# Cloud Computing & Big Data

PARALLEL & SCALABLE MACHINE LEARNING & DEEP LEARNING

## Prof. Dr. – Ing. Morris Riedel

Associated Professor

School of Engineering and Natural Sciences, University of Iceland, Reykjavik, Iceland

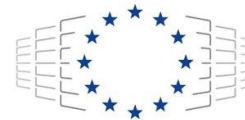Research Group Leader, Juelich Supercomputing Centre, Forschungszentrum Juelich, Germany

@Morris Riedel    @MorrisRiedel    @MorrisRiedel

# Using Deep Learning Techniques in Clouds

October 27, 2020
Online Lecture

# Review of Lecture 7 – Deep Learning Applications in Clouds

■ Limited Perceptron Learning Model



*[30] F. Rosenblatt, 1957*

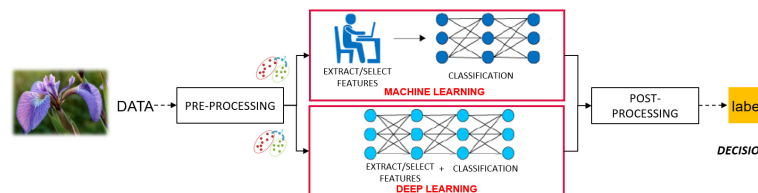*[29] XOR Problem*

| | Input 1 | Input 2 | Output |
|---|---|---|---|
| | 0 | 0 | 0 |
| | 0 | 1 | 1 |
| | 1 | 1 | 0 |
| | 1 | 0 | 1 |

■ Artificial Neural Networks (ANNs)



**(logical combination of two linear classifiers solves XOR problem)**

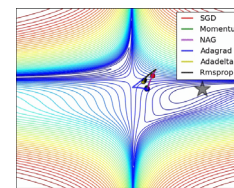■ Feature Engineering vs. Feature Learning



■ Training via Optimization & Backpropagation

**(stochastic gradient & mini-batches)**



*[31] MIT Deep Learning*     *[32] Optimizers*

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence
3. **Pick batch of B data points**
4. Compute gradient $\frac{\partial \mathcal{L}_i(W)}{\partial W} = \frac{1}{B}\sum_{k=1}^{B} \frac{\partial \mathcal{L}_i(W)}{\partial W}$
5. Update weights $W := W - \eta \frac{\partial \mathcal{L}(W)}{\partial W}$
6. Return weights

**(training has here effect!)**

■ Distributed Deep Learning



**(data parallel)**     **(model parallel)**     **(pipelining)**

*[22] Horovod*

*[23] T. Ben-Nun & T. Hoefler*

# Outline of the Course

1. Cloud Computing & Big Data Introduction

2. Machine Learning Models in Clouds

3. Apache Spark for Cloud Applications

4. Virtualization & Data Center Design

5. Map-Reduce Computing Paradigm

6. Deep Learning driven by Big Data

7. Deep Learning Applications in Clouds

8. Infrastructure-As-A-Service (IAAS)

9. Platform-As-A-Service (PAAS)
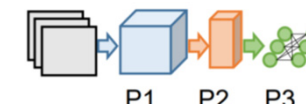
10. Software-As-A-Service (SAAS)

11. Big Data Analytics & Cloud Data Mining

12. Docker & Container Management

13. OpenStack Cloud Operating System

14. Online Social Networking & Graph Databases

15. Big Data Streaming Tools & Applications

16. Epilogue

+ additional practical lectures & Webinars for our hands-on assignments in context

- Practical Topics
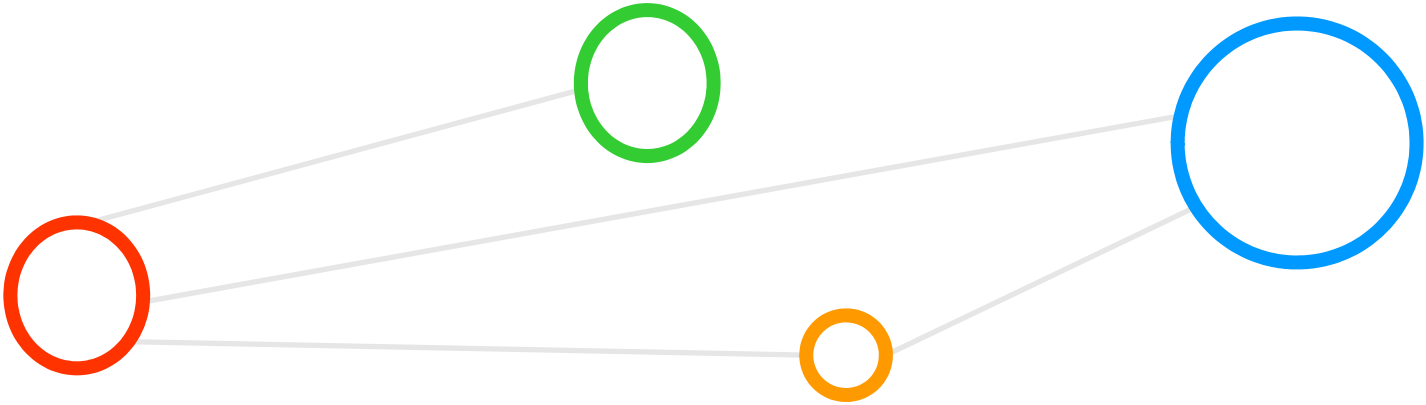
- Theoretical / Conceptual Topics

# Outline

- **Using Artificial Neural Network (ANN) & CPUs in Clouds**
  - Handwritten Character Recognition MNIST Dataset Revisited
  - AWS Elastic Compute Cloud (EC2) & Virtual Server Cloud Instances
  - Using EC2 Amazon Machine Images (AMIs) for Machine Learning
  - Limitations of Free Usage Tiers & Review Challenges of Deploying
  - Observe Growth of Trainable Parameter & Understanding Overfitting

- **Using Convolutional Neural Network (CNN) & GPUs in Clouds**
  - Using EC2 Amazon Machine Images (AMIs) for Deep Learning via CPU
  - Growth of Trainable Parameters & Hyperparameter Complexity
  - Understanding difference between CPUs & GPUs in Training
  - Using Google Colaboratory 'Colab' Cloud Service for Deep Learning
  - Neural Architecture Search and Auto-ML & Resource Requirements

- Promises from previous lecture(s):
- *Practical Lecture 0.1:* Lecture 6 & 7 will provide more insights into deep learning algorithms and networks including the use of TensorFlow and Keras libraries
- *Practical Lecture 0.1:* Lecture 6 & 7 will provide more details on how artificial neural networks (ANNs) and deep learning networks can be used with this data
- *Lecture 2:* Lectures 6 & 7 offer more details on feature selection concepts including working with spatial aspects in image recognition tasks
- *Lecture 3:* Lecture 6 & 7 offer insights of how to use deep learning with cutting-edge GPUs via Google 'colab' notebooks within the Google Cloud

# Using Artificial Neural Network (ANN) & CPUs in Clouds

# Handwritten Character Recognition MNIST Dataset – Preprocessing with Python

- Metadata (cf. Practical Lecture 0.1 )
  - Not very challenging dataset, but good for benchmarks & tutorials

- When working with the dataset
  - Dataset is not in any standard image format like jpg, bmp, or gif (i.e. file format not known to a graphics viewer)
  - Data samples are stored in a simple file format that is designed for storing vectors and multidimensional matrices (i.e. numpy arrays)
  - The pixels of the handwritten digit images are organized row-wise with pixel values ranging from 0 (white background) to 255 (black foreground)
  - Images contain grey levels as a result of an anti-aliasing technique used by the normalization algorithm that generated this dataset
  - Initially input for an Artificial Neural Network (ANN)  *[33] www.big-data.tips, 'MNIST Database'*
  - Afterwards input for a deep learning network   *[36] www.big-data.tips, 'MNIST Dataset'*

