

Deep Water

GPU Deep Learning for H2O

Arno Candel, PhD
Chief Architect, Physicist & Hacker, H2O.ai
@ArnoCandel

The Open Tour, Dallas, TX
Oct 26 2016

Computer Science (CS)

Artificial Intelligence (A.I.)

Machine Learning (ML)

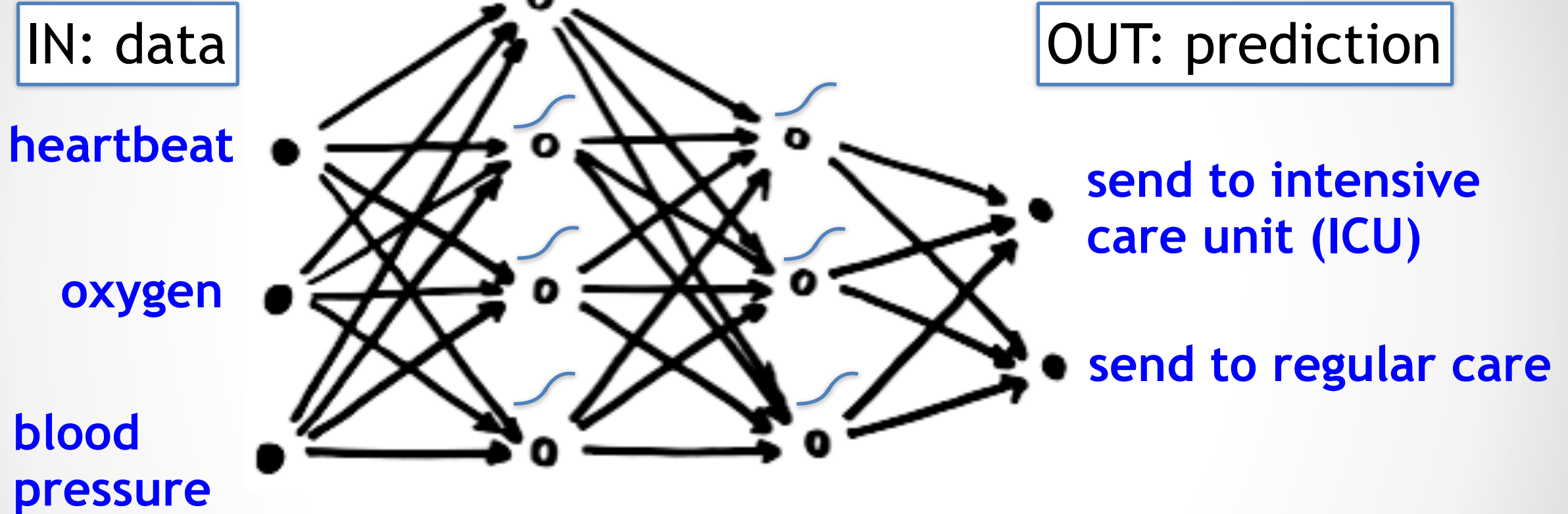
Deep Learning (DL)

hot hot hot hot hot

H2O.ai

A Simple Deep Learning Model: Artificial Neural Network

from 1970s, now rebranded as DL



nodes : neuron activations (real numbers) – represent features
arrows : connecting weights (real numbers) – learned during training
~ : non-linearity $x \rightarrow f(x)$ – adds model complexity

A step back: A.I. was coined over 60 years ago



John McCarthy

Princeton, Bell Labs, Dartmouth, later: MIT, Stanford

1955: “A proposal for the Dartmouth summer research project on Artificial Intelligence”

with Marvin Minsky (MIT), Claude Shannon (Bell Labs) and Nathaniel Rochester (IBM)

http://www.asiapacific-mathnews.com/04/0403/0015_0020.pdf

1955 proposal for the Dartmouth summer research project on A.I.

“We propose that a 2-month, 10-man study of artificial intelligence be carried out during the summer of 1956 at Dartmouth College in Hanover, New Hampshire. The study is to proceed on the basis of the conjecture that every aspect of learning and any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it. An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves. We think that a significant advance can be made in one or more of these problems if a carefully selected group of scientists work on it together for one summer.”

Step 1: Great Algorithms + Fast Computers

“No computer will ever beat me at playing chess.”



1997: Playing Chess
(IBM Deep Blue beats Kasparov)

Computer Science

30 custom CPUs, 60 billion moves in 3 mins

<http://nautil.us/issue/18/genius/why-the-chess-computer-deep-blue-played-like-a-human>

Step 2: More Data + Real-Time Processing

“It takes a human to drive a car!”



2005: Self-driving Cars
DARPA Grand Challenge, 132 miles
(won by Stanford A.I. lab*)

Sensors & Computer Science
video, radar, laser, GPS, 7 Pentium computers

<http://cs.stanford.edu/group/roadrunner/old/presskit.html>

*A.I. lab was established by McCarthy et al. in the early '60s

Step 3: Big Data + In-Memory Clusters

“Computers can’t answer arbitrary questions!”



