

# Train Deep Learning Models Using H2O Deep Water

Arno Candel, Magnus Stensmo, Wen Phan

GTC 2017

May 11, 2017, San Jose, CA

# LAB CONNECTION INSTRUCTIONS - Part 1

Go to [nvlabs.qwiklab.com](https://nvlabs.qwiklab.com)

Sign in or create an account

Check for Access Codes (each day):

- Click **My Account**
- Click **Credits & Subscriptions**

If no Access Codes, ask for paper one from TA.

Please tear in half once used

An Access Code is needed to start the lab

The screenshot shows the 'My Account' page on the Qwiklabs website. At the top, there is a navigation bar with the Qwiklabs logo and a user profile icon labeled 'JOHN S.'. Below the navigation bar, there are three main menu items: 'CREDITS & SUBSCRIPTIONS' (which is highlighted in blue), 'PROFILE INFO', and 'SECURITY'. The main content area is titled 'Credits & Subscriptions'. It displays a message 'You have 5 Access Codes' with a 'More Details' link. A green rectangular box highlights this message. In the bottom right corner of the content area, there is a button labeled 'BUY CREDITS OR SUBSCRIPTION'. At the very bottom of the page, there is a link 'Share Your Credits with Friends'.

WIFI SSID: **GTC\_Hands\_On**

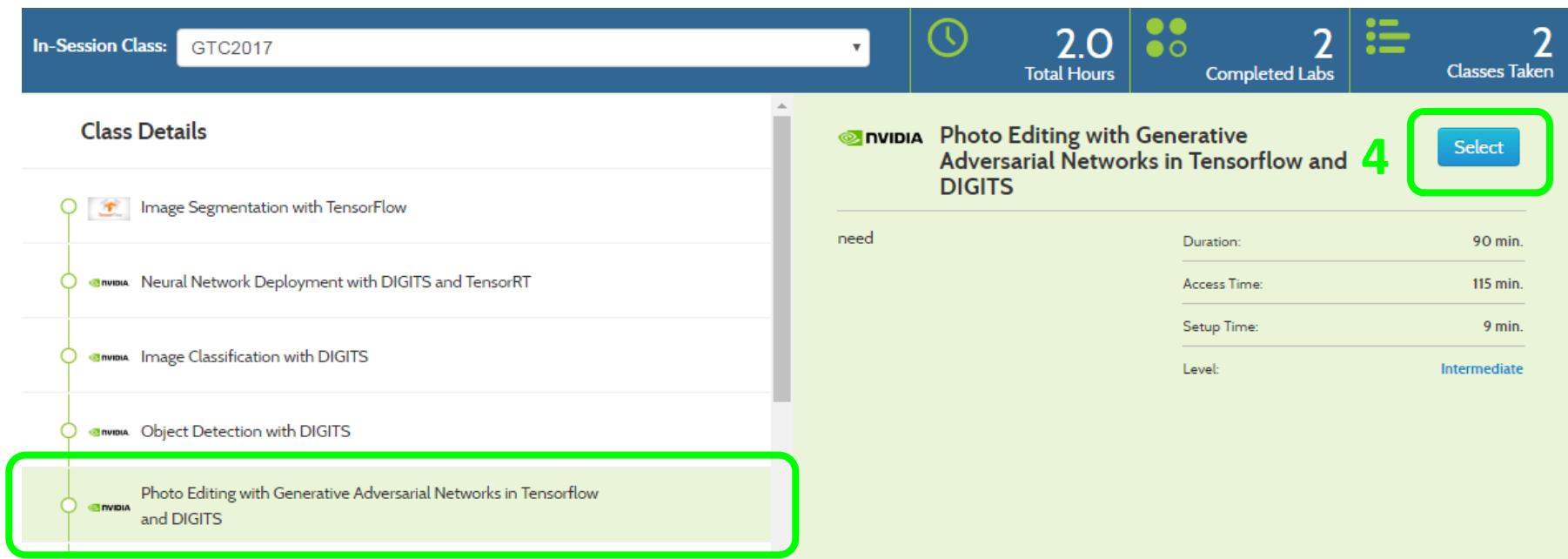
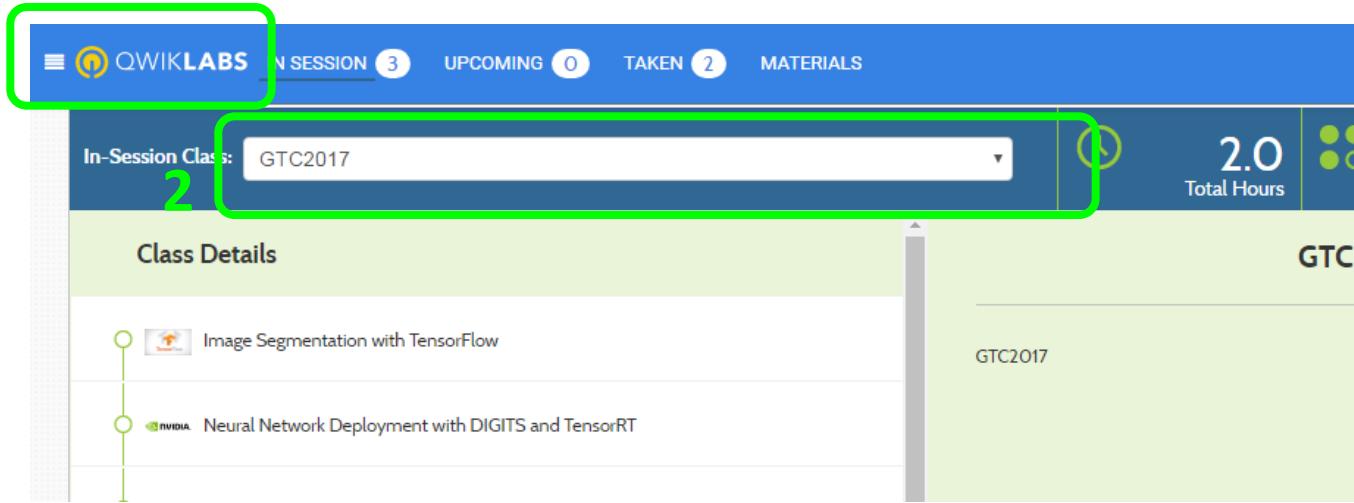
Password: **HandsOnGpu**

# LAB CONNECTION INSTRUCTIONS - Part 2

1. Click **Qwiklabs** in upper-left
2. Select GTC2017 Class
3. Find lab and click on it
4. Click on Select
5. Click Start Lab

WIFI SSID:  
**GTC\_Hands\_On**

Password:  
**HandsOnGpu**



# Agenda

- Introduction to H2O
- H2O Deep Learning and Deep Water
- Image Classification
- Custom and User-Defined Networks
- Pre-Trained Networks
- Ensembles
- Deep Features
- Deployment

# Introduction to H2O

# Company Overview

Founded	2011 Venture-backed, debuted in 2012
Products	<ul style="list-style-type: none"><li>• H2O Open Source In-Memory AI Prediction Engine</li><li>• Sparkling Water</li><li>• Enterprise Steam</li><li>• Deep Water</li></ul>
Mission	Operationalize data science and provide a platform for users to build beautiful data products
Team	<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>
Headquarters	Mountain View, CA



# Leading Users Across Industries



Hospital Corporation of America<sup>SM</sup>



Financial

Telecom

Insurance

Healthcare

Marketing



# Gartner MQ: Data Science Platform



***“Overall customer satisfaction is very high”***

***“H2O is especially suited to IoT edge and device scenarios”***

***“H2O had the highest reference customer analytics support score of all the vendors”***

# Market Validation



Gartner Magic Quadrant for Data

Science Platforms

**H2O.ai Selected** in Top 16 MQ Vendors  
for 2017



Forrester Wave for Predictive Analytics

**H2O.ai Selected** in Top Wave Vendor  
for 2017



**H2O.ai Selected** in AI Top 100

**46** PRESS  
MENTIONS



THE  
HUFFINGTON  
POST



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

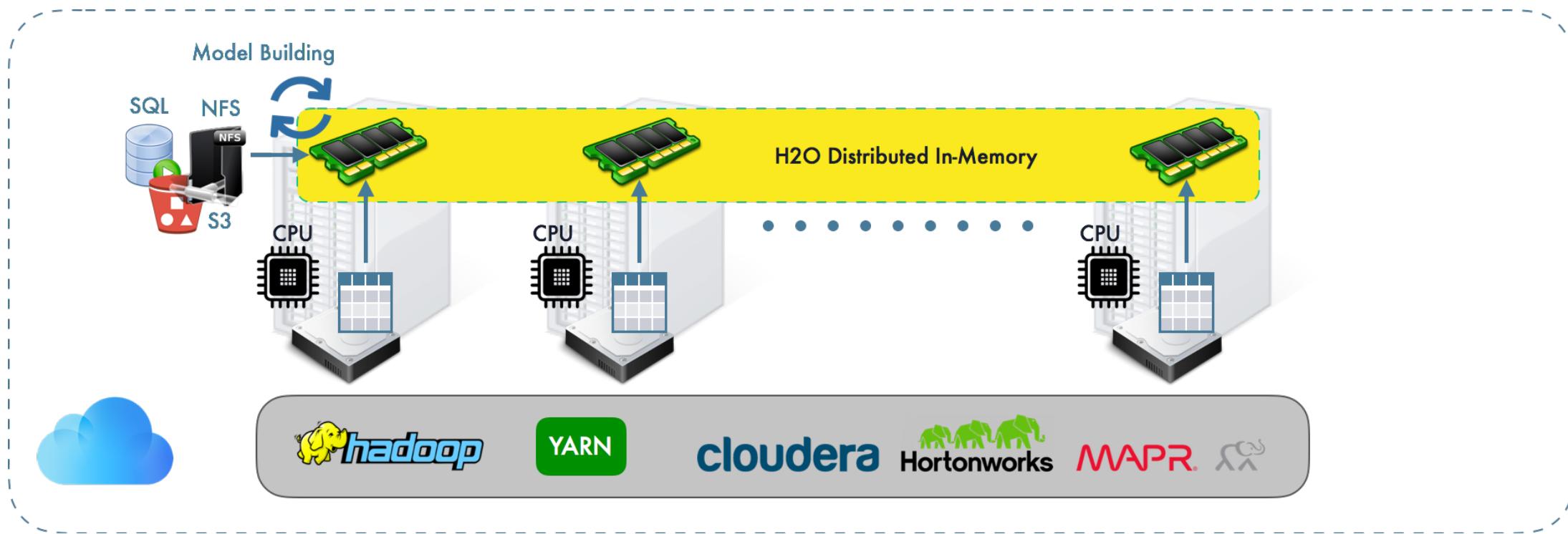
**107** OF FORTUNE  
THE 500  
 H<sub>2</sub>O

**8** OF TOP 10  
BANKS

**7** OF TOP 10  
INSURANCE COMPANIES

**4** OF TOP 10  
HEALTHCARE COMPANIES

# H2O Core



# Algorithm Overview

## Statistical Analysis

---

- Linear Models (GLM)
- Naïve Bayes

## Ensembles

---

- Random Forest
- Distributed Trees
- Gradient Boosting Machine
- Stacking / Super Learner

## Deep Neural Networks

---

- MLP
- Autoencoder
  - Anomaly Detection
  - Deep Features
- CNN, RNN (Deep Water)

## Clustering

---

- K-Means (Auto-K)

## Dimension Reduction

---

- Principal Component Analysis
- Generalized Low Rank Models

## Word Embedding

---

- Word2Vec

## Time Series

---

- iSAX

## Machine Learning Tuning

---

- Hyperparameter Search
- Early Stopping

# User Interfaces and APIs



H<sub>2</sub>O FLOW

The Python logo consists of two interlocking snakes, one blue and one yellow, with the word "python" in a brown serif font to the right.

# H2O Flow

The screenshot shows the H2O Flow web application running at [10.0.1.13:54321/flow/index.html#](http://10.0.1.13:54321/flow/index.html#). The title bar says "H2O Flow". The main header includes "Flow", "Cell", "Data", "Model", "Score", "Admin", and "Help". Below the header is a toolbar with various icons for file operations like opening, saving, and deleting. The main content area is titled "Untitled Flow" and contains a search bar with "assist" and a timestamp of "131ms". To the right of the search bar is a "HELP" tab. The "Assistance" section lists several routines with their descriptions:

Routine	Description
<code>importFiles</code>	Import file(s) into H <sub>2</sub> O
<code>getFrames</code>	Get a list of frames in H <sub>2</sub> O
<code>splitFrame</code>	Split a frame into two or more frames
<code>getModels</code>	Get a list of models in H <sub>2</sub> O
<code>getGrids</code>	Get a list of grid search results in H <sub>2</sub> O
<code>getPredictions</code>	Get a list of predictions in H <sub>2</sub> O
<code>getJobs</code>	Get a list of jobs running in H <sub>2</sub> O
<code>buildModel</code>	Build a model
<code>importModel</code>	Import a saved model
<code>predict</code>	Make a prediction

On the right side, there's a "Help" section with links to "Quickstart Videos" and "example Flows". It also features a "Star" button with 1,238 stars and a "GENERAL" section with a list of links:

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

At the bottom left is a "Ready" status indicator, and at the bottom right are "Connections: 0" and the H<sub>2</sub>O logo.

