



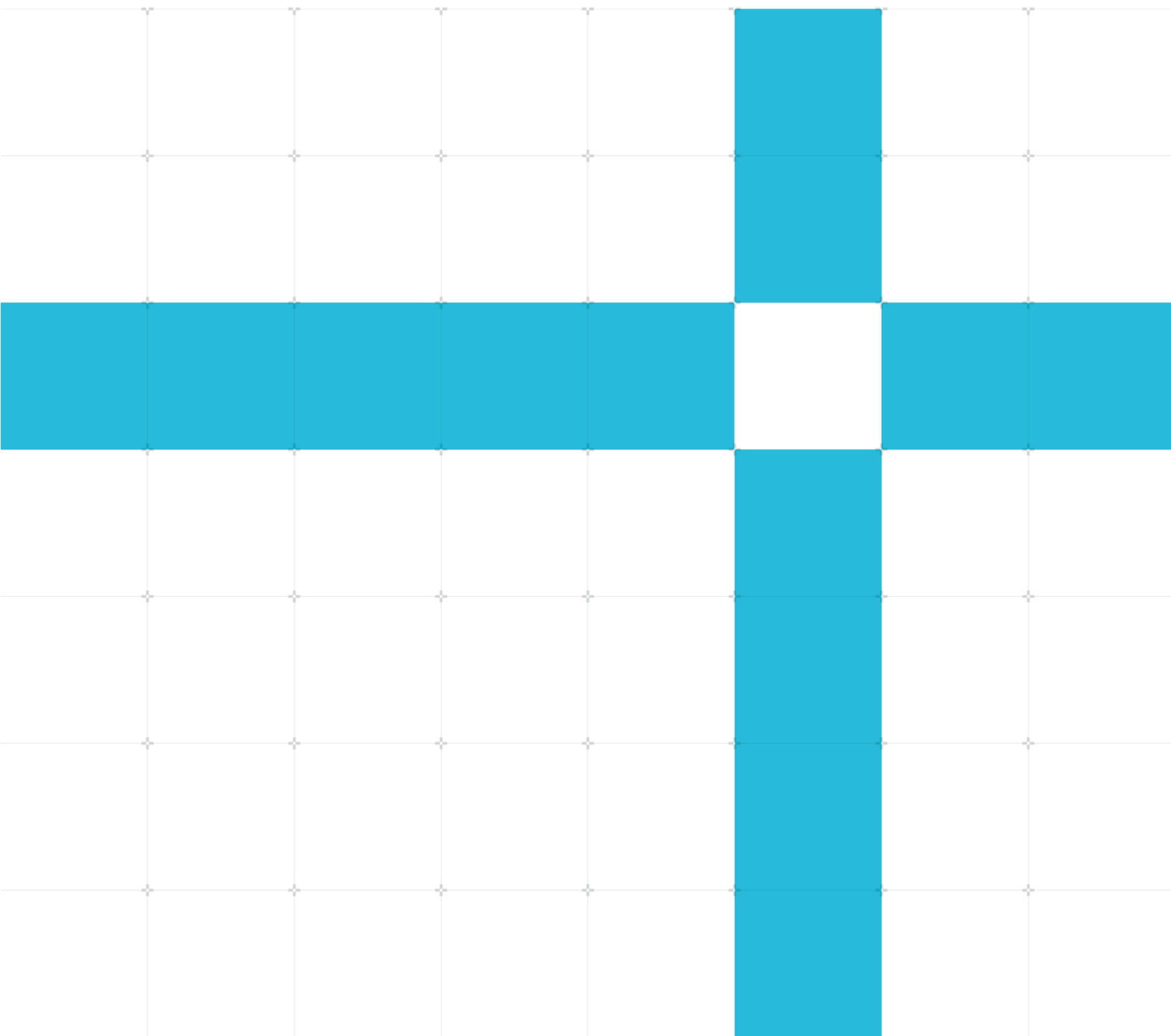
# Running AlexNet on Raspberry Pi with Arm Compute Library

Non-Confidential

Copyright © 2018, 2020 Arm Limited (or its affiliates).  
All rights reserved.

