

CertusPro-NX MobileNet Object Classification on VVML Board

Reference Design



Disclaimers

Lattice makes no warranty, representation, or guarantee regarding the accuracy of information contained in this document or the suitability of its products for any particular purpose. All information herein is provided AS IS, with all faults and associated risk the responsibility entirely of the Buyer. Buyer shall not rely on any data and performance specifications or parameters provided herein. Products sold by Lattice have been subject to limited testing and it is the Buyer's responsibility to independently determine the suitability of any products and to test and verify the same. No Lattice products should be used in conjunction with mission- or safety-critical or any other application in which the failure of Lattice's product could create a situation where personal injury, death, severe property or environmental damage may occur. The information provided in this document is proprietary to Lattice Semiconductor, and Lattice reserves the right to make any changes to the information in this document or to any products at any time without notice.



Contents

Acronyms in This Document	8
1. Introduction	9
1.1. Design Process Overview	9
2. Setting Up the Basic Environment	10
2.1. Software and Hardware Requirements	10
2.1.1. Lattice Software	10
2.1.2. Hardware	10
2.2. Setting Up the Linux Environment for Machine Training	10
2.2.1. Installing the CUDA Toolkit	
2.2.2. Installing the cuDNN	
2.2.3. Installing Anaconda and Python 3	
3. Preparing the Dataset	
3.1. Preparing Dataset for Training	
3.2. Generating Dataset from Board Camera	
3.3. Generating Video for Demo	
4. Training the Machine	
4.1. Training Code Structure	
4.2. Neural Network Architecture	
4.2.1. Object Detection Training Network Layers	
4.2.2. Object Detection Network Output	
4.2.3. Training Code Overview	
4.2.3.1. Model Configuration	
4.2.3.2. Model Building	
4.2.3.3. Training	
4.3. Training from Scratch and/or Transfer Learning	
5. Testing the Trained Model	
5.1. Running the Test Shell Script	
5.2. Checking the Mean Average Precision (mAP) of the Model	
5.3. Testing the Centroid Tracking Algorithm	
5.3.1. Testing the Images Captured using the Board Camera	
5.3.2. Live Testing of the Centroid Tracking using Real-time Data from the Board	
6. Creating Frozen File	
6.1. Generating the Frozen .pb File	
7. Creating Binary File with Lattice sensAl	
Hardware Implementation	
8.1. Top Level Information	
8.1.1. Block Diagram	
•	
8.1.2. Operational Flow	
8.2. Architecture Details	
8.2.1. SPI Flash Operation	
8.2.2. Pre-processing CNN	
8.2.2.1. Pre-processing Flow	
8.2.3. HyperRAM Operations	
8.2.4. Post-processing CNN	
8.2.4.1. Confidence Sorting	
8.2.4.2. Bounding Box Calculation	
8.2.4.3. NMS – Non Max Suppression	
8.2.4.4. Bounding Box Upscaling	
8.2.4.5. Centroid Tracker	
8.2.4.6. OSD Text Display	
8.2.4.7. USB Wrapper	52



9. Creating FPGA Bitstream File	53
9.1. Bitstream Generation Using Lattice Radiant Software	53
9.2. Configuring the IP in Lattice Radiant Software	56
10. Programming the Demo	
10.1. Programming the CertusPro-NX SPI Flash	
10.1.1. Erasing the CertusPro-NX SRAM Prior to Reprogramming	
10.1.2. Programming the CertusPro-NX Board	59
10.1.3. Programming sensAl Firmware Binary to the CertusPro-NX SPI Flash	
10.1.3.1. Convert Flash sensAl Firmware Hex to CertusPro-NX SPI Flash	
11. Running the Demo	64
Appendix A. Other Labelling Tools	
References	66
Technical Support Assistance	67
Revision History	68