---

- Handwritten Character Recognition MNIST dataset is a subset of a larger dataset from US National Institute of Standards (NIST)
- MNIST handwritten digits includes corresponding labels with values 0-9 and is therefore a labeled dataset
- MNIST digits have been size-normalized to 28 * 28 pixels & are centered in a fixed-size image for direct processing
- Two separate files for training & test: 60000 training samples (~47 MB) & 10000 test samples (~7.8 MB)

(10 class classification problem)

# AWS Educate Starter Account – Account Status in Classrooms

- Workbench & Example Classroom
  - Cloud Computing & Big Data – Parallel and Scalable Machine Learning and Deep Learning
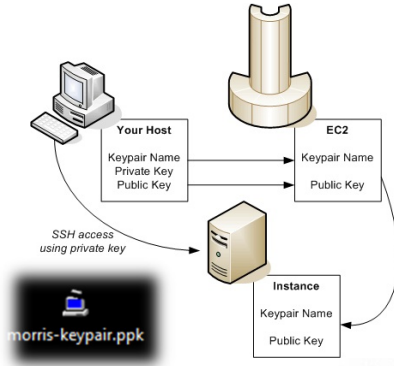


[5] Amazon Web Services

[4] AWS Educate Web page

# AWS Elastic Compute Cloud (EC2) Virtual Servers & Using Key Pairs – Revisited
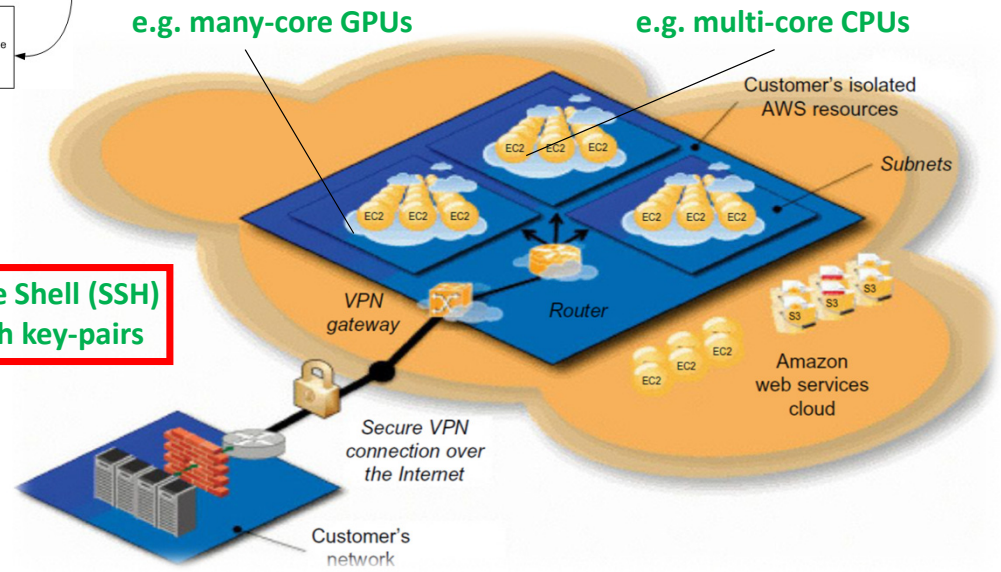
- Secure Shell (SSH)
  - Universal technique to securely access remote clusters & HPC machines

- **The Secure Shell (SSH) is a technique to securely access remote AWS computing instances (e.g., AWS EC2) using a named key pair**
- **An SSH key pair consists of a public key that is known by the Amazon Cloud and a private key that remains only on the laptop of cloud users**
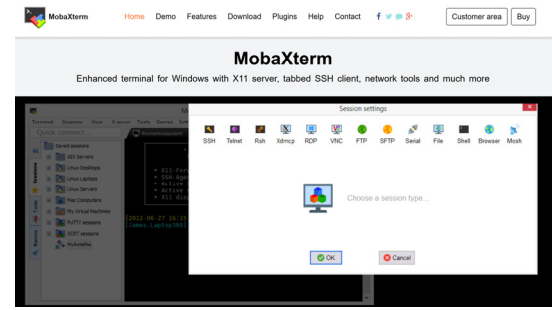
- **Generated AWS key pairs are created per region (e.g., Virginia) in the AWS Cloud**
- **Switching regions means new and/or other SSH keys needs to be used as before**

*[5] Amazon Web Services*

e.g. many-core GPUs          e.g. multi-core CPUs

e.g. Secure Shell (SSH) access with key-pairs
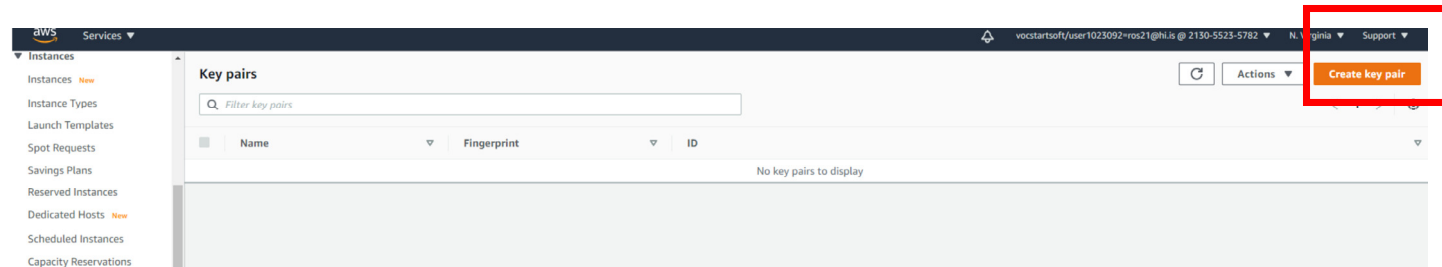
(SSH client is necessary)

*[7] MobaXterm Web page*

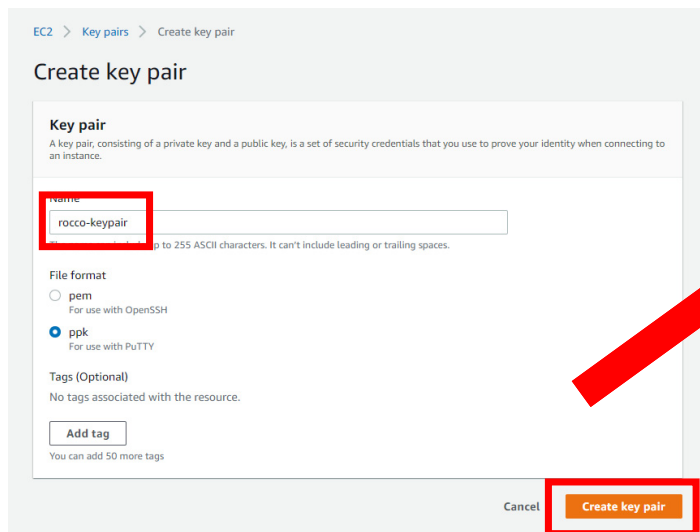*[6] Key Concepts from the AWS Cloud*

# AWS Key Pair – Key Pair Generation (cf. Practical Lecture 5.1)

- Usage
  - Public Key remains in Cloud
  - **Private Key on Laptop**
  - Use SSH Client tool with private key to access remote cloud with matching public key



- After the AWS Key Pair generation, the name of the key is known in many AWS service configuration deployment options such as within the Elastic Compute Cloud (EC2) service or Elastic Map-Reduce (EMR) service

Creating an SSH key pair and keeping private key in pem can be more convenient in certain connections with SSH

# AWS Elastic Compute Cloud (EC2) & Launch Virtual Server Cloud Instances

# Popular Deep Learning Frameworks used with Python in Cloud Computing

- ## TensorFlow

  - One of the most popular deep learning frameworks available today

  - Execution on multi-core CPUs or many-core GPUs

  - **Tensorflow is an open source library for deep learning models using a flow graph approach**
  - **Tensorflow nodes model mathematical operations and graph edges between the nodes are so-called tensors (also known as multi-dimensional arrays)**
  - **The Tensorflow tool supports the use of CPUs and GPUs (much more faster than CPUs)**
  - **Tensorflow work with the high-level deep learning tool Keras in order to create models fast**
  - **New versions of Tensorflow have Keras shipped with it as well & many further tools**

