



Running AlexNet on Raspberry Pi with Compute Library

Document number: ARM-ECM-074668I Version: 1.0

Date of Issue: 05/02/2018

© Copyright ARM Limited 2018. All rights reserved.

Abstract

This guide demonstrates how to run the AlexNet Convolutional Neural Network (CNN) using the Compute Library and a Raspberry Pi for simple image category classification.

Keywords

Machine learning, artificial intelligence, neural network, AlexNet, Raspberry Pi, Compute Library.

Contents

1	OVERVIEW	2
2	PREREQUISITES	2
2.1	Raspberry Pi and host machine requirements	3
3	INTRODUCING THE GRAPH API	3
4.1	What does AlexNet consist of?	5
4.2	Grouping	6
5.1	Header files	6
5.2	Mean subtraction pre-processing	7
5.3	Network description	7

7	COMPILE THE COMPUTE LIBRARY	8
8	RUN THE CLASSIFIER	9
9	DEVELOP YOUR OWN NETWORK USING THE COMPUTE LIBRARY	10

I Overview

This guide follows on from [our blog post](#) that introduces the Compute Library.

The instructions given here show you how to develop a Convolutional Neural Network (CNN) called AlexNet using just the Compute Library and a Raspberry Pi. Links are provided to all of the software tools that you need to get up and running.

The guide starts by introducing Compute library's graph API and AlexNet, two tools that help simplify developing neural networks on a Raspberry Pi. An example of AlexNet using the graph API, in C++, is explained in detail to help you get started running your own and other classifiers.

By following the steps in this guide, you'll be up and running with AlexNet, one of the first Deep Convolutional Neural Networks (CNN) designed to recognize 1000 different object categories within images. You will use AlexNet to classify an image of a go-kart with the neural network returning some predictions based on the image content. The network can only be used to predict the 1000 object categories that it has been trained with. If you want to use the CNN for a different task, then the network has to be re-trained.



2 Prerequisites

This guide is a sequel to our blog post on [how to apply a cartoon effect with the Compute Library](#). The blog post introduces the Compute Library and provides a simple example of how to use Raspberry with SSH. The

post explains how to compile or cross-compile the Compute Library for Raspberry Pi, but this is also covered here.

In addition to some basic knowledge of the Compute Library, this guide assumes some knowledge of a CNN. You don't need to be an expert, just have an idea of the main functions.

Further ahead in this guide, you will download and unzip [the tutorial .zip file](#) on the Raspberry Pi. This file contains:

- The AlexNet model. This is the same one as in the Caffe Model Zoo.
- A text file, containing the ImageNet labels that are required to map the predicted objects to the name of the classes.
- Several images in the ppm file format that are ready to be used with the network.

2.1 Raspberry Pi and host machine requirements

1. Raspberry Pi 2 or 3 with [Ubuntu Mate 16.04.02](#).
2. A blank Micro SD card. We recommend an 8 GB (minimum 6GB) Class 6 or Class 10 microSDHC card for installing the Raspberry Pi OS and storing the CNN model.
3. Router and ethernet cable. This is to connect to the Raspberry Pi using SSH.

3 Introducing the Graph API

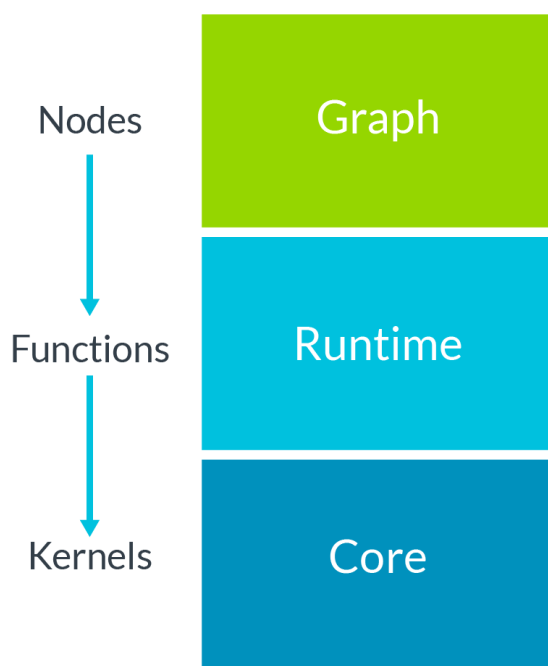
In release 17.09 of the Compute Library, we introduced the Graph API which is an important feature to make life easier for developers, and anyone else benchmarking the library.