Figures

Figure 1.1. Lattice Machine Learning Design Flow	9
Figure 2.1. Lattice CertusPro-NX Voice and Vision Machine Learning Board	10
Figure 2.2. CUDA Repo Download	11
Figure 2.3. CUDA Repo Installation	11
Figure 2.4. Fetch Keys	11
Figure 2.5. Update Ubuntu Packages Repositories	11
Figure 2.6. CUDA Installation Completed	
Figure 2.7. cuDNN Library Installation	12
Figure 2.8. Anaconda Installation	
Figure 2.9. Accept License Terms	12
Figure 2.10. Confirm/Edit Installation Location 6	
Figure 2.11. Launch/Initialize Anaconda Environment on Installation Completion	
Figure 3.1. Directory After Unzipping	
Figure 3.2. CUDA Version 10.0	
Figure 3.3. Directory Path given as Input	
Figure 3.4. Object Name to be Labeled	
Figure 3.5. Selected Object	
Figure 3.6. Text Label in the _Annotations_KITTI Directory	
Figure 3.7. Generated Media Directory	
Figure 4.1. Training Code Directory Structure	
Figure 4.2. Training Code Flow Diagram	
Figure 4.3. Code Snippet – Input Image Size Config	
Figure 4.4. Code Snippet – Anchors Per Grid Config #1 (Grid Sizes)	
Figure 4.5. Code Snippet – Anchors Per Grid Config #2	
Figure 4.6. Code Snippet – Anchors Per Grid Config #3	
Figure 4.7. Code Snippet – Training Parameters	
Figure 4.8. Code Snippet – Filter Values	
Figure 4.9. Code Snippet – Forward Graph Fire Layers for MobileNet v1	25
Figure 4.10. Code Snippet – Forward Graph Last Convolution Layer	
Figure 4.11. Grid Output Visualization #1	
Figure 4.12. Grid Output Visualization #2	
Figure 4.13. Code Snippet – Interpret Output Graph	
Figure 4.14. Code Snippet – Bbox Loss	
Figure 4.15. Code Snippet – Confidence Loss	
Figure 4.16. Code Snippet – Class Loss	
Figure 4.17. Code Snippet – Training	
Figure 4.18. Training Code Snippet for Mean and Scale	
Figure 4.19. Training Code Snippet for Dataset Path	
Figure 4.20. Training Input Parameter	
Figure 4.21. Execute train.sh Script	
Figure 4.22. TensorBoard	
Figure 4.23. Example of Checkpoint Data Files at weights/Object_Detector/train_weights/ folder	
Figure 5.1. Test Script Last Log	
Figure 5.2. Test Result in images_out_v1 Directory	
Figure 5.3. Running map.py	
Figure 5.4. map.png in MAP/results directory	
Figure 5.5. crop.py Terminal Output	
Figure 5.6. Input Parameter for the Script	
Figure 5.7. Output Window with Overlay on Images	
Figure 5.8. Object Detector on PC	
Figure 6.1. pb File Generation from Checkpoint	
Figure 6.2. Frozen .pb File	
· · · · · · · · · · · · · · · · · · ·	



Figure 7.1. sensAl Home Screen	37
Figure 7.2. sensAl – Network File Selection	38
Figure 7.3. sensAl – Input Image File Selection	38
Figure 7.4. sensAl – Project Settings	39
Figure 7.5. sensAl – Analyze Project	39
Figure 7.6. Q Format Settings for Each Layer	40
Figure 7.7. Compile Project	40
Figure 8.1. RTL Top Level Block Diagram	41
Figure 8.2. SPI Read Command Sequence	43
Figure 8.3. Masking	44
Figure 8.4. Downscaling	45
Figure 8.5. HyperRAM Memory Addressing	46
Figure 8.6. HyperRAM Access Block Diagram	47
Figure 8.7. CNN Output Data Format	48
Figure 8.8. Confidence Sorting	49
Figure 8.9. Intersection-Union Area NMS	50
Figure 9.1. Radiant – Default Screen	
Figure 9.2. Radiant – Open CertusPro-NX Project File (.rdf)	54
Figure 9.3. Radiant – Design Load Check After Opening the Project File	54
Figure 9.4. Radiant – Trigger Bitstream Generation	55
Figure 9.5. Radiant – Radiant Bit File Generation Report Window	55
Figure 9.6. Radiant – Uninstall Old IP	
Figure 9.7. Radiant – Install New IP	56
Figure 9.8. Radiant – Select User IP Package to Install	57
Figure 10.1. Radiant Programmer Default Screen	58
Figure 10.2. Radiant programmer- Device selection	58
Figure 10.3. Radiant Programmer – Device Operation	59
Figure 10.4. Radiant Programmer – Selecting Device Properties Options for CertusPro-NX Flashing	60
Figure 10.5. CertusPro-NX Flashing Switch – SW4 Push Button	61
Figure 10.6. Radiant Programmer – Output Console	61
Figure 10.7. Radiant Programmer – Selecting Device Properties Options for CertusPro-NX Flashing	62
Figure 10.8. Radiant Programmer – Output Console	63
Figure 11.1. Running the Demo	64



