



Deep Learning Workshop II 2018

Organizing Committee:

Maggi Zhu

Aly El Gamal

Stanley Chan

Charles Bouman

Greg Buzzard

Dong Hye Ye

Amir Ziahari

Sri Yarlagada

Diyu Yang



Part 1:
*GPU Implementation of Deep Neural
Networks Using Tensorflow and Keras*

Amirkoushyar Ziabari

aziabari@purdue.edu

Aug 08 2018



- Deep Learning Frameworks
- Tensors
- What is Tensorflow (TF)?
- Tensorflow Structure
- CNN to classify MNIST
- GPU implementation
- Keras



Deep Learning Frameworks



Framework	Distributed Execution	Architecture Optimization	Visualization	Community Support	Portability	Notes
Tensorflow	✓ ✓	✓ ✓	✓ ✓	✓ ✓	✓ ✓	<ul style="list-style-type: none">• Google,• wide usage, ecosystem and community support• Visualization is superior
Pytorch	✓ ✓	✓ ✓	✓ ✓	✓ ✓	✓ ✓	<ul style="list-style-type: none">• Facebook• Easy to use & technically impressive
CNTK	✓ ✓	✓ ✓	✓	-	✓ ✓	<ul style="list-style-type: none">• Microsoft• Licensig issues
MXNet	✓	✓ ✓	✓	-	✓ ✓	<ul style="list-style-type: none">• Behind other frameworks in DE, and lacks details
Caffe2	✓ ✓	✓ ✓	-	-	✓ ✓	<ul style="list-style-type: none">• Facebook• Still moving
Caffe	-	✓ ✓	✓	✓	✓	<ul style="list-style-type: none">• Facebook• Lacks future community support
Theano	✓	✓ ✓	✓	✓	✓	<ul style="list-style-type: none">• University of Montreal• Lacks future community support

□ <https://agi.io/2018/02/09/survey-machine-learning-frameworks/>

- ❖ Tensors are the standard way of representing data in Tensorflow (deep learning)
- ❖ Tensors are multidimensional arrays, an extension of matrices to data with higher dimensions

't'
'e'
'n'
's'
'o'
'r'

Tensor of dimension[1]

3	1	4	1
5	9	2	6
5	3	5	8
9	7	9	3
2	3	8	4
6	2	6	4

Tensor of dimensions[2]

2	1	8	8	1	8		
2	8	4	5	0	4	5	
2	3	5	3	0	2	8	
7	4	7	1	3	5	2	6

Tensor of dimensions[3]



Tensor Rank



Rank	Math Entity	Python Example
0	Scalar (magnitude only)	<code>s = 483</code>
1	Vector (magnitude and direction)	<code>v = [1.1, 2.2, 3.3]</code>
2	Matrix (table of numbers)	<code>m = [[1, 2, 3], [4, 5, 6], [7, 8, 9]]</code>
3	3-Tensor (cube of numbers)	<code>t = [[[2], [4], [6]], [[8], [10], [12]], [[14], [16], [18]]]</code>
n	n-Tensor (you get the idea)	<code>....</code>



