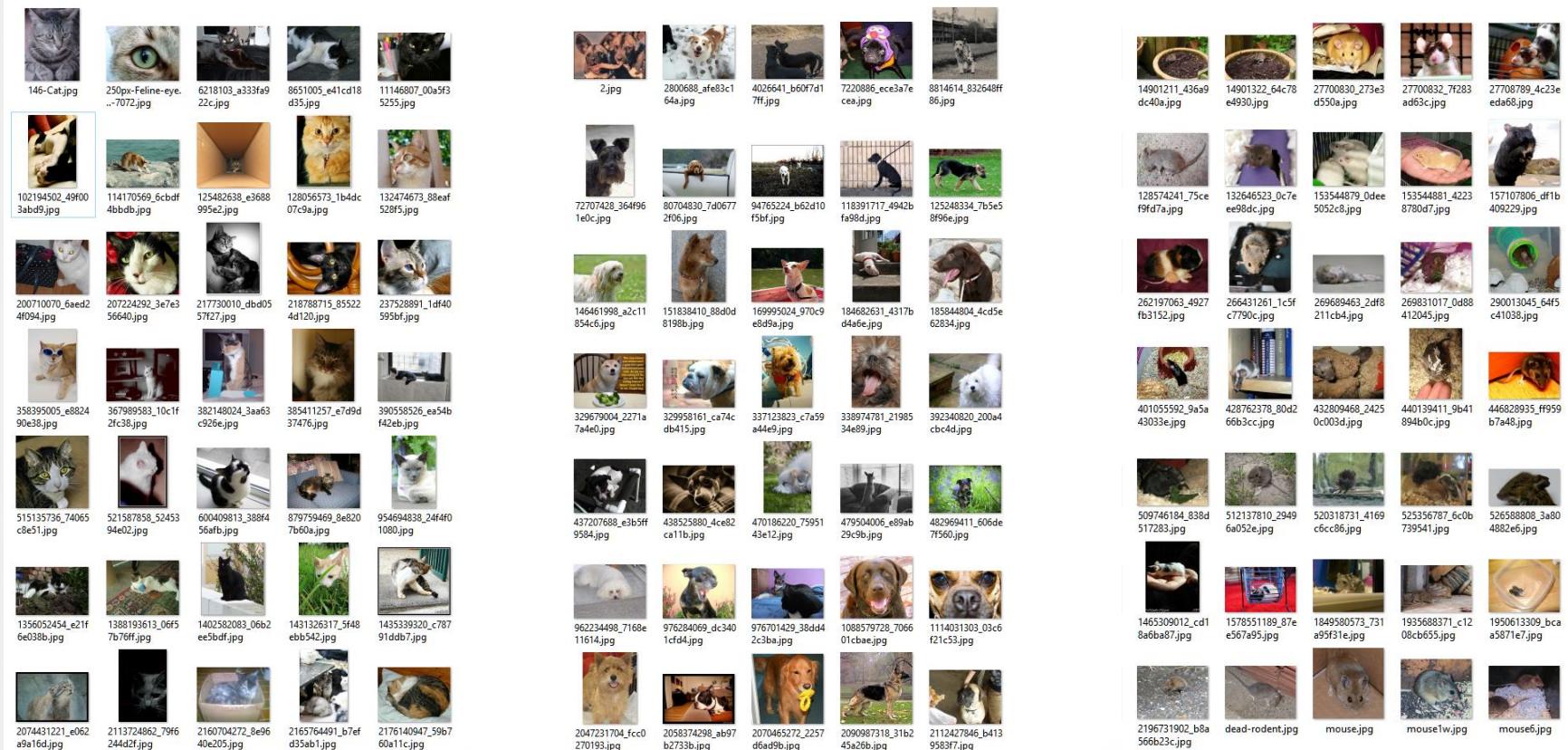
# Deep Water Demo



# Deep Water Demo

- H2O + mxnet
  - Dataset:
    - Cat / Dog / Mouse
  - H2O Python interface
  - mxnet GPU backend
  - Train a LeNet (CNN) model
  - Explore model in Flow
- Code and Data
  - [bit.ly/h2o\\_milan\\_1](http://bit.ly/h2o_milan_1)
  - subfolder
    - deep\_water\_demo