**Issue 1.2**

ARM-ECM-0746681



## Running AlexNet on Raspberry Pi with Arm Compute Library

Copyright © 2018, 2020 Arm Limited (or its affiliates). All rights reserved.

### Release information

#### Document history

Issue	Date	Confidentiality	Change
1.1	05/02/2018	Non-Confidential	First release
1.2	07/09/2020	Non-Confidential	General editorial updates

## Non-Confidential Proprietary Notice

This document is protected by copyright and other related rights and the practice or implementation of the information contained in this document may be protected by one or more patents or pending patent applications. No part of this document may be reproduced in any form by any means without the express prior written permission of Arm. No license, express or implied, by estoppel or otherwise to any intellectual property rights is granted by this document unless specifically stated.

Your access to the information in this document is conditional upon your acceptance that you will not use or permit others to use the information for the purposes of determining whether implementations infringe any third party patents.

THIS DOCUMENT IS PROVIDED "AS IS". ARM PROVIDES NO REPRESENTATIONS AND NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, INCLUDING, WITHOUT LIMITATION, THE IMPLIED WARRANTIES OF MERCHANTABILITY, SATISFACTORY QUALITY, NON-INFRINGEMENT OR FITNESS FOR A PARTICULAR PURPOSE WITH RESPECT TO THE DOCUMENT. For the avoidance of doubt, Arm makes no representation with respect to, and has undertaken no analysis to identify or understand the scope and content of, patents, copyrights, trade secrets, or other rights.

This document may include technical inaccuracies or typographical errors.

TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL ARM BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF ARM HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

This document consists solely of commercial items. You shall be responsible for ensuring that any use, duplication or disclosure of this document complies fully with any relevant export laws and regulations to assure that this document or any portion thereof is not exported, directly or indirectly, in violation of such export laws. Use of the word "partner" in reference to Arm's customers is not intended to create or refer to any partnership relationship with any other company. Arm may make changes to this document at any time and without notice.

If any of the provisions contained in these terms conflict with any of the provisions of any click through or signed written agreement covering this document with Arm, then the click through or signed written agreement prevails over and supersedes the conflicting provisions of these terms. This document may be translated into other languages for convenience, and you agree that if there is any conflict between the English version of this document and any translation, the terms of the English version of the Agreement shall prevail.

The Arm corporate logo and words marked with ® or ™ are registered trademarks or trademarks of Arm Limited (or its affiliates) in the US and/or elsewhere. All rights reserved. Other brands and names mentioned in this document may be the trademarks of their respective owners. Please follow Arm's trademark usage guidelines at <http://www.arm.com/company/policies/trademarks>.

Copyright © 2018, 2020 Arm Limited (or its affiliates). All rights reserved.

Arm Limited. Company 02557590 registered in England.

110 Fulbourn Road, Cambridge, England CB1 9NJ.

(LES-PRE-20349)

## Confidentiality Status

This document is Non-Confidential. The right to use, copy and disclose this document may be subject to license restrictions in accordance with the terms of the agreement entered into by Arm and the party that Arm delivered this document to.

Unrestricted Access is an Arm internal classification.

## Product Status

The information in this document is Final, that is for a developed product.

## Web Address

[www.arm.com](http://www.arm.com)