Tables

Table 4.1. Object Detection MobileNet v1 Training Topology	19
Table 8.1. Core Parameter	
Table 8.2. Data Parameters of CNN Output	47
Table 8.3. Pre-Selected Width and Height of Anchor Boxes	
Table 8.4. Grid Center Values (X, Y) for Anchor Boxes	49
Table A.1. Other Labelling Tools	65



Acronyms in This Document

A list of acronyms used in this document.

Acronym	Definition	
AXI	Advanced Extensible Interface	
CNN	Convolutional Neural Network	
FPGA	Field-Programmable Gate Array	
NN	Neural Network	
SD	Secure Digital	
SPI	Serial Peripheral Interface	
SRAM	Static Random Access Memory	
VVML	Voice and Vision Machine Learning	



1. Introduction

This document describes the Object detection Design process using the CertusPro™-NX platform.

Generic Object Classification base design is used. This document covers MobileNet v1 version of the reference design.

1.1. Design Process Overview

The design process involves the following steps:

- 1. Training the model
 - Setting up the basic environment
 - Preparing the dataset
 - Training the machine
 - Creating the checkpoint data
 - Creating the frozen file (*.pb)
- 2. Compiling Neural Network
 - Creating the binary file with Lattice Neural Network Compiler Software 4.0 program
- 3. FPGA Design
 - Creating the FPGA bitstream file
- 4. FPGA Bitstream and Quantized Weights and Instructions
 - Flashing the binary and bitstream files
 - Binary File to Flash Memory on CertusPro-NX board
 - Bitstream to Flash Memory on CertusPro-NX board

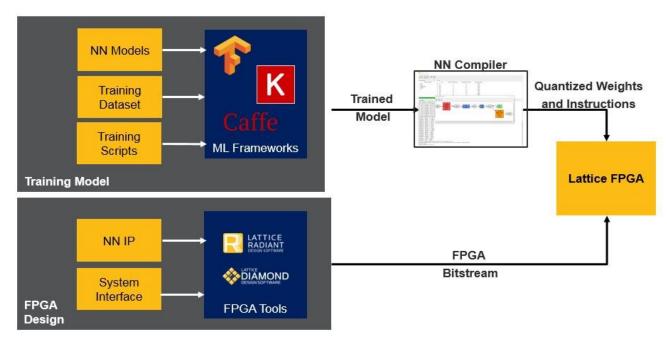


Figure 1.1. Lattice Machine Learning Design Flow



2. Setting Up the Basic Environment

2.1. Software and Hardware Requirements

This section describes the required tools and environment setup for the FPGA bitstream and flashing.

2.1.1. Lattice Software

- Lattice Radiant™ Tool Refer to http://www.latticesemi.com/LatticeRadiant.
- Lattice Radiant Programmer Refer to http://www.latticesemi.com/programmer.
- Lattice Neural Network Compiler Software 4.0

 Refer to

 https://www.latticesemi.com/Products/DesignSoftwareAndIP/AIML/NeuralNetworkCompiler.

2.1.2. Hardware

This design uses the CertusPro-NX Voice and Vision Machine Learning Board as shown in Figure 2.1.