- ## Keras

  - Often used in combination with low-level frameworks like Tensorflow

  - **Keras is a high-level deep learning library implemented in Python that works on top of existing other rather low-level deep learning frameworks like Tensorflow, CNTK, or Theano**
  - **Created deep learning models with Keras run seamlessly on CPU and GPU via low-level deep learning frameworks**
  - **The key idea behind the Keras tool is to enable faster experimentation with deep networks**

# AWS Elastic Compute Cloud (EC2) & Amazon Machine Images (AMIs)



- **AWS Amazon Machine Images (AMIs) are templates that contain the software configuration (e.g., operating system, libraries, application server, and applications) required to launch a EC2 virtual server instance for a specific purpose**
- **AWS EC2 AMI offers solutions that enormously simplify the deployment of required machine learning and deep learning stacks that can be complicated to make work together due to the many different software versions and fast moving technologies (e.g., different NVIDIA GPUs)**
- **AWS EC2 AMI The AMIs are independent from the underlying hardware infrastructure (i.e. concrete CPUs) and can be easily migrated (cf. Lecture 4) to other hardware – be aware of different hardware costs here**
- **Amazon offers pre-configured AMIs for deep learning consisting of preinstalled deep learning packages such as MXNet, TensorFlow, PyTorch, Keras, etc.**
- **Pre-configured AMIs for deep learning feature preinstalled GPU NVIDIA CUDA, cuDNN, and NCCL libraries that usually requires a lot of efforts in installation and version checks with deep learning packages**

# Choose EC2 Instance for AMI & Review Costs Using Free Tier Eligible CPUs



*[8] AWS EC2 Pricing*

# Review & Launching EC2 Instance with AMI – Problems with Free Usage Tier?!



e.g. Secure Shell (SSH) access with key-pairs

# Using SSH Client to Connect to AWS EC2 AMI Instance & Jupyter Notebooks



- Jupyter Notebook
  - Use a browser with http://localhost:8888
  - Use forwarding ssh connection (-L) to connect to localhost although we actually connect to the Amazon AMI instance
  - E.g., ssh -L 127.0.0.1:8888:127.0.0.1:8888 -i ~/Desktop/morris-key-pair-3.ppk ec2-user@ec2-54-92-173-254.compute-1.amazonaws.com
  - It appears we work local, but indeed we work remotely in the cloud

[9] Jupyter

# Using Jupyter with a Kernel & Machine & Deep Learning Software Configuration



*[9] Jupyter*

# Lessons Learned – Dead Environments in the Cloud & Reboot Cloud Instance



- **Despite the fact that Clouds are often stable and production ready they show sometimes still faults and errors that are partly also related to the way of using them, e.g., not terminating properly the Jupyter environment**

- **Because Clouds run remotely on computing systems they can still continue to run even if the local Laptop has no SSH connection open nor is there an active browser window, i.e. remember that this still costs money even if you do not actively use the cloud resources**

# MNIST Dataset – Training/Testing Datasets & One Character Encoding

<div style="border: 2px solid black; background: yellow;">

- **Different phases in machine learning**
- **Training phases is a hypothesis search**
- **Testing phase checks if we are on the right track once the hypothesis is clear**
- **Validation phane for model selection (set fixed parameters and set model types)**

</div>

- Work on two disjoint datasets
  - One for training only (i.e. training set)
  - One for testing only (i.e. test set)
  - Exact seperation is rule of thumb per use case (e.g. 10 % training, 90% test)
  - Practice: If you get a dataset take immediately test data away ('throw it into the corner and forget about it during modelling')
  - Once we learned from training data it has an 'optimistic bias'
  - Usually start by exploring the dataset and its format & labels



'training set'          'test set'

**Training Examples**

$$(\mathbf{x}_1, y_1), ..., (\mathbf{x}_N, y_N)$$

(historical records, groundtruth data, examples)

# MNIST Dataset – Data Exploration Script Training Data – Revisited

```python
import numpy as np
from keras.datasets import mnist

# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()
```

```python
# function to explore one hand-written character
def character_show(character):
    for y in character:
        row = ""
        for x in y:
            row += '{0: <4}'.format(x)
        print(row)
```

```python
# view first 10 hand-written characters
for i in range (0,9):
    character_show(X_train[i])
    print("\n")
    print("Label:")
    print(y_train[i])
    print("\n")
```

- **Loading MNIST training datasets (X) with labels (Y) stored in a binary numpy format**
- **Format is 28 x 28 pixel values with grey level from 0 (white background) to 255 (black foreground)**

- **Small helper function that prints row-wise one 'hand-written' character with the grey levels stored in training dataset**
- **Should reveal the nature of the number (aka label)**

- **Example: loop of the training dataset (e.g. first 10 characters as shown here)**
- **At each loop interval the 'hand-written' character (X) is printed in 'matrix notation' & label (Y)**



*[1] Jupyter @ JSC*

# Data Inspection using Keras Dataset MNIST with Visualization in Jupyter

# MNIST Dataset with Perceptron Learning Model – Need for Reshape

- Two dimensional dataset (28 x 28)
  - Does not fit well with input to Perceptron Model
  - Need to prepare the data even more
  - Reshape data → we need one long vector

- Note that the reshape from two dimensional MNIST data to one long vector means that we loose the surrounding context
- Loosing the surrounding context is one factor why later in this lecture deep learning networks achieving essentially better performance by, e.g., keeping the surrounding context



$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

# MNIST Dataset – Reshape & Normalization – Example

**(one long input vector with length 784)**

**(two dimensional original input)**

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

**(numbers are between 0 and 1)**

```
784 input pixel values per train samples
784 input pixel values per test samples
[0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.01176471 0.07058824 0.07058824 0.07058824
 0.49411765 0.53333336 0.6862745  0.10196079 0.6509804  1.
 0.96862745 0.49803922 0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.11764706 0.14117648 0.36862746 0.6039216
 0.6666667  0.99215686 0.99215686 0.99215686 0.99215686 0.99215686
 0.88235295 0.6745098  0.99215686 0.9490196  0.7647059  0.2509804
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.         0.         0.19215687
 0.93333334 0.99215686 0.99215686 0.99215686 0.99215686 0.99215686
 0.99215686 0.99215686 0.99215686 0.9843137  0.3647059  0.32156864
 0.32156864 0.21960784 0.15294118 0.         0.         0.
 0.         0.         0.         0.         0.         0.
 0.         0.         0.         0.07058824 0.85882354 0.99215686
 0.99215686 0.99215686 0.99215686 0.99215686 0.7764706  0.7137255
```

```
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 3 18 18 18 126 136 175 26 166 255 247 127 0 0 0 0 0
0 0 0 0 0 0 0 30 36 94 154 170 253 253 253 253 253 225 172 253 242 195 64 0 0 0 0 0
0 0 0 0 0 0 49 238 253 253 253 253 253 253 253 253 251 93 82 82 56 39 0 0 0 0 0 0
0 0 0 0 0 0 18 219 253 253 253 253 253 198 182 247 241 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 80 156 107 253 253 205 11 0 43 154 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 14 1 154 253 90 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 139 253 190 2 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 11 190 253 70 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 35 241 225 160 108 1 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 81 240 253 253 119 25 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 45 186 253 253 150 27 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 16 93 252 253 187 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 249 253 249 64 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 46 130 183 253 253 207 2 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 39 148 229 253 253 253 250 182 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 24 114 221 253 253 253 253 201 78 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 23 66 213 253 253 253 253 198 81 2 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 18 171 219 253 253 253 253 195 80 9 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 55 172 226 253 253 253 253 244 133 11 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 136 253 253 253 212 135 132 16 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

Label:
5
```