# Data – Cat/Dog/Mouse Images



# Data - CSV

	A	B
1	bigdata/laptop/deepwater/imagenet/cat/102194502_49f003abd9.jpg	cat
2	bigdata/laptop/deepwater/imagenet/cat/11146807_00a5f35255.jpg	cat
3	bigdata/laptop/deepwater/imagenet/cat/1140846215_70e326f868.jpg	cat
4	bigdata/laptop/deepwater/imagenet/cat/114170569_6cbdf4bbdb.jpg	cat
5	bigdata/laptop/deepwater/imagenet/cat/1217664848_de4c7fc296.jpg	cat
6	bigdata/laptop/deepwater/imagenet/cat/1241603780_5e8c8f1ced.jpg	cat
7	bigdata/laptop/deepwater/imagenet/cat/1241612072_27ececbeef.jpg	cat
8	bigdata/laptop/deepwater/imagenet/cat/1241613138_ef1d82973f.jpg	cat
9	bigdata/laptop/deepwater/imagenet/cat/1244562192_35becd66bd.jpg	cat
10	bigdata/laptop/deepwater/imagenet/cat/125482638_e3688995e2.jpg	cat
11	bigdata/laptop/deepwater/imagenet/cat/128056573_1b4dc07c9a.jpg	cat
12	bigdata/laptop/deepwater/imagenet/cat/12945197_75e607e355.jpg	cat
13	bigdata/laptop/deepwater/imagenet/cat/132474673_88eaf528f5.jpg	cat
14	bigdata/laptop/deepwater/imagenet/cat/1350530984_ecf3039cf0.jpg	cat
15	bigdata/laptop/deepwater/imagenet/cat/1351606235_c9fbebf634.jpg	cat
16	bigdata/laptop/deepwater/imagenet/cat/1356052454_e21f6e038b.jpg	cat
17	bigdata/laptop/deepwater/imagenet/cat/1388193613_06f57b76ff.jpg	cat

# H2O + mxnet Demo



## Deep Water Demo: H2O + mxnet GPU backend

Original reference: ([Link](#))

### Introduction - What can Deep Water do?

- Train user-defined or pre-defined deeplearning models for image classification from Flow, R, Python, Java, Scala or REST API
- Train on a single GPU (requires CUDA) or CPU (requires BLAS)
- Uses the MXNet backend transparently
- Behave just like any other H2O model (Flow, cross-validation, early stopping, hyper-parameter search, etc.)

### Prerequisite

- Install Ubuntu 16.04 LTS
- Install the latest NVIDIA Display driver
- Install CUDA 8 (latest available) in /usr/local/cuda
- Install CUDNN 5 (to lib and include directories in /usr/local/cuda/)
- Obtain GPU-enabled h2o.jar (preview: <https://slack-files.com/T0329MHH6-F2GQ0B72S-bb15ff7626>) - not strictly necessary, as h2o.jar is also in the python module below, but done here for simplicity (manual launch below)
- Obtain Deep Water edition of H2O's python module (preview: <https://slack-files.com/T0329MHH6-F2GQH4D34-8d9295e775>), install with `sudo pip install h2o*.whl`
- Optional (only for custom networks) - Obtain mxnet python egg (preview: <https://slack-files.com/T0329MHH6-F2C7LQWMR-6b78dfab1a>), install with `sudo easy_install <egg-file>`
- Set environment variables: `export CUDA_PATH=/usr/local/cuda` and `export LD_LIBRARY_PATH=$CUDA_PATH/lib64:$LD_LIBRARY_PATH`
- Run `java -jar h2o.jar`
- Download dataset ([https://h2o-public-test-data.s3.amazonaws.com/bigdata/laptop/deepwater/imagenet/cat\\_dog\\_mouse.tgz](https://h2o-public-test-data.s3.amazonaws.com/bigdata/laptop/deepwater/imagenet/cat_dog_mouse.tgz), unpack contents into directory `./bigdata/laptop/deepwater/imagenet/`, relative to where h2o was launched)