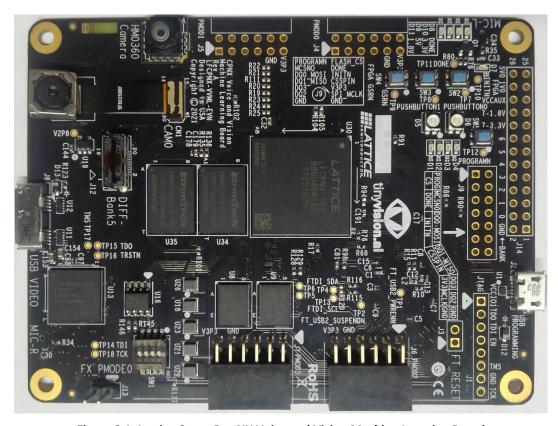


Figure 2.1. Lattice CertusPro-NX Voice and Vision Machine Learning Board

2.2. Setting Up the Linux Environment for Machine Training

2.2.1. Installing the CUDA Toolkit

To install the CUDA toolkit, run the following commands in the order specified below:

\$ curl -0

 $\label{lem:https://developer.download.nvidia.com/compute/cuda/repos/ubuntu1604/x86_64/cudarepo-ubuntu1604_10.1.105-1_amd64.deb$

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



```
$ curl -0 https://developer.download.nvidla.com/compute/cuda/repos/ubuntu1604/x86_64/cuda-repo-ubuntu1604_10.1.105-1_amd64.deb

% Total % Received % Xferd Average Speed Time Time Time Current

Dload Upload Total Spent Left Speed

100 2832 100 2832 0 0 2204 0 0:00:01 0:00:01 -:--;-- 2205
```

Figure 2.2. CUDA Repo Download

\$ sudo dpkg -I ./cuda-repo-ubuntu1604 10.1.105-1 amd64.deb

```
$ sudo dpkg -i ./cuda-repo-ubuntu1604_10.1.105-1_amd64.deb
Selecting previously unselected package cuda-repo-ubuntu1604.
(Reading database ... 5287 files and directories currently installed.)
Preparing to unpack .../cuda-repo-ubuntu1604_10.1.105-1_amd64.deb ...
Unpacking cuda-repo-ubuntu1604 (10.1.105-1) ...
Setting up cuda-repo-ubuntu1604 (10.1.105-1) ...
The public CUDA GPG key does not appear to be installed.
To install the key, run this command:
sudo apt-key adv --fetch-keys http://developer.download.nvidia.com/compute/cuda/repos/ubuntu1604/x86_64/7fa2af80.pub
```

Figure 2.3. CUDA Repo Installation

\$ sudo apt-key adv --fetch-keys
http://developer.download.nvidia.com/compute/cuda/repos/ubuntu1604/x86_64/7fa2af80.
pub

Figure 2.4. Fetch Keys

\$ sudo apt-get update

FPGA-RD-02242-1 0

```
$ sudo apt-get update
Ign:1 http://developer.download.nvidia.com/compute/cuda/repos/ubuntu1604/x86_64 InRelease
Hit:2 http://archive.ubuntu.com/ubuntu xenial InRelease
Get:3 http://developer.download.nvidia.com/compute/cuda/repos/ubuntu1604/x86_64 Release [697 B]
Get:4 http://developer.download.nvidia.com/compute/cuda/repos/ubuntu1604/x86_64 Release.gpg [836 B]
Hit:5 http://archive.ubuntu.com/ubuntu xenial-updates InRelease
Hit:6 http://security.ubuntu.com/ubuntu xenial-security InRelease
Hit:7 http://archive.ubuntu.com/ubuntu xenial-backports InRelease
Ign:8 http://developer.download.nvidia.com/compute/cuda/repos/ubuntu1604/x86_64 Packages
Get:8 http://developer.download.nvidia.com/compute/cuda/repos/ubuntu1604/x86_64 Packages [428 kB]
Fetched 429 kB in 1s (386 kB/s)
Reading package lists... Done
```

