

Deep Learning with OpenCL Caffe

A Guide to Real-Time Inference using Python

Version 1.5.0

Deep Learning with OpenCL Caffe

A Guide to Real-Time Inference using Python

Version 1.5.0

Rahul Atlury and Fabian David Tschopp

Media Press

Atlanta, GA

Copyright © 2017 Rahul Atlury and Fabian David Tschopp.

Media Press
Checkpoint Ave
GA 30008

Background for Rahul Atlury:

Rahul Atlury holds a Master's degree in IT. He is a member of the Intel Channel Board of Advisors (BOA), IoT for the Asian region. He contributes to open source and helps in OS releases for the Intel Edison platforms. Previously he worked as a research associate at the University of Hyderabad and helped with setup and maintenance of the PRAGMA Grid from India. He has worked as a consultant with various companies including the Indian Defense agencies.

Background for Fabian Tschopp:

Fabian Tschopp holds a Bachelor's degree in computer science and is currently a neuroinformatics Master student at the University and ETH Zürich. He is the developer and maintainer of the official BVLC Caffe OpenCL deep learning framework, in collaboration with HHMI Janelia, Intel and AMD. He has in-depth experience with visual computing, deep neural networks, 3D scanning and high performance computing. Fabian Tschopp is an Intel ISEF 2013 alumni and 2nd award winner.

Contents

Preface	xi
1 Python programming basics	1
1.1 Variables and types	1
1.2 Sequences	3
1.2.1 Tuples and lists	4
1.2.2 Bytes and bytearray	5
1.2.3 Slicing	6
1.2.4 Multi value assignments	7
1.3 Sets	8
1.4 Mappings	9
1.5 Conditions	10
1.5.1 If statement	10
1.5.2 Conditional expression selector	11
1.6 Loops	11
1.6.1 While statement	11
1.6.2 For statement	12
1.6.3 Enumerate	12
1.6.4 Comprehensions	13
1.7 Error handling	13

1.8	Context management	14
1.9	Input output	14
1.10	Functions	16
1.10.1	Generator function	16
1.10.2	Lambda function	17
1.11	Classes and objects	18
1.11.1	Instance objects and instance methods	18
1.11.2	Class variables and instance variables	19
1.11.3	Private variables	20
1.11.4	Special methods	20
1.11.5	Inheritance	21
1.11.6	Operator overloading	22
1.11.7	Multiple inheritance	23
1.11.8	Abstract methods	25
1.11.9	Decorators	25
1.11.10	Iterators	26
1.12	Modules and packages	26
2	Python scientific computing	29
2.1	The nature of data	29
2.2	Numpy data representation	29
2.3	Numpy array creation	32
2.4	Numpy data types	33
2.5	Numpy operations	33
2.5.1	Indexing	33
2.5.2	Slicing	33

2.5.3	Iterating arrays	35
2.5.4	Arithmetic Operators	36
2.5.5	Matrix operator	36
2.5.6	Increment decrement operators	36
2.5.7	Universal operations	37
2.5.8	Aggregate operations	37
2.5.9	Conditions and boolean arrays	37
2.5.10	Shape manipulation	38
2.5.11	Joining arrays	38
2.5.12	Splitting arrays	38
2.5.13	Copy/Views	39
2.5.14	Vectorization	40
2.5.15	Broadcasting	40
2.5.16	Structured arrays	40
2.5.17	Reading/Writing array data	40
2.5.18	Recording files with tabular data	40
2.6	Matplotlib	41
3	Installation and configuration	43
3.1	Introduction	43
3.2	BLAS libraries	44
3.3	OpenCL backend	45
3.4	Installation prerequisites	47
3.4.1	Standard requirements	47
3.4.2	OpenCL requirements	48
3.5	Device specific considerations	48