# Contents

<b>1 Introduction</b> .....	<b>5</b>
1.1 Product revision status.....	5
1.2 Intended audience.....	5
1.3 Conventions .....	5
1.3.1 Glossary.....	5
1.3.2 Typographical conventions .....	6
1.4 Additional reading.....	6
1.5 Feedback.....	7
1.5.1 Feedback on this product.....	7
1.5.2 Feedback on content .....	7
<b>2 Overview</b> .....	<b>8</b>
<b>3 Before you begin</b> .....	<b>9</b>
<b>4 Introducing the Graph API</b> .....	<b>10</b>
<b>5 Introducing AlexNet</b> .....	<b>11</b>
5.1 What is AlexNet consist of?.....	11
5.2 Grouping.....	13
<b>6 Evaluate the example code</b> .....	<b>14</b>
6.1 Mean subtraction pre-processing.....	14
6.2 Network description.....	15
<b>7 Download and install the tutorial ZIP file</b> .....	<b>16</b>
<b>8 Compile the Arm Compute Library</b> .....	<b>17</b>
<b>9 Run the classifier</b> .....	<b>18</b>
<b>10 Develop your own network using the Arm Compute Library</b> .....	<b>19</b>
<b>11 Related information</b> .....	<b>20</b>
<b>12 Next steps</b> .....	<b>21</b>

# 1 Introduction

## 1.1 Product revision status

The *mpn* identifier indicates the revision status of the product described in this book, for example, *r1p2*, where:

*rm*

Identifies the major revision of the product, for example, *r1*.

*pn*

Identifies the minor revision or modification status of the product, for example, *p2*.

## 1.2 Intended audience

This guide is for software developers who want to know how to run AlexNet on Raspberry Pi with the Arm Compute Library.

## 1.3 Conventions






The following subsections describe conventions used in Arm documents.

### 1.3.1 Glossary

The Arm Glossary is a list of terms used in Arm documentation, together with definitions for those terms. The Arm Glossary does not contain terms that are industry standard unless the Arm meaning differs from the generally accepted meaning.

See the Arm Glossary for more information: <https://developer.arm.com/glossary>.

## 1.3.2 Typographical conventions

Convention	Use
<i>italic</i>	Introduces citations.
<b>bold</b>	Highlights interface elements, such as menu names. Denotes signal names. Also used for terms in descriptive lists, where appropriate.
monospace	Denotes text that you can enter at the keyboard, such as commands, file and program names, and source code.
monospace <b>bold</b>	Denotes language keywords when used outside example code.
monospace <u>underline</u>	Denotes a permitted abbreviation for a command or option. You can enter the underlined text instead of the full command or option name.
<and>	Encloses replaceable terms for assembler syntax where they appear in code or code fragments. For example: <pre>MRC p15, 0, &lt;Rd&gt;, &lt;CRn&gt;, &lt;CRm&gt;, &lt;Opcode_2&gt;</pre>
SMALL CAPITALS	Used in body text for a few terms that have specific technical meanings, that are defined in the Arm® Glossary. For example, IMPLEMENTATION DEFINED, IMPLEMENTATION SPECIFIC, UNKNOWN, and UNPREDICTABLE.
 Caution	This represents a recommendation which, if not followed, might lead to system failure or damage.
 Warning	This represents a requirement for the system that, if not followed, might result in system failure or damage.
 Danger	This represents a requirement for the system that, if not followed, will result in system failure or damage.
 Note	This represents an important piece of information that needs your attention.
 Tip	This represents a useful tip that might make it easier, better or faster to perform a task.
 Remember	This is a reminder of something important that relates to the information you are reading.

## 1.4 Additional reading

This document contains information that is specific to this product. See the following documents for other relevant information:

## 1.5 Feedback

Arm welcomes feedback on this product and its documentation.

### 1.5.1 Feedback on this product

If you have any comments or suggestions about this product, contact your supplier and give:

- The product name.
- The product revision or version.
- An explanation with as much information as you can provide. Include symptoms and diagnostic procedures if appropriate.

### 1.5.2 Feedback on content

If you have comments on content, send an email to [errata@arm.com](mailto:errata@arm.com) and give:

- The title [Product Name] [Document Title].
- The number [Document ID].
- If applicable, the page number(s) to which your comments refer.
- A concise explanation of your comments.

Arm also welcomes general suggestions for additions and improvements.



Arm tests the PDF only in Adobe Acrobat and Acrobat Reader and cannot guarantee the quality of the represented document when used with any other PDF reader.