The Graph API's primary function is to reduce the boilerplate code, but it can also reduce errors in your code and improve its readability. It is simple and easy-to-use, with a stream interface that is designed to be similar to other C++ objects.

At the current stage, the Graph API only supports the ML functions, such as convolution, fully connected, activation, pooling, and so on. To use the Graph API, you must compile the library with both NEON and OpenCL enabled by setting `neon=1` and `opencl=1`.

Note: As Raspberry Pi does not have OpenCL, the Graph API will automatically fall back to using NEON. This is why you need to compile the Compute Library with both NEON and OpenCL enabled.

In terms of building blocks, the Graph API represents the third computation block, together with core and runtime. In terms of hierarchy, the Graph API lies just above the runtime, which in turn lies above the core block as this image shows.



4 Introducing AlexNet

AlexNet is a convolutional neural network (CNN) that rose to prominence when it won [the ImageNet Large Scale Visual Recognition Challenge](#) (ILSVRC), an annual challenge that aims to evaluate algorithms for object detection and image classification. The model is trained on more than a million images and can classify images into 1000 object categories.

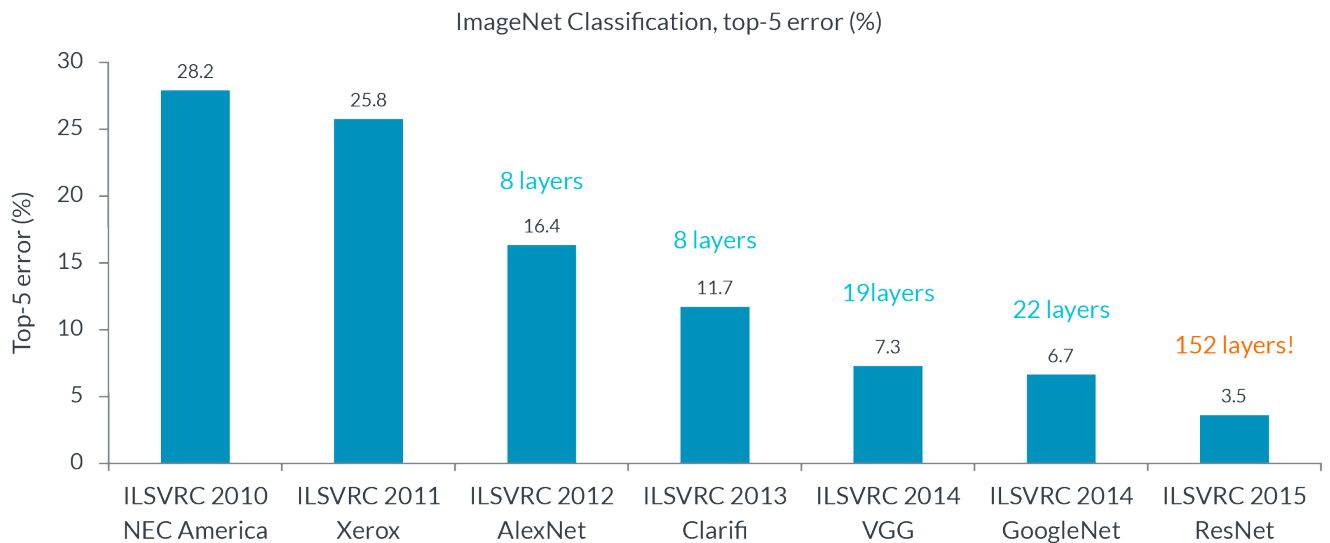
The ILSVRC evaluates the success of image classification solutions by using two important metrics, the top-5 and top-1 errors. When given a set of N images, often called test images, and mapped to a target class for each metric:

- top-1 error checks if the top predicted class is the same as the target class.
- top-5 error checks if the target class is one of the top five predictions.
-

For both metrics, the top error is calculated as, "the number of times the predicted class does not match the target class, divided by the total number of test images". In simpler terms, a lower score is better.