2011: Jeopardy (IBM Watson)

In-Memory Analytics/ML

4 TB of data (incl. wikipedia), 90 servers,
16 TB RAM, Hadoop, 6 million logic rules

<https://www.youtube.com/watch?v=P18EdAKuC1U>

[https://en.wikipedia.org/wiki/Watson_\(computer\)](https://en.wikipedia.org/wiki/Watson_(computer))

Note: IBM Watson received the question in electronic written form, and was often able to press the answer button faster than the competing humans.

Step 4: Deep Learning

“Computers don’t understand our language!”



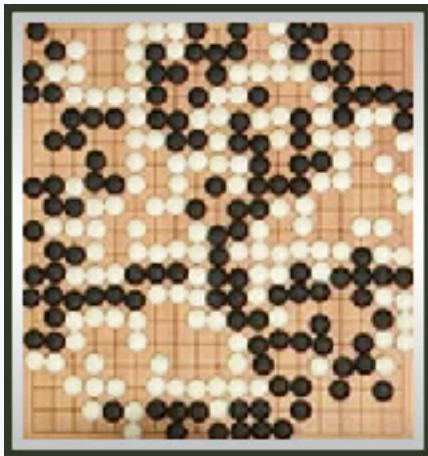
2014: Google
(acquired Quest Visual)

Deep Learning
Convolutional and Recurrent
Neural Networks,
with training data from users

- Translate between 103 languages by typing
- Instant camera translation: Use your camera to translate text instantly in 29 languages
- Camera Mode: Take pictures of text for higher-quality translations in 37 languages
- Conversation Mode: Two-way instant speech translation in 32 languages
- Handwriting: Draw characters instead of using the keyboard in 93 languages

Step 5: Augmented Deep Learning

“Go is too complex for computers to master!”



Go board has approx.

[illegible]

2014: Atari Games (DeepMind)

trained from raw pixel values, no human rules

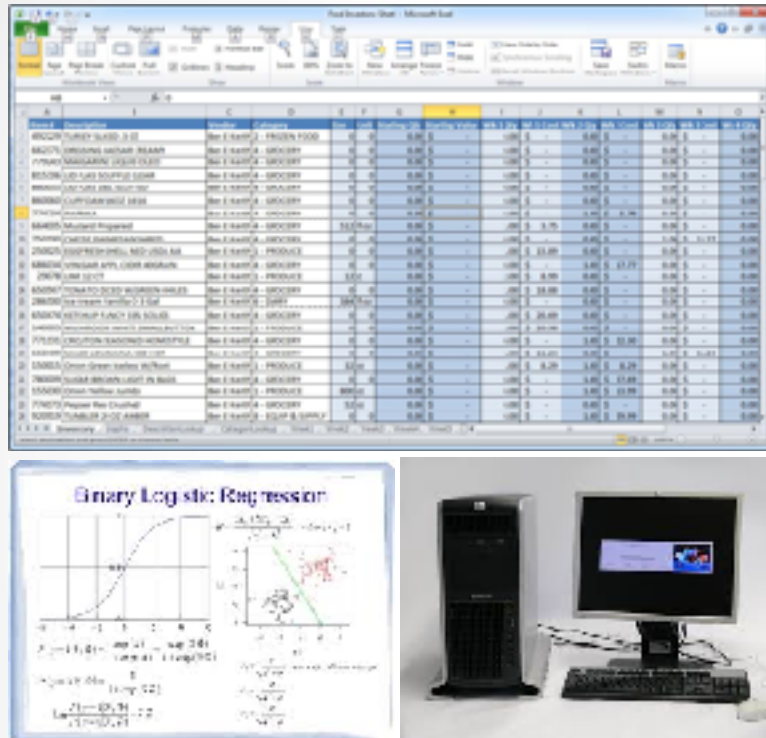
Deep Learning

**+ reinforcement learning, tree search,
Monte Carlo, GPUs, playing against itself, ...**

2016: AlphaGo (Google DeepMind)



Things are Changing Quickly



Yesterday: Small Data (<GB)

Data + Skills

H₂O.ai are good for business



Today: Big Data (TeraBytes, ExaBytes)

Data + Machine Learning

ARE the business



Why Deep Learning?

Deep Learning is in the Center of This Revolution

- conceptually simple
- solves many problems
- benefits from big data
- super-human results

Why Deep Learning?