---

## 2 Overview

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

The instructions given here show you how to develop a Convolutional Neural Network (CNN) called AlexNet using just the Arm 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 Arm 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 will be up and running with AlexNet. AlexNet is one of the first Deep Convolutional Neural Networks (CNN) designed to recognize 1000 different object categories within images. You 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 retrained.





## 3 Before you begin

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 Arm Compute Library and provides a simple example of how to use Raspberry with SSH. The post explains how to compile or cross-compile the Arm Compute Library for Raspberry Pi, but this is also covered here.

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

Further ahead in this guide, you can 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.

### 4 Raspberry Pi and host machine requirements

- Raspberry Pi 2 or 3 with [Ubuntu Mate 16.04.02](#) or Raspberry Pi 4 with Raspbian 32-bit.
- A blank Micro SD card. We recommend an 8GB (minimum 6GB) Class 6 or Class 10 microSDHC card for installing the Raspberry Pi OS and storing the CNN model.
- Router and ethernet cable. This is to connect to the Raspberry Pi using SSH.

## 4 Introducing the Graph API

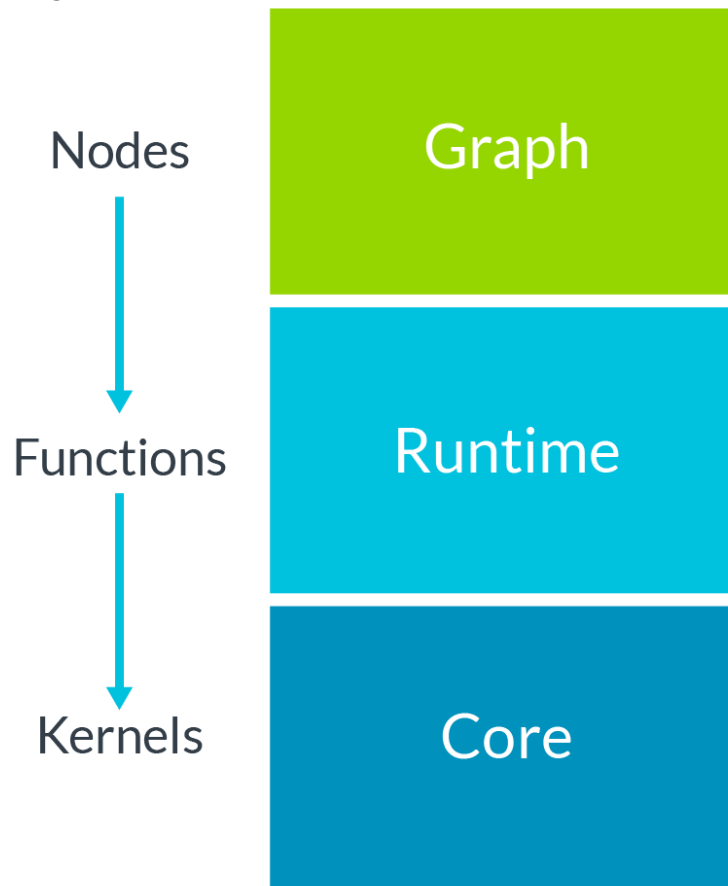
In release 17.09 of the Arm Compute Library, we introduced the Graph API. This is an important feature to make life easier for developers, and anyone else benchmarking the library.

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 similar to other C++ objects.

At the current stage, the Graph API only supports the ML functions, such as convolution, fully connected, activation, pooling. 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 Arm 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.