# H2O Clients and API

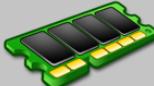


```
library(h2o)
h2o.init(ip="localhost", port=54321)
h2o.ls()
df_allyears2k <- h2o.getFrame("allyears2k.hex")
deeplearning_model <- h2o.getModel("deeplearning_model")
summary(df_allyears2k)
```

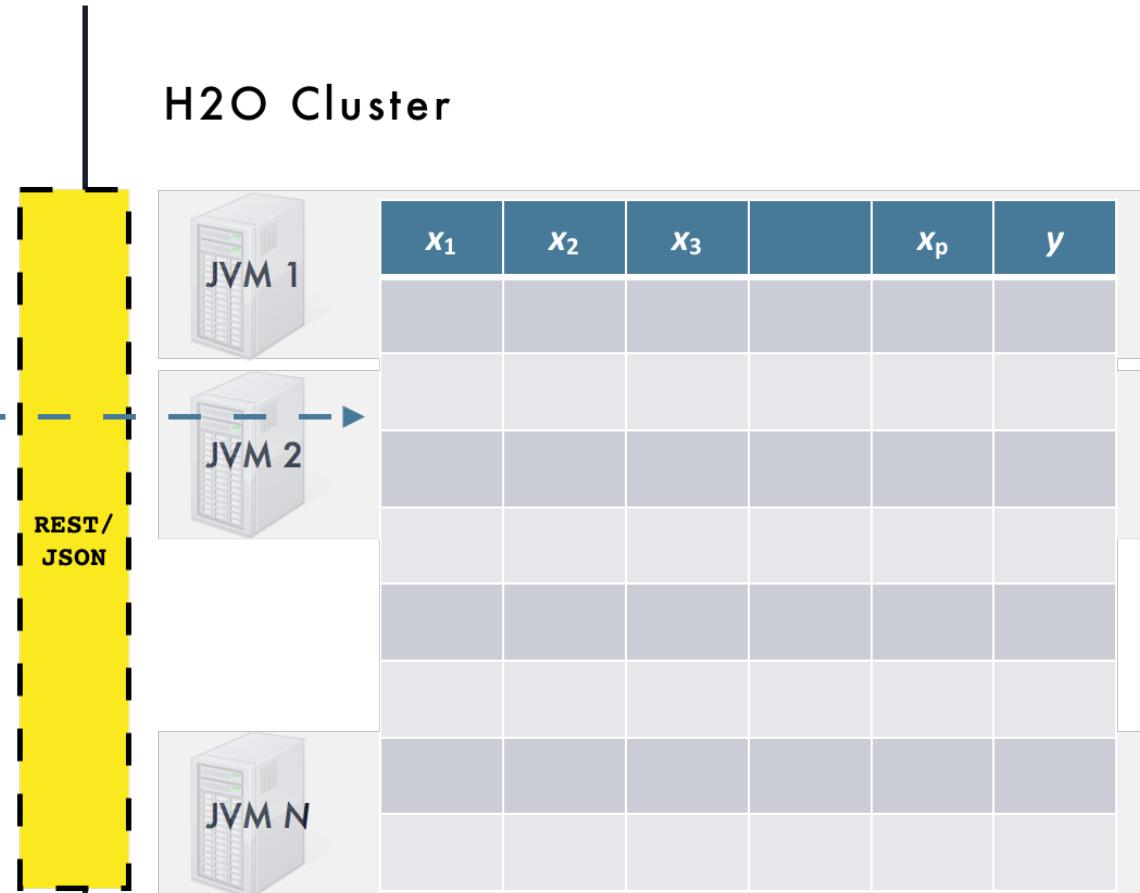
Console ~ /Library/Mobile Documents/com~apple~CloudDocs/\_h2o\_docs/demo

```
SecurityDelay LateAircraftDelay IsArrDelayed IsDepDelayed
Min. : 0.0000 Min. : 0.00 YES:24441 YES:23091
1st Qu.: 0.0000 1st Qu.: 0.00 NO :19537 NO :20887
Median : 0.0000 Median : 0.00
Mean : 0.01702 Mean : 7.62
3rd Qu.: 0.0000 3rd Qu.: 0.00
Max. :14.00000 Max. :373.00
NA's :35045 NA's :35045
```

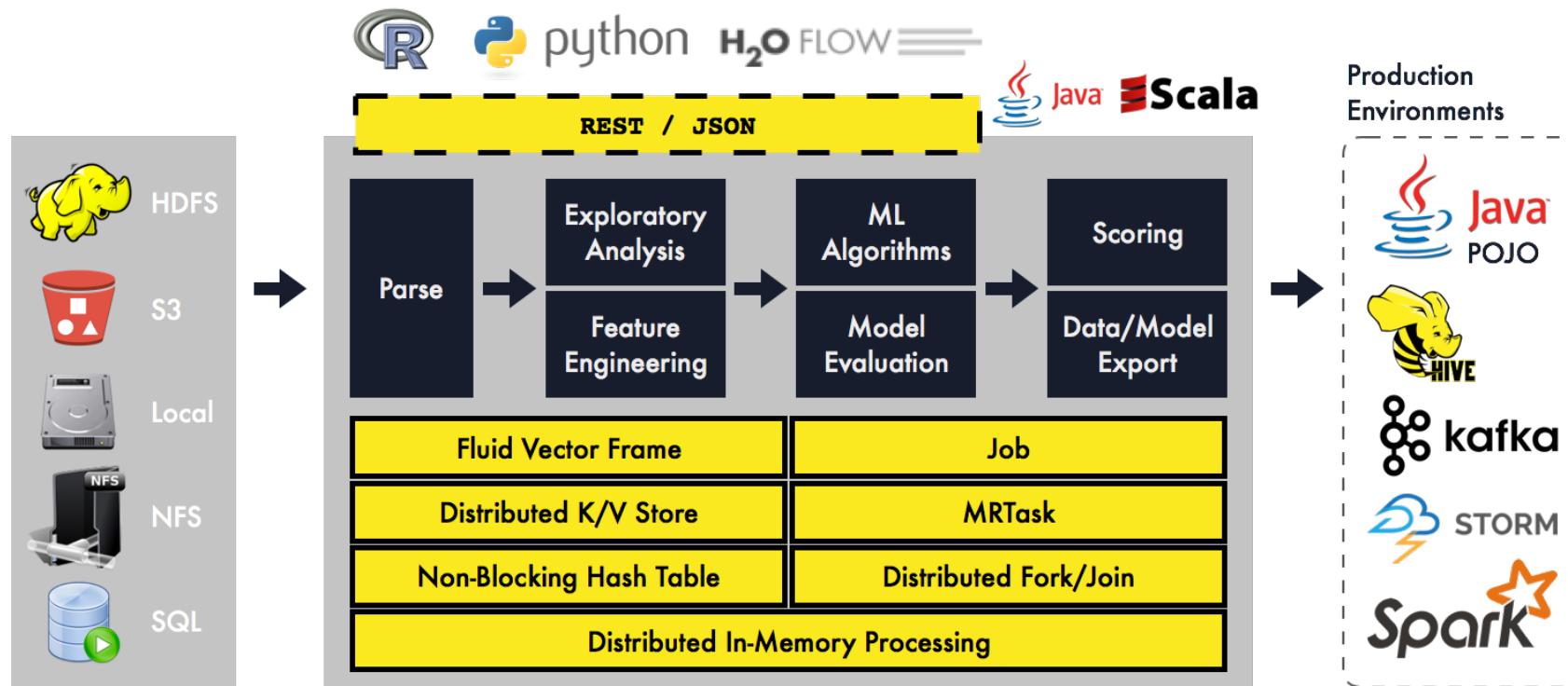
Local  
Machine



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



# Core H2O Architecture



cloudera  
MAPR



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

Spark + H<sub>2</sub>O  
SPARKLING  
WATER

# Other Solutions

- Sparkling Water (H<sub>2</sub>O + Spark)
- Enterprise Steam

# Upcoming and Ongoing

- Automatic Machine Learning (AutoML)
- Machine Learning Interpretability

# H2O Deep Learning and Deep Water

# H2O Deep Learning

- **“Best in Class” Multi-Layer Feed-Forward Deep Neural Networks**
- **Modern training options:** specifications for distributions (Bernoulli, Multinomial, Poisson, Gamma, Tweedie, Laplace, Huber, Quantile, Gaussian), loss functions (cross entropy, quadratic, absolute, Huber), learning rate, annealing, momentum, mini-batch size, and initialization
- **Deep autoencoders for unsupervised learning:** deep features and anomaly detection

# H2O Deep Learning (cont.)

- **Automatic and flexible data handling to maximize productivity:** standardization, one-hot encoding, observation weights and offsets, class balancing, sampling factors, ignoring constant columns, sparse data handling, and input layer constraints
- **Tuning parameters to prevent model overfitting and efficient model development:** cross-validation, regularization, drop out, early stopping, model checkpointing, and hyperparameter search

# World-Record MNIST Performance

```
> library(h2o)
> h2oServer <- h2o.init(ip="mr-0xd1", port=53322)
> train_hex <- h2o.importFile(h2oServer, "mnist/train.csv.gz")
> test_hex <- h2o.importFile(h2oServer, "mnist/test.csv.gz")
> train_hex[,785] <- as.factor(train_hex[,785])
> test_hex[,785] <- as.factor(test_hex[,785])
> dl_model <- h2o.deeplearning(x=c(1:784), y=785, training_frame=train_hex, l1=1e-5,
                                 activation="RectifierWithDropout", input_dropout_ratio=0.2,
                                 classification_stop=-1, hidden=c(1024,1024,2048), epochs=8000)
|=====
=====| 100%
> h2o.confusionMatrix(dl_model, test_hex)
Confusion Matrix - (vertical: actual; across: predicted):
      0   1   2   3   4   5   6   7   8   9   Error   Rate
0   974   1   1   0   0   0   2   1   1   0  0.00612 = 6 / 980
1    0 1135   0   1   0   0   0   0   0   0  0.00088 = 1 / 1,135
2    0   0 1028   0   1   0   0   3   0   0  0.00388 = 4 / 1,032
3    0   0   1 1003   0   0   0   3   2   1  0.00693 = 7 / 1,010
4    0   0   1   0 971   0   4   0   0   6  0.01120 = 11 / 982
5    2   0   0   5   0 882   1   1   1   0  0.01121 = 10 / 892
6    2   3   0   1   1   2 949   0   0   0  0.00939 = 9 / 958
7    1   2   6   0   0   0   0 1019   0   0  0.00875 = 9 / 1,028
8    1   0   1   3   0   4   0   2 960   3  0.01437 = 14 / 974
9    1   2   0   0   4   3   0   2   0  997 0.01189 = 12 / 1,009
Totals 981 1142 1038 1013 977 891 956 1031 964 1007 0.00830 = 83 / 10,000
```

- No data augmentation
- No distortions
- No convolutions
- No pre-training
- No ensemble

Table 1: Classification error rate comparison: DBN vs. DCN

DBN [3] (Hinton's)	DBN (MSR's)	DCN (Fine-tuning)	DCN (no Fine-tuning)	Shallow (D)CN (Fine-tuned single layer)
1.20%	1.06%	0.83%	0.95%	1.10%

Current World Record:  
H2O Deep Learning  
0.83% Test Set Error

Full Run: 10 Hours on 10-Node Cluster

Desktop: 2 Hours to 0.9% Test Set Error

# H2O Deep Water

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

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



Caffe

mxnet

H<sub>2</sub>O Deep Learning Algo

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**

# H2O Deep Water (cont.)

- **Deep learning framework integration:** TensorFlow, MXNet, and Caffe
  - GPU-accelerated training, modern deep learning architectures
- **Machine learning platform:** Other algos. Ensemble.
- **Ease of use and APIs:** Web UI (Flow), R, Python, Java
- **Deployment**

# MXNet

*dmlc*  
**mxnet** for Deep Learning

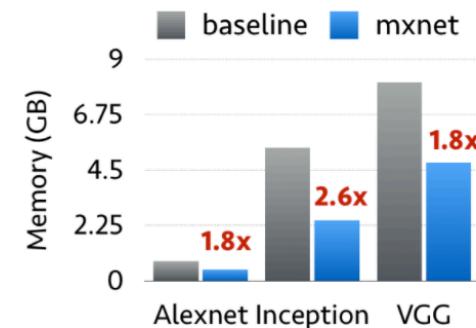
[build](#) passing [docs](#) [latest](#) [license](#) Apache 2.0



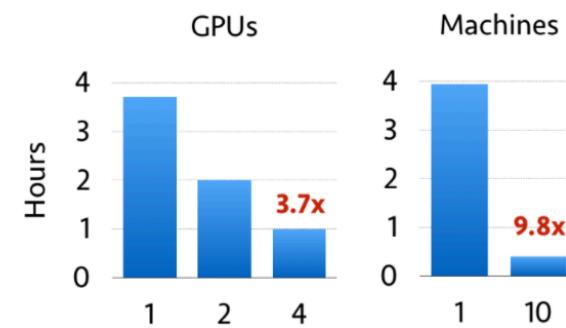
## Portable



## Efficient



## Scalable



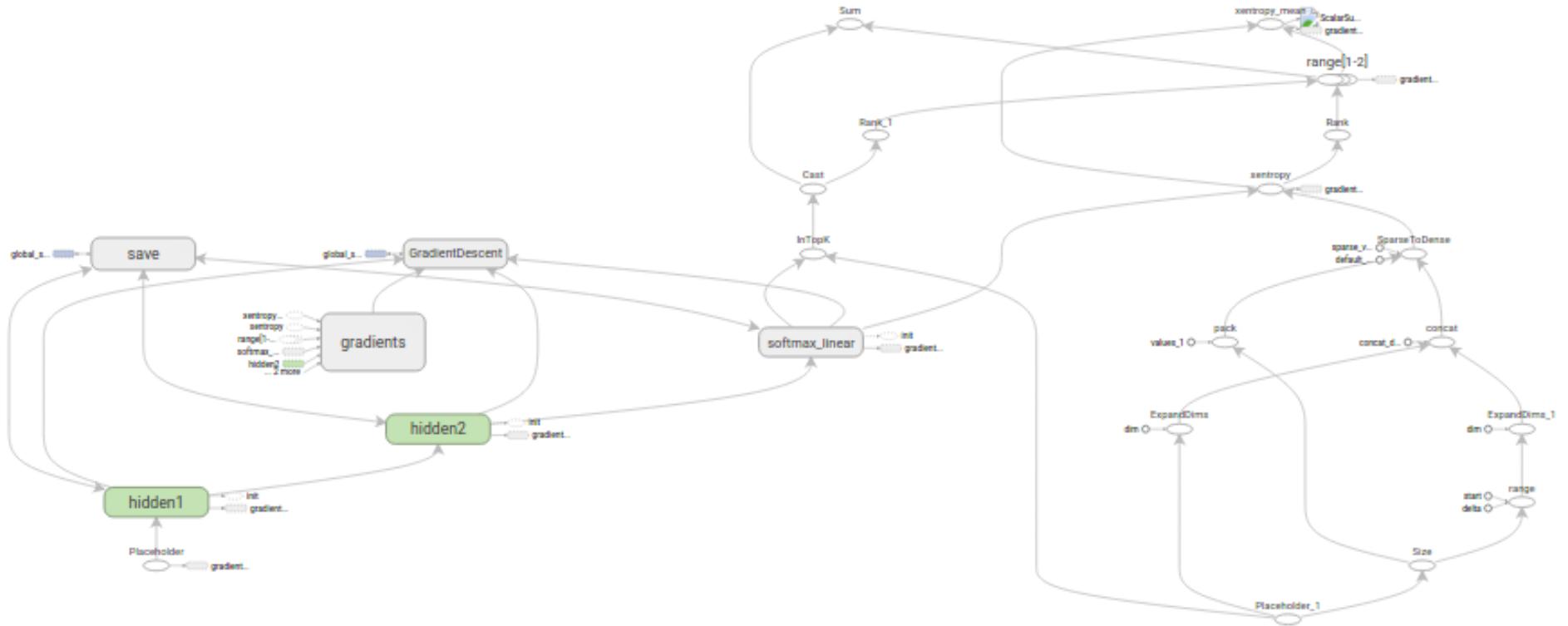
# TensorFlow



TensorFlow

Google

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



Caffe

Caffe



BERKELEY ARTIFICIAL INTELLIGENCE RESEARCH



# Setup

# Setup

- Extract data.
  - cd /gtc-2017/data
  - tar zxvf mnist\_png.tar.gz
- Start H2O
  - cd /gtc-2017
  - java -Xmx10g -jar /opt/h2o.jar &

# Setup (cont.)

- Start Jupyter notebook
  - cd /gtc-2017
  - jupyter notebook --allow-root --ip=\* &
    - This will give you a link with localhost. Copy that link.
- Go to browser
  - H2O: <http://you-ip-address:54321>
  - Jupyter: Paste link into browser. Replace localhost with your ip address.

# MNIST Data

# MNIST Database

MNIST handwritten digit database X yann.lecun.com/exdb/mnist/

## THE MNIST DATABASE of handwritten digits

Yann LeCun, Courant Institute, NYU  
Corinna Cortes, Google Labs, New York  
Christopher J.C. Burges, Microsoft Research, Redmond

The MNIST database of handwritten digits, available from this page, has a training set of 60,000 examples, and a test set of 10,000 examples. It is a subset of a larger set available from NIST. The digits have been size-normalized and centered in a fixed-size image.