Figure 2.5. Update Ubuntu Packages Repositories

\$ sudo apt-get install cuda-10-0

```
$ sudo apt-get install cuda-9-0
Reading package lists... Done
Building dependency tree
Reading state information... Done
```

Figure 2.6. CUDA Installation Completed

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



2.2.2. Installing the cuDNN

To install the cuDNN:

- 1. Create an NVIDIA developer account in https://developer.nvidia.com.
- Download cuDNN lib in https://developer.nvidia.com/compute/machinelearning/cudnn/secure/v7.1.4/prod/9.0 20180516/cudnn-9.0-linux-x64-v7.1.
- 3. Execute the commands below to install cuDNN.

```
$ tar xvf cudnn-9.0-linux-x64-v7.1.tgz
$ sudo cp cuda/include/cudnn.h /usr/local/cuda/include
$ sudo cp cuda/lib64/libcudnn* /usr/local/cuda/lib64
$ sudo chmod a+r /usr/local/cuda/include/cudnn.h /usr/local/cuda/lib64/libcudnn*
```

```
$ tar xcf cudnn-9.0-linux-x64-v7.1.tgz
cuda/include/cudnn.h
cuda/NVIDIA_SLA_cuDNN_Support.txt
cuda/lib64/libcudnn.so
cuda/lib64/libcudnn.so.7
cuda/lib64/libcudnn.so.7.1.4
cuda/lib64/libcudnn_static.a
$ sudo cp cuda/include/cudnn.h /usr/local/cuda/include
$ sudo cp cuda/lib64/libcudnn* /usr/local/cuda/lib64
$ sudo chmod a+r /usr/local/cuda/include/cudnn.h /usr/local/cuda/lib64
```

Figure 2.7. cuDNN Library Installation

2.2.3. Installing Anaconda and Python 3

To install Anaconda and Python 3:

- 1. Go to https://www.anaconda.com/products/individual#download-section.
- 2. Download Python 3 version of Anaconda for Linux.
- 3. Install the Anaconda environment by running the command below:

```
$ sh Anaconda3-2019.03-Linux-x86 64.sh
```

Note: Anaconda3-<version>-Linux-x86_64.sh version may vary based on the release.

```
(base) $ sh Anaconda3-2020.07-Linux-x86_64.sh

Welcome to Anaconda3 2020.07

In order to continue the installation process, please review the license agreement.

Please, press ENTER to continue

>>>
```

Figure 2.8. Anaconda Installation

4. Accept the license.

```
Do you accept the license terms? [yes|no] [no] >>> yes
```

Figure 2.9. Accept License Terms

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



5. Confirm the installation path. Follow the instruction onscreen if you want to change the default path.

```
[no] >>> yes
Anaconda3 will now be installed into this location:
/home/user/anaconda3
   - Press ENTER to confirm the location
   - Press CTRL-C to abort the installation
   - Or specify a different location below
[/home/user/anaconda3] >>> /home/user/anaconda3
```

Figure 2.10. Confirm/Edit Installation Location 6

6. After installation, enter **No**, as shown in Figure 2.11.

```
Preparing transaction: done
Executing transaction: done
installation finished.
Do you wish the installer to initialize Anaconda3
by running conda init? [yes|no]
[no] >>> no
```

Figure 2.11. Launch/Initialize Anaconda Environment on Installation Completion



3. Preparing the Dataset

This section describes how to create a synthetic dataset using python script.

The dataset generated has five classes. The Python script generates the images of resolution 1920×1080p and the labels are also generated in parallel.

3.1. Preparing Dataset for Training

To prepare the dataset:

1. Unzip the package by running the following commands:

```
$ unzip synthetic_mobilenetV1.zip
$ cd Synthetic data test
```

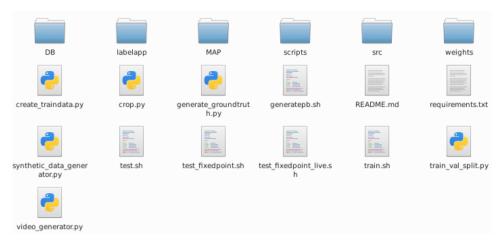


Figure 3.1. Directory After Unzipping

2. Run the script below to create a virtual python environment:

```
$ virtualenv lattice_env --python=python3
```

3. Run the script below to activate the environment:

```
$ source lattice env/bin/activate
```

Ensure that CUDA version is 10.0.