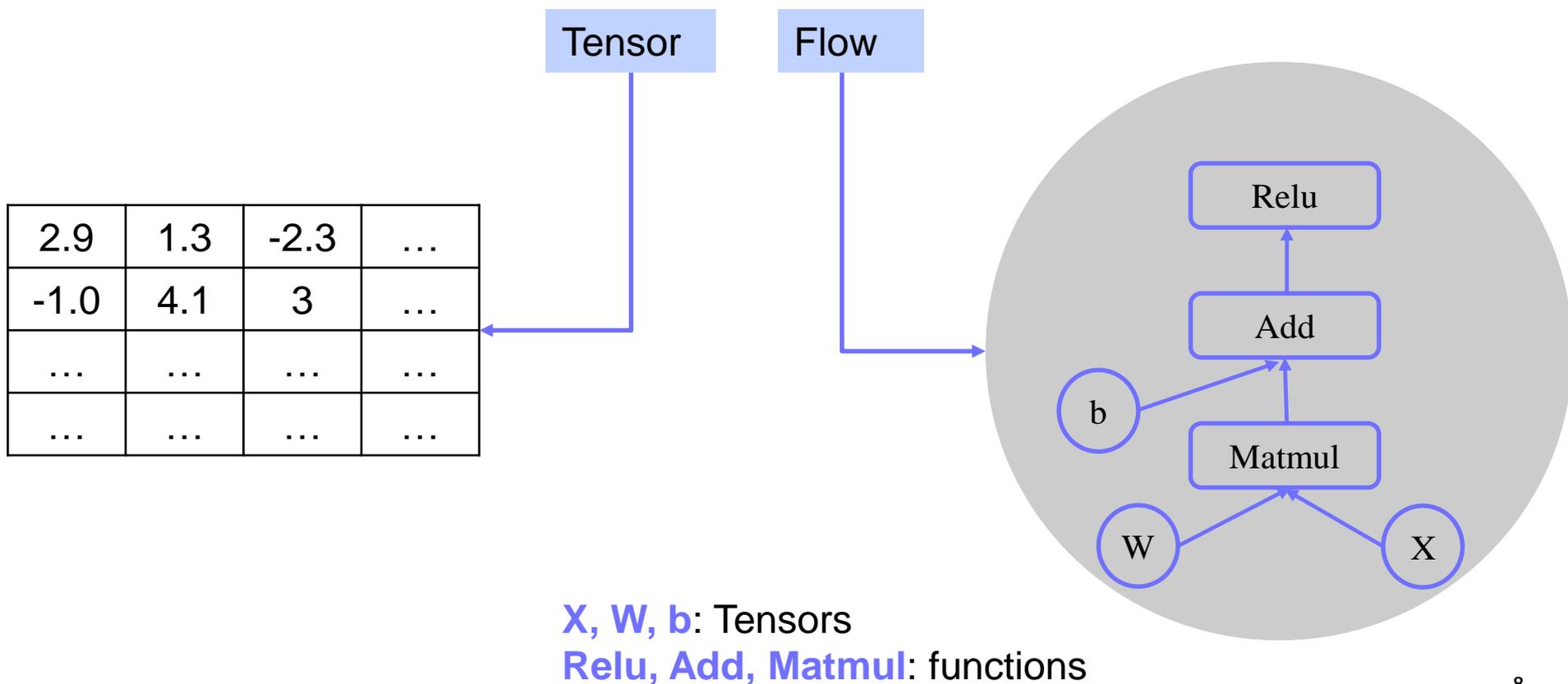
Tensor Data Types



- ❖ In addition to dimensionality Tensors have different data types

Data type	Python type	Description
DT_FLOAT	<code>tf.float32</code>	32 bits floating point.
DT_DOUBLE	<code>tf.float64</code>	64 bits floating point.
DT_INT8	<code>tf.int8</code>	8 bits signed integer.
DT_INT16	<code>tf.int16</code>	16 bits signed integer.
DT_INT32	<code>tf.int32</code>	32 bits signed integer.
DT_INT64	<code>tf.int64</code>	64 bits signed integer.
DT_UINT8	<code>tf.uint8</code>	8 bits unsigned integer.
DT_STRING	<code>tf.string</code>	Variable length byte arrays. Each element of a tensor is a byte array.
DT_BOOL	<code>tf.bool</code>	Boolean.

- ❖ Tensorflow (TF): a python library to implement deep networks
 - ❖ Very simple to install on all operating systems ([tensorflow.org/install](https://www.tensorflow.org/install))
 - ❖ Pycharm, ipython, etc can be used to run TF on windows
- ❖ In Tensorflow, computation is approached as a dataflow graph



TensorFlow core programs consists of two discrete sections:

Building a computational graph

Running a computational graph

A computational graph is a series of TensorFlow operations arranged into a graph of nodes

```
import tensorflow as tf
```

```
a = tf.constant(5.0, tf.float32)
```

```
b = tf.constant(6.0)
```

```
c = a*b
```