It is a good database for people who want to try learning techniques and pattern recognition methods on real-world data while spending minimal efforts on preprocessing and formatting.

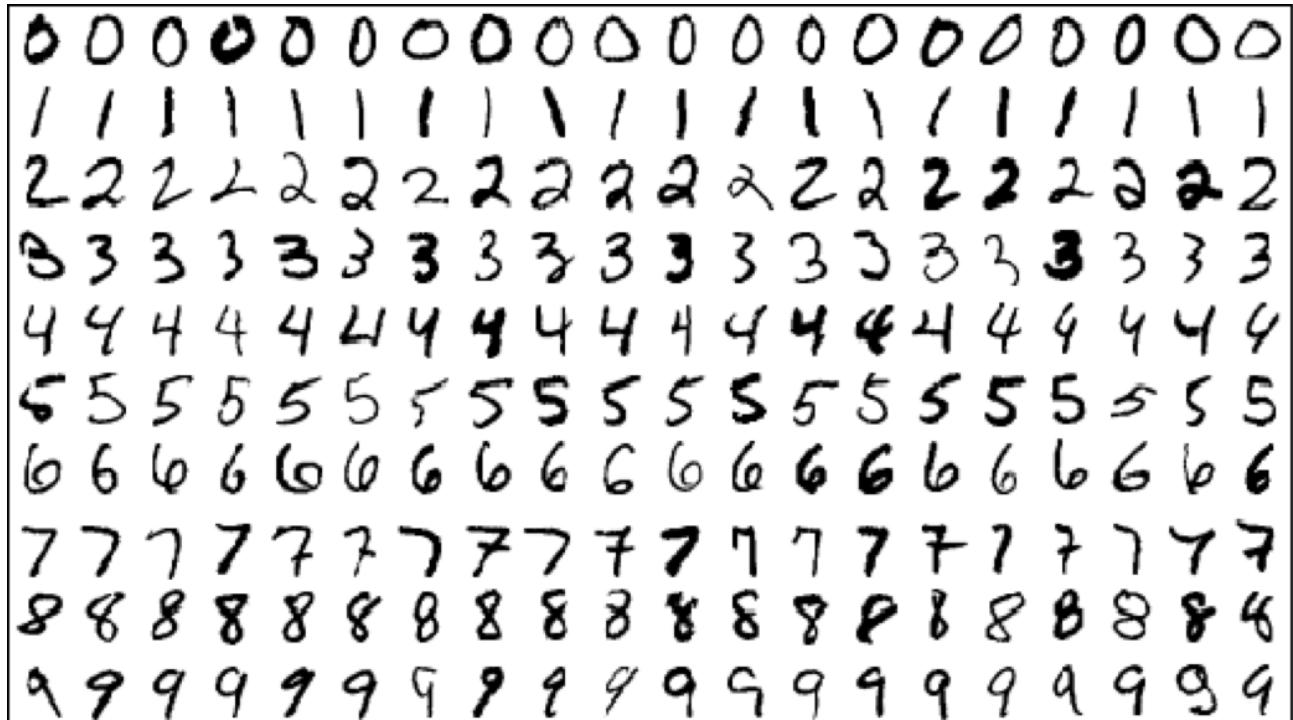
Four files are available on this site:

[train-images-idx3-ubyte.gz](#): training set images (9912422 bytes)  
[train-labels-idx1-ubyte.gz](#): training set labels (28881 bytes)  
[t10k-images-idx3-ubyte.gz](#): test set images (1648877 bytes)  
[t10k-labels-idx1-ubyte.gz](#): test set labels (4542 bytes)

please note that your browser may uncompress these files without telling you. If the files you downloaded have a larger size than the above, they have been uncompressed by your browser. Simply rename them to remove the .gz extension. Some people have asked me "my application can't open your image files". These files are not in any standard image format. You have to write your own (very simple) program to read them. The file format is described at the bottom of this page.

The original black and white (bilevel) images from NIST were size normalized to fit in a 20x20 pixel box while preserving their aspect ratio. The resulting images contain grey levels as a result of the anti-aliasing technique used by the normalization algorithm. the images were centered in a 28x28 image by computing the center of mass of the pixels, and translating the image so as to position this point at the center of the 28x28 field.

With some classification methods (particularly template-based methods, such as SVM and K-nearest neighbors), the error rate improves when the digits are centered by bounding box rather than center of mass. If you do this kind of pre-processing, you should report it in your publications.



# MNIST Digit

- 28 x 28 Pixels

Shape: (28, 28)

# Image as Pixel Vector

- 784 features (28x28)

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_P \end{bmatrix}$$

# MNIST Data Set

- Training Set: 60,000
- Testing/Validation: 10,000

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & \dots & x_{1,P} \\ \vdots & \ddots & \vdots \\ x_{N,1} & \dots & x_{N,P} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}$$

$$\mathbf{x}_i = \begin{bmatrix} x_{i,1} \\ \vdots \\ x_{i,P} \end{bmatrix}$$

where  $\mathbf{x}_i$  is  $i$ -th data record

# mnist\_training.hex

Actions: [View Data](#) [Split...](#) [Build Model...](#) [Predict](#) [Download](#) [Export](#)

Rows	Columns	Compressed Size
60000	2	5MB

## ▼ COLUMN SUMMARIES

label	type	Missing	Zeros	+Inf	-Inf	min	max	mean	sigma	cardinality	Actions
uri	string	0	0	0	0	.	.	.	.	.	<a href="#">Convert to enum</a>
label	enum	0	5923	0	0	0	9.0	.	.	10	<a href="#">Convert to numeric</a>

# mnist\_testing.hex

Actions: [View Data](#) [Split...](#) [Build Model...](#) [Predict](#) [Download](#) [Export](#)

Rows	Columns	Compressed Size
10000	2	807KB

## ▼ COLUMN SUMMARIES

label	type	Missing	Zeros	+Inf	-Inf	min	max	mean	sigma	cardinality	Actions
uri	string	0	0	0	0	.	.	.	.	.	<a href="#">Convert to enum</a>
label	enum	0	980	0	0	0	9.0	.	.	10	<a href="#">Convert to numeric</a>

# mnist\_training.hex

## DATA

[Previous 20 Columns](#)[Next 20 Columns](#)

Row	uri	label
1	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/6453.png	6
2	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/13970.png	6
3	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/21967.png	6
4	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/11277.png	6
5	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/27350.png	6
6	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/32119.png	6
7	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/44201.png	6
8	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/58770.png	6
9	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/39440.png	6
10	/home/ubuntu/sales-engineering/wen/gtc2017/data/mnist_png/training/6/41988.png	6

# Image Classification

# LeNet

Yann LeCun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. **Gradient-Based Learning Applied to Document Recognition.** In *Proceedings of the IEEE*, Nov. 1998. URL <http://yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf>

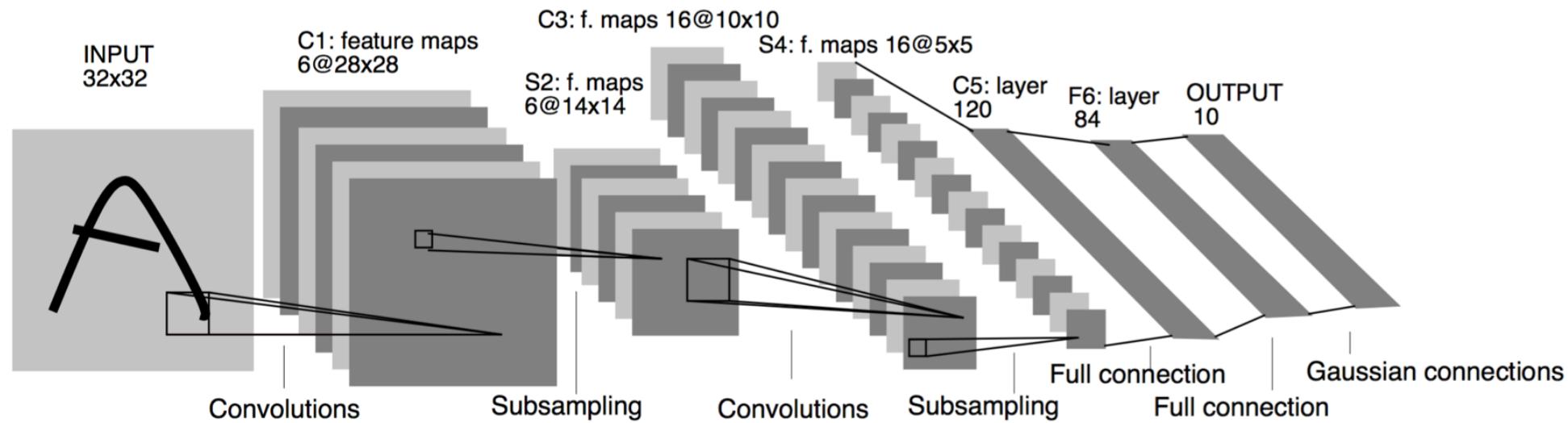
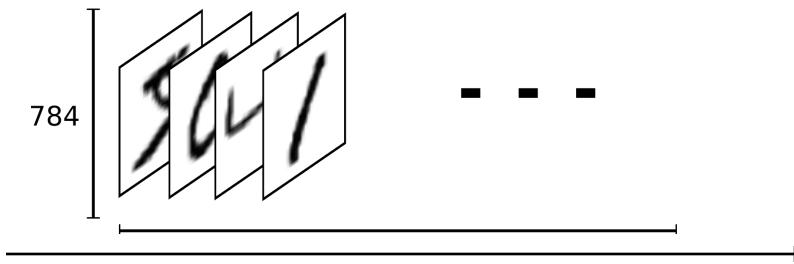


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

# Supervised Learning

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & \dots & x_{1,P} \\ \vdots & \ddots & \vdots \\ x_{N,1} & \dots & x_{N,P} \end{bmatrix}$$

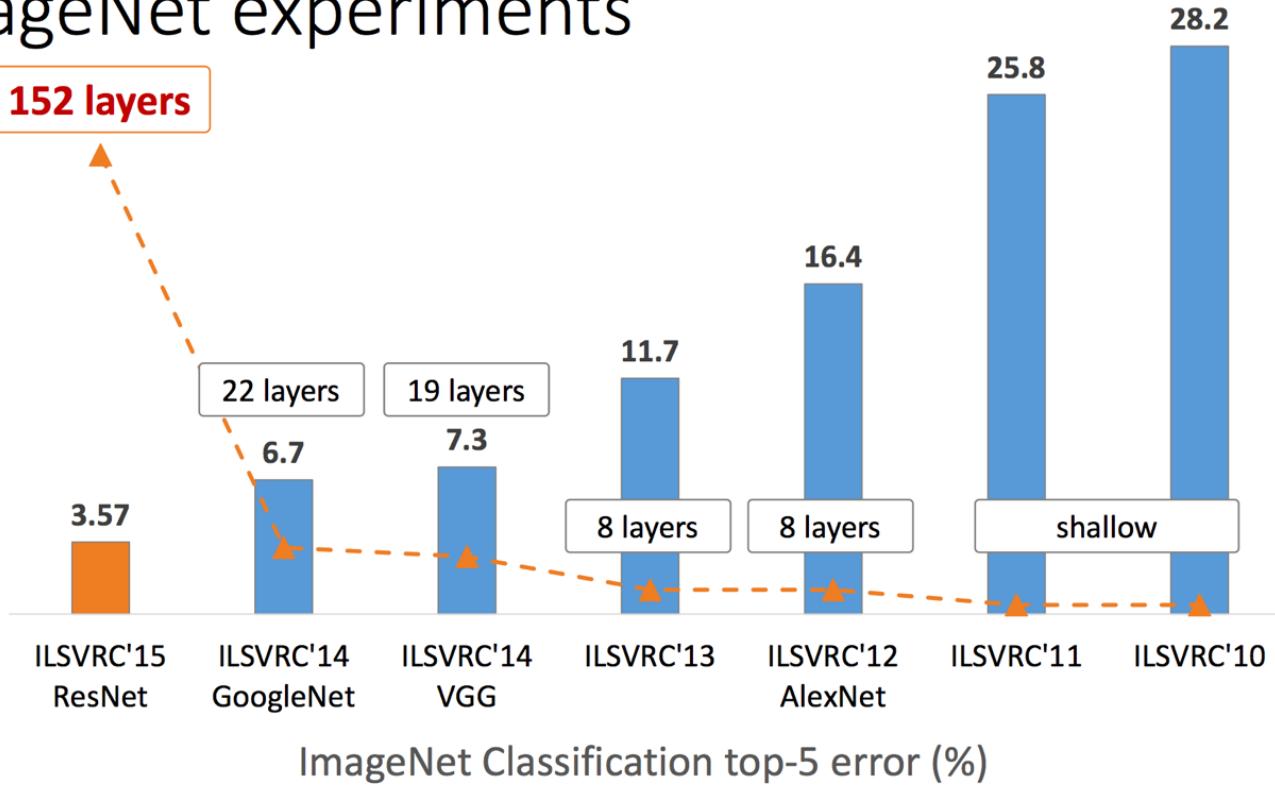


$$\mathbf{Y} = \begin{bmatrix} y_{1,1} & \dots & y_{1,C} \\ \vdots & \ddots & \vdots \\ y_{N,1} & \dots & y_{N,C} \end{bmatrix}$$