bit.ly/h2o\_milan\_1  
Subfolder:  
**deep\_water\_demo**

### Run Demo Python Scripts

- Run demo 1 `python demo_01_lenet.py`
- (Optional) Try to inspect and/or build Deep Water model from Flow (i.e. `localhost:54321`)



```
demo_01_lenet.py x
1 from __future__ import print_function
2 import sys, os
3 sys.path.insert(1, os.path.join("../", "..", ".."))
4
5 # H2O
6 import h2o
7 from h2o.estimators.deepwater import H2ODeepWaterEstimator
8
9 # Start and connect to H2O local cluster
10 h2o.init()
11
12 # Import CSV
13 frame = h2o.import_file("bigdata/laptop/deepwater/imagenet/cat_dog_mouse.csv")
14 print(frame.head(5))
15
16 # Define LeNet model
17 model = H2ODeepWaterEstimator(epochs=300, rate=1e-3, network='lenet', score_interval=0, train_samples_per_iteration=1000)
18
19 # Train LeNet model on GPU
20 model.train(x=[0], y=1, training_frame=frame)
21 model.show()
```

H2O's Python Module

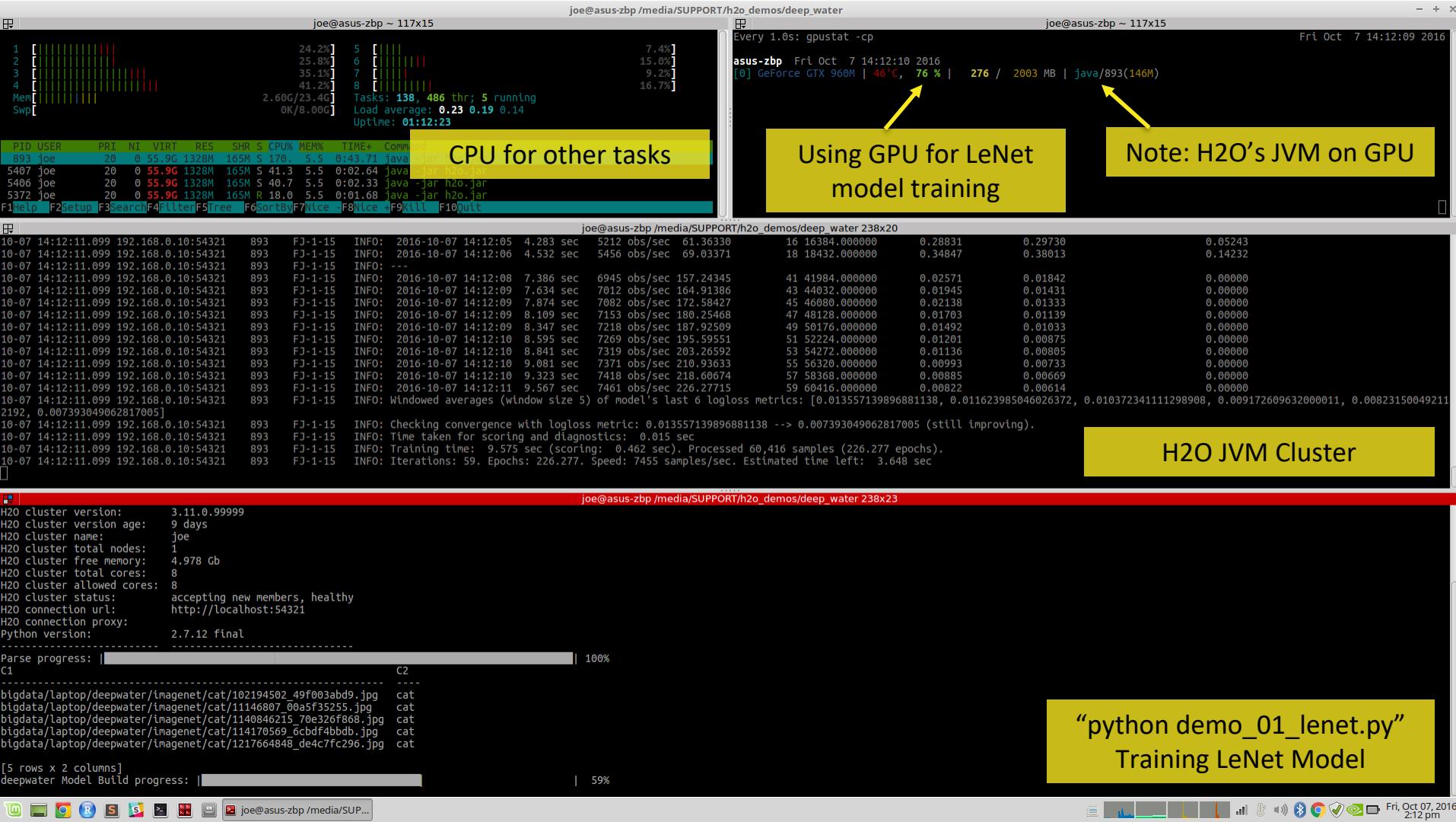
Deep Water module

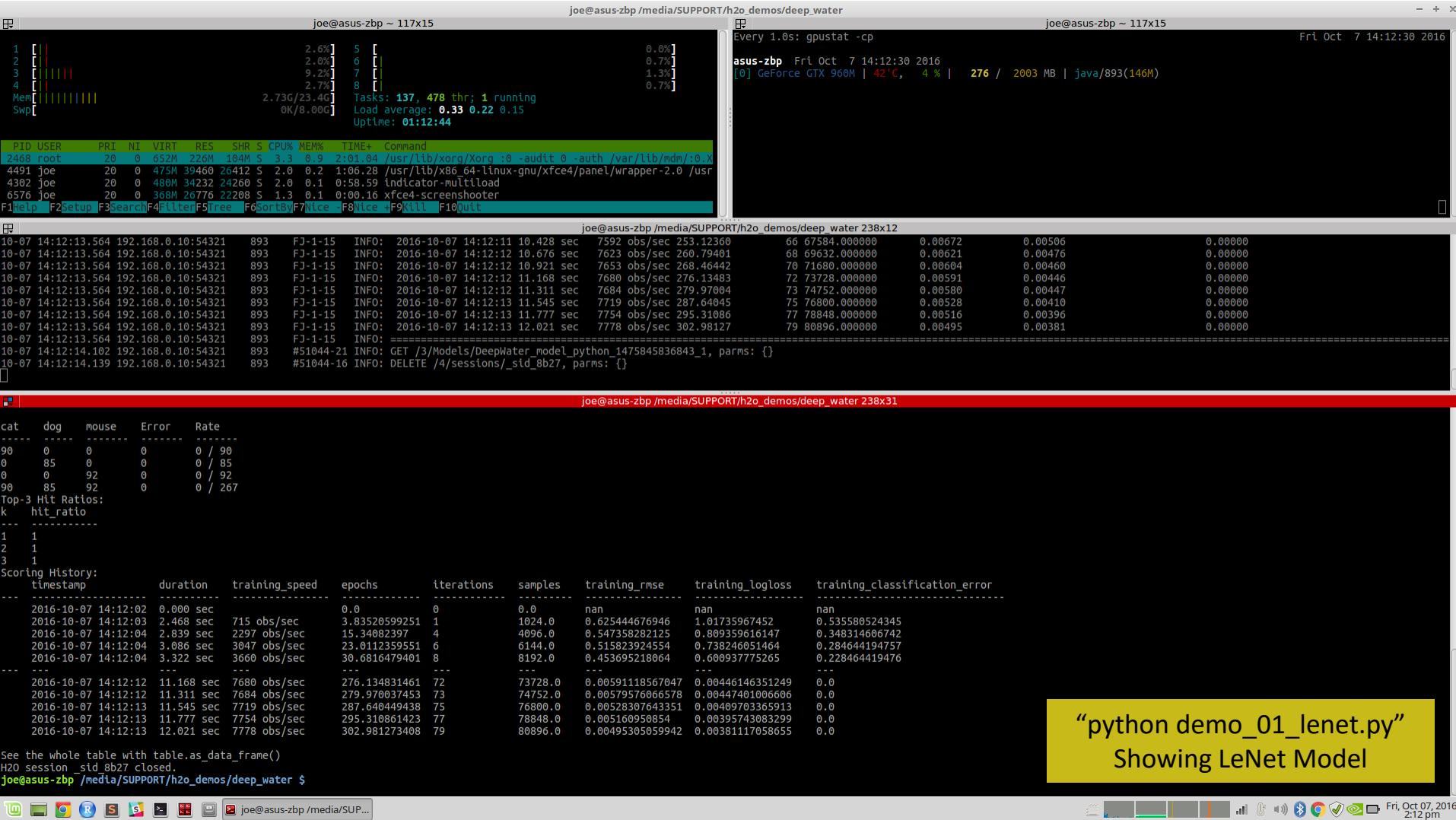
Connect to H2O Cluster

Import CSV

Define LeNet model in Deep Water

Train and show model