3.6	Compilation	49
3.6.1	Windows	49
3.6.2	Linux	49
3.7	Testing	55
4	Basics of Deep Learning	57
4.1	Introduction to Neural Networks	57
4.1.1	Weights	58
4.1.2	Neurons	60
4.1.3	Layers of Neurons	62
4.2	Designing Neural Networks	65
4.2.1	Training Neurons	67
4.2.2	Inference	75
4.3	Introduction to Caffe Framework	75
4.3.1	Parameters	75
4.3.2	Training Loop	75
4.4	Modeling of Neural Networks	76
4.4.1	Single neuron model	76
4.4.2	Tiny 3-layer model	76
4.4.3	Model for Numeral Classifier	77
4.5	Modeling of Deep Neural Networks	77
4.6	Modeling of Convolutional Neural Networks	77
4.7	Architectures	77
4.8	Current Trends	77
5	Understanding Caffe Architecture	79
5.1	OpenCL implementation	79

5.1.1	OpenCL Hybrid	80
5.2	Convolution Methods	81
5.3	The Python Interface	83
5.4	Useful Python Tools for Caffe	83
5.5	The Workflow of Caffe	83
5.6	The Caffe layers	83
6	Exploring and Preparing Data	85
6.1	Storage Formats	85
6.2	Loading Data with Python	85
6.3	Data Augmentation, Filtering and Balancing	85
7	Real-Time Inference	87
7.1	Classification	88
7.1.1	Introduction	88
7.1.2	BVLC GoogleNet	88
7.1.3	NIN-Imagenet	88
7.1.4	VGG CNN-S	88
7.1.5	VGG 16-Layer/19-Layer	88
7.2	Detection using Machine Learning	88
7.2.1	Introduction	88
7.2.2	Gaussian EM Segmentation	88
7.2.3	Salient Region Detection	88
7.2.4	Graph based Image Segmentation	88
7.2.5	Object Proposals using Filters	88
7.3	Detection using Deep Learning	88
7.3.1	Towards Real-Time Object Detection	88

7.3.2	CNN Model for Salient Object Subitizing	88
7.3.3	Fully Convolutional Networks for Semantic Segmentation	88
	U-Net Segmentation Architectures	88
	SK-Net Tricks for Segmentation with Existing Classifiers	88
8	Training using Transfer Learning	89
8.1	Understanding the trained data	89
8.2	Fixed feature extractors	89
8.3	Fine-tuning weights	89
9	Quantization and Compression of Networks	91
9.1	FP16, INT8 and Binary Quantization	91
10	Caffe in Production	93
10.1	Exploring notMNIST and Handwritten Alphabets	93
10.2	Identify Pets, Beer and Traffic Signs	93
10.3	Sound Classification for Explosion Detection	93
10.4	Generating Mountain Terrain, Paintings and Music	93
10.5	Human Action Recognition	93
10.6	Voice to Text Translation	93
10.7	Intelligent Chat-bot	93
11	Event and Stream Processing	95
11.1	Introduction to CEP	95
11.2	Using CEP with deep learning	95
12	References	97
	Index	101

Chapter 3

Installation and configuration

3.1 Introduction

Basic Linear Algebra Subprograms (BLAS) are the workhorses behind deep learning applications. They are low-level routines that provide the building blocks for performing common linear algebra operations such as, scalar multiplication, vector addition, dot products, linear combinations, and matrix multiplication.

BLAS libraries are highly optimized for specific hardware architectures. The routines in BLAS libraries consist of kernels written in assembly, or C. Kernels include a small-dimension matrix or vector operations also known as "**kernel operations**". The kernel operations are loops that break down large matrix multiplications into smaller sequences. The loops determine the overall performance of the matrix-matrix product since most of the time is spent there.

Large (2000x2000) matrix multiplication operations are broken into a sequence of smaller (100x100) matrix multiplications. These fixed-size small-dimension operations are hard coded and optimized using intrinsics and specific assembly code to improve performance. The specific assembly code makes use of special floating point units found in devices consisting of features like:

- SIMD-type instructions
- Instruction level parallelism

- Cache-awareness

Different levels of specifications for these kernel operations exist that define the BLAS implementations. The different levels also correlate to both the chronological order and publication, as well as the algorithm complexities.