# MNIST Dataset & Multi Output Perceptron Model

- 10 Class Classification Problem
  - Use 10 Perceptrons for 10 outputs with softmax activation function (enables probabilities for 10 classes)

**(Dense Layer)**   **(Softmax Layer)**   **(output probabilities)**

$x_1$   $\Sigma$   $\int g$   $\hat{y}_1$

$x_{...}$   $\Sigma$   $\hat{y}_{...}$

$x_m$   $\Sigma$   $\hat{y}_{10}$

**(input m = 784)**   **(10 neurons sum with 10 bias)**   **(softmax activation)**   **(NB_CLASSES = 10)**

```python
from keras.models import Sequential
from keras.layers.core import Dense, Activation

# model Keras sequential
model = Sequential()

# add fully connected layer – input with output
model.add(Dense(NB_CLASSES, input_shape=(RESHAPED,)))

# add activation function layer to get class probabilities
model.add(Activation('softmax'))

# printout a summary of the model to understand model complexity
model.summary()
```

**K**

**(parameters = 784 * 10 + 10 bias = 7850)**

```
Layer (type)                 Output Shape              Param #
=================================================================
dense_1 (Dense)              (None, 10)                7850
_____
activation_1 (Activation)    (None, 10)                0
=================================================================
Total params: 7,850
Trainable params: 7,850
Non-trainable params: 0
```

- Note that the output units are independent among each other in contrast to neural networks with one hidden layer
- The output of softmax gives class probabilities
- The non-linear Activation function 'softmax' represents a generalization of the sigmoid function – it squashes an n-dimensional vector of arbitrary real values into a n-dimensional vector of real values in the range of 0 and 1 – here it aggregates 10 answers provided by the Dense layer with 10 neurons

# MNIST Dataset & Compile Multi Output Perceptron Model

- ## Compile the model
  - Optimizer as algorithm used to update weights while training the model
  - Specify loss function (i.e. objective function) that is used by the optimizer to navigate the space of weights
  - (note: process of optimization is also called loss minimization, cf. Invited lecture Gabriele Cavallaro)
  - Indicate metric for model evaluation (e.g., accuracy)

- ## Specify loss function
  - Compare prediction vs. given class label
  - E.g. categorical crossentropy



```
from keras.optimizers import SGD

OPTIMIZER = SGD()  # optimization technique
```

- Compile the model to be executed by the Keras backend (e.g. TensorFlow)
- Optimizer Gradient Descent (GD) uses all the training samples available for a step within a iteration
- Optimizer Stochastic Gradient Descent (SGD) converges faster: only one training samples used per iteration
- Loss function is a multi-class logarithmic loss: target is $t_{i,j}$ and prediction is $p_{i,j}$
- Categorical crossentropy is suitable for multiclass label predictions (default with softmax)

$$L_i = -\sum_j t_{i,j} \log(p_{i,j})$$

*[10] Big Data Tips, Gradient Descent*

```
# specify loss, optimizer and metric
model.compile(loss='categorical_crossentropy', optimizer=OPTIMIZER, metrics=['accuracy'])
```

# Full Script: MNIST Dataset – Model Parameters & Data Normalization

```python
import numpy as np
from keras.datasets import mnist
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.optimizers import SGD
from keras.utils import np_utils

# parameter setup
NB_EPOCH = 20
BATCH_SIZE = 128
NB_CLASSES = 10 # number of outputs = number of digits
OPTIMIZER = SGD() # optimization technique
VERBOSE = 1

# download and shuffled as training and testing set
(X_train, y_train), (X_test, y_test) = mnist.load_data()

# X_train is 60000 rows of 28x28 values --> reshaped in 60000 x 784
RESHAPED = 784
X_train = X_train.reshape(60000, RESHAPED)
X_test = X_test.reshape(10000, RESHAPED)
X_train = X_train.astype('float32')
X_test = X_test.astype('float32')

# normalize
X_train /= 255
X_test /= 255

# output number of samples
print(X_train.shape[0], 'train samples')
print(X_test.shape[0], 'test samples')
```

- **NB_CLASSES: 10 Class Problem**
- **NB_EPOCH: number of times the model is exposed to the overall training set – at each iteration the optimizer adjusts the weights so that the objective function is minimized – <u>increasing leads to better accuracy, but also to overfitting (cf. Lecture 7)</u>**
- **BATCH_SIZE: number of training instances taken into account before the optimizer performs a weight update to the model**
- **OPTIMIZER: Stochastic Gradient Descent ('SGD') – only one training sample/iteration**

- **Data load shuffled between training and testing set in files**
- **Data preparation, e.g. X_train is 60000 samples / rows of 28 x 28 pixel values that are reshaped in 60000 x 784 including type specification (i.e. float32)**
- **Data normalization: divide by 255 – the max intensity value to obtain values in range [0,1]**

*'training set'*    *'test set'*

**Training Examples**
$$(\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_N, y_N)$$

(historical records, groundtruth data, examples)

➤ **Assignment #2 will explore the change of parameters in context of changes in running time when training models on GPUs vs. CPUs**

# Full Script: MNIST Dataset – Fitting a Multi Output Perceptron Model

**(full script continued from previous slide)**

```python
# convert class label vectors using one hot encoding
Y_train = np_utils.to_categorical(y_train, NB_CLASSES)
Y_test = np_utils.to_categorical(y_test, NB_CLASSES)

# model Keras sequential
model = Sequential()

# add fully connected layer - input with output
model.add(Dense(NB_CLASSES, input_shape=(RESHAPED,)))

# add activation function layer to get class probabilities
model.add(Activation('softmax'))

# printout a summary of the model to understand model complexity
model.summary()

# specify loss, optimizer and metric
model.compile(loss='categorical_crossentropy', optimizer=OPTIMIZER, metrics=['accuracy'])

# model training
history = model.fit(X_train, Y_train, batch_size=BATCH_SIZE, epochs=NB_EPOCH, verbose=VERBOSE)

# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

- The Sequential() Keras model is a linear pipeline (aka 'a stack') of various neural network layers including Activation functions of different types (e.g. softmax)
- Dense() represents a fully connected layer used in ANNs that means that each neuron in a layer is connected to all neurons located in the previous layer

- The non-linear activation function 'softmax' is a generalization of the sigmoid function – it squashes an n-dimensional vector of arbitrary real values into a n-dimenensional vector of real values in the range of 0 and 1 – here it aggregates 10 answers provided by the Dense layer with 10 neurons

- Loss function is a multi-class logarithmic loss: target is $t_{i,j}$ and the prediction is $p_{i,j}$

$$L_i = -\Sigma_j t_{i,j} \log(p_{i,j})$$

- Train the model ('fit') using selected batch & epoch sizes on training & test data

➤ **Assignment #2 will explore the change of parameters in context of changes in running time when training models on GPUs vs. CPUs**

# Running a Simple ANN with no hidden layers – Multi-Output-Perceptron



Note that the outcome of the training process is the result of optimization techniques like SGD that tend to vary 'a bit'

Note that the outcome of the training process can be dependent on the length of training increasing accuracy to a certain point when overfitting starts

Overfitting can be controlled with validation and regularization techniques that belong to advanced machine learning methods to be studied in full university machine learning course in detail

**cf. Lecture 6 and 7**

# MNIST Dataset – A Multi Output Perceptron Model – Output & Evaluation