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

Flow ▾ Cell ▾ Data ▾ Model ▾ Score ▾ Admin ▾ Help ▾

### Untitled Flow



Expression...

getFrames

60ms

### Frames

- Type ID
- cat\_dog\_mouse.hex

Build Model... Predict... Inspect

Predict on selected frames... Delete selected frames

Using Flow (localhost:54321) to explore data frame and model

Rows 267 Columns 2 Size 18KB

CS

getModels

24ms

### Models

- Key
- DeepWater\_model\_python\_1475845836843\_1

Inspect Delete selected

Algorithm

Deep Water

Actions

Predict... Inspect

Ready

Connections: 0



## H2O FLOW

Flow ▾ Cell ▾ Data ▾ Model ▾ Score ▾ Admin ▾ Help ▾

Untitled Flow



### Model

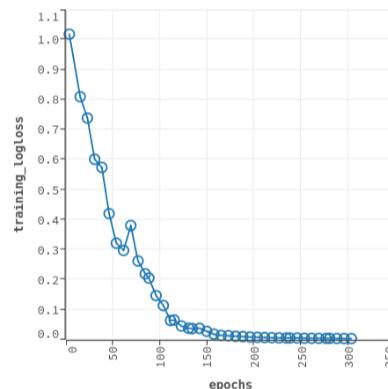
Model ID: DeepWater\_model\_python\_1475845836843\_1

Algorithm: Deep Water

Actions: Refresh Predict... Download POJO Export Inspect Delete

#### MODEL PARAMETERS

#### SCORING HISTORY - LOGLOSS



Using Flow (localhost:54321) to explore data frame and model

#### TRAINING METRICS - CONFUSION MATRIX VERTICAL: ACTUAL; ACROSS: PREDICTED

	cat	dog	mouse	Error	Rate