Level 1: Set of linear algebra routines that operate on vectors only.

Level 2: Set of matrix-vector operations including, a generalized matrix-vector multiplication (**GEMV**).

Level 3: Set of matrix-matrix product operations including general matrix multiplication (**GEMM**).

3.2 BLAS libraries

BLAS libraries can be **tuned** and **fused**.

Tuning is the process where parameters are changed to improve the kernel performance according to the hardware.

Fusing combines two kernels to improve performance by exchanging data of the algorithms in cache or local memory instead of CPU or GPU global memory.

Many open-source BLAS libraries exist today, some of the common ones include the following:

Library	Description
ViennaCL	OpenCL Context + BLAS
ISAAC	Tuned OpenCL BLAS
clBLAS	AMD OpenCL BLAS
CLBlast	C. Nugteren OpenCL BLAS

Table 3.1: List of few BLAS libraries

Essentially all of the above are just BLAS libraries. ViennaCL BLAS is very basic: easy to install, but a bit older and slower than state-of-the-art BLAS

libraries. CLBlast has been added to favor i.e. ARM chips (experimental). clBLAS is preferentially used on high-end desktop GPUs by AMD and nVidia.

LibDNN is a fused im2col+GEMM convolution and pooling library with tuning and JIT/RTC support. It has an extensive collection of fast convolution implementations. It is the core of the OpenCL Caffe framework. It handles kernel generation, caching, loading and launching.

im2col and GEMM are usually two separate kernels, and LibDNN fuses them. im2col is an image to column transfer kernel function, which reformats a tensor to a 2D matrix to use GEMM on it. The reformatting of the tensor to a 2D matrix structure takes place within local GPU memory and GPU registers rather than global memory, giving a large speedup.

LibDNN uses its GEMM implementation, and im2col+GEMM fused kernels require an OpenCL kernel which has both the im2col and GEMM code in it. Specific high-performance OpenCL kernels are available for various devices. The Intel spatial kernels are one such breed of OpenCL kernels that are heavily optimized for Intel HD graphics.

Intel spatial kernels support auto tuning which takes care of tuning the parameters and selecting the fastest kernel within LibDNN for a particular hardware.

LibDNN has im2col+GEMM fused and implemented as JIT/RTC code so that the code complexity can be reduced based on known parameters of the convolution (stride, dilation, pad, kernel size, image sizes). im2col requires memory, and this memory is not needed when fusing GEMM+im2col.

3.3 OpenCL backend

Kernels written for GPUs are based on different ISA than x86, but are generally based on C.

In the OpenCL Caffe library, a versatile backend for various compute devices, based on OpenCL and ViennaCL, is available. The backend is called *Greentea* and is part of the *Project Greentea* consisting of frontend, models and modified Caffe library (see Figure 3.1). LibDNN is one of the many OpenCL toolsets that are part of the greentea project.