However:

- hard to interpret — but solutions exist
- lots of architectural choices — require lucky PhDs, open-source helps
- lots of hyper-parameters — AutoML can do the tuning for you
- slow to train on big data — dedicated hardware helps (GPU clusters)
- rapidly changing landscape — Deep Water unifies open-source APIs

H2O pioneered Easy-To-Use Open Source Deep Learning

The screenshot displays the H2O FLOW web interface. At the top, the navigation bar includes 'H2O FLOW' and several dropdown menus: 'Flow', 'Cell', 'Data', 'Model', 'Score', 'Admin', and 'Help'. The 'Model' dropdown menu is currently open, showing a list of machine learning models: 'Aggregator...', 'Deep Learning...', 'Distributed Random Forest...', 'Gradient Boosting Machine...', 'Generalized Linear Modeling...', 'Generalized Low Rank Modeling...', 'K-means...', 'Naive Bayes...', and 'Principal Components Analysis...'. Below this list are links for 'List All Models', 'List Grid Search Results', 'Import Model...', and 'Export Model...'. On the left side, the 'Assistance' sidebar is visible, featuring a search bar with the text 'assist' and a table of routines. The table has two columns: 'Routine' and 'Description'. The routines listed are: 'importFiles' (Import file(s) into H2O), 'getFrames' (Get a list of frames in H2O), 'splitFrame' (Split a frame into two or more), 'mergeFrames' (Merge two frames into one), 'getModels' (Get a list of models in H2O), 'getGrids' (Get a list of grid search results), 'getPredictions' (Get a list of predictions in H2O), 'getJobs' (Get a list of jobs running in H2O), 'buildModel' (Build a model), 'importModel' (Import a saved model), and 'predict' (Make a prediction). On the right side, there is a 'CLIPS' button and a 'HELP' button. Below these, there is a section titled 'for the first time?' with a 'Start Videos' button. At the bottom right, there is a 'GENERAL' section with a list of links: 'Flow Web UI ...', '... Importing Data', '... Building Models', '... Making Predictions', '... Using Flows', and '... Troubleshooting Flow'. At the bottom left, there is a logo for 'H2O.ai'. At the bottom right, there is a circular logo for 'H2O.ai' with the text 'THE OPEN-TOUR' and 'H2O.ai' in the center.

H2O FLOW

Flow Cell Data Model Score Admin Help

Untitled Flow

assist

CS

Assistance

Routine	Description
importFiles	Import file(s) into H2O
getFrames	Get a list of frames in H2O
splitFrame	Split a frame into two or more
mergeFrames	Merge two frames into one
getModels	Get a list of models in H2O
getGrids	Get a list of grid search results
getPredictions	Get a list of predictions in H2O
getJobs	Get a list of jobs running in H2O
buildModel	Build a model
importModel	Import a saved model
predict	Make a prediction

Aggregator...

Deep Learning...

Distributed Random Forest...

Gradient Boosting Machine...

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...

CLIPS

HELP

for the first time?

Start Videos

Flows to explore and learn

GENERAL

- Flow Web UI ...
- ... Importing Data
- ... Building Models
- ... Making Predictions
- ... Using Flows
- ... Troubleshooting Flow