```
Epoch 7/20
60000/60000 [==============================] - 2s 26us/step - loss: 0.4419 - acc: 0.8838
Epoch 8/20
60000/60000 [==============================] - 2s 26us/step - loss: 0.4271 - acc: 0.8866
Epoch 9/20
60000/60000 [==============================] - 2s 25us/step - loss: 0.4151 - acc: 0.8888
Epoch 10/20
60000/60000 [==============================] - 2s 26us/step - loss: 0.4052 - acc: 0.8910
Epoch 11/20
60000/60000 [==============================] - 2s 26us/step - loss: 0.3968 - acc: 0.8924
Epoch 12/20
60000/60000 [==============================] - 2s 25us/step - loss: 0.3896 - acc: 0.8944
Epoch 13/20
60000/60000 [==============================] - 2s 26us/step - loss: 0.3832 - acc: 0.8956
Epoch 14/20
60000/60000 [==============================] - 2s 25us/step - loss: 0.3777 - acc: 0.8969
Epoch 15/20
60000/60000 [==============================] - 2s 25us/step - loss: 0.3727 - acc: 0.8982
Epoch 16/20
60000/60000 [==============================] - 1s 24us/step - loss: 0.3682 - acc: 0.8989
Epoch 17/20
60000/60000 [==============================] - 1s 25us/step - loss: 0.3641 - acc: 0.9001
Epoch 18/20
60000/60000 [==============================] - 1s 25us/step - loss: 0.3604 - acc: 0.9007
Epoch 19/20
60000/60000 [==============================] - 2s 25us/step - loss: 0.3570 - acc: 0.9016
Epoch 20/20
60000/60000 [==============================] - 1s 24us/step - loss: 0.3538 - acc: 0.9023
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

```
10000/10000 [==============================] - 0s 41us/step
Test score: 0.33423959468007086
Test accuracy: 0.9101
```



(Dense Layer)  (Softmax Layer)  (output probabilities)

$x_1$   $\sum$   $\hat{y}_1$

$x_{...}$   $\sum$   $\int g$   $\hat{y}_{...}$

$x_m$   $\sum$   $\hat{y}_{10}$

(input m = 784)  (10 neurons sum with 10 bias)  (softmax activation)  (NB_CLASSES = 10)

- **How to improve the model design by extending the neural network topology?**
- **Which layers are required?**
- **Think about input layer need to match the data – what data we had?**
- **Maybe hidden layers?**
- **How many hidden layers?**
- **What activation function for which layer (e.g. maybe ReLU)?**
- **Think Dense layer – Keras?**
- **Think about final Activation as Softmay (cf. Day One) → output probability**

# MNIST Dataset – Add Two Hidden Layers for Artificial Neural Network (ANN)

- **All parameter value remain the same as before**
- **We add N_HIDDEN as parameter in order to set 128 neurons in one hidden layer – this number is a hyperparameter that is not directly defined and needs to be find with parameter search**

[34] big-data.tips, 'Relu Neural Network'

[35] big-data.tips, 'tanh'

Input      Hidden Layers      Output

$x_1$, $x_2$, $x_n$ → ... → $y_1$, $y_c$

```
# parameter setup
NB_EPOCH = 20
BATCH_SIZE = 128
NB_CLASSES = 10 # number of outputs = number of digits
OPTIMIZER = SGD() # optimization technique
VERBOSE = 1
N_HIDDEN = 128 # number of neurons in one hidden layer
```

**(activation functions ReLU & Tanh)**

```
# model Keras sequential
model = Sequential()
```

**K**

```
model.add(Dense(N_HIDDEN))
model.add(Activation('relu'))
```

```
model.add(Dense(N_HIDDEN))
model.add(Activation('tanh'))
```

```
# modeling step
# 2 hidden layers each N_HIDDEN neurons
model.add(Dense(N_HIDDEN, input_shape=(RESHAPED,)))
model.add(Activation('relu'))
model.add(Dense(N_HIDDEN))
model.add(Activation('relu'))
model.add(Dense(NB_CLASSES))
```

```
# add activation function layer to get class probabilities
model.add(Activation('softmax'))
```

- **The non-linear Activation function 'relu' represents a so-called Rectified Linear Unit (ReLU) that only recently became very popular because it generates good experimental results in ANNs and more recent deep learning models – it just returns 0 for negative values and grows linearly for only positive values**
- **A hidden layer in an ANN can be represented by a fully connected Dense layer in Keras by just specifying the number of hidden neurons in the hidden layer**

➤ **Assignment #2 will explore the change of parameters in context of changes in running time when training models on GPUs vs. CPUs**

# Running a Simple ANN with two hidden layers

# MNIST Dataset – ANN Model Parameters & Output Evaluation

```
Epoch 7/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.2743 - acc: 0.9223
Epoch 8/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.2601 - acc: 0.9266
Epoch 9/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.2477 - acc: 0.9301
Epoch 10/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.2365 - acc: 0.9329
Epoch 11/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.2264 - acc: 0.9356
Epoch 12/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.2175 - acc: 0.9386
Epoch 13/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.2092 - acc: 0.9412
Epoch 14/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.2013 - acc: 0.9432
Epoch 15/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.1942 - acc: 0.9454
Epoch 16/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.1876 - acc: 0.9472
Epoch 17/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.1813 - acc: 0.9487
Epoch 18/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.1754 - acc: 0.9502
Epoch 19/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.1700 - acc: 0.9522
Epoch 20/20
60000/60000 [==============================] - 1s 18us/step - loss: 0.1647 - acc: 0.9536
```

```
# model evaluation
score = model.evaluate(X_test, Y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])

10000/10000 [==============================] - 0s 33us/step
Test score: 0.16286438911408185
Test accuracy: 0.9514
```

✓ **Multi Output Perceptron: ~91,01% (20 Epochs)**

✓ **ANN 2 Hidden Layers: ~95,14 % (20 Epochs)**

```
# printout a summary of the model to understand model complexity
model.summary()
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| dense_1 (Dense) | (None, 128) | 100480 |
| activation_1 (Activation) | (None, 128) | 0 |
| dense_2 (Dense) | (None, 128) | 16512 |
| activation_2 (Activation) | (None, 128) | 0 |
| dense_3 (Dense) | (None, 10) | 1290 |
| activation_3 (Activation) | (None, 10) | 0 |

```
Total params: 118,282
Trainable params: 118,282
Non-trainable params: 0
```

Input    Hidden Layers    Output

$x_1$

$x_2$

$x_n$

$y_1$

$y_c$

- **Dense Layer connects every neuron in this dense layer to the next dense layer with each of its neuron also called a fully connected network element with weights as trainiable parameters**
- **Choosing a model with different layers is a model selection that directly also influences the number of parameters (e.g. add Dense layer from Keras means new weights)**
- **Adding a layer with these new weights means much more computational complexity since each of the weights must be trained in each epoch (depending on #neurons in layer)**

# Using Convolutional Neural Network (CNN) & GPUs in Clouds

# Innovative Deep Learning Techniques – Revisited (cf. Lecture 6 & 7)



[11] M. Riedel, 'Deep Learning - Using a Convolutional Neural Network', Invited YouTube Lecture, six lectures, University of Ghent, 2017



[12] M. Riedel et al., 'Introduction to Deep Learning Models', JSC Tutorial, three days, JSC, 2019



[13] H. Lee et al., 'Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical

- ▪ **Innovation via specific layers and architecture types**



[14] Neural Network 3D Simulation

[15] A. Rosebrock

# Complex Relationships: ML & DL vs. HPC/Clouds & Big Data (cf. Lecture 0)



Model Performance / Accuracy

'small datasets'

*manual feature engineering' changes the ordering*

Large Deep Learning Networks

Medium Deep Learning Networks

Small Neural Networks

Traditional Learning Models

*Computing*

*Training Time*

**High Performance Computing & Cloud Computing**

*Random Forests*

*SVMs*

*Statistical Computing with R*

*MatLab*

*scikit-learn*   *Weka*   *Octave*

Dataset Volume

→ 'Big Data'      [16] www.big-data.tips

# Understanding Feature Maps & Convolutions – Online Web Tool



*[17] Harley, A.W., An Interactive Node-Link Visualization of Convolutional Neural Networks*

# MNIST Dataset – Convolutional Neural Network (CNN) Model

```python
from keras import backend as K
from keras.models import Sequential
from keras.layers.convolutional import Conv2D
from keras.layers.convolutional import MaxPooling2D
from keras.layers.core import Activation
from keras.layers.core import Flatten
from keras.layers.core import Dense
from keras.datasets import mnist
from keras.utils import np_utils
from keras.optimizers import SGD, RMSprop, Adam
import numpy as np
import matplotlib.pyplot as plt
```

- Increasing the number of filters learned to 50 in the next layer from 20 in the first layer
- Increasing the number of filters in deeper layers is a common technique in deep learning architecture modeling
- Flattening the output as input for a Dense layer (fully connected layer)
- Fully connected / Dense layer responsible with softmay activation for classification based on learned filters and features

*[18] A. Gulli et al.*