# 5 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.

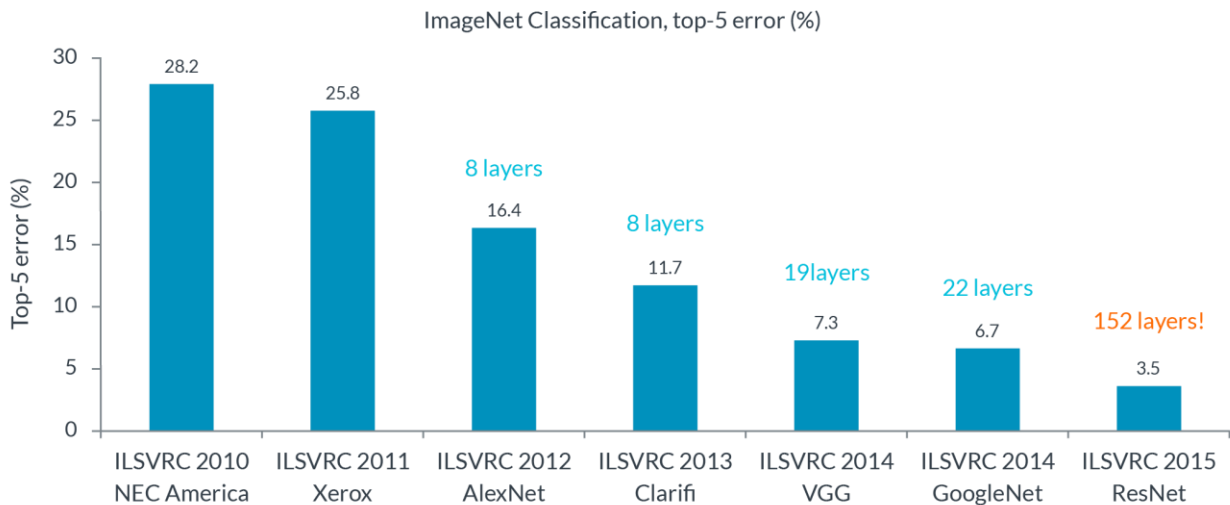
The ILSVRC evaluates the success of image classification solutions by using two important metrics, the top-5 and top-1 errors. When given a set of N images, often called test images, and mapped to a target class for each metric:

- 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 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.



## 5.1 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:

- Three pooling layers.

- Two normalization layers.
- One dropout layer. This is only used for training to reduce the overfitting.

This table shows the layers and their details:

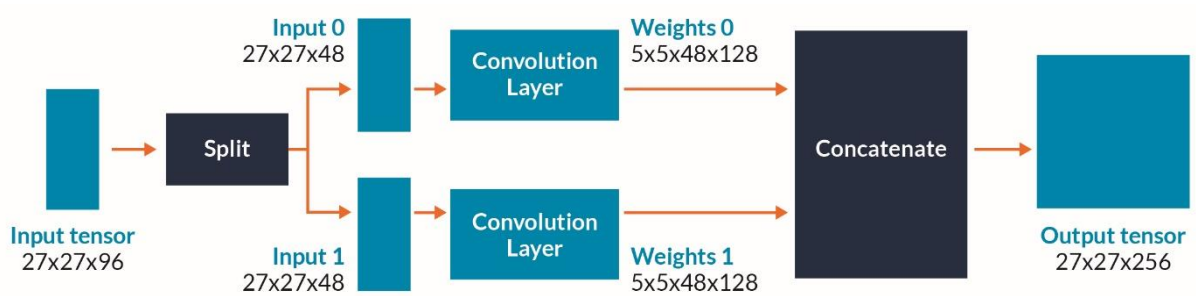
n.	Layer	Info
1	Convolution	11x11x3x96 - (Stride(4,4) - Pad(0,0))
2	Activation	ReLU
3	Normalization	Cross Map, 5,0.0001,0.75
4	Pooling	3x3 - Stride(2,2)
5	Grouping Convolution	5x5x96x256 - Stride(1,1) - Pad(2,2)
6	Activation	ReLU
7	Normalization	Cross Map, 5,0.0001,0.75
8	Pooling	3x3 - Stride(2,2)
9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

n.	Layer	Info
20	Fully connected	1000x4096 - Stride(1,1) - Pad(1,1)
21	Softmax	

## 5.2 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 is connected to the first half of the input feature maps. The second half is connected to the second half of the input feature maps, as this image shows.



The grouping convolution not only allows you to spread the workload over multiple GPUs, it also reduces the number of MACs needed for the layer by half.