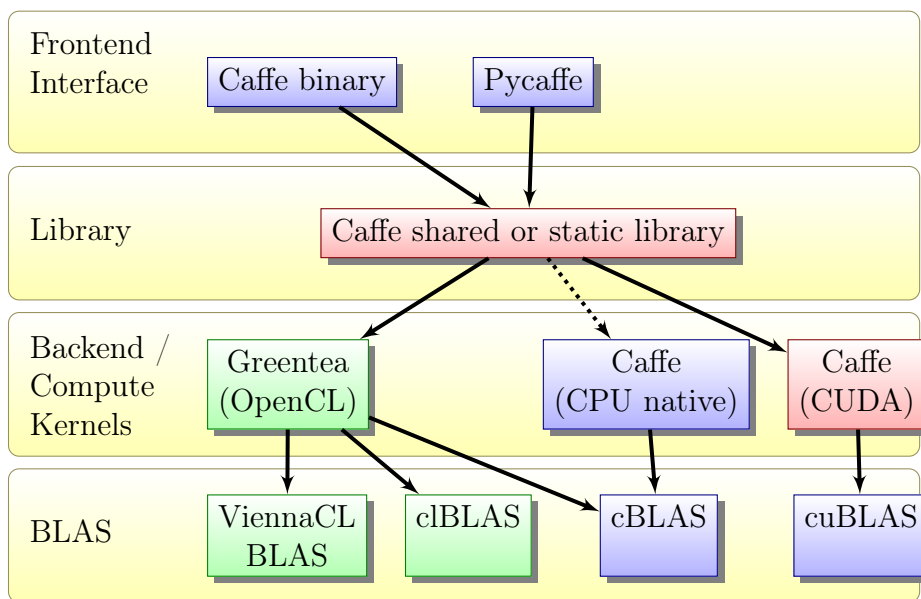


Figure 3.1: *Project Greentea* overview.

Green boxes denote entirely new additions to the Caffe landscape. Red boxes are parts that have been re-implemented or adapted from existing work to fit the needs of this project. Blue parts are mostly unchanged from existing work. Dashed arrows denote deprecated and only partially supported combinations.

Caffe supports CUDA backend. In CUDA Caffe, cuBLAS takes care of the BLAS, cuDNN does convolutions, and the rest is implemented as GPU kernels in CUDA code.

A key feature is that the OpenCL backend in Caffe is **feature equivalent** to the CUDA backend. All GPU layers can be used on both backends. The OpenCL backend is also unit test verified and passes all test cases of the original Caffe library. The tests can be invoked by executing "make runtest" on the source code folder.

It remains possible to compile the Caffe library with support for all backends at once. The compute kernel and BLAS calls can be dispatched dynamically at runtime, depending on what kind of device is selected.

Each of the devices (GPU/CPU) available in the hardware is registered in a new *DeviceContext* object that stores the device and back-end type.

```
$ caffe device_query

CUDA devices: 0
OpenCL devices: 2

Device id: 0
Device backend: OpenCL
Backend details: Intel(R) Corporation: OpenCL 1.2
Device vendor: Intel(R) Corporation
Name: Intel(R) HD Graphics 4400
Total global memory: 1708759450

Device id: 1
Device backend: OpenCL
Backend details: Intel(R) Corporation: OpenCL 1.2
Device vendor: Intel(R) Corporation
Intel(R) Core(TM) i5-4210U CPU @ 1.70GHz
Total global memory: 8513761280
```

After running the above command, the Caffe library enumerates all devices on all enabled backends, starting with CUDA devices and displays them. The selected GPU number determines which *DeviceContext* is set as the default.

3.4 Installation prerequisites

3.4.1 Standard requirements

There are standard requirements for installing/compiling Caffe framework.

- BLAS: Caffe framework requires BLAS for performing backend operations such as matrix and vector computations. BLAS implementations are provided by Intel Math Kernel Library (Intel MKL), ATLAS, open-BLAS, ViennaCL and so forth.
- Boost: Caffe framework makes use of the extensive collection of high-quality C++ libraries available from Boost for its math and shared pointer functions.

- leveldb, lmdb, HDF5: They are fast, high-performance storage libraries used by Caffe framework for preparing, storing and managing data.
- protobuf: Caffe framework uses the protobuf library to define the structure and format of data.
- glog, gflags: Caffe framework uses them as essential utilities for logging and for debugging.

3.4.2 OpenCL requirements

LibDNN DNN is integrated into OpenCL Caffe framework, and separate installation is not required. There are two types of **OpenCL modes** when compiling and using the Caffe framework.

- CPU + OpenCL: Faster than traditional Caffe (CPU_ONLY compile) mode but slower than GPU + OpenCL mode. BLAS can still be OpenBLAS and not an OpenCL BLAS.
- GPU + OpenCL: Faster than CPU+OpenCL mode and choice of BLAS does not matter.

The second mode is preferred for real-time performance. The default convolution engine for OpenCL GPUs is LibDNN and is a replacement for cuDNN. There are also new Intel kernels for Intel GPUs available and enabled by default.

Intel spatial kernels included in OpenCL Caffe do not support efficient back propagation and not all shapes of convolutions, but yes, they are the fastest forward propagation on Intel iGPUs.