```python
#define the CNN model
class CNN:
  @staticmethod
  def build(input_shape, classes):
    model = Sequential()
    # CONV => RELU => POOL
    model.add(Conv2D(20, kernel_size=5, padding="same",
    input_shape=input_shape))
    model.add(Activation("relu"))
    model.add(MaxPooling2D(pool_size=(2, 2), strides=(2, 2)))
    # CONV => RELU => POOL
    model.add(Conv2D(50, kernel_size=5, border_mode="same"))
    model.add(Activation("relu"))
    model.add(MaxPooling2D(pool_size=(2, 2), strides=(2, 2)))
    # Flatten => RELU layers
    model.add(Flatten())
    model.add(Dense(500))
    model.add(Activation("relu"))
    # a softmax classifier
    model.add(Dense(classes))
    model.add(Activation("softmax"))
    return model
```



```python
# initialize the optimizer and model
model = CNN.build(input_shape=INPUT_SHAPE, classes=NB_CLASSES)
model.compile(loss="categorical_crossentropy", optimizer=OPTIMIZER,
metrics=["accuracy"])
```

```
# printout a summary of the model to understand model complexity
model.summary()
_____
Layer (type)                 Output Shape              Param #
================================================================
conv2d_1 (Conv2D)            (None, 20, 28, 28)        520
_____
activation_1 (Activation)    (None, 20, 28, 28)        0
_____
max_pooling2d_1 (MaxPooling2 (None, 20, 14, 14)        0
_____
conv2d_2 (Conv2D)            (None, 50, 14, 14)        25050
_____
activation_2 (Activation)    (None, 50, 14, 14)        0
_____
max_pooling2d_2 (MaxPooling2 (None, 50, 7, 7)          0
_____
flatten_1 (Flatten)          (None, 2450)              0
_____
dense_1 (Dense)              (None, 500)               1225500
_____
activation_3 (Activation)    (None, 500)               0
_____
dense_2 (Dense)              (None, 10)                5010
_____
activation_4 (Activation)    (None, 10)                0
================================================================
Total params: 1,256,080
Trainable params: 1,256,080
Non-trainable params: 0
_____
```

# MNIST Dataset – Model Parameters & 2D Input Data

```python
# parameter setup
NB_EPOCH = 20
BATCH_SIZE = 128
VERBOSE = 1
OPTIMIZER = Adam()
VALIDATION_SPLIT=0.2
IMG_ROWS, IMG_COLS = 28, 28 # input image dimensions
NB_CLASSES = 10 # number of outputs = number of digits
INPUT_SHAPE = (1, IMG_ROWS, IMG_COLS)
```

```python
# data: shuffled and split between train and test sets
(X_train, y_train), (X_test, y_test) = mnist.load_data()
K.set_image_dim_ordering("th")
# consider them as float and normalize
X_train = X_train.astype('float32')
X_test = X_test.astype('float32')
X_train /= 255
X_test /= 255
# we need a 60K x [1 x 28 x 28] shape as input to the CONVNET
X_train = X_train[:, np.newaxis, :, :]
X_test = X_test[:, np.newaxis, :, :]
print(X_train.shape[0], 'train samples')
print(X_test.shape[0], 'test samples')
# convert class vectors to binary class matrices
y_train = np_utils.to_categorical(y_train, NB_CLASSES)
y_test = np_utils.to_categorical(y_test, NB_CLASSES)
```

- **OPTIMIZER: Adam - advanced optimization technique that includes the concept of a momentum (a certain velocity component) in addition to the acceleration component of Stochastic Gradient Descent (SGD)**
- **Adam computes individual adaptive learning rates for different parameters from estimates of first and second moments of the gradients**
- **Adam enables faster convergence at the cost of more computation and is currently recommended as the default algorithm to use (or SGD + Nesterov Momentum)**

*[19] D. Kingma et al., 'Adam: A Method for Stochastic Optimization'*

- **Compared to the Multi-Output Perceptron and Artificial Neural Networks (ANN) model, the input dataset remains as 2d matricew with 1 x 28 x 28 per image, including also the class vectors that are converted to binary class matrices**

➢ **Assignment #2 will explore the change of parameters in context of changes in running time when training models on GPUs vs. CPUs**

# MNIST Dataset – CNN Model Output & Evaluation

```
Epoch 14/20
48000/48000 [==============================] - 4s 88us/step - loss: 0.0065 - acc: 0.9980 - val_loss: 0.0346 - val_acc: 0.9921
Epoch 15/20
48000/48000 [==============================] - 4s 89us/step - loss: 0.0030 - acc: 0.9990 - val_loss: 0.0418 - val_acc: 0.9903
Epoch 16/20
48000/48000 [==============================] - 4s 88us/step - loss: 0.0057 - acc: 0.9980 - val_loss: 0.0470 - val_acc: 0.9910
Epoch 17/20
48000/48000 [==============================] - 4s 88us/step - loss: 0.0043 - acc: 0.9985 - val_loss: 0.0440 - val_acc: 0.9906
Epoch 18/20
48000/48000 [==============================] - 4s 88us/step - loss: 0.0046 - acc: 0.9985 - val_loss: 0.0474 - val_acc: 0.9891
Epoch 19/20
48000/48000 [==============================] - 4s 88us/step - loss: 0.0047 - acc: 0.9986 - val_loss: 0.0353 - val_acc: 0.9928
Epoch 20/20
48000/48000 [==============================] - 4s 88us/step - loss: 3.4055e-04 - acc: 1.0000 - val_loss: 0.0374 - val_acc: 0.9927
```

```python
# model evaluation
score = model.evaluate(X_test, y_test, verbose=VERBOSE)
print("Test score:", score[0])
print('Test accuracy:', score[1])
```