AlexNet achieved a top-5 error around 16%, which was an extremely good result back in 2012. To put it into context, until that year no other classifier had been able to achieve results under 20%. AlexNet was also more than 10% more accurate than the runner up.

Since 2012, other CNNs, such as VGG and ResNet, have improved on AlexNet's performance, as illustrated in this graph.



4.1 What does AlexNet consist of?

AlexNet is made up of eight trainable layers, five convolution layers and three fully connected layers. All the trainable layers are followed by a ReLu activation function, except for the last fully connected layer, where the Softmax function is used.

Besides the trainable layers, the network also has:

1. Three pooling layers.
2. Two normalization layers.
3. One dropout layer. This is only used for training to reduce the overfitting.

This table shows the layers and their details:

n.	Layer	Info
1	Convolution	11x11x3x96 - (Stride(4,4) - Pad(0,0))
2	Activation	ReLu
3	Normalization	Cross Map,5,0.0001,0.75
4	Pooling	3x3 - Stride(2,2)
5	Grouping Convolution	5x5x96x256 - Stride(1,1) - Pad(2,2)
6	Activation	ReLu
7	Normalization	Cross Map,5,0.0001,0.75
8	Pooling	3x3 - Stride(2,2)
9	Convolution	3x3x256x384 - Stride(1,1) - Pad(1,1)
10	Activation	ReLu
11	Grouping Convolution	3x3x384x384 - Stride(1,2) - Pad(1,1)
12	Activation	ReLu
13	Grouping Convolution	3x3x384x256 - Stride(1,1) - Pad(1,1)

n.	Layer	Info
14	Activation	ReLu
15	Pooling	3x3 - Stride(2,2)
16	Fully connected	4096x9216 - Stride(1,1) - Pad(1,1)
17	Activation	ReLu
18	Fully connected	4096x9216 - Stride(1,1) - Pad(1,1)
19	Activation	ReLu
20	Fully connected	1000x4096 - Stride(1,1) - Pad(1,1)
21	Softmax	

4.2 Grouping

In the table, there are some convolution layers that are actually grouping convolutions. This is an efficient engineering trick that allows the acceleration of the network over two GPUs without sacrificing accuracy.

If the group size is set to two, then the first half of the filters will be connected to the first half of the input feature maps and the second half will be connected to the second half of the input feature maps, as this image shows.

5 Evaluate the example code

A C++ implementation of AlexNet using the Graph API is proposed in `examples/graph_alexnet.cpp`. Here, we will look at some key elements of the code and what they do.

To run the AlexNet example you need to include these four command line arguments:

```
./graph_alexnet <target> <path_cnn_data> <input_image> <labels>
```

Where:

1. Target is the the type of acceleration (NEON=0 or OpenCL=1).
2. Path to `cnn_data`.
3. Path to your input image. Note that only .ppm files are supported).
4. Path to your ImageNet labels.

These subsections describe the key aspects of the example:

5.1 Header files

In order to use the Graph API you need to include these three header files:

```
// Contains the definitions for the graph
#include "arm_compute/graph/Graph.h"
```

```
// Contains the definitions for the nodes (convolution, pooling, fully
connected)
#include "arm_compute/graph/Nodes.h"

// Contains the utility functions for the graph such as the accessors for
the input, trainable and output nodes. The accessors will be presented when
we are going to talk about the graph.
#include "utils/GraphUtils.h"
```