Build a computational graph

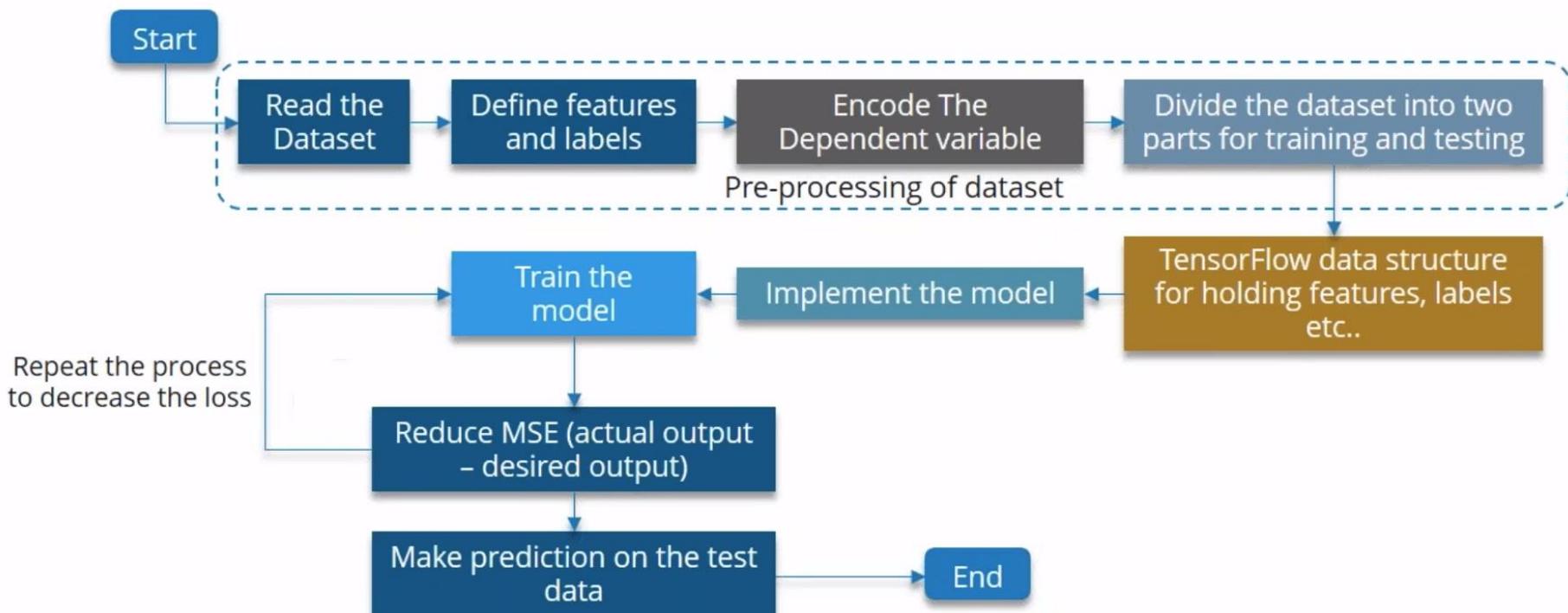
Run the computational graph



Goal



- Develop a Convolutional Neural Network To Classify MNIST DATA**



Build the computational Graph

0. Import the necessary libraries
1. Read the input data; define parameters, constants..
2. Define input/target size, type
Assign space for input/target
3. Define weights and biases
4. Define and construct the model
(e.g. Convolutional Neural Net)
5. Define loss function
6. Choose optimization technique
7. Define Training operation
8. Define Initialization operation
9. Define events logs and saving operations