# ILSVRC Winners

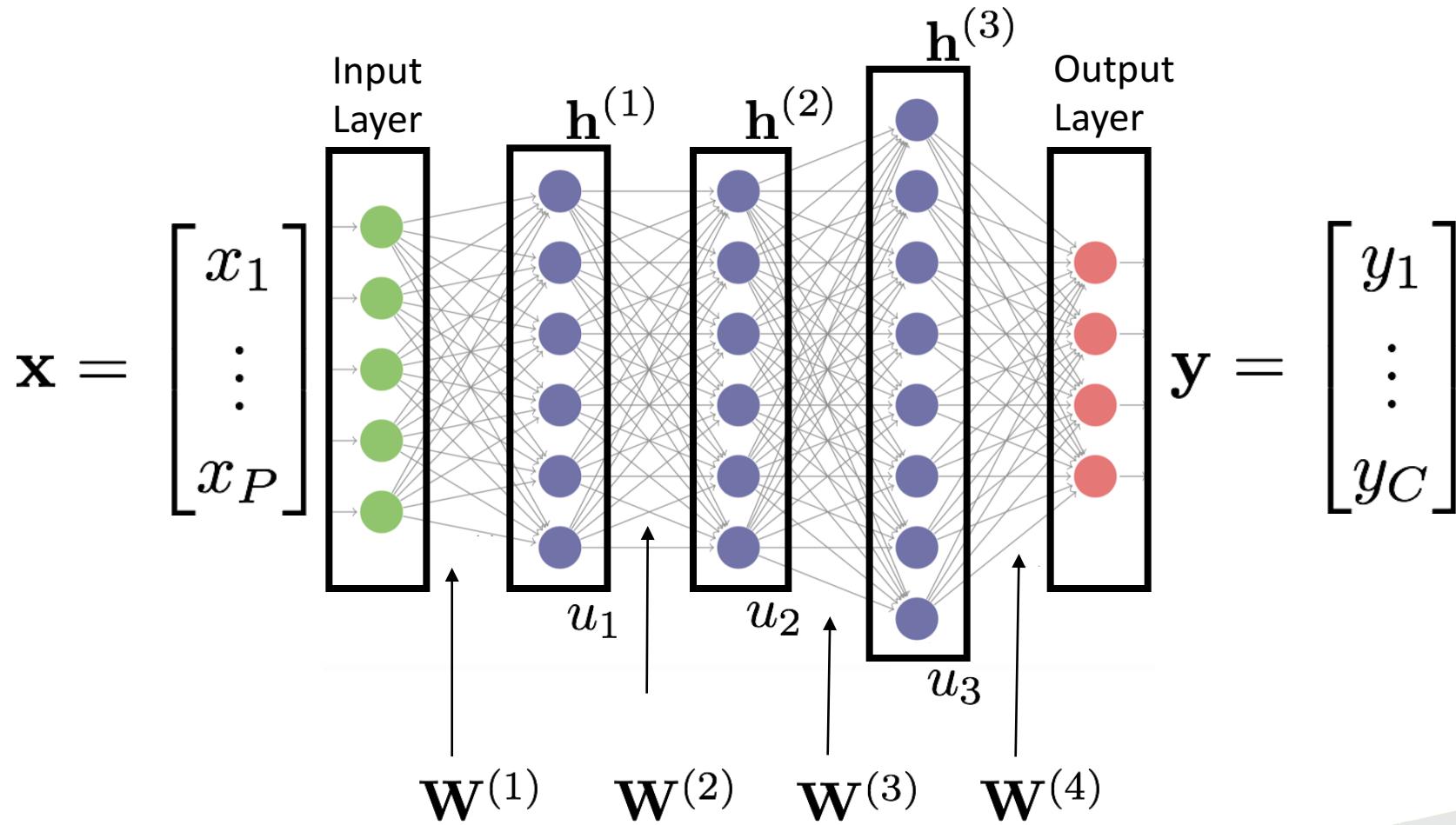
Microsoft  
Research

## ImageNet experiments

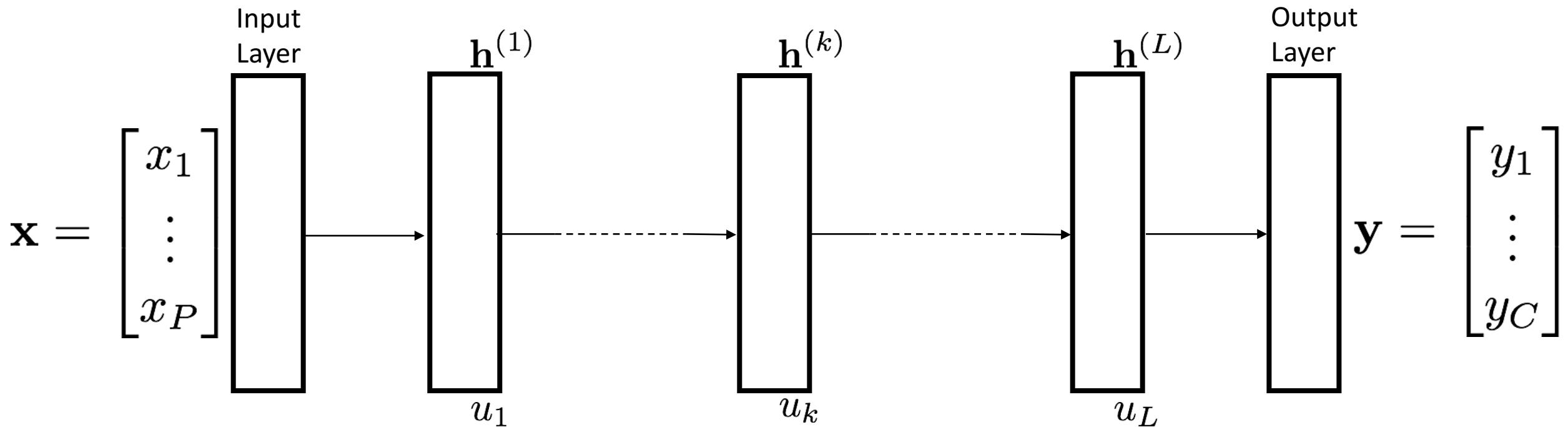


# Custom and User-Defined Networks

# Neural Network



# Graph



# MXNet

```
net = mx.sym.Variable('data')
net = mx.sym.FullyConnected(data=net, name='fc1', num_hidden=128)
net = mx.sym.Activation(data=net, name='relu1', act_type="relu")
net = mx.sym.FullyConnected(data=net, name='fc2', num_hidden=10)
net = mx.sym.SoftmaxOutput(data=net, name='out')
mx.viz.plot_network(net, shape={'data':(100,200)})
```

**Symbol.save**

Saves symbol to a file.

**BaseModule.save\_params**

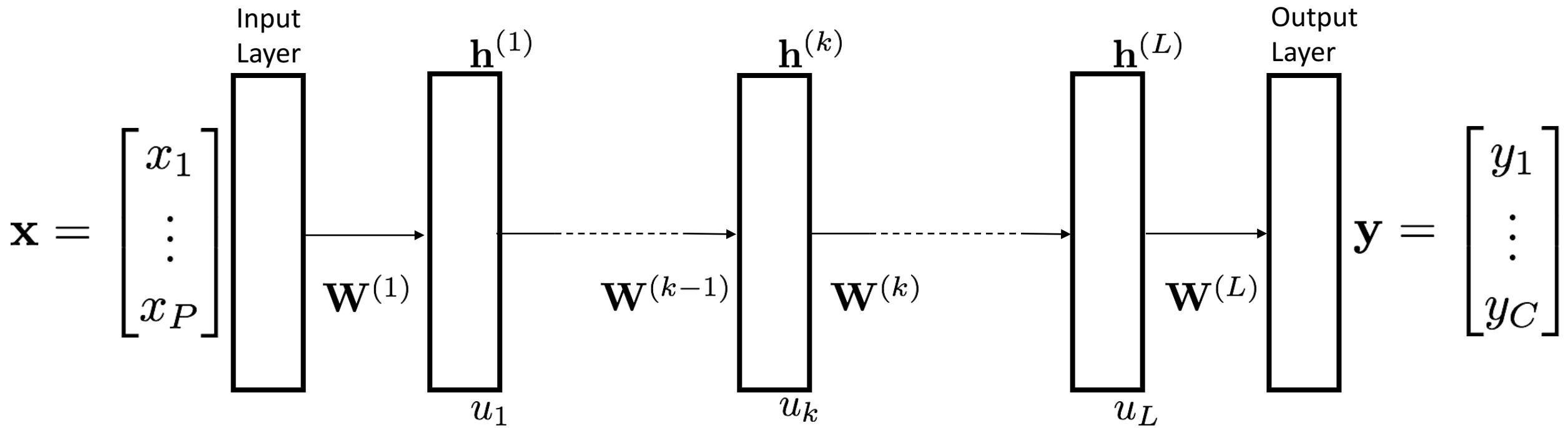
Save model parameters to file.

# TensorFlow

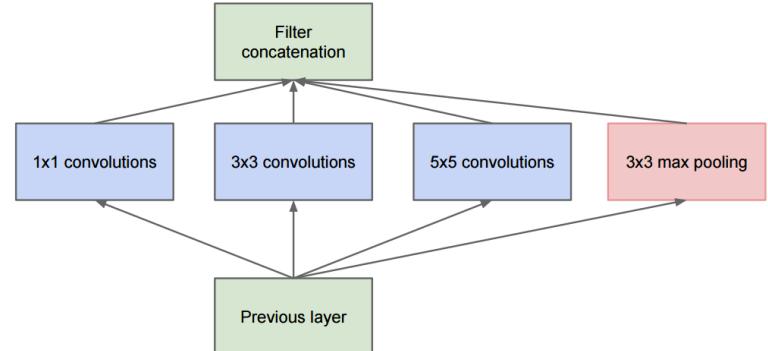
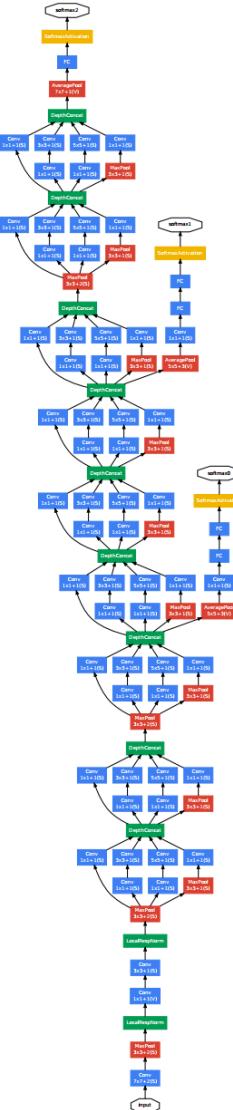
```
with tf.Session() as sess:  
    # Initializes all the variables.  
    sess.run(init_all_op)  
    # Runs to logit.  
    sess.run(logits)  
    # Creates a saver.  
    saver0 = tf.train.Saver()  
    saver0.save(sess, 'my-save-dir/my-model-10000')  
    # Generates MetaGraphDef.  
    saver0.export_meta_graph('my-save-dir/my-model-10000.meta')
```

# Pre-Trained Networks

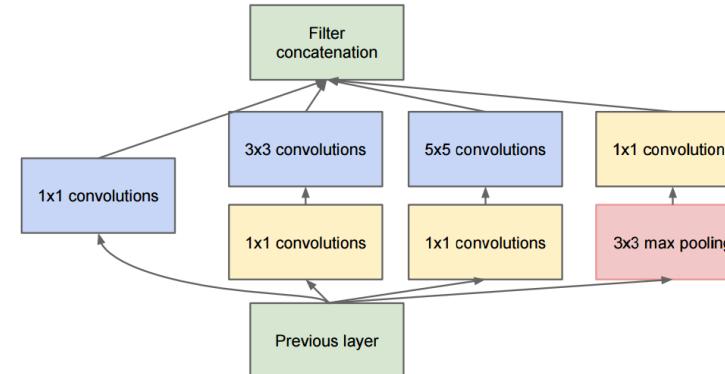
# Graph and Weights



# GoogLeNet and the Inception Module



(a) Inception module, naïve version



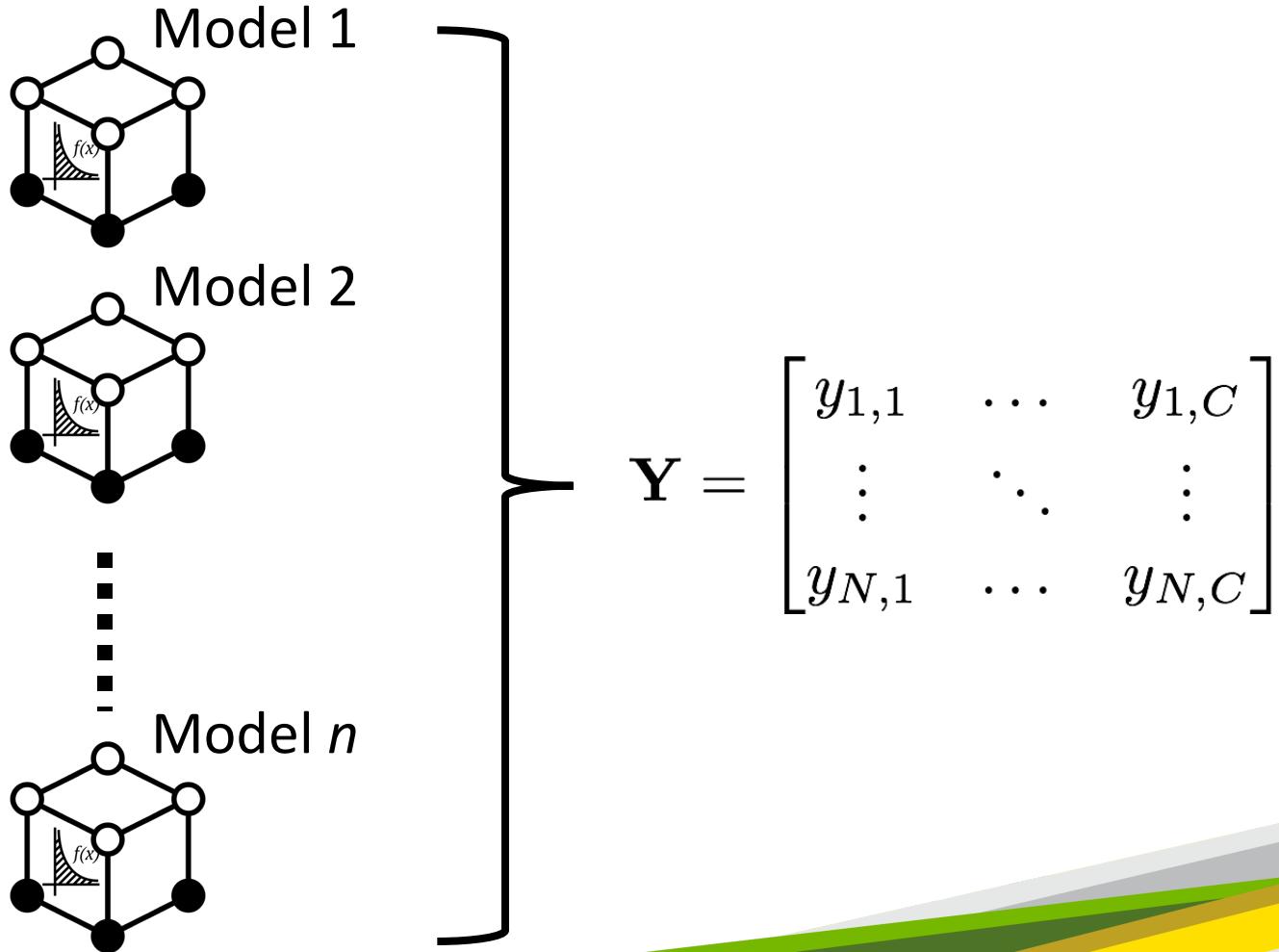
(b) Inception module with dimension reductions

Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott E. Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. **Going Deeper with Convolutions**. 2014. URL <https://arxiv.org/pdf/1409.4842.pdf>

# Ensembles

# Ensembles

$$\mathbf{X} = \begin{bmatrix} x_{1,1} & \dots & x_{1,P} \\ \vdots & \ddots & \vdots \\ x_{N,1} & \dots & x_{N,P} \end{bmatrix}$$



# Stacking / Super Learner

- Train a second machine learning model (called a “meta-learner”) to learn the optimal combination of the base learners

# “Level-Zero” Data

$$n \left\{ \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix} \right| \begin{bmatrix} y \\ \vdots \\ y \end{bmatrix} \right\}^m$$

1. Start with design matrix,  $X$ , and response,  $y$
2. Specify  $L$  base learners
3. Specify a meta-learner
4. Perform  $k$ -fold cross-validation (CV) on each of the  $L$  learners