cat	90	0	0	0	0 / 90
dog	0	85	0	0	0 / 85
mouse	0	0	92	0	0 / 92
Total	90	85	92	0	0 / 267

Ready



Connections: 0 H2O

# H2O FLOW

Flow ▾ Cell ▾ Data ▾ Model ▾ Score ▾ Admin ▾ Help ▾

Untitled Flow



## SplitOptions

Frame:

Splits: Ratio

0.9

Key

cat\_dog\_mouse.hex\_0.90



0.10

cat\_dog\_mouse.hex\_0.10

Add a new split

Seed: 123456

Using Flow (localhost:54321) to split data and train Deep Water model

SplitOptions Create

CS

```
splitFrame "cat_dog_mouse.hex", [0.9], ["cat_dog_mouse.hex_0.90","cat_dog_mouse.hex_0.10"], 123456
```

53ms

## Split Frames

Type Key

cat\_dog\_mouse.hex\_0.90

cat\_dog\_mouse.hex\_0.10

Ratio

0.0999999999999999



Ready

Connections: 0 H2O



Flow ▾ Cell ▾ Data ▾

Model ▾

Score ▾

Admin ▾

Help ▾

Deep Learning...

Deep Water...

Distributed Random Forest...

Gradient Boosting Method...

Generalized Linear Modeling...

Generalized Low Rank Modeling...

K-means...

Naive Bayes...

Principal Components Analysis...

List All Models

List Grid Search Results

Import Model...

Export Model...

## Untitled Flow



CS

Expression...