Launch the computational Graph

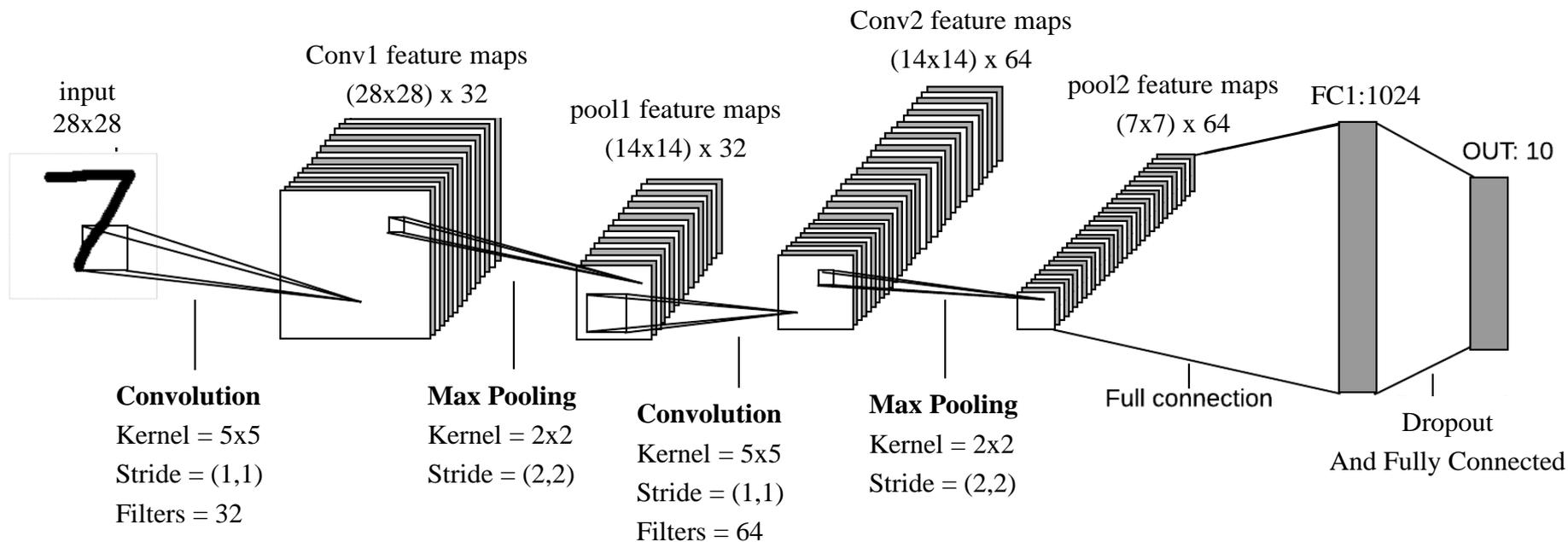
10. Define Session and run initialization
11. Train the model (run the training op.)
Print outputs, Save (or restore) model and events logs

MNIST Dataset Overview

This example is using MNIST handwritten digits. The dataset contains 60,000 examples for training and 10,000 examples for testing. The digits have been size-normalized and centered in a fixed-size image (28x28 pixels) with values from 0 to 1. For simplicity, each image has been flattened and converted to a 1-D numpy array of 784 features (28*28).



More info: <http://yann.lecun.com/exdb/mnist/>





Review the code



Supported devices

On a typical system, there are multiple computing devices. In TensorFlow, the supported device types are `CPU` and `GPU`. They are represented as `strings`. For example:

- `"/cpu:0"`: The CPU of your machine.
- `"/device:GPU:0"`: The GPU of your machine, if you have one.
- `"/device:GPU:1"`: The second GPU of your machine, etc.

If a TensorFlow operation has both CPU and GPU implementations, the GPU devices will be given priority when the operation is assigned to a device. For example, `matmul` has both CPU and GPU kernels. On a system with devices `cpu:0` and `gpu:0`, `gpu:0` will be selected to run `matmul`.



- 1. To perform particular operation on a device of your choice**
- 2. with `tf.device ()`**
 1. Creates a context such that all operations with the context will have the same device assignment



- ❑ By default, TF maps nearly all of the GPU memory of all GPUs visible to the process
- ❑ In some cases it is desirable for the process to only allocate a subset of the available memory, or to only grow the memory usage as it is needed by the process

```
config = tf.ConfigProto()  
config.gpu_options.allow_growth = True  
session = tf.Session(config=config, ...)
```



Allow soft placement



If you would like TensorFlow to automatically choose an existing and supported device to run the operations in case the specified one doesn't exist, you can set **allow_soft_placement** to **True** in the configuration option when creating the session.



Websites



- https://www.tensorflow.org/api_docs/python/tf/ConfigProto
- https://www.tensorflow.org/guide/using_gpu
- https://www.tensorflow.org/api_docs/python/tf/device



Assignment 1



- Change code 1 to run on a gpu**
 - Set allow growth and soft placement to True
 - Set the graph to be on the cpu, while optimization and loss calculation to be on the gpu
 - If you confront memory error on google colab when calculating the test accuracy, then please write the code to input the test data in batches

What is Keras?