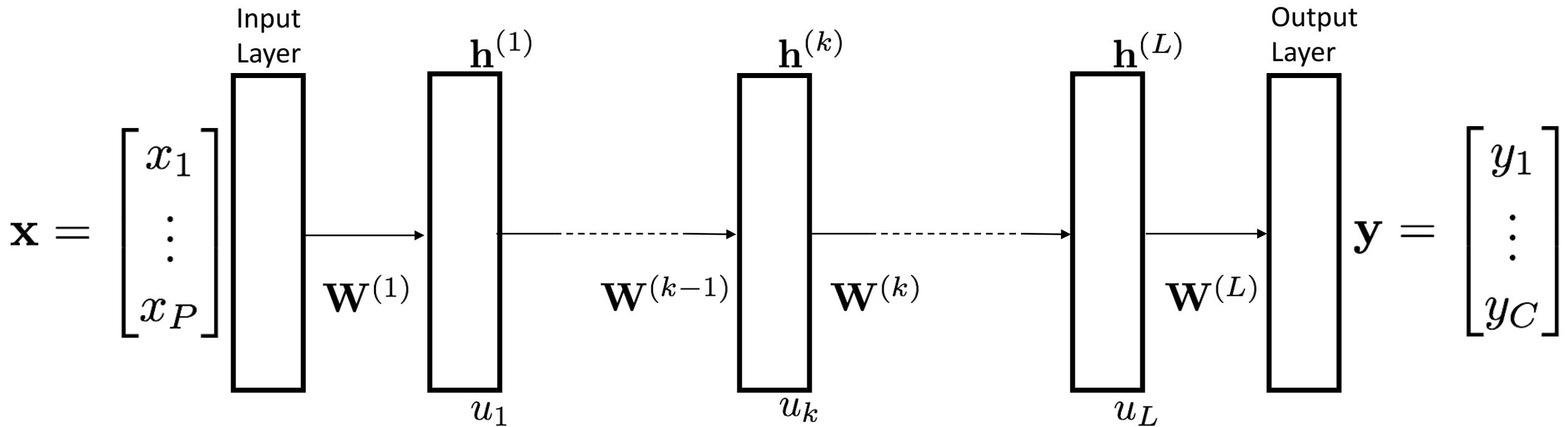
# “Level-One” Data

$$n \left\{ \begin{bmatrix} p_1 \\ \vdots \\ p_L \end{bmatrix} \quad \begin{bmatrix} y \end{bmatrix} \right\} \rightarrow n \left\{ \underbrace{\begin{bmatrix} \quad & \quad & \quad \\ \quad & \quad & \quad \\ \quad & \quad & \quad \end{bmatrix}}_L \quad \begin{bmatrix} z \\ y \end{bmatrix} \right\}$$

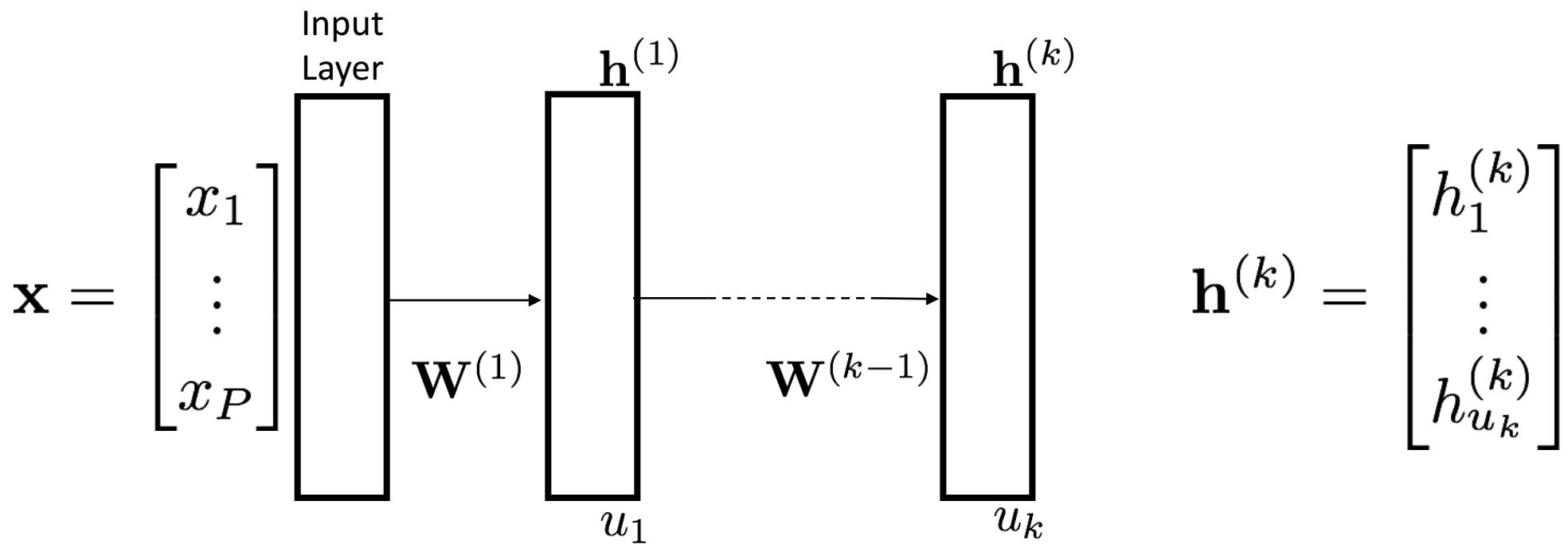
1. Collect the predicted values from  $k$ -fold CV that was performed on each of the  $L$  base learners
2. Column-bind these prediction vectors together to form a new design matrix,  $Z$
3. Train the metalearner using  $Z, y$

# Deep Features

# Trained Network



# Hidden Layer Features



# Matrix of Hidden Layer (Deep) Features

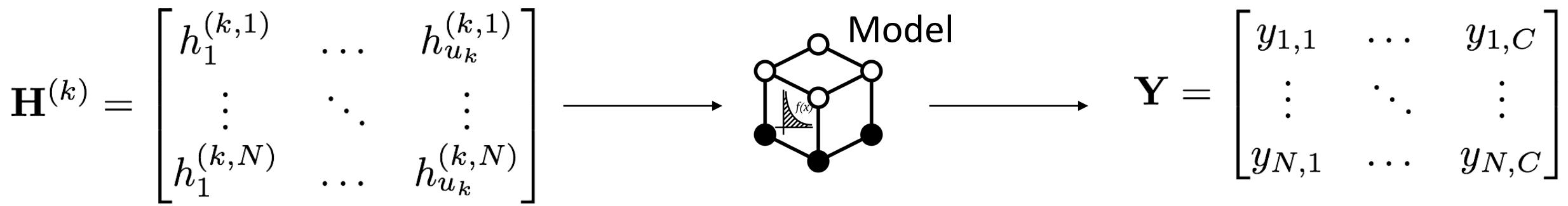
$$\mathbf{X} = \begin{bmatrix} x_{1,1} & \dots & x_{1,P} \\ \vdots & \ddots & \vdots \\ x_{N,1} & \dots & x_{N,P} \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}$$

$$\mathbf{H}^{(k)} = \begin{bmatrix} h_1^{(k,1)} & \dots & h_{u_k}^{(k,1)} \\ \vdots & \ddots & \vdots \\ h_1^{(k,N)} & \dots & h_{u_k}^{(k,N)} \end{bmatrix} = \begin{bmatrix} [\mathbf{h}^{(k,1)}]^T \\ \vdots \\ [\mathbf{h}^{(k,N)}]^T \end{bmatrix}$$

# Deep Features

- Can use the deep features as input features for another model (algorithm)
- Can use the deep features as data representation for similarity analysis

# Deep Features Classifier



# Distance

$$\mathbf{A} = \begin{bmatrix} a_{1,1} & \dots & a_{1,P} \\ \vdots & \ddots & \vdots \\ a_{M,1} & \dots & a_{M,P} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_M^T \end{bmatrix} \quad \text{where } \mathbf{a}_i = \begin{bmatrix} a_{i,1} \\ \vdots \\ a_{i,P} \end{bmatrix}$$

$$\mathbf{B} = \begin{bmatrix} b_{1,1} & \dots & b_{1,P} \\ \vdots & \ddots & \vdots \\ b_{O,1} & \dots & b_{O,P} \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1^T \\ \vdots \\ \mathbf{b}_O^T \end{bmatrix} \quad \text{where } \mathbf{b}_i = \begin{bmatrix} b_{i,1} \\ \vdots \\ b_{i,P} \end{bmatrix}$$

# “Distance” Matrix

**distance** ( $\mathbf{A}, \mathbf{B}$ ) =  $\mathbf{Z}$  :  $z_{i,j} = \text{similarity} (\mathbf{a}_i, \mathbf{b}_j) = \text{sim} (\mathbf{a}_i, \mathbf{b}_j)$   
, where  $\mathbf{A} \in \mathbb{R}^{M \times P}$ ,  $\mathbf{B} \in \mathbb{R}^{O \times P}$ ,  $\mathbf{Z} \in \mathbb{R}^{M \times O}$

$$\mathbf{Z} = \begin{bmatrix} \text{sim}(\mathbf{a}_1, \mathbf{b}_1) & \dots & \text{sim}(\mathbf{a}_1, \mathbf{b}_O) \\ \vdots & \ddots & \vdots \\ \text{sim}(\mathbf{a}_M, \mathbf{b}_1) & \dots & \text{sim}(\mathbf{a}_M, \mathbf{b}_O) \end{bmatrix}$$

ai.stanford.edu/~jkrause/cars/ai.stanford.edu/~jkrause/cars/car\_dataset.html

## Cars Dataset



### Overview

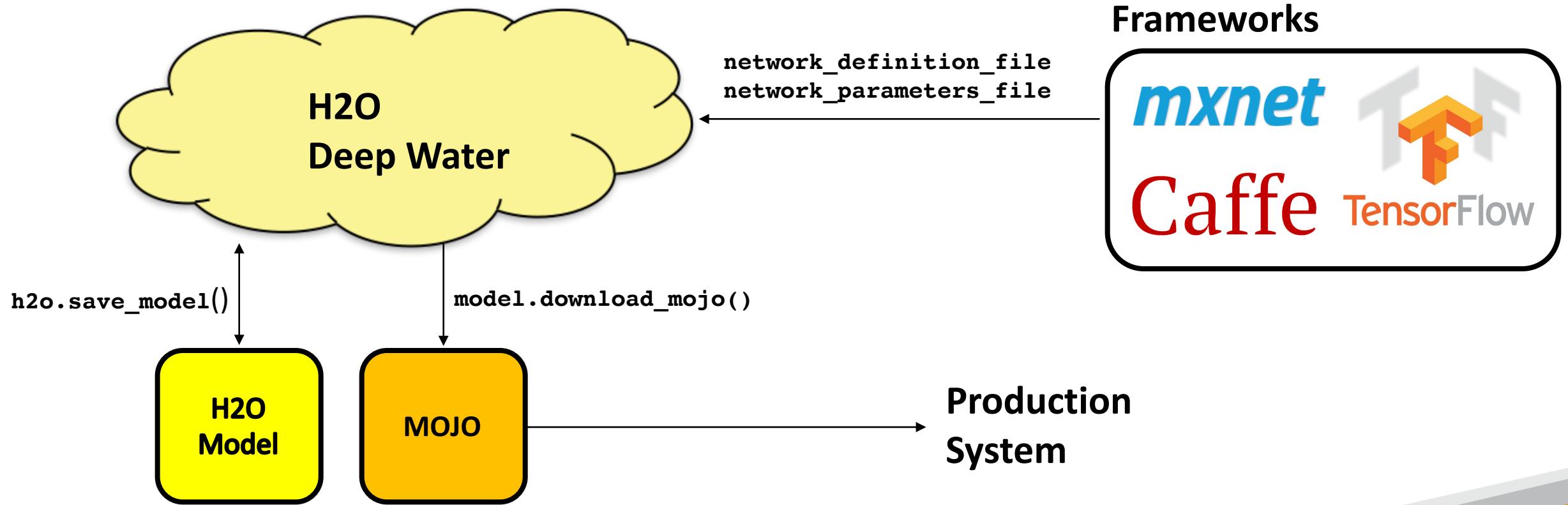
The *Cars* dataset contains 16,185 images of 196 classes of cars. The data is split into 8,144 training images and 8,041 testing images, where each class has been split roughly in a 50-50 split. Classes are typically at the level of *Make, Model, Year*, e.g. 2012 Tesla Model S or 2012 BMW M3 coupe.



[http://ai.stanford.edu/~jkrause/cars/car\\_dataset.html](http://ai.stanford.edu/~jkrause/cars/car_dataset.html)