Using Flow to train Deep Water Model



Ready

# H2O FLOW

Flow ▾ Cell ▾ Data ▾ Model ▾ Score ▾ Admin ▾ Help ▾

## Untitled Flow



CS

```
buildModel "deepwater"
```

251ms

### Build a Model

Select an algorithm: Deep Water ▾

#### PARAMETERS

GRID ?

model\_id deepwater-flow Destination id for this model; auto-generated if not specified.

training\_frame cat\_dog\_mouse.hex\_0.90 ▾ Id of the training data frame (Not required, to allow initial validation of model parameters).

validation\_frame cat\_dog\_mouse.hex\_0.10 ▾ Id of the validation data frame.

nfolds 0 Number of folds for N-fold cross-validation (0 to disable or >= 2).

response\_column C2 Response variable column.

ignored\_columns Search...

Showing page 1 of 1.

C1 STRING

C2 ENUM(3)

Ready

Connections: 0 H2O



**H<sub>2</sub>O FLOW**

Flow ▾ Cell ▾ Data ▾ Model ▾ Score ▾ Admin ▾ Help ▾

## Untitled Flow



epochs	300	How many times the dataset should be iterated (streamed), can be fractional.
network	auto ▾ (Choose...)	Network architecture.
activation	auto	Activation function.
hidden	user lenet alexnet vgg googlenet inception_bn resnet	Hidden layer sizes (e.g. [200, 200]).
ADVANCED		<b>Choosing Different Network Structure</b>
checkpoint		Model checkpoint to resume training with.
fold_column		Column with cross-validation fold index assignment per observation.
score_each_iteration	<input checked="" type="checkbox"/>	Whether to score during each iteration of model training.
categorical_encoding	OneHotExplicit ▾	Encoding scheme for categorical features
train_samples_per_iteration	-2	Number of training samples (globally) per MapReduce iteration. Special values are 0: one epoch, -1: all available data (e.g., replicated training data), -2: automatic.
distribution	AUTO ▾	Distribution function
score_interval	5	Shortest time interval (in seconds) between model scoring.
score_training_samples	10000	Number of training set samples for scoring (0 for all).
score_validation_samples	0	Number of validation set samples for scoring (0 for all).
score_duty_cycle	0.1	Maximum duty cycle fraction for scoring (lower: more training, higher: more scoring).
stopping_rounds	5	Early stopping based on convergence of stopping_metric. Stop if simple moving average of length k of the stopping_metric

Ready

Connections: 0 H<sub>2</sub>O

# H2O FLOW

Flow ▾ Cell ▾ Data ▾ Model ▾ Score ▾ Admin ▾ Help ▾

## Untitled Flow



stopping_metric	AUTO	does not improve for k:=stopping_rounds scoring events (0 to disable) Metric to use for early stopping (AUTO: logloss for classification, deviance for regression)
stopping_tolerance	0	Relative tolerance for metric-based stopping criterion (stop if relative improvement is not at least this much)
max_runtime_secs	0	Maximum allowed runtime in seconds for model training. Use 0 to disable.
backend	mxnet ▾	Deep Learning Backend.
image_shape	(Choose...) auto	Width and height of image.
channels	mxnet	Number of (color) channels.
network_definition_file	tensorflow	Path of file containing network definition (graph, architecture).
network_parameters_file		Path of file containing network (initial) parameters (weights, biases).
mean_image_file		Path of file containing the mean image data for data normalization.
export_native_model_prefix		Path (prefix) where to export the native model after every iteration.
input_dropout_ratio	0	Input layer dropout ratio (can improve generalization, try 0.1 or 0.2).
hidden_dropout_ratios		Hidden layer dropout ratios (can improve generalization), specify one value per hidden layer, defaults to 0.5.

### Choosing Different Backend

#### EXPERT

overwrite_with_best_model	<input checked="" type="checkbox"/>	If enabled, override the final model with the best model found during training.
target_ratio_comm_to_comp	0.05	Target ratio of communication overhead to computation. Only for multi-node operation and train_samples_per_iteration = -2 (auto-tuning).