5.2 Mean subtraction pre-processing

A pre-processing stage is needed for preparing the input RGB image before feeding the network, so we are going to subtract the channel means from each individual color channel. This operation centres the red, green, and blue channels around the origin:

$$R_{\text{norm}}(x,y) = R(x,y) - R_{\text{mean}}(x,y)$$

$$G_{\text{norm}}(x,y) = G(x,y) - G_{\text{mean}}(x,y)$$

$$B_{\text{norm}}(x,y) = B(x,y) - B_{\text{mean}}(x,y)$$

Where:

- $r_{\text{norm}}(x,y)$, $g_{\text{norm}}(x,y)$, $b_{\text{norm}}(x,y)$ are the RGB values at coordinates x,y after the mean subtraction.
- $r(x,y)$, $g(x,y)$, $b(x,y)$ are the RGB values at coordinates x,y before the mean subtraction.
- r_{mean} , g_{mean} , and b_{mean} are the mean values to use for the RGB channels.

For simplicity, the mean values for the examples are already hard-coded as:

```
constexpr float mean_r = 122.68f; /* Mean value to subtract from red channel
*/
constexpr float mean_g = 116.67f; /* Mean value to subtract from green
channel */
constexpr float mean_b = 104.01f; /* Mean value to subtract from blue
channel */
```

If you are not familiar with mean subtraction pre-processing before, the [Compute Image Mean section on the Caffe website](#) provides a useful explanation.

5.3 Network description

The body of the network is described through the Graph API.

The graph consists of three main parts:

1. Mandatory: One input - Tensor object. This layer describes the geometry of the input data along with the data type to use. In this case, we'll have a 3D input image with shape 227x227x3, using the FP32 data type.
2. The Convolution Neural Network layers, or 'nodes' in the graph's terminology. These are needed for the network.
3. Mandatory: One output - Tensor object. This is used to get the result back from the network.

As you can see in the example, the Tensor objects (input and output) and all of the trainable layers accept an input function called accessor.

Important: The accessor is the only way to access the internal Tensors.

- The accessor used by the input Tensor object can initialize the input Tensor of the network. This function can also be responsible for the mean subtraction pre-processing and reading the input image from a file or camera.
- The accessor used by the trainable layers, such as convolution, fully connected, and so on, can initialize the weights and the biases reading. For example, the values from a NumPy file.
- The accessor used by the output Tensor object can return the result of the classification, along with the score.

If you want to understand how the accessor works, see the the utils/GraphUtils.h file. This has some ready-to-use accessors for your Tensor objects and trainable layers.

6 Download and install the tutorial ZIP file

Here, you are going to turn on your Raspberry Pi and test AlexNet with some images.

First, turn on your Raspberry Pi and navigating to the home directory of your Raspberry Pi or host machine.

Then, on your Raspberry Pi enter the following commands:

```
# Install unzip
sudo apt-get install unzip
# Download the zip file with the AlexNet model, input images and labels
wget https://developer.arm.com/-/media/developer/technologies/Machine
learning/on_Arm/Tutorials/Running_AlexNet_on_Pi_with_Compute
Library/compute_library_alexnet.zip?revision=c1a232fa-f328-451f-9bd6-
250b83511e01
# Create a new folder
mkdir assets_alexnet
# Unzip
unzip compute_library_alexnet.zip -d assets_alexnet
```

The contents of the ZIP file are now extracted into the assets_alexnet folder.

7 Compile the Compute Library

To compile natively on your Raspberry Pi, enter the following on the command line:


```
# Clone Compute Library
git clone https://github.com/Arm-software/ComputeLibrary.git
# Enter ComputeLibrary folder
cd ComputeLibrary
# Native Build the library and the examples
scons Werror=1 debug=0 asserts=0 neon=1 opencl=1 examples=1 build=native -j2
Or, to cross-compile, enter the following on the command line:
# Clone Compute Library
git clone https://github.com/Arm-software/ComputeLibrary.git
# Enter ComputeLibrary folder
cd ComputeLibrary
# Build the library and the examples scons Werror=1 debug=0 asserts=0 neon=1
opencl=1 examples=1 os=linux arch=armv7a -j4 # Copy the example and dynamic
libraries on the Raspberry Pi scp build/example/graph_alexnet
build/libarm_compute.so build/libarm_compute_core.so
build/libarm_compute_graph.so @:Desktop
```

Where:

- <username_raspberrypi> is the username used on your Raspberry Pi.
- <ip_addr_raspberrypi> is the IP address of your Raspberry Pi.