ISAAC, clBLAS, and CLBlast are strictly optional and can be compiled if required.

3.5 Device specific considerations

INCOMPLETE (this section will show how to get higher performance with example configurations of the Makefile. It is reserved for future, and depends if FP16 code can make it by July/August because in that case, LibDNN will have it is own BLAS as there is currently no good FP16 BLAS)

3.6 Compilation

An Intel device with HD graphics 515 with up to 384 gigaflops of GPU performance is used to demonstrate the compilation process. However, any new device that has Intel HD graphics or discrete AMD graphics can be utilized.

3.6.1 Windows

Installation in Windows is fairly simple and an easy process.

Step 1: General

Visual Studio 2015 x64 community edition or higher is needed along with the latest versions of Intel OpenCL SDK and Intel OpenCL run-time. The graphics driver update is done using the Intel update utility.

Install Anaconda version of Python distribution as well Git for Windows. Make sure the latest 64-bit version of CMake for Windows is installed. Create a folder called "projects." Clone or download the OpenCL Caffe framework from URL <https://github.com/bvlc/caffe/tree/opencv>.

```
C:\> cd projects
C:\Projects> git clone -b opencv https://github.com/BVLC/caffe
C:\Projects> git clone https://github.com/viennacl/viennacl-dev.git
C:\Projects> cd caffe
```

Set the following variables in build_win.cmd file.

```
C:\projects\caffe\scripts\build_win.cmd.

if NOT DEFINED WITH_NINJA set WITH_NINJA=0
if NOT DEFINED PYTHON_VERSION set PYTHON_VERSION=3

C:\Projects\caffe> scripts\build_win.cmd
```

3.6.2 Linux

Step 1: General OS

Use the URL <http://cdimage.ubuntu.com/daily-live/current/> and install Ubuntu 17.0.4 64-bit (at the time of writing).

```
$ cd /home/oreilly
$ sudo su
$ apt-get update
$ apt-get install git
$ apt-get install xz-utils

$ updatedb; locate libicui18n
$ apt-get install dkms
$ apt-get install mono-devel
$ apt-get install lsb-core

$ add-apt-repository ppa:webupd8team/java
$ apt-get update; apt-get install oracle-java8-installer
$ apt-get install oracle-java8-set-default
```

Step 2: Intel OpenCL SDK and runtime

Download OpenCL SDK from the URL <https://software.intel.com/en-us/intel-opencl/download> and run the following commands.

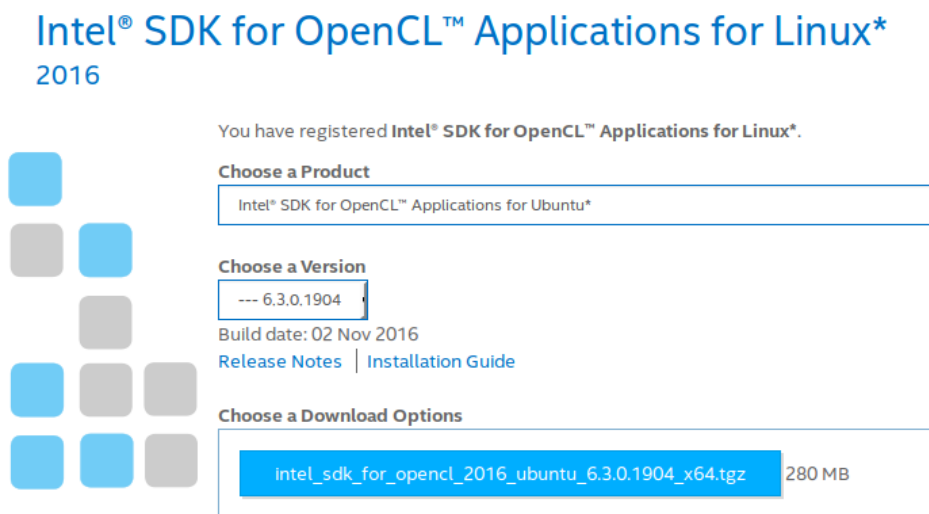


Figure 3.2: opencl.

```
$ tar -xvzf /home/oreilly/intel_sdk_for_openc1_2016_ubuntu_6.3.0.1904_x64.tg
$ cd /home/oreilly/intel_sdk_for_openc1_2016_ubuntu_6.3.0.1904_x64
$ ./install_GUI.sh
```

You can ignore missing prerequisite files `libcui18n.so.52()(64bit)`, and `libcuc.so.52()(64bit)` as newer versions are already present.