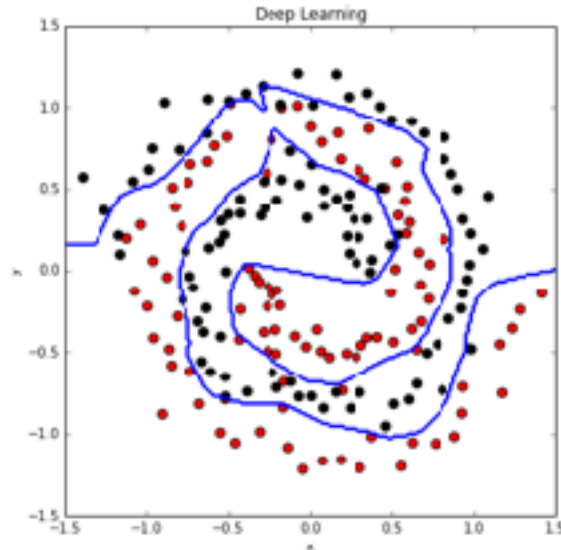
EXAMPLES

H2O.ai

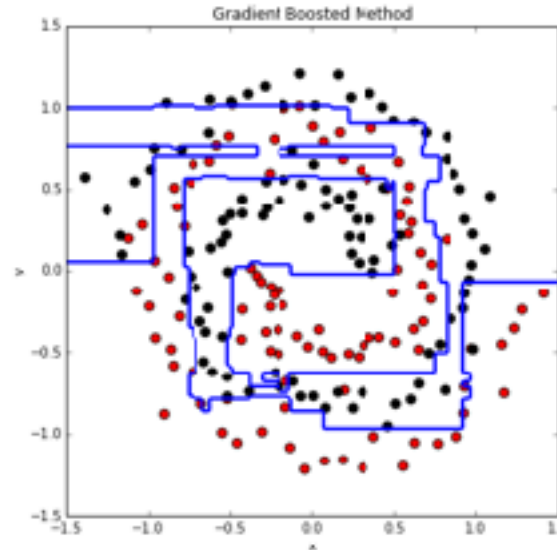
H2O.ai THE OPEN-TOUR

H2O pioneered Easy-To-Use Open Source Deep Learning

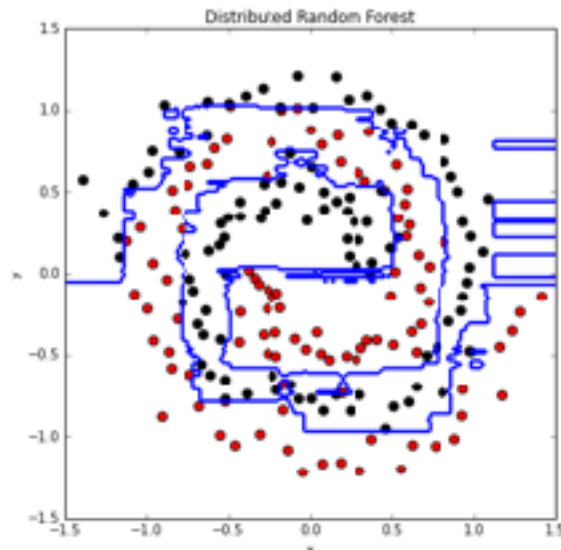
Deep Learning



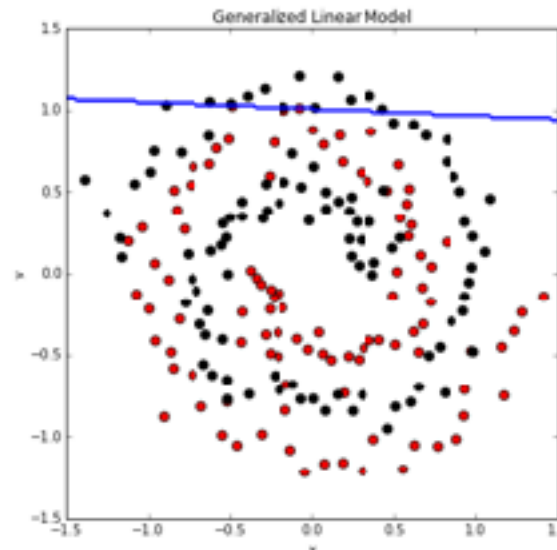
Gradient Boosting Machine



Distributed Random Forest



Generalized Linear Modeling



All algorithms are distributed and scalable

H2O Deep Learning Is 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>

H2O and TensorFlow are tied



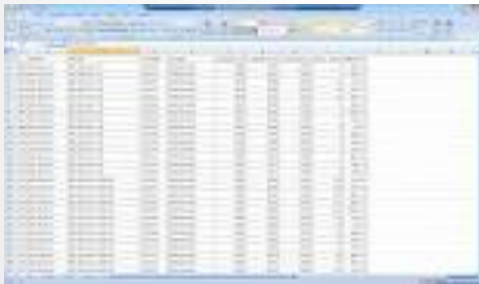
Deep Learning Tools & Platforms: Caffe (63)	0%
Cuda-convnet (22)	0%
Convnet.js (8)	0%
carch (2)	0%
Deeplearning4j (46)	0%
H2O (190)	1%
MATLAB Deep Learning Toolbox (57)	0%
Microsoft CNTK (25)	0%
mxnet (16)	0%
Nervana (2)	0%
Tensorflow (190)	1%
Theano ecosystem including Keras, Lasagne, Pylearn2 (140)	1%
Torch (28)	0%
Veles (2)	0%
Other Deep Learning Tools (104)	1%

usage of Deep Learning tools in past year

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

H2O Deep Learning: simple multi-layer networks, CPUs

1-5 layers
MBs/GBs of data

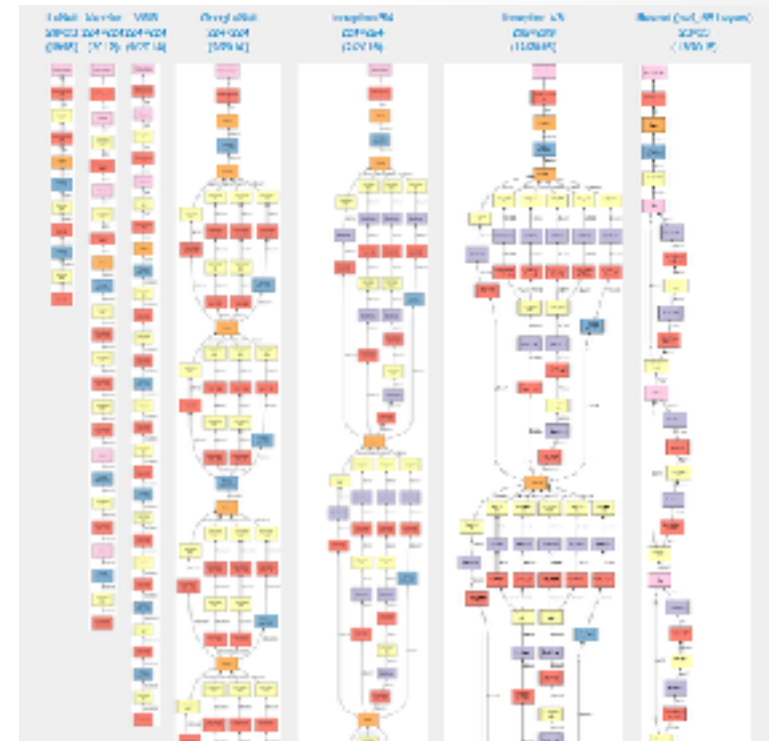
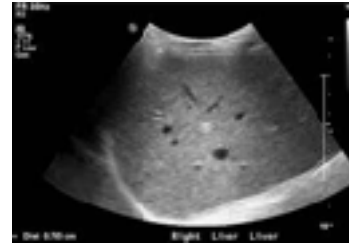