```
10000/10000 [==============================] - 1s 70us/step
Test score: 0.0303058747581508
Test accuracy: 0.9936
```

- ✓ **Multi Output Perceptron: ~91,01% (20 Epochs)**
- ✓ **ANN 2 Hidden Layers: ~95,14 % (20 Epochs)**
- ✓ **CNN Deep Learning Model: ~99,36 % (20 Epochs)**

*[18] A. Gulli et al.*



Why not 100% **?**

some samples even for a human unrecognizable

# Running a Deep Learning Model with Convolutional Neural Network (CNN)



- **Using Deep Learning Techniques such as Convolutional Neural Networks (CNNs) in clouds can lead to significant improvements in accuracy, but also to significant longer run-times than traditional Artificial Neural Networks (ANNs) and are thus much more costly in clouds**

- **Using CPU resources for deep learning techniques is usually not recommended**

# More Computation: Deep Learning via RESNET-50 Architecture (cf. Lecture 6 & 7)

- Application Example: Classification of land cover in scenes on remote sensing datasets
  - Very suitable for parallelization via distributed training on multi GPUs

*[20] RESNET*



*[22] Horovod*

MPI_Allreduce()

- **RESNET-50 is a known neural network architecture that has established a strong baseline in terms of accuracy**
- **The computational complexity of training the RESNET-50 architecture relies in the fact that is has ~ 25.6 millions of trainable parameters**
- **RESNET-50 still represents a good trade-off between accuracy, depth and number of parameters**
- **The setups of RESNET-50 makes it very suitable for parallelization via distributed training on multi GPUs**

Time per epoch [sec]

A partition of the JUWELS system has 56 compute nodes, each with 4 NVIDIA V100 GPUs (equipped with 16 GB of memory)

24 nodes x 4 GPUs = 96 GPUs

*[21] R. Sedona & M. Riedel et al., 2019*

# Cloud Computing & HPC using GPUs for Deep Learning – Hardware Complexity



**Facts:** GPUs are mostly used today for deep learning compared to CPUs, FPGA, and specialized hardware

**Facts:** ~55% of all users that use deep learning use it with multiple nodes instead of just a single node

**Facts:** The communication layer MPI is mostly used as communication layer for distributed training compared to Apache Spark, Remote Procedure Calls, Apache Hadoop MapReduce, or traditional Sockets

Most users use deep learning today with minibatches that are selected numbers of samples for performing the optimization (e.g. SGD on minibatches)

Minibatches should be not too small to increate performance, but also not too large to increase validation error

*[23] T. Ben-Nun & T. Hoefler*

> **Complementary High Performance Computing course offers insights into parallel programming models such as MPI & hardware impact**

# AWS Amazon Sagemaker – SAAS Service to Abstract from Hardware Complexity

- **AWS Cloud – Amazon Sagemaker**
  - **Fully managed service** that enables quick & easy machine & deep learning applications
  - **Avoids time-consuming manual installation** of many required software frameworks
  - Builds on-top of various IAAS & PAAS services



**Build**
Connect to other AWS services and transform data in SageMaker notebooks

**Train**
Use SageMaker's algorithms and frameworks, or bring your own, for distributed training

**Tune**
SageMaker automatically tunes your model by adjusting multiple combinations of algorithm parameters

**Deploy**
Once training is completed, models can be deployed to SageMaker endpoints, for real-time predictions

**(SAAS solutions often abstracts away completely underlying resources)**

MACHINE LEARNING

Free Tier          FREE TRIAL

Amazon SageMaker
**250 Hours**
per month of t2.medium notebook usage for the first two months

Fully managed platform to build, train, and deploy machine learning models.

250 hours per month of t2.medium notebook usage for the first two months

50 hours per month of m4.xlarge for training for the first two months

125 hours per month of m4.xlarge for hosting for the first two months



- **AWS Amazon Sagemaker is a SAAS oriented service that provides fully managed instances running Jupyter notebooks that include examples training & tuning various machine and deep learning models**
- **Offers Amazon SageMaker Studioas a fully integrated development environment (IDE) for machine learning in the AWS cloud**
- **SAAS services are usually not free and often require a subscription**

*[9] Jupyter Web page*

*[24] AWS – Amazon Sagemaker*

> **Lecture 10 provides more details about AWS Cloud services and its Software-as-a-Service (SAAS) models & other SAAS cloud services**

# Using Google Colaboratory Cloud Infrastructure for Deep Learning with GPUs

- **Google Colaboratory** (free & pro version for 9.99 $ / month)
  - **'Colab' notebooks are Jupyter notebooks** that run in the Google cloud
  - Possible to run Apache Spark via PySpark Jupyter notebooks in Colab (cf. Lecture 3)
  - Possible to train Deep Learning networks via GPUs & Jupyter notebooks in Colab
  - Highly integrated with other Google services (e.g., Google Drive for data)
  - Access to vendor-specific Tensor Processing Units (TPUs)

*[27] Google Colaboratory*

**(for international students: watch out – it uses the browser language automatically)**

- **Google Colaboratory offers 'Colab' notebooks that are implemented with Jupyter notebooks that in turn run in the Google cloud and are highly integrated with other Google cloud services such as Google Drive thus making 'Colab' notebooks easy to set up, access, and share with others**

```
[ ]  # initialize the optimizer and model
     model = CNN.build(input_shape=INPUT_SHAPE, classes=NB_CLASSES)
     model.compile(loss="categorical_crossentropy", optimizer=OPTIMIZER,
     metrics=["accuracy"])

     ---------------------------------------------------------------------
     TypeError                         Traceback (most recent call last)
     <ipython-input-11-13ca928a46c2> in <module>()
           1 # initialize the optimizer and model
     ----> 2 model = CNN.build(input_shape=INPUT_SHAPE, classes=NB_CLASSES)
           3 model.compile(loss="categorical_crossentropy", optimizer=OPTIMIZER,
           4 metrics=["accuracy"])

                          ↕ 5 frames
     /usr/local/lib/python3.6/dist-packages/tensorflow/python/keras/utils/generic_utils.py in validate_kwargs(kwargs, allowed_kwargs, error_message)
         776    for kwarg in kwargs:
         777       if kwarg not in allowed_kwargs:
     --> 778          raise TypeError(error_message, kwarg)
         779
         780

     TypeError: ('Keyword argument not understood:', 'border_mode')