\$ nvcc -V

```
(lattice_env) sidhanatht@gamora:~/projects/lattice/Package_test/synthetic_data_object_detector-dev1$ nvcc -V
nvcc: NVIDIA (R) Cuda compiler driver
Copyright (c) 2005-2018 NVIDIA Corporation
Built on Sat_Aug_25_21:08:01_CDT_2018
Cuda compilation tools, release 10.0, V10.0.130
```

Figure 3.2. CUDA Version 10.0

4. Run the script below to install the required packages:

```
$ pip install -r requirements.txt
```

Synthetic data generation takes 5 to 10 minutes to complete. Image count can be given as an argument as shown in the script below.

```
$ python synthetic datagenerator.py 5000
```

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.
All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



5. Split the generated dataset into train and validation set by running the script below:

```
$ python train val split.py
```

6. Crop the images and labels from 1920×1080p to 896×896p by running the script below. This also creates the val and train text files in the ImageSet directory and takes a few minutes to complete.

```
$ python src/crop data.py
```

3.2. Generating Dataset from Board Camera

After capturing the images from the board camera, use the annotate-to-KITTI software to label the data.

To generate the dataset from the board camera:

- 1. Copy folder containing captured images to labelapp folder
- 2. Run annotate-folder.py followed by directory name

```
$ cd labelapp/
$ python annotate-folder.py
```

3. Provide the directory path containing the captured images as shown in Figure 3.3Figure 3.3.

```
(lattice_env) sidhanatht@gamora:~/projects/lattice/synthetic/Synthetic_data_test
/labelapp$ python annotate-folder.py
Enter the path to dataset: /home/sidhanatht/projects/lattice/synthetic/Synthetic
_data_test/labelapp/captured_images/
```

Figure 3.3. Directory Path given as Input

4. Provide the label name as shown in Figure 3.4. The image from the directory opens in a window.

```
Enter the path to dataset: /home/sidhanatht/Lattice_release/labelapp/capture_images
/home/sidhanatht/Lattice_release/labelapp/capture_images_Images_KITTI
/home/sidhanatht/Lattice_release/labelapp/capture_images_Annotations_KITTI
Enter default object label: circle
```

Figure 3.4. Object Name to be Labeled

5. Left-click and drag to draw the box around the area of interest.



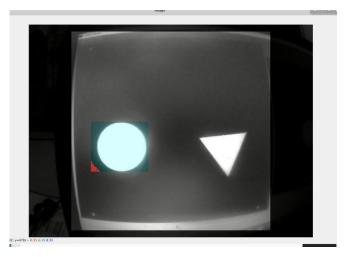


Figure 3.5. Selected Object

6. Press **Q** to save and exit. The Label file is generated in the _Annotations_KITTI directory.

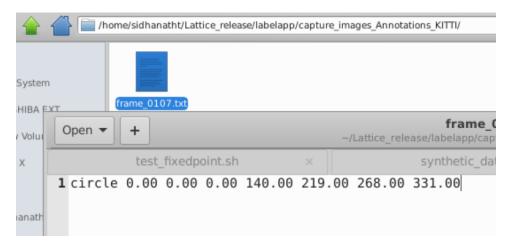


Figure 3.6. Text Label in the _Annotations_KITTI Directory

For more detailed description of the annotate-to-KITTI, go to https://github.com/SaiPrajwal95/annotate-to-KITTI.



3.3. Generating Video for Demo

To generate a video for demo:

- 1. Run the following script in the same virtual environment to generate the video.
 - \$ python video generator.py
- 2. The media folder is generated and contains the items shown in Figure 3.7.