Limited to business analytics,
statistical models (CSV data)

H₂O.ai

H2O Deep Water: arbitrary networks, CPUs or GPUs

1-1000 layers
GBs/TBs of data






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

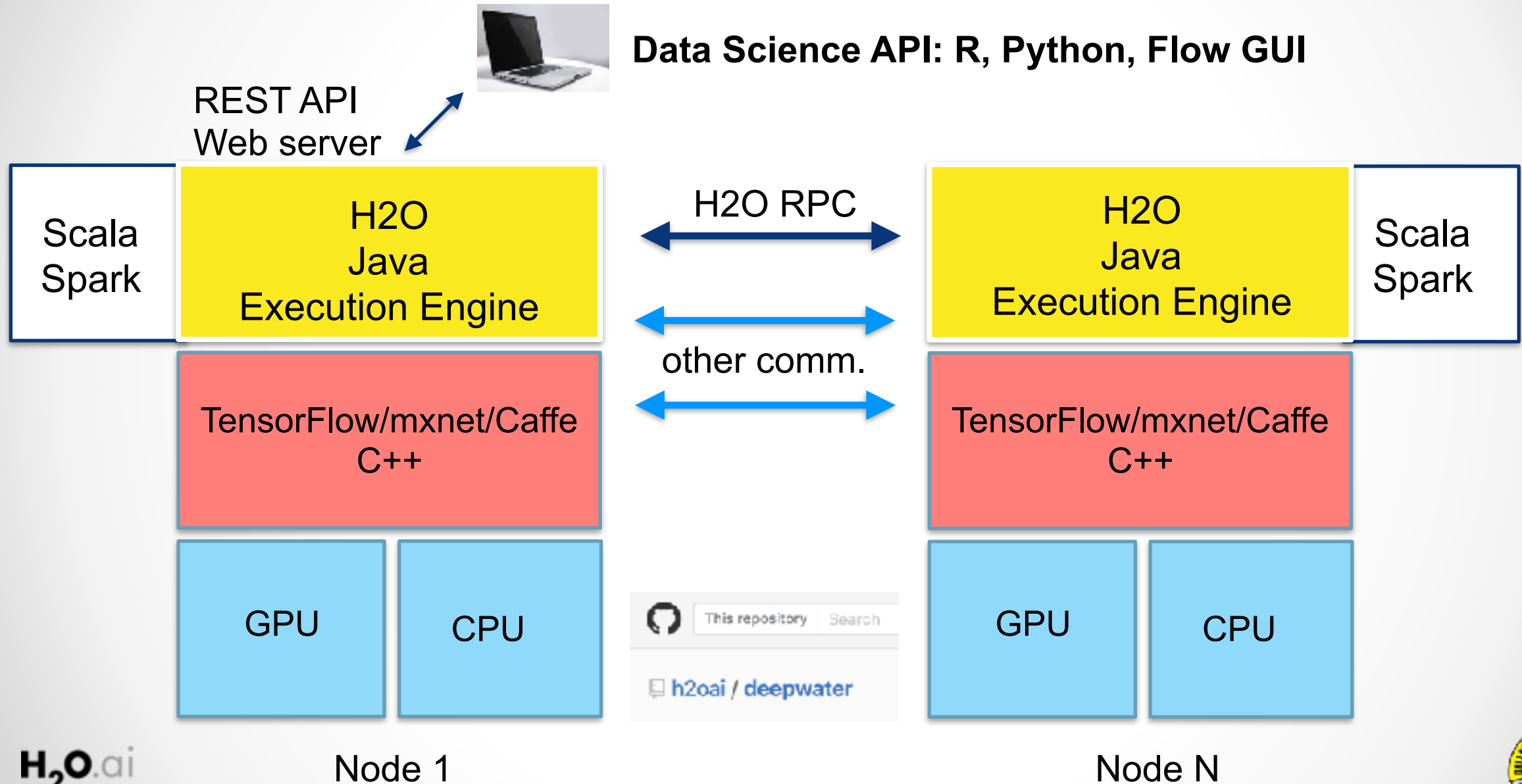


Deep Water: Best Open-Source Deep Learning

Enterprise Deep Learning for Business Transformation

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

Deep Water Architecture



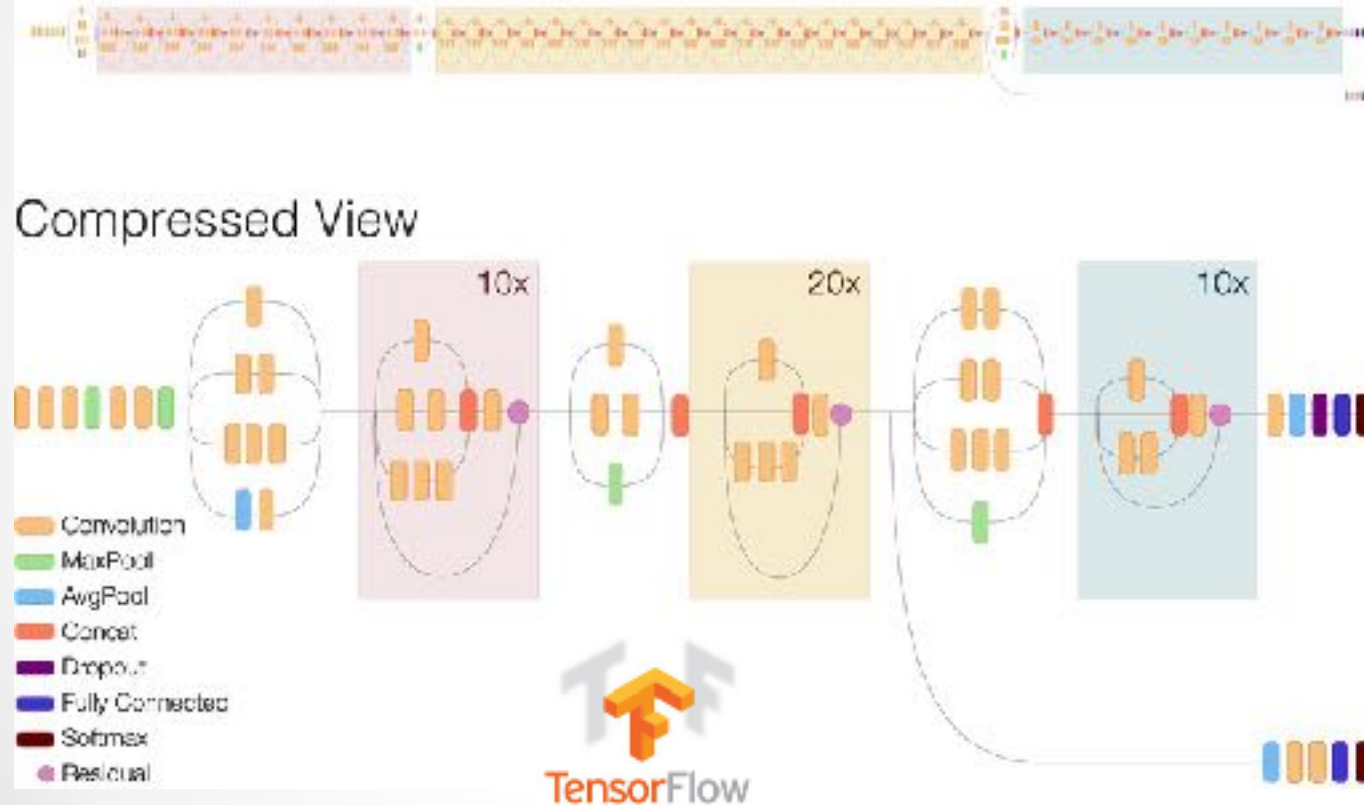
Step 1: Leverage Research Community Code, Data and Models

World's best Image Classifier (Google + Microsoft, Aug 2016)



open-source implementation

Inception Resnet V2 Network



<https://research.googleblog.com/2016/08/improving-inception-and-image.html>

```
stem_2_1x1 = Conv(data=concat1, num_filter=num_4_1,
stem_2_7x1 = Conv(data=stem_2_1x1, num_filter=num_4_1,
                    suffix='_conv_2')
stem_2_1x7 = Conv(data=stem_2_7x1, num_filter=num_4_1,
                    suffix='_conv_3')
stem_2_3x3 = Conv(data=stem_2_1x7, num_filter=num_4_1,