# Deployment

# Deployment



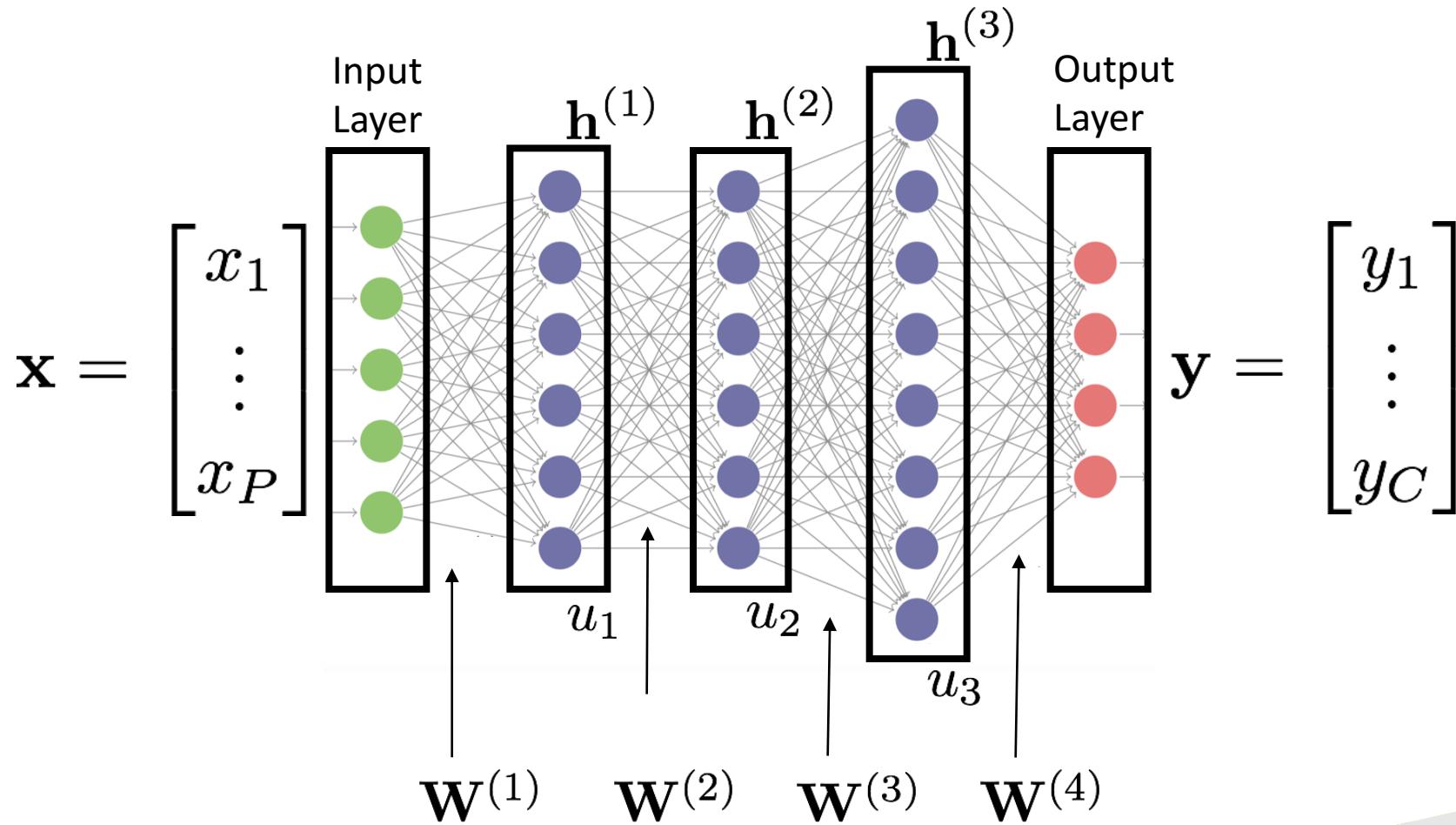
# Slides and Materials

<https://github.com/h2oai/h2o-tutorials/gtc-2017-deep-water>

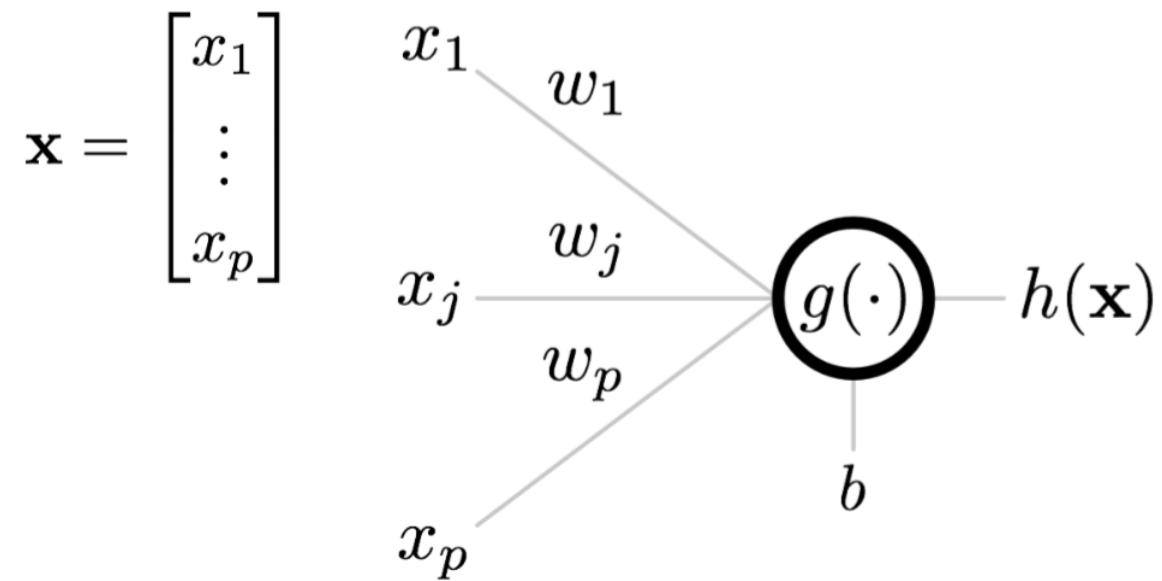


# Appendix: Deep Learning Review

# Neural Network

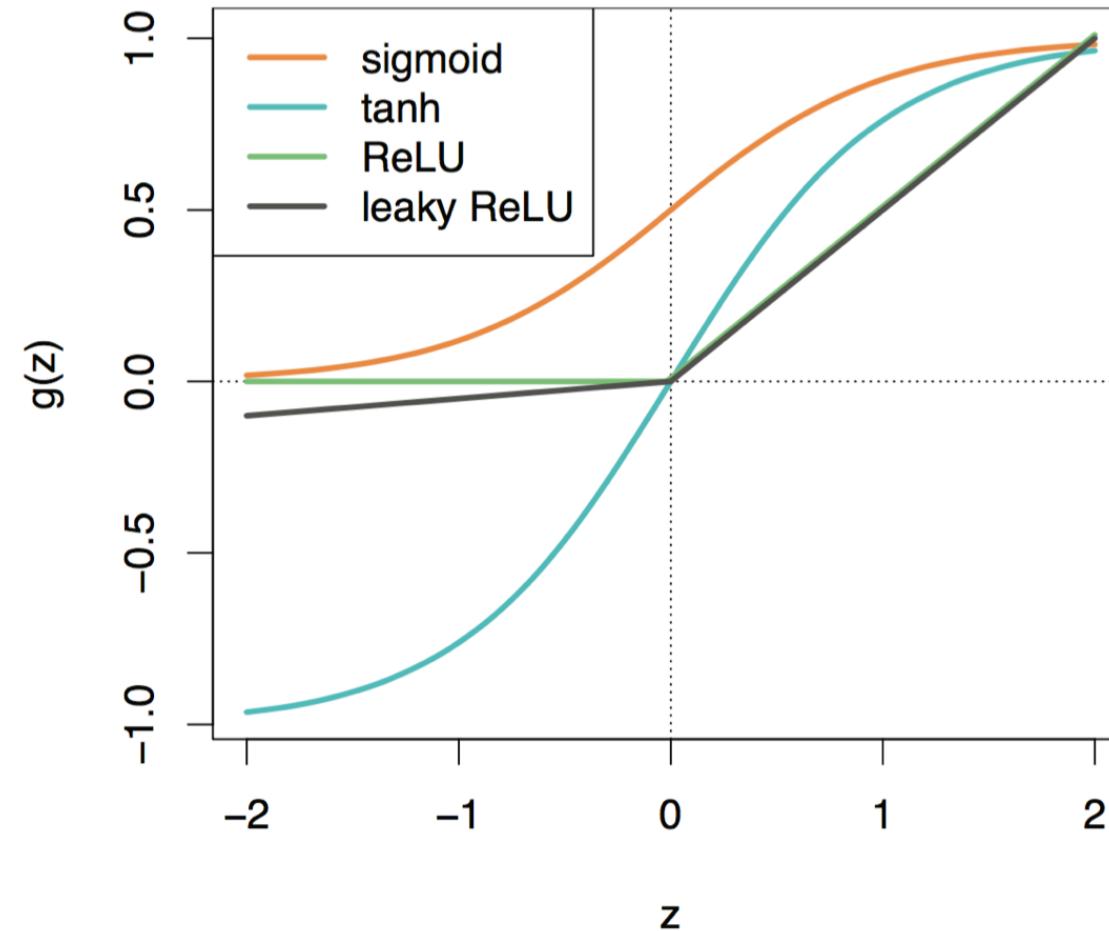


# Neuron

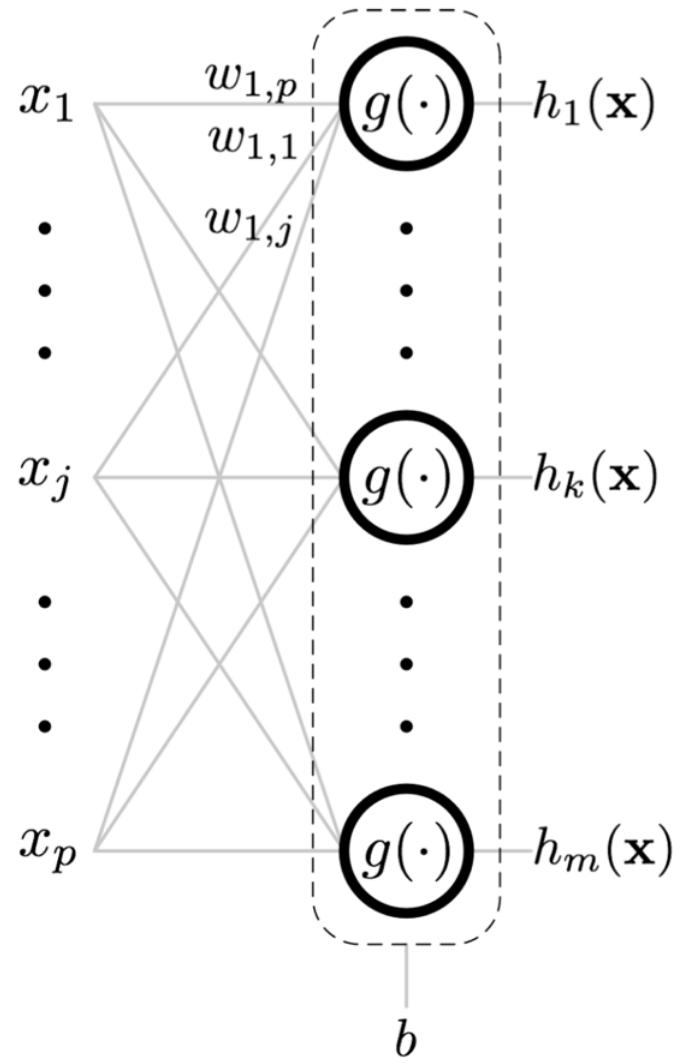


$$\begin{aligned} h(\mathbf{x}) &= g \left( \sum_{i=1}^p w_i x_i + b \right) \\ &= g(\mathbf{w}^T \mathbf{x} + b) \end{aligned}$$

# Activation Functions

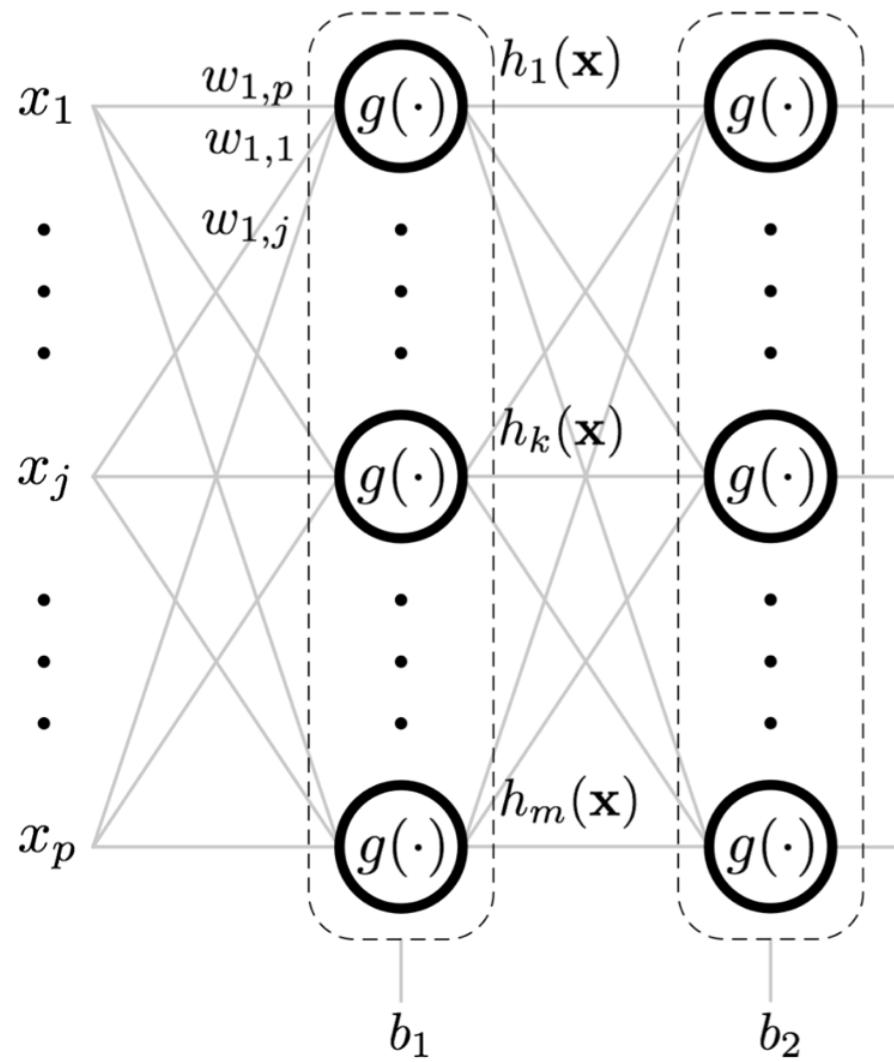


# Layer of Neurons



$$h_m(\mathbf{x}) = g \left( \sum_{i=1}^p w_{m,i} x_i + b \right)$$

# Hidden Layers



- ▶ Output from one layer is input to next layer
- ▶  $L$  hidden layers
- ▶ Layer pre-activation,  $k > 0$   
$$a^{(k)}(\mathbf{x}) = \mathbf{W}^{(k)}\mathbf{h}^{(k-1)} + \mathbf{b}^{(k)}$$
$$\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x}$$
- ▶ Hidden layer activation,  
 $k \in [1, L]$   
$$\mathbf{h}^{(k)}(\mathbf{x}) = g((a^{(k)}(\mathbf{x}))$$

# Output Layer

$$0 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, 1 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \dots, 9 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$$\Pr(y = k|\mathbf{x}) = \begin{bmatrix} \Pr(y = 0|\mathbf{x}) \\ \Pr(y = 1|\mathbf{x}) \\ \vdots \\ \Pr(y = 9|\mathbf{x}) \end{bmatrix}$$

# Output Layer Requirements

- ▶ All probabilities must sum up to one.

$$\sum_{k \in K} \Pr(y = k | \mathbf{x}) = 1$$

- ▶ Probabilities must be between 0 and 1.

$$0 \leq \Pr(y = k | \mathbf{x}) \leq 1, k \in K$$

# Softmax

$$\begin{aligned}\text{softmax}_k(a_1, a_2, \dots, a_{|K|}) &= \frac{\exp(a_k)}{\sum_{j \in K} \exp(a_j)} \\ &= \frac{\text{evidence for class } k}{\text{total evidence for all classes}}\end{aligned}$$

# Cross Entropy

$$H(p, q) = - \sum_i p_i \log q_i$$

$$= - \sum_i p_i \log_e q_i$$

$$= - \sum_i p_i \ln q_i$$

# Cross-Entropy (cont.)

- Binary Classification:  $p \in \{y, 1 - y\}$ ,  $q \in \{\hat{y}, 1 - \hat{y}\}$

$$\begin{aligned} H(p, q) &= - \sum_{i=1}^2 p_i \ln q_i \\ &= -p_1 \ln q_1 - p_2 \ln q_2 \\ H(y, \hat{y}) &= -y \ln \hat{y} - (1 - y) \ln(1 - \hat{y}) \end{aligned}$$

- $y = 1$ :  $H(1, \hat{y}) = \ln \hat{y}$
- $y = 0$ :  $H(0, \hat{y}) = \ln(1 - \hat{y})$

# Cross-Entropy Loss Function

$$\begin{aligned} L(\hat{\mathbf{y}}, \mathbf{y}) &= \frac{1}{N} \sum_{i=1}^N H(p_j, q_j) \\ &= -\frac{1}{N} \sum_{i=1}^N \left[ \sum_{j=1}^K y_j^{(i)} \log \hat{y}_j^{(i)} \right] \\ &= -\frac{1}{N} \sum_{i=1}^N \left[ \sum_{j=1}^{10} y_j^{(i)} \log \hat{y}_j^{(i)} \right] \end{aligned}$$