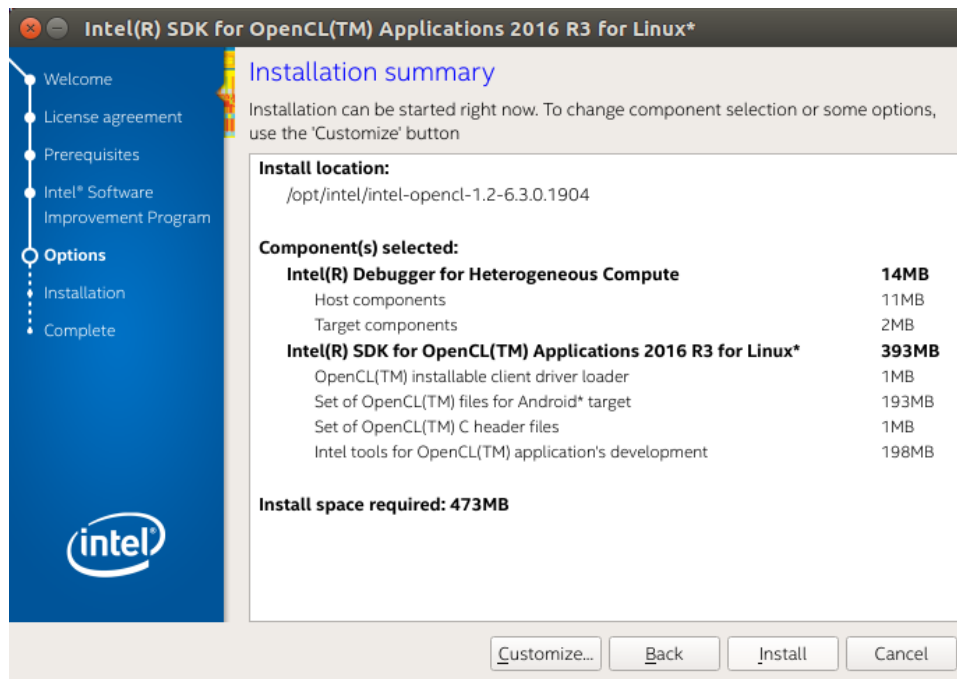


Figure 3.3: install.

Download OpenCL runtime from URL <https://software.intel.com/en-us/articles/openc1-drivers> and run the following commands.

```
$ unzip /home/oreilly/SRB4.1_linux64.zip
$ cd /home/oreilly/
$ tar -xJf intel-openc1-cpu-r4.1-61547.x86_64.tar.xz
$ tar -xJf intel-openc1-devel-r4.1-61547.x86_64.tar.xz
$ tar -xJf intel-openc1-r4.1-61547.x86_64.tar.xz
```

The extracted `opt`, etc, directories need to be copied to the root of the file system, and `clinfo` utility is required to query and display OpenCL information.

```
$ cp -avr opt/ /
$ cp -avr etc/ /

$ apt-get install cmake-qt-gui
$ apt-get install clinfo vim
```

Step 3: Clone ViennaCL, OpenCL Caffe

Clone the OpenCL branch of Caffe framework, and clone or download the latest ViennaCL from the URL <http://viennacl.sourforge.net/viennacl-download.html> and extract.

```
$ cd /home/oreilly
$ git clone https://github.com/viennacl/viennacl-dev.git
$ cd /home/oreilly/viennacl-dev
$ mkdir build

$ cd /home/oreilly/viennacl-dev/build
$ cmake -DBUILD_TESTING=OFF -DBUILD_EXAMPLES=OFF
-DMAKE_INSTALL_PREFIX=$HOME/local
-DOPENCL_LIBRARY=/opt/intel/opencl/libOpenCL.so ..

$ make -j8
$ make install
$ git clone -b opencl https://github.com/BVLC/caffe
```