concat2 = mx.sym.Concat(*[stem_1_3x3, stem_2_3x3], :

pool2 = mx.sym.Pooling(data=concat2, kernel=[3, 3],
                        name='%s_pool2' % ('max',
stem_3_3x3 = Conv(data=concat2, num_filter=num_5_1,
                    suffix='_conv_1', withRelu=False)

concat3 = mx.sym.Concat(*[pool2, stem_3_3x3], name=
bn1 = mx.sym.BatchNorm(data=concat3, name='%s_bn1'
act1 = mx.sym.Activation(data=bn1, act_type='relu',

return act1

def InceptionResnetV2A(data,
                        num_1_1,
                        num_2_1, num_2_2,
                        num_3_1, num_3_2, num_3_3,
                        proj,
                        name,
                        scaleResidual=True):
    import mxnet as mx
    init = data

    a1 = Conv(data=data, num_filter=num_1_1, name='%s_1'
    a2 = Conv(data=a1, num_filter=num_2_1, name='%s_2'
    a2 = Conv(data=a2, num_filter=num_2_2, kernel=[3, 3]
    a3 = Conv(data=a2, num_filter=num_3_1, name='%s_3'
    a3 = Conv(data=a3, num_filter=num_3_2, kernel=[3, 3]
    a3 = Conv(data=a3, num_filter=num_3_3, kernel=[3, 3]
```