Ready

Connections: 0 H2O

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

Flow ▾ Cell ▾ Data ▾ Model ▾ Score ▾ Admin ▾ Help ▾

## Untitled Flow



cs

```
buildModel 'deepwater', {"model_id":"deepwater-
flow","training_frame":"cat_dog_mouse.hex_0.90","validation_frame":"cat_dog_mouse.hex_0.10","nfolds":0,"response_column":"C2","ignored_c
olumns":[],"epochs":300,"network":"lenet","hidden":
[],"checkpoint":"","score_each_iteration":true,"categorical_encoding":"OneHotExplicit","train_samples_per_iteration":-2,"distribution":"
AUTO","score_interval":0,"score_training_samples":10000,"score_validation_samples":0,"score_duty_cycle":0.1,"stopping_rounds":5,"stop
ping_metric":AUTO,"stopping_tolerance":0,"max_runtime_secs":0,"backend":"mxnet","image_shape":
[0,0],"channels":3,"network_definition_file":"","network_parameters_file":"","mean_image_file":"","export_native_model_prefix":"",
"input
_dropout_ratio":0,"hidden_dropout_ratios":
[],"overwrite_with_best_model":true,"target_ratio_comm_to_comp":0.05,"seed":-1,"rate":0.005,"rate_annealing":0.000001,"momentum_start":0
.9,"momentum_ramp":10000,"momentum_stable":0.99,"single_node_mode":false,"shuffle_training_data":true,"mini_batch_size":32,"clip_gradien
t":10,"gpu":true,"device_id":0}
```

4.1s

## Job

Run Time 00:00:03.374

Remaining Time 00:00:00.0

Type Model

Key deepwater-flow

Description DeepWater

Status DONE

Progress 100%

Ready

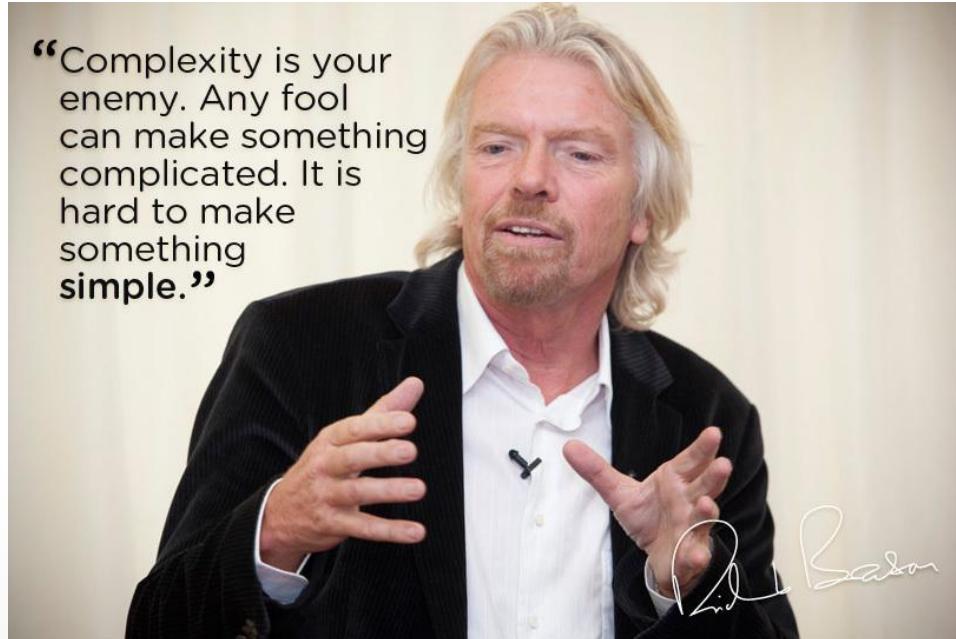
Training Deep Water Models without Programming



Connections: 0 H<sub>2</sub>O

# H2O's Mission

Making Machine Learning Accessible to Everyone



*Photo credit: Virgin Media*

# Grazie mille!

- Data Science Milan
- Gianmario Spacagna
- Politecnico di Milano



- Resources
  - [bit.ly/h2o\\_milan\\_1](https://bit.ly/h2o_milan_1)
  - [www.h2o.ai](http://www.h2o.ai)
  - [docs.h2o.ai](http://docs.h2o.ai)
- Contact
  - [joe@h2o.ai](mailto:joe@h2o.ai)
  - [@matlabulous](https://twitter.com/matlabulous)
  - [github.com/woobe](https://github.com/woobe)