Step 5: Install dependencies

```
$ apt-get install libprotobuf-dev
$ apt-get install libleveldb-dev

$ apt-get install libsnappy-dev
$ apt-get install libhdf5-serial-dev

$ apt-get install protobuf-compiler
$ apt-get install --no-install-recommends libboost-all-dev

$ apt-get install libgflags-dev
$ apt-get install libgoogle-glog-dev

$ apt-get install liblmdb-dev
$ apt-get install libopenblas-dev
```

```
$ apt-get install -y cython
$ apt-get install python3-dev python3-skimage
$ apt-get install python3-pip
$ pip3 install pyzmq jsonschema pillow numpy scipy ipython jupyter pyyaml
```

```
$ apt-get install build-essential
$ apt-get install cmake git libgtk2.0-dev pkg-config
$ apt-get install libavcodec-dev libavformat-dev libswscale-dev
$ apt-get install libtbb2 libtbb-dev libjpeg-dev
$ apt-get install libpng-dev libtiff5-dev libdc1394-22-dev
$ apt-get install libjasper-dev
$ apt-get install doxygen
```

Step 6: Compile OpenCV 3

Compile and install OpenCV 3.

```
$ cd /home/oreilly
$ git clone https://github.com/opencv/opencv.git
$ git clone https://github.com/opencv/opencv_contrib.git
$ cd /home/oreilly/opencv
$ mkdir build
$ cmake-gui
```

The source code and build folders must be filled as shown in figure 3.4 and upon pressing the configure button, a pop-up window enables the selection of Unix Makefiles as the generator for the project.

Post configure, the OpenCL and OpenCV paths must be set, and the generate button generates the makefiles. The make command compiles the OpenCV package.

```
$ cd /home/oreilly/opencv/build
$ make -j8
$ make install
$ ldconfig
```

Make install command will copy the compiled libraries to system directories such as:


```

USE_GREENTEA := 1
DISABLE_DEVICE_HOST_UNIFIED_MEMORY := 0
USE_LIBDNN := 1
USE_INTEL_SPATIAL := 1
VIENNAACL_DIR = ../viennacl-dev
OPENCV_VERSION := 3
BLAS := open
PYTHON_LIBRARIES := boost_python3 python3.5m
PYTHON_INCLUDE := /usr/include/python3.5m \
    /usr/lib/pthong3.5/dist-packages/numpy/core/include
WITH_PYTHON_LAYER := 1
INCLUDE_DIRS := $(PYTHON_INCLUDE) /usr/local/include \
    /usr/include/hdf5/serial

```

Determine the correction version present for boost python and libhdf5 to symlink taking into account their respective versions.

```

$ ln -s /usr/lib/x86_64-linux-gnu/libboost_python-py35.so
  /usr/lib/x86_64-linux-gnu/libboost_python3.so
$ ln -s /usr/lib/x86_64-linux-gnu/libhdf5_serial.so.100
  /usr/lib/x86_64-linux-gnu/libhdf5.so
$ ln -s /usr/lib/x86_64-linux-gnu/libhdf5_serial_hl.so.100
  /usr/lib/x86_64-linux-gnu/libhdf5_hl.so

```

```

$ make all -j8
$ make pycaffe -j8
$ make runtest -j8
$ make install

```

The Caffe module for Python `_caffe.so` should now be available in the directory `/home/oreilly/caffe/python/caffe`.

The entire caffe folder is itself the module. It can just be added to the python path or added to the python script itself.

```
$ echo $PYTHONPATH=/home/oreilly/caffe/python:$PYTHONPATH
```

3.7 Testing

```
$ /home/oreilly/caffe/tools/caffe device_query -gpu all
```