Step 2: Train Models with Familiar APIs

```
frame = h2o.import_file(PATH+"/bigdata/laptop/deepwater/imagenet/cat_dog_mouse.csv")
print(frame.head(5))
```

Parse progress: 100%

C1	C2
bigdata/laptop/deepwater/imagenet/cat/102194502_49f003abcd9.jpg	cat
bigdata/laptop/deepwater/imagenet/cat/11146807_00a5f35255.jpg	cat
bigdata/laptop/deepwater/imagenet/cat/1140846215_70e326f868.jpg	cat

```
nclasses = frame[1].nlevels()[0]
get_symbol(nclasses).save("/tmp/symbol_inception_resnet_v2-py.json")
model = H2ODeepWaterEstimator(epochs=20, learning_rate=1e-3, learning_rate_annealing=1e-5,
                              mini_batch_size=16,
                              network_definition_file="/tmp/symbol_inception_resnet_v2-py.json",
                              image_shape=[299,299],
                              channels=3)
model.train(x=[0],y=1, training_frame=frame)
model.show()
```

deepwater Model Build progress: 100%

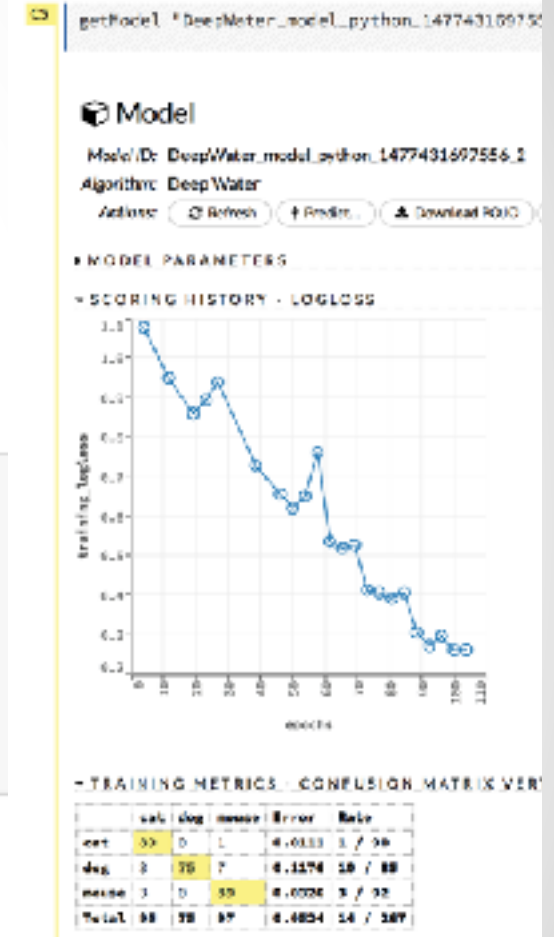
Model Details

=====

H2ODeepWaterEstimator : Deep Water

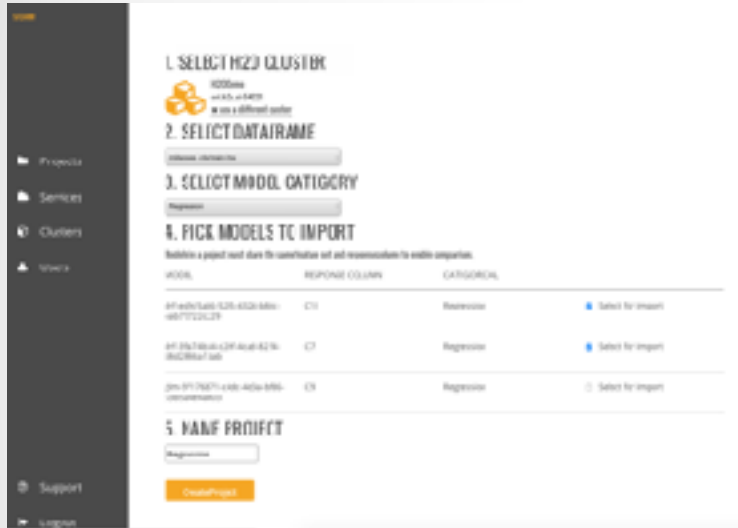
Model Key: DeepWater_model_python_1477179782032_5

Status of Deep Learning Model: user, 116.1 MB, predicting C2, 3-class classification, 5,632 training samples, mini-batch size 16