## 6 Evaluate the example code

A C++ implementation of AlexNet using the Graph API is proposed in `examples/graph_alexnet.cpp`. Here, we 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 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:

### 6.1 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 centers 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_norm(x,y)`, `g_norm(x,y)`, `b_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:

```
const std::array<float, 3=""> mean_rgb{ { 122.68f, 116.67f, 104.01f } }; </float,>
```

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

## 6.2 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 will 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 `utils/GraphUtils.h` file. This has some ready-to-use accessors for your Tensor objects and trainable layers.

# 7 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/43359E999DEF433BAF63523C529D21AD.ashx?revision=cla232fa-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.



## 8 Compile the Arm 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.

## 9 Run the classifier

The CNN has been trained to recognize 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:

```
./build/examples/graph_alexnet --data=$PATH_ASSETS --image=$PATH_ASSETS/go_kart.ppm --labels=$PATH_ASSETS/labels.txt --target=neon -type=f32 -threads=4
```

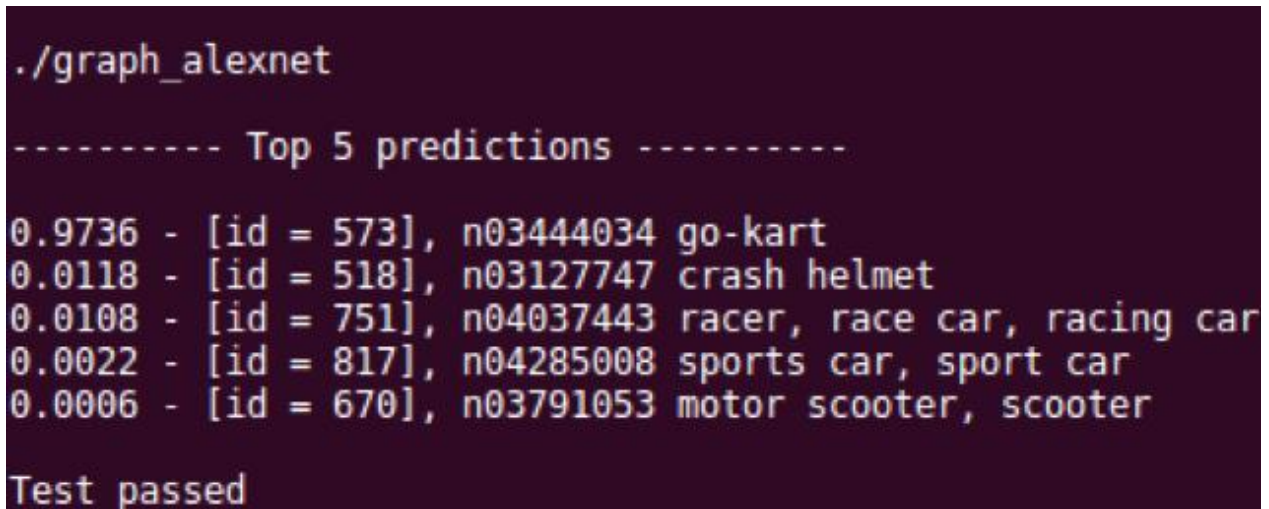
Or, if you cross-compiled, on your host machine open an SSH session by entering:

```
ssh  
<username_raspberrypi>@<ip_addr_raspberrypi>&gt;?</ip_addr_raspberrypi></username_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 screenshot 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 screenshot, then it is highly likely that the assets have not been correctly copied to the SD card.

# 10 Develop your own network using the Arm 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 Arm Compute Library contains other CNNs that are capable of recognizing 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 Arm Compute Library.

# 11 Related information

Arm Compute Library documentation

- <https://github.com/ARM-software/ComputeLibrary/wiki/Documentation>
- <https://community.arm.com/developer/tools-software/graphics/b/blog/posts/cartoonifying-images-on-raspberry-pi-with-the-compute-library>

## 12 Next steps

In this tutorial, we have learned how to run a model for image classification. Run the following command when to complete the guide.

- Run model for object detection (e.g. `graph_ssd_mobilenet`, `graph_yolov3`)