# Training and Optimization

$$\underset{\mathbf{W}, \mathbf{b}}{\text{minimize}} -\frac{1}{N} \sum_{i=1}^N \left[ \sum_{j=1}^K y_j \log \hat{y}_j \right]$$

# Appendix: Convolutional Neural Networks

# Frobenius Inner Product

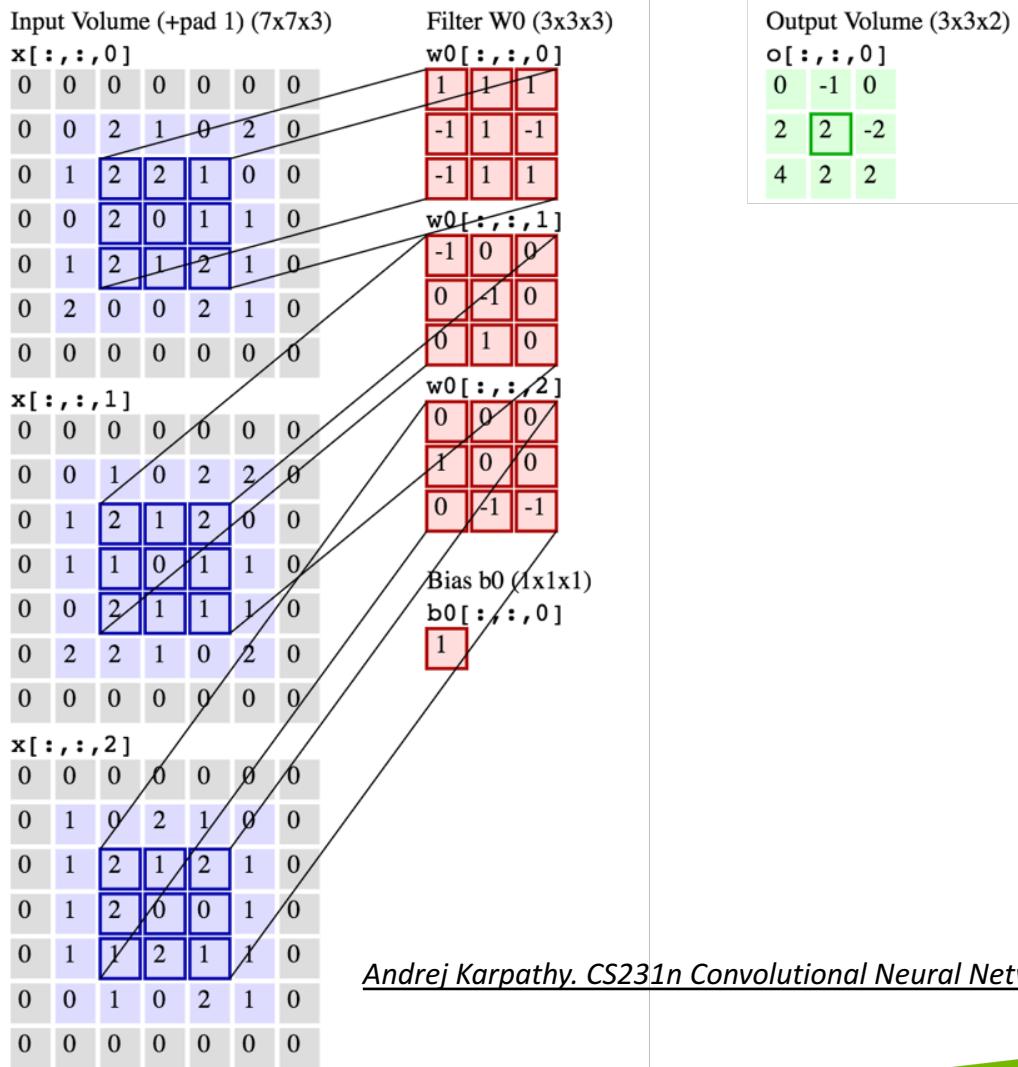
$$\mathbf{X} = \begin{bmatrix} 2 & 2 & 1 \\ 2 & 0 & 1 \\ 2 & 1 & 2 \end{bmatrix}, \mathbf{K} = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 1 & -1 \\ -1 & 1 & 1 \end{bmatrix}$$

$$\begin{aligned}\langle \mathbf{X}, \mathbf{K} \rangle_F &= \sum_{i,j} x_{i,j} k_{i,j} \\&= (2 * 1) + (2 * 1) + (1 * 1) + (2 * -1) + (0 * 1) + (1 * -1) + (2 * -1) + (1 * 1) + (2 * 1) \\&= 2 + 2 + 1 - 2 + 0 - 1 - 2 + 1 + 2 \\&= 3\end{aligned}$$

# Convolution

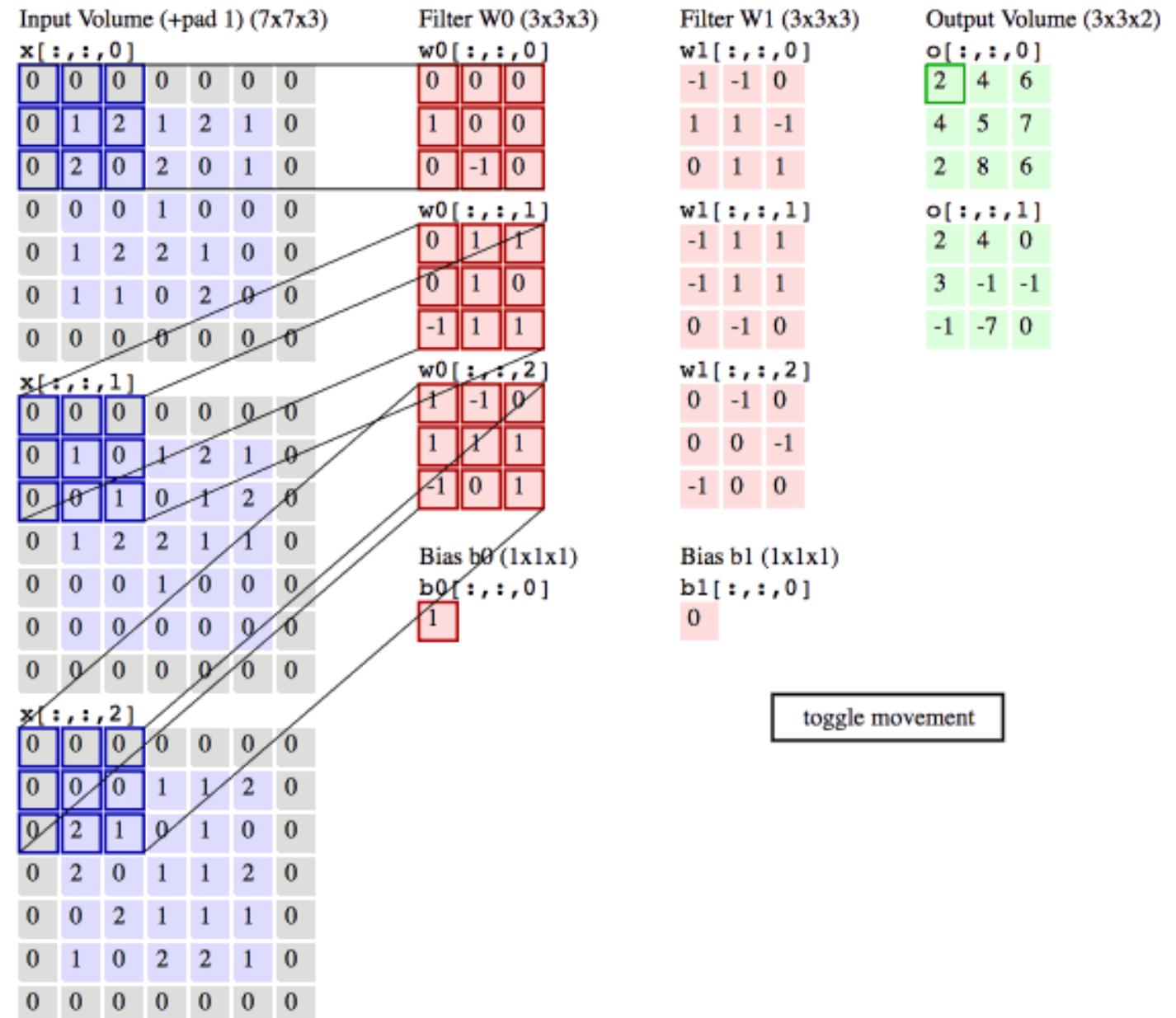
$$\mathbf{X} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 1 & 0 & 2 & 0 \\ 0 & 1 & 2 & 2 & 1 & 0 & 0 \\ 0 & 0 & 2 & 0 & 1 & 1 & 0 \\ 0 & 1 & 2 & 1 & 2 & 1 & 0 \\ 0 & 2 & 0 & 0 & 2 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$
$$\mathbf{K} = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 1 & -1 \\ -1 & 1 & 1 \end{bmatrix}$$
$$\mathbf{A} = \begin{bmatrix} 1 & 4 & 0 & -4 & 1 \\ 3 & 4 & 1 & 4 & 1 \\ 4 & 9 & 3 & 5 & 0 \\ 3 & 0 & 2 & 5 & 0 \\ 5 & 2 & 3 & 5 & 2 \end{bmatrix}$$

# Convolution



$$\sum_d \sum_{i,j} x_{i,j,d} k_{i,j,d}$$

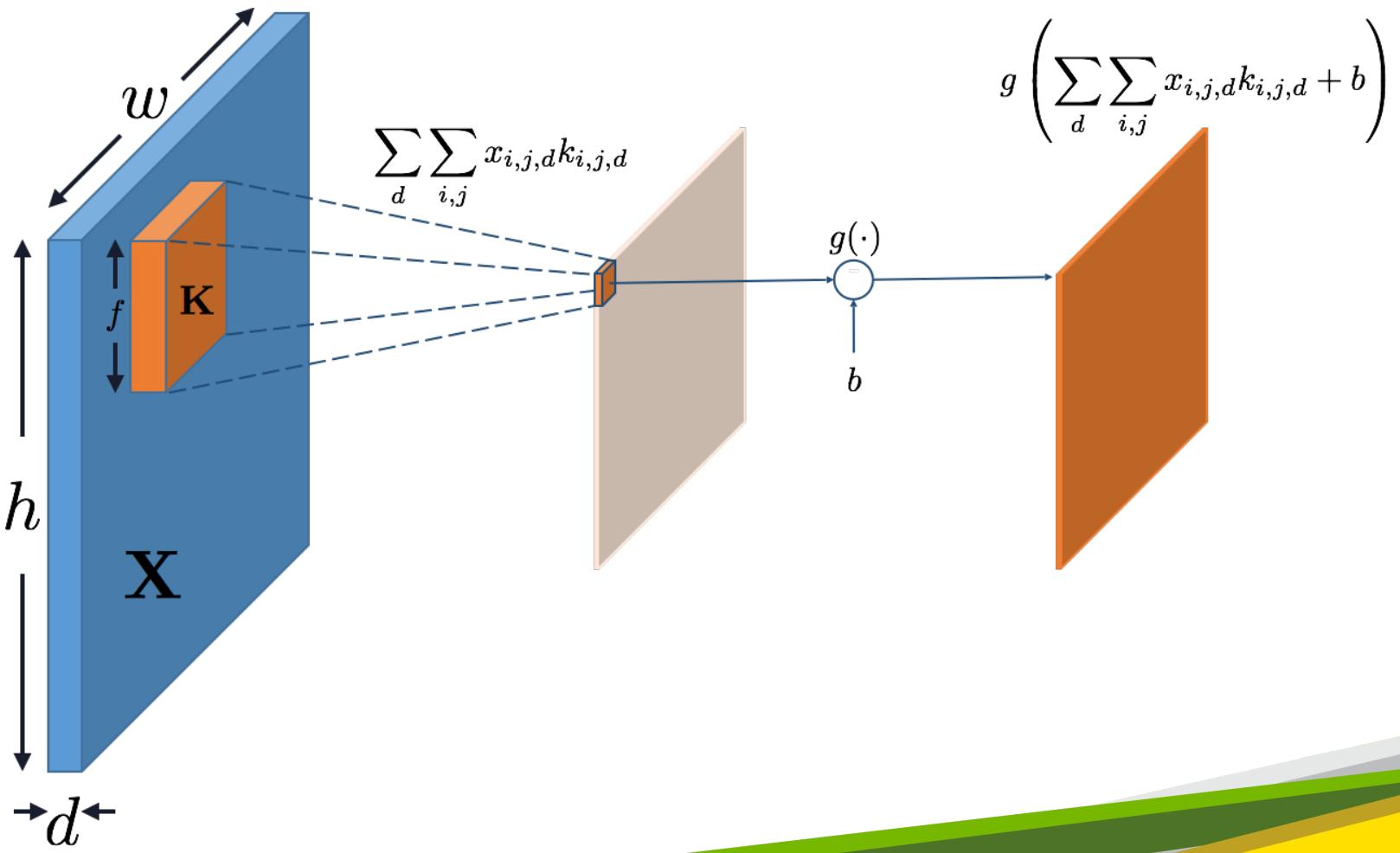
# Convolution Layer



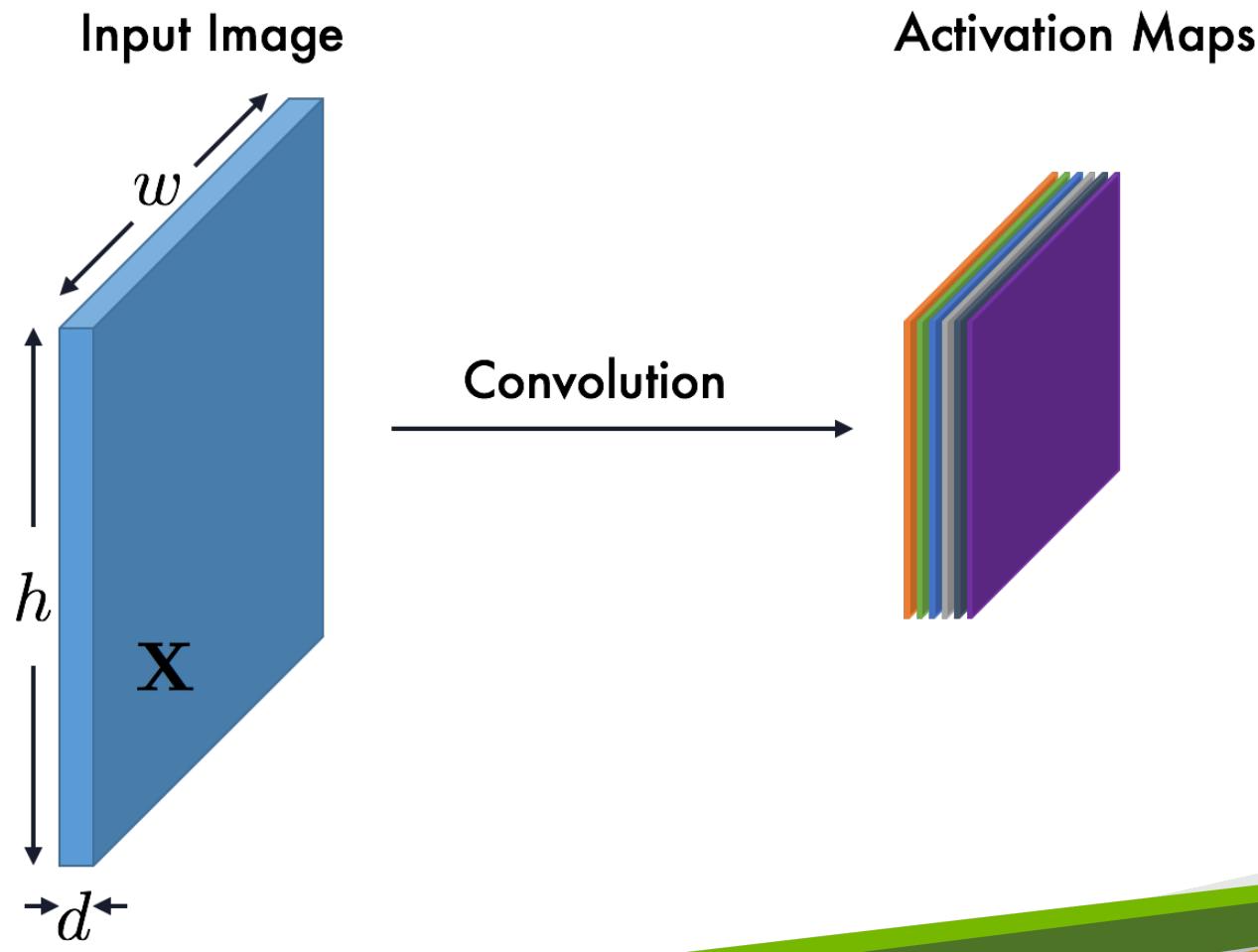
[Andrey Karpathy. CS231n Convolutional Neural Networks for Visual Recognition.](#)