     SEARCH STACK OVERFLOW
```

**(Keras API update now available in Google 'Colab' creates an issue to port our CNN model directly from our Amazon EC2 AMI example because it is based on previous versions of Keras)**

**(tutorials & codes need updates)**

- Update Oct/2016: Updated for Keras 1.1.0, TensorFlow 0.10.0 and scikit-learn v0.18.
- Update Mar/2017: Updated for Keras 2.0.2, TensorFlow 1.0.1 and Theano 0.9.0.
- Update Sep/2019: Updated for Keras 2.2.5 API.

- **The portability of deep learning codes is hindered by the frequent updates of the different APIs of deep learning frameworks like Keras, Tensorflow, etc. (cf. different AWS EC2 AMI versions)**

**(Clouds also face this update problem)**

*[28] Machine Learning Mastery MNIST Tutorial*

**Notebook settings**

Hardware accelerator
GPU
None
GPU
TPU

☐ Omit code cell output when saving this notebook

CANCEL     SAVE

# Massive Requirement for Cloud Resources: Neural Architecture Search (NAS)

- **Often a Recurrent Neural Network (RNN) technique that performs the agent steps**



Search Space $\mathcal{A}$ → Search Strategy → Architecture $A \in \mathcal{A}$ → Performance Estimation Strategy

performance estimate of A

**Child Architectures**

**Controller**

**Task-Dependent Objectives** $\mathbb{O}_T$

**Architecture-Dependent Objectives** $\mathbb{O}_A$

Reward

*[25] A.C. Cheng et al., 'InstaNAS: Instance-aware Neural Architecture Search', 2018*

- **Employed neural networks architectures are often developed manually by human experts that is time-consuming and error-prone**
- **Deep learning success has been accompanied by a rising demand for architecture engineering, where increasingly more complex neural architectures are designed manually**
- **Neural Architecture Search (NAS) methods can be categorized in (a) search space, (b) search strategy, and (c) performance estimation strategy**
- **Automated Neural Architecture (NAS) search methods aim to solve this problem as a process of automating Architecture engineering**

*[26] M. Riedel, 'NAS with Reinforcement Learning'*

- **Derived specific architectures that perform good for specific dataset samples**

- **E.g. what is the accuracy or error rate we obtain as metric to guide the search for specific architectures for specific dataset samples**

- **E.g. what is the latency of the network for a given dataset sample to guide the search for specific architectures that offer better latency by keeping accuracy(!)**

# Lecture Bibliography

# Lecture Bibliography (1)

- [1] Jupyter @ Juelich Supercomputing Centre, Online:
  https://jupyter-jsc.fz-juelich.de/
- [2] Tensorflow Deep Learning Framework, Online:
  https://www.tensorflow.org/
- [3] Keras Python Deep Learning Library, Online:
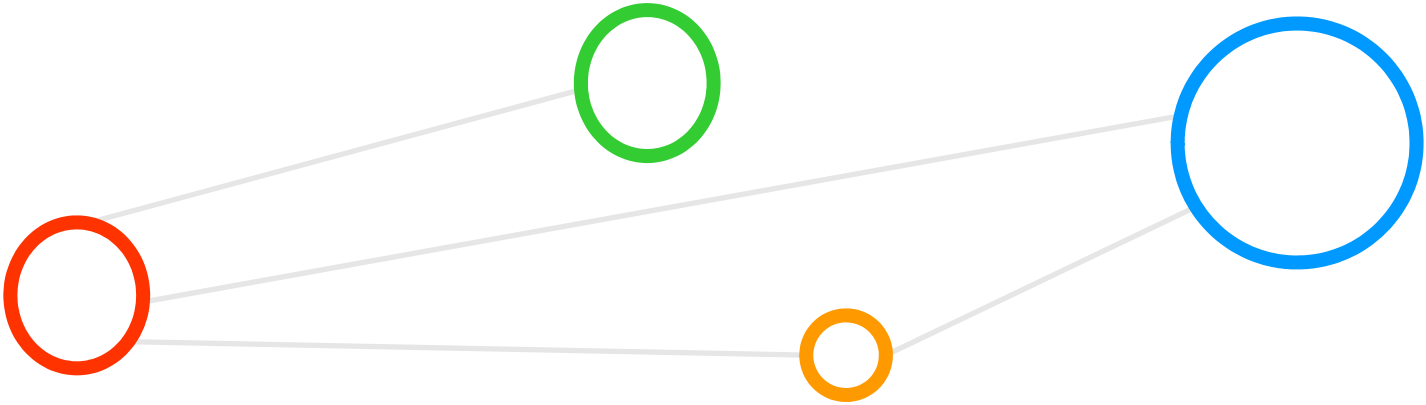  https://keras.io/
- [4] Amazon Web Services Educate Web Page, Online:
  https://aws.amazon.com/education/awseducate/
- [5] Amazon Web Services Web Page, Online:
  https://aws.amazon.com
- [6] Key Concepts from the AWS Cloud, Online:
  https://blogs.sap.com/2010/01/28/key-concepts-from-the-aws-cloud/
- [7] SSH Client MobaXterm, Online:
  https://mobaxterm.mobatek.net/
- [8] Amazon Web Services EC2 On-Demand Pricing models, Online:
  https://aws.amazon.com/ec2/pricing/on-demand/
- [9] Jupyter Web page, Online:
  http://jupyter.org/
- [10] Big Data Tips, 'Gradient Descent', Online:
  http://www.big-data.tips/gradient-descent
- [11] Morris Riedel, 'Deep Learning - Using a Convolutional Neural Network', Invited YouTube Lecture, University of Ghent, 2017, Online:
  https://www.youtube.com/watch?v=gOL1_YIosYk&list=PLrmNhuZo9sgZUdaZ-f6OHK2yFW1kTS2qF
- [12] M. Riedel et al., 'Introduction to Deep Learning Models', JSC Tutorial, three days, JSC, 2019, Online:
  http://www.morrisriedel.de/introduction-to-deep-learning-models
- [13] H. Lee et al., 'Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations', Online:
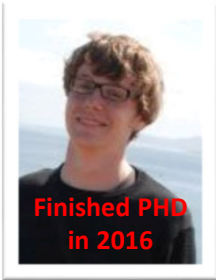  http://doi.acm.org/10.1145/1553374.1553453

# Lecture Bibliography (2)

- [13] H. Lee et al., 'Convolutional Deep Belief Networks for Scalable Unsupervised Learning of Hierarchical Representations', Online:
  http://doi.acm.org/10.1145/1553374.1553453
- [14] YouTube Video, 'Neural Network 3D Simulation', Online:
  https://www.youtube.com/watch?v=3JQ3hYko51Y
- [15] A. Rosebrock, 'Get off the deep learning bandwagon and get some perspective', Online:
  http://www.pyimagesearch.com/2014/06/09/get-deep-learning-bandwagon-get-perspective/
- [16] Big Data Tips – Big Data Mining & Machine Learning, Online:
  http://www.big-data.tips/
- [17] Harley, A.W., An Interactive Node-Link Visualization of Convolutional Neural Networks, Online:
  http://scs.ryerson.ca/~aharley/vis/conv/flat.html
- [18] A. Gulli and S. Pal, 'Deep Learning with Keras' Book, ISBN-13 9781787128422, 318 pages, Online:
  https://www.packtpub.com/big-data-and-business-intelligence/deep-learning-keras
- [19] D. Kingma and Jimmy Ba, 'Adam: A Method for Stochastic Optimization', Online:
  https://arxiv.org/abs/1412.6980
- [20] Kaiming He et al., 'Deep Residual Learning for Image Recognition', Online:
  https://arxiv.org/pdf/1512.03385.pdf
- [21] Sedona, R., Cavallaro, G., Riedel, M., Benediktsson, J.A. et. al.: Remote Sensing Big Data Classification with High Performance Distributed Deep Learning, Journal of Remote Sensing, Multidisciplinary Digital Publishing Institute (MDPI), Special Issue on Analysis of Big Data in Remote Sensing, 2019, Online:
  https://www.researchgate.net/publication/338077024_Remote_Sensing_Big_Data_Classification_with_High_Performance_Distributed_Deep_Learning
- [22] Horovod: Uber's Open Source Distributed Deep Learning Framework for TensorFlow, Online:
  https://www.slideshare.net/databricks/horovod-ubers-open-source-distributed-deep-learning-framework-for-tensorflow
- [23] T. Ben-Nun & T. Hoefler, 'Demystifying Parallel and Distributed Deep Learning: An In-depth Concurrency Analysis', Online:
  http://doi.acm.org/10.1145/3320060
- [24] AWS Amazon Sagemaker Service, Online:
  https://aws.amazon.com/sagemaker

# Lecture Bibliography (3)

- [25] Cheng, A.C, Lin, C.H., Juan, D.C., InstaNAS: Instance-aware Neural Architecture Search, Online:
  https://arxiv.org/abs/1811.10201
- [26] M. Riedel, 'Neural Architecture Search with Reinforcement Learning', Online:
  http://www.morrisriedel.de/neural-architecture-search-with-reinforcement-learning
- [27] Google Colaboratory, Online:
  https://colab.research.google.com
- [28] Machine Learning Mastery MNIST Tutorial, Online:
  https://machinelearningmastery.com/handwritten-digit-recognition-using-convolutional-neural-networks-python-keras/
- [29] The XOR Problem in Neural Networks, Online:
  https://medium.com/@jayeshbahire/the-xor-problem-in-neural-networks-50006411840b
- [30] F. Rosenblatt, 'The Perceptron - a perceiving and recognizing automaton', Report 85-460-1, Cornell Aeronautical Laboratory, 1957, Online:
  https://blogs.umass.edu/brain-wars/files/2016/03/rosenblatt-1957.pdf
- [31] MIT 6.S191: Introduction to Deep Learning, Online:
  http://introtodeeplearning.com/
- [32] An overview of gradient descent optimization algorithm, Online:
  http://ruder.io/optimizing-gradient-descent/index.html#visualizationofalgorithms
- [33] www.big-data.tips, 'MNIST Database', Online:
  http://www.big-data.tips/mnist-database
- [34] www.big-data.tips, 'Relu Neural Network', Online:
  http://www.big-data.tips/relu-neural-network
- [35] www.big-data.tips, 'tanh', Online:
  http://www.big-data.tips/tanh
- [36] www.big-data.tips,'MNIST dataset', Online:
  http://www.big-data.tips/mnist-dataset

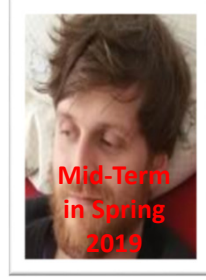# Acknowledgements – High Productivity Data Processing Research Group
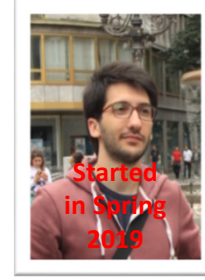


**PD Dr.**
**G. Cavallaro**

**Senior PhD**
**Student A.S. Memon**

**Senior PhD**
**Student M.S. Memon**

**PhD Student**
**E. Erlingsson**

**PhD Student**
**S. Bakarat**

**PhD Student**
**R. Sedona**

**Dr. M. Goetz**
**(now KIT)**

**MSc M.**
**Richerzhagen**
**(now other division)**

**MSc**
**P. Glock**
**(now INM-1)**

**MSc**
**C. Bodenstein**
**(now**
**Soccerwatch.tv)**

**MSc Student**
**G.S. Guðmundsson**
**(Landsverkjun)**