Step 3: Model Comparison and Rapid Deployment

Pretrained Image Classifier - Standalone Scoring Server Deployed in Seconds



Steam
Compare Models


<http://deepwater.h2o.ai/classy-demo/>


Prediction Service

Select input parameters, OR enter your own custom query string to predict

MODEL INPUT PARAMETERS

Parameters

1. C1 



Predicting **n02108000 EntleBucher**

Index	Labels	Probability
241	n02108000 EntleBucher	0.8957
238	n02107574 Greater Swiss Mountain dog	0.0701
240	n02107908 Appenzeller	0.0341
239	n02107683 Bernese mountain dog	0.0001

Scoring Server

Build smarter applications and data products

High-Performance Multi-GPU Training with Caffe

Training on a 16-GPU EC2 instance with Caffe

<https://aws.amazon.com/blogs/aws/new-p2-instance-type-for-amazon-ec2-up-to-16-gpus/>

Instance Name	GPU Count	vCPU Count	Memory	Parallel Processing Cores	GPU Memory	Network Performance
p2.xlarge	1	4	61 GiB	2,495	12 GiB	High
p2.8xlarge	8	32	488 GiB	19,988	96 GiB	10 Gigabit
p2.16xlarge	16	64	732 GiB	39,935	192 GiB	20 Gigabit

Caffe launched from Java:

```
java -cp target/dependency/*:target/classes deepwater.examples.ImageNet
```

H2O integration is in progress



NVIDIA-SMI 367.48					Driver Version: 367.48						
GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile Uncorr. ECC	GPU	Name	Persistence-M	Bus-Id	Disp.A	Volatile Uncorr. ECC
Id	Temp	Mode	Path	Mode	Mode	Id	Temp	Mode	Path	Mode	Mode
0	N/A	07C	P0	134W / 149W	1246MiB / 11439MiB	0	65C	P0	134W / 149W	1246MiB / 11439MiB	55%
1	N/A	07C	P0	121W / 119W	1750MiB / 11439MiB	1	65C	P0	121W / 119W	1750MiB / 11439MiB	71%
2	N/A	07C	P0	144W / 149W	2861MiB / 11439MiB	2	65C	P0	144W / 149W	2861MiB / 11439MiB	58%
3	N/A	07C	P0	147W / 149W	2750MiB / 11439MiB	3	65C	P0	147W / 149W	2750MiB / 11439MiB	92%
4	N/A	07C	P0	147W / 149W	1852MiB / 11439MiB	4	65C	P0	147W / 149W	1852MiB / 11439MiB	89%
5	N/A	07C	P0	143W / 149W	2750MiB / 11439MiB	5	65C	P0	143W / 149W	2750MiB / 11439MiB	93%
6	N/A	07C	P0	159W / 149W	2861MiB / 11439MiB	6	65C	P0	159W / 149W	2861MiB / 11439MiB	85%
7	N/A	07C	P0	154W / 149W	1750MiB / 11439MiB	7	65C	P0	154W / 149W	1750MiB / 11439MiB	92%
8	N/A	07C	P0	147W / 149W	2965MiB / 11439MiB	8	65C	P0	147W / 149W	2965MiB / 11439MiB	92%
9	N/A	07C	P0	157W / 149W	2750MiB / 11439MiB	9	65C	P0	157W / 149W	2750MiB / 11439MiB	93%
10	N/A	07C	P0	141W / 149W	2861MiB / 11439MiB	10	65C	P0	141W / 149W	2861MiB / 11439MiB	47%
11	N/A	07C	P0	155W / 149W	2750MiB / 11439MiB	11	65C	P0	155W / 149W	2750MiB / 11439MiB	97%
12	N/A	07C	P0	119W / 119W	1852MiB / 11439MiB	12	65C	P0	119W / 119W	1852MiB / 11439MiB	75%
13	N/A	07C	P0	147W / 149W	2750MiB / 11439MiB	13	65C	P0	147W / 149W	2750MiB / 11439MiB	97%
14	N/A	07C	P0	124W / 149W	2861MiB / 11439MiB	14	65C	P0	124W / 149W	2861MiB / 11439MiB	85%
15	N/A	07C	P0	122W / 119W	1750MiB / 11439MiB	15	65C	P0	122W / 119W	1750MiB / 11439MiB	97%

Caffe

All 16 GPUs are busy training



Roadmap for Deep Water (Q4 2016):



**Finish TensorFlow integration (C++/Python/Java):
Package Python on the backend to create trainable graphs**



**Finish Caffe integration (pure C++/Java):
Optimized Multi-GPU training (NVIDIA NCCL)**



Add multi-GPU support for mxnet



**Add more capabilities to H2O Deep Water:
Text/NLP, Time Series, LSTM, AutoEncoder,
Feature Extraction, Input/Output shape mapping, etc.**

Breakout Tracks on Deep Water & GPU Backends



11:05-11:30 Dmitry Larko — Credit Card Default Prediction



2:40-3:00 Fabrizio Milo — TensorFlow internals, Wide & Deep models



3:00-4:00 Arno Candel — Hands-On Workshop with Image Classification, Credit Card Default Prediction, Benchmarking and Hyper-Parameter Search in Flow, Python and R

Live Demo