8 Run the classifier

The CNN has been trained to recognise 1000 object categories. A go-kart is one of these. This step demonstrates how the object is recognized when the CNN is passed an image of a go-kart. In contrast, if you defined a random image which is not part of the 1000 categories then the CNN will not be able to recognize it.

If you compiled natively on your Raspberry Pi, enter the following on the command line to run the classifier against the go_kart.ppm image:

```
export LD_LIBRARY_PATH=build/
PATH_ASSETS=../assets_alexnet
./build/examples/graph_alexnet 0 $PATH_ASSETS $PATH_ASSETS/go_kart.ppm
$PATH_ASSETS/labels.txt
```

Or, if you cross-compiled, on your host machine open an SSH session by entering:

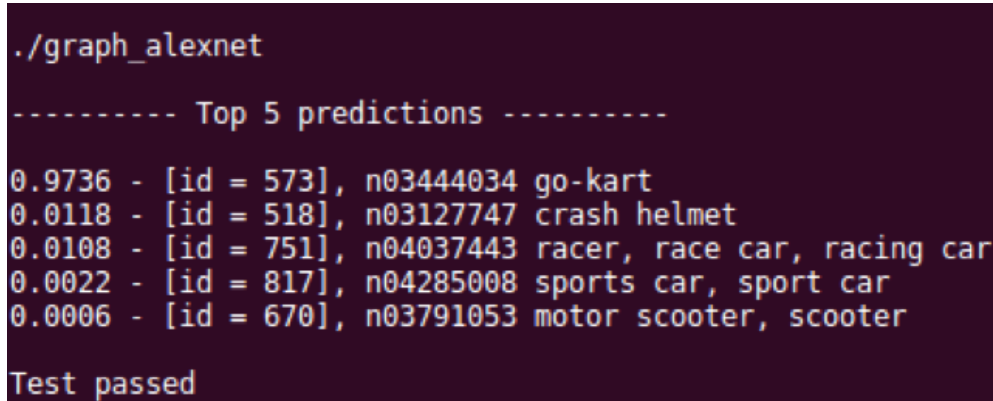
```
ssh
<username_raspberrypi>@<ip_addr_raspberrypi>>?</ip_addr_raspberrypi></use
rname_raspberrypi>
```

And in the SSH session, enter the Desktop folder and run the classifier against the go-kart .ppm image:

```
cd Desktop
export LD_LIBRARY_PATH=build/
```

```
PATH_ASSETS=../assets_alexnet
./build/examples/graph_alexnet 0 $PATH_ASSETS $PATH_ASSETS/go_kart.ppm
$PATH_ASSETS/labels.txt
```

Whether or not you are building the library natively, the output should look like this if a successful classification has been performed:



```
./graph_alexnet
----- Top 5 predictions -----
0.9736 - [id = 573], n03444034 go-kart
0.0118 - [id = 518], n03127747 crash helmet
0.0108 - [id = 751], n04037443 racer, race car, racing car
0.0022 - [id = 817], n04285008 sports car, sport car
0.0006 - [id = 670], n03791053 motor scooter, scooter
Test passed
```

This screen-shot shows that the classifier has provided five predictions of the content of the image against the object categories that the CNN has been trained with. If the output does not look like the screen-shot, then it is highly likely that the assets have not been correctly copied to the SD card.

9 Develop your own network using the Compute Library

If your test passed then you have successfully used AlexNet to classify the characteristics of an image.

You can now go on to use what you have learned here to develop even more exciting and performant intelligent vision solutions on Arm-based platforms.

The Compute Library contains other CNNs that are capable of recognising objects with greater accuracy than AlexNet. These include MobileNet, VGG-16, and GoogleNet. Try experimenting with these as you begin to develop your own network using the Compute Library.