Keras is a high-level neural networks API, written in Python and capable of running on top of [TensorFlow](#), [CNTK](#), or [Theano](#). It was developed with a focus on enabling fast experimentation. *Being able to go from idea to result with the least possible delay is key to doing good research.*

Why Keras?

- User friendly
- Modularity
- Easy Extensibility
- Work with python

0. Import the necessary libraries

1. Read the input data; define parameters, constants..

2. Define model structure. (**Sequential** model: linear stack of layers, or more complicated **Keras functional API** user defined model)

3. Add layers of the model

4. Compile the model

5. Fit the model

6. Evaluate the model



<https://keras.io/>



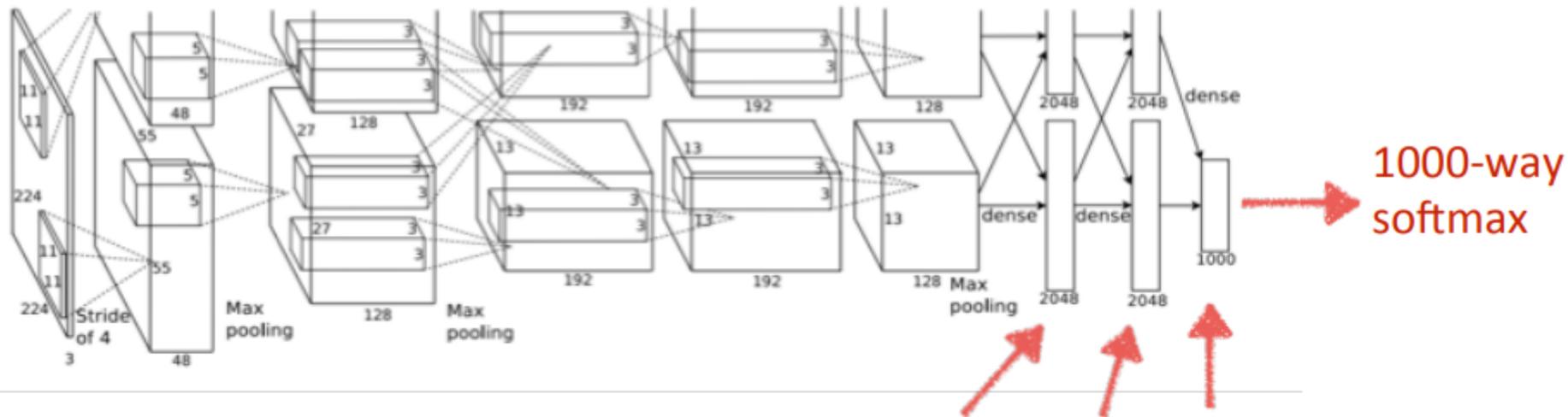
Let's convert our TF code to Keras



Show how to run Keras on GPU and Multi-GPUs

□ Convert the tensorflow code for AlexNet into Keras

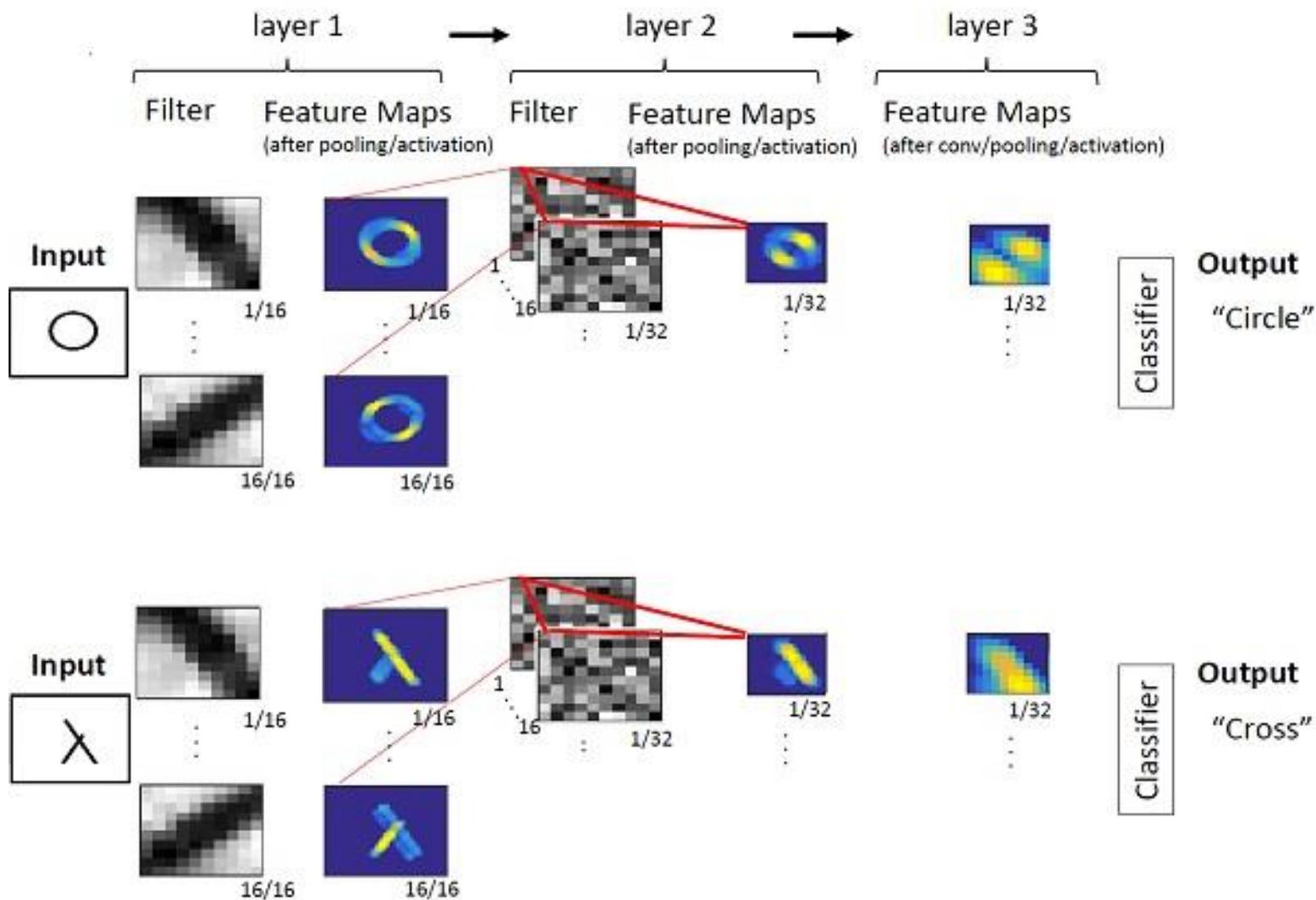
5 Convolutional Layers

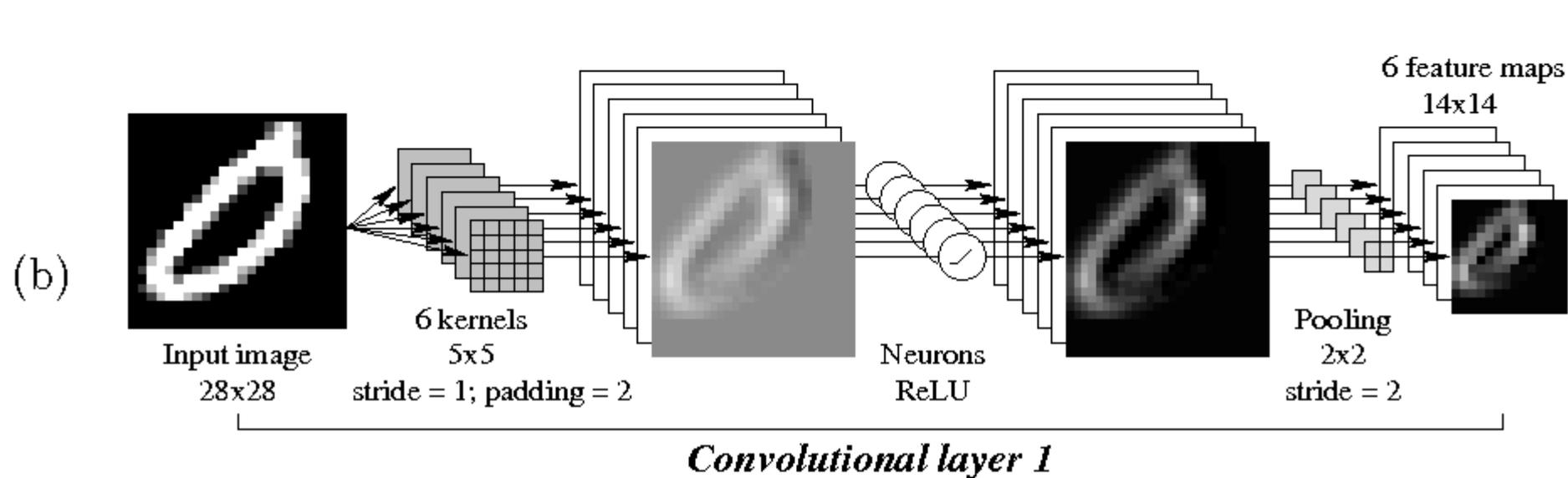
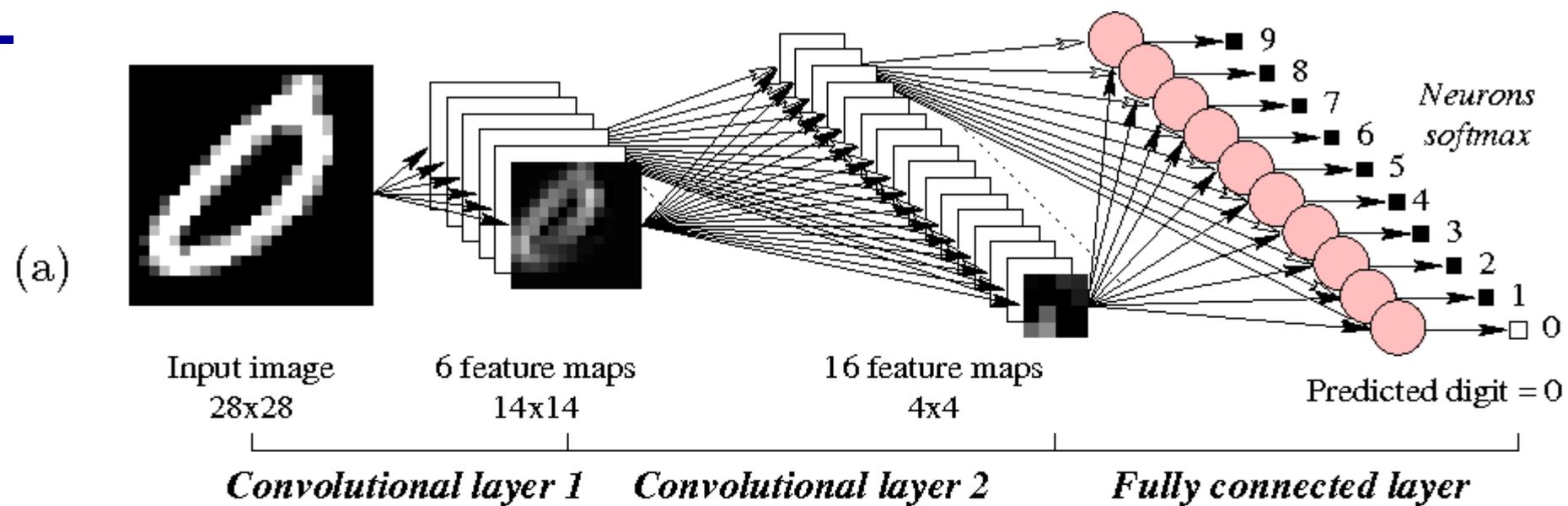


3 Fully Connected Layers



Backup







- **training**: Either a Python boolean, or a TensorFlow boolean scalar tensor (e.g. a placeholder). Whether to return the output in training mode (normalized with statistics of the current batch) or in inference mode (normalized with moving statistics). **NOTE**: make sure to set this parameter correctly, or else your training/inference will not work properly.

□ Training (mini-batch)

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$s^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

• Testing (one-example)

$$z_{norm}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{norm}^{(i)} + \beta$$

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \quad // \text{ scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.