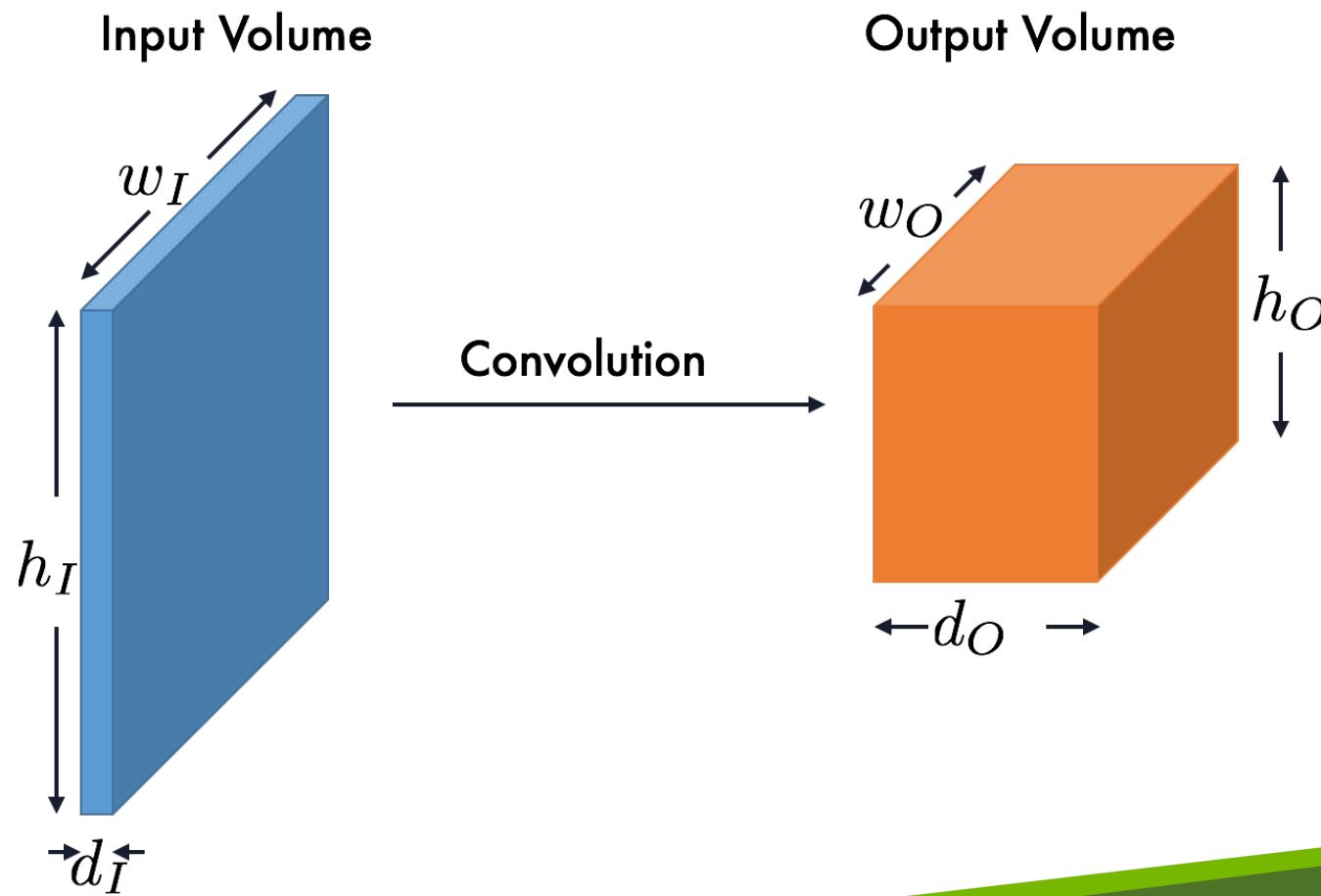
# Convolution Layer



# Convolution Layer



# Convolution Layer

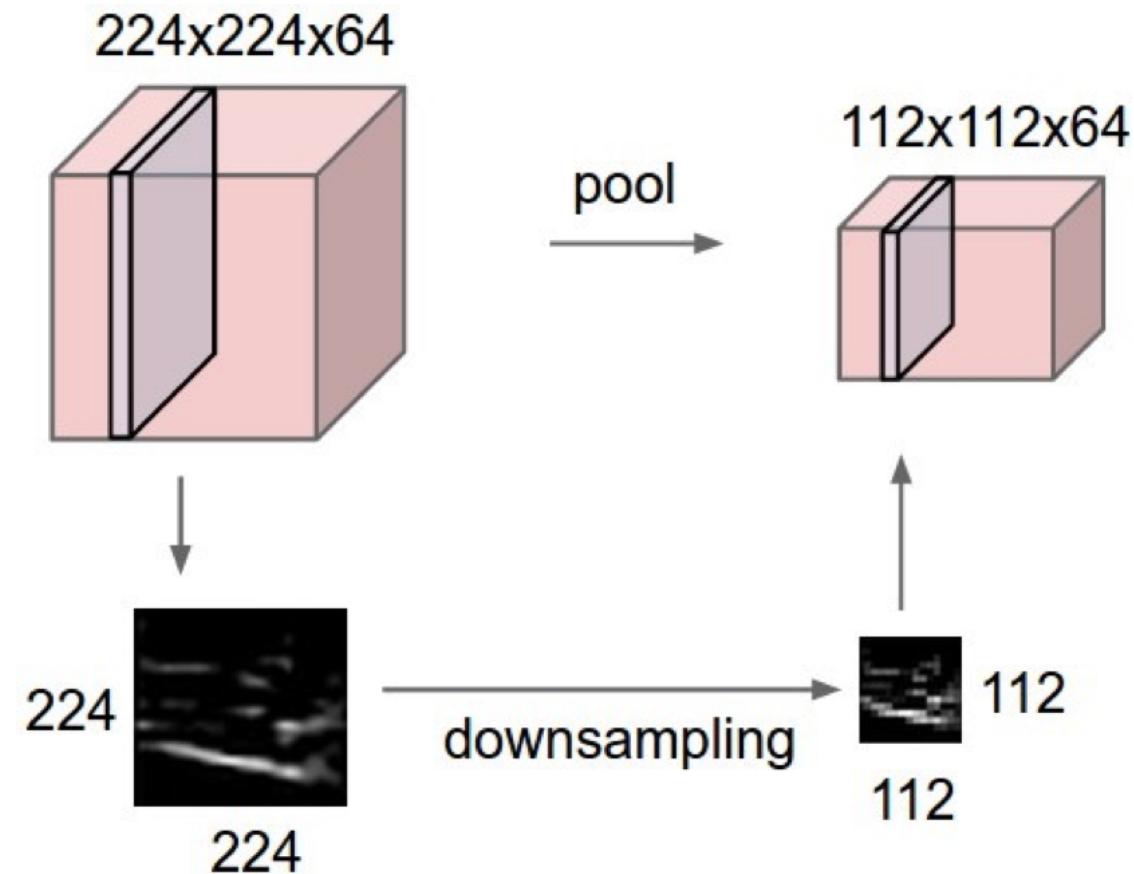


$$w_O = \frac{w_I - f + 2p}{s} + 1$$

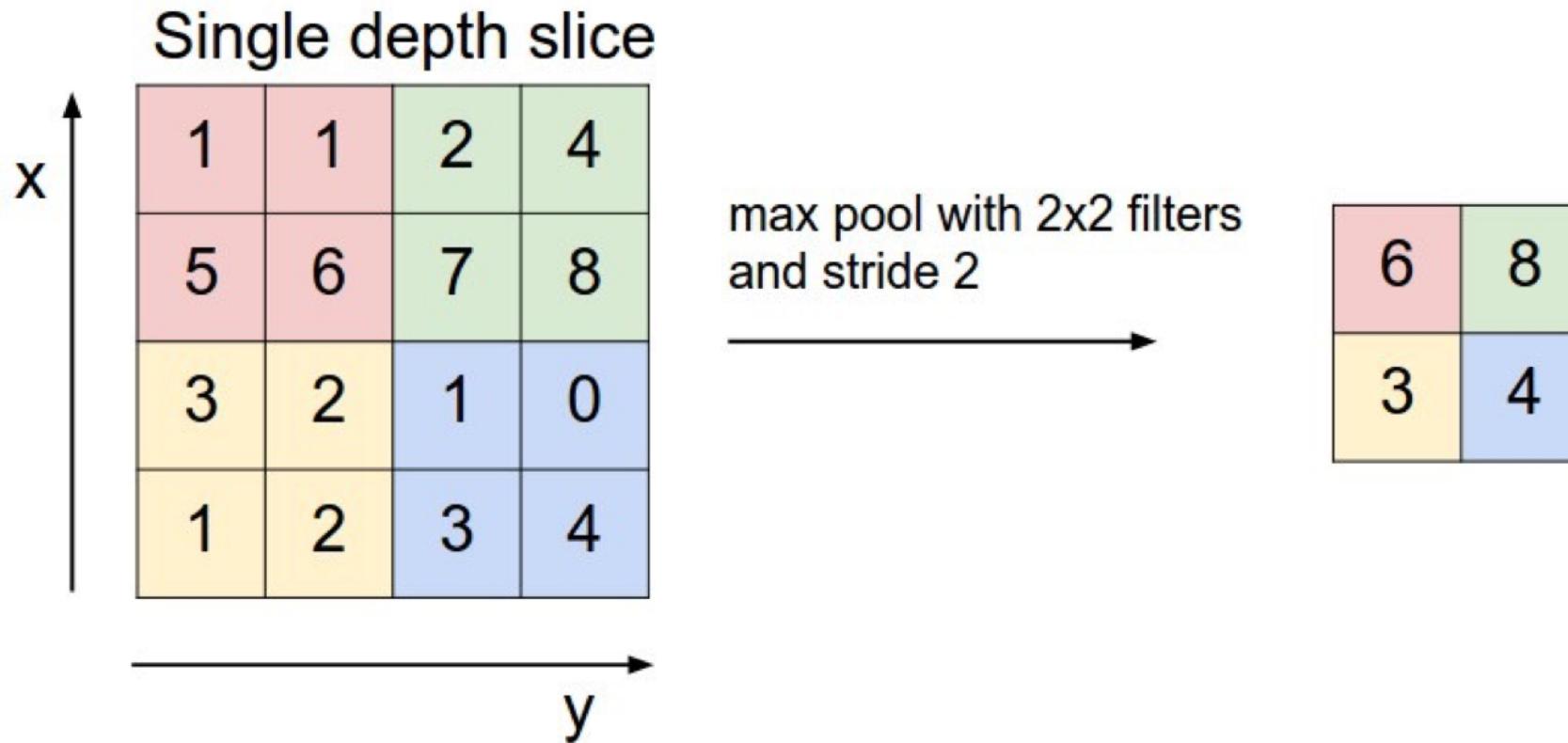
$$h_O = \frac{h_I - f + 2p}{s} + 1$$

$$d_O = m$$

# Pooling



# Max Pooling



# LeNet

Yann LeCun, Leon Bottou, Yoshua Bengio, and Patrick Haffner. **Gradient-Based Learning Applied to Document Recognition.** In *Proceedings of the IEEE*, Nov. 1998. URL <http://yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf>

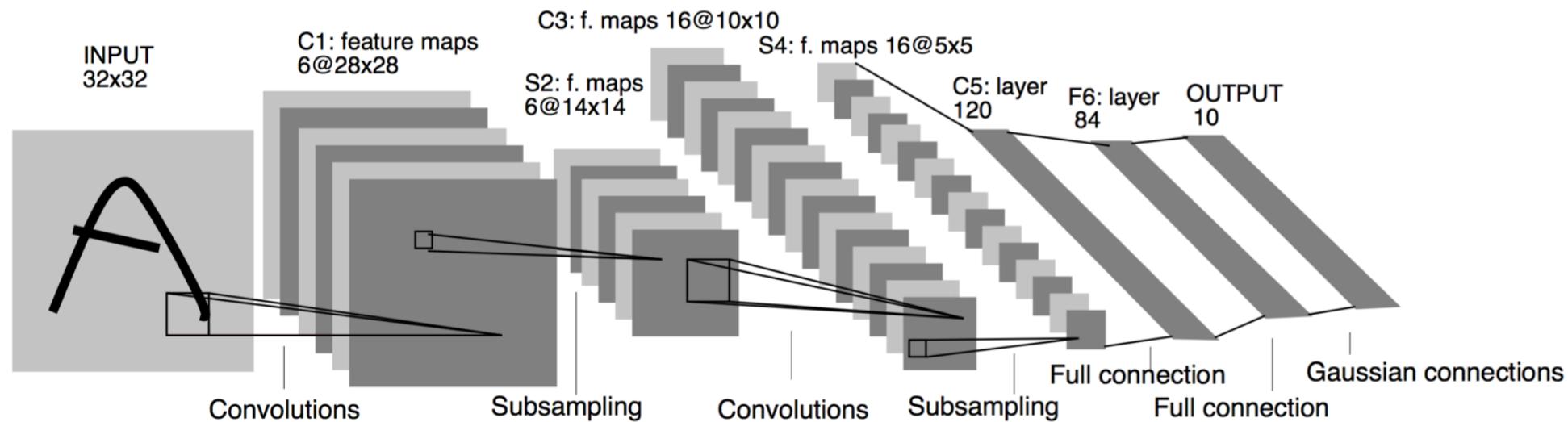


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.