Figure 3.7. Generated Media Directory



4. Training the Machine

4.1. Training Code Structure

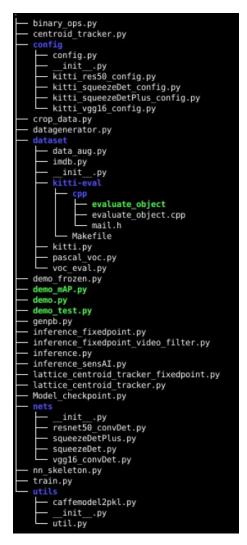


Figure 4.1. Training Code Directory Structure

4.2. Neural Network Architecture

4.2.1. Object Detection Training Network Layers

This section provides information on the Convolution Network Configuration of the Object Detection design. Table 4.1 shows the MobileNet v1 structure.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



Table 4.1. Object Detection MobileNet v1 Training Topology

T		out (224 × 224 × 1)
	Conv3-32	Conv3 – # where:
Fire 1	BN	• Conv3 = 3 × 3 Convolution filter Kernel size
	ReLU	• # = The number of filter
	MaxPool	DWConv3-32 - # where:
	DWConv3-32	 DWConv3 = Depthwise convolution filter with 3 × 3 size # = The number of filter
	BN	Conv1–32 – # where:
Fire 2	ReLU	• Conv1 = 1 × 1 Convolution filter Kernel size
	Conv1-24	# = The number of filter
	BN	
	ReLU	For example, Conv3–16 = 16 3 × 3 convolution filters
	DWConv3-24	
	BN	BN – Batch Normalization
Fine 2	ReLU	
Fire 3	Conv1-24	
	BN	
	ReLU	
	DWConv3-24	
	BN	
	ReLU	
Fire 4	MaxPool	
	Conv1–48	
	BN	
	ReLU	
	DWConv3-48	
	BN	
	ReLU	
Fire 5	Conv1–56	
	BN	
	ReLU	
	DWConv3-56	
	BN	
	ReLU	
Fire 6	MaxPool	
	Conv1-100	
_	BN	
	ReLU	
Fire 7	DWConv3-100	
	BN	
	ReLU	
	Conv1-104	
	BN	
	ReLU	
Conv12	Conv3-42	



Conv3-#where:

- Conv3 = 3 × 3 Convolution filter Kernel size
- # = The number of filter

DWConv3-32-# where:

DWConv3 = Depthwise convolution filter with 3 × 3 size # = The number of filter

Conv1-32-# where:

- Conv1 = 1 × 1 Convolution filter Kernel size
- # = The number of filter

For example, Conv3–16 = 16.3×3 convolution filters.

- The Object Detection network structure consists of seven fire layers followed by one convolution layer. A fire layer contains Convolutional, Batch Normalization, and ReLU (Rectified Linear Unit). Pooling layers are only in Fire 1, Fire 3, Fire 5, and Fire 7. Fire 4, and Fire 6 do not contain pooling layers.
- In Table 4.1, the layer contains Convolution (Conv), Batch Normalization (BN), and ReLU layers.
- Layer information:
 - Convolutional Layer

In general, the first layer in a CNN is always a convolutional layer. Each layer consists of a number of filters (sometimes referred as kernels), which convolves with the input layer/image and generates an activation map (that is, feature map). This filter is an array of numbers (called weights or parameters). Each of these filters can be thought of as feature identifiers, such as straight edges, simple colors, curves, and other high-level features. For example, the filters on the first layer convolve around the input image and *activate* (or compute high values) when the specific feature it is looking for (such as curve) is in the input volume.

• ReLU (Activation Layer)

It is the convention to apply a nonlinear layer (or activation layer) immediately after each conv layer. The purpose of this layer is to introduce nonlinearity to a system that is basically computing linear operations during the conv layers (element wise multiplications and summations). In the past, nonlinear functions such as tanh and sigmoid were used, but researchers found out that ReLU layers work far better because the network is able to train a lot faster (because of the computational efficiency) without making a significant difference in accuracy. The ReLU layer applies the function f(x) = max(0, x) to all of the values in the input volume. In basic terms, this layer changes all the negative activations to 0. This layer increases the nonlinear properties of the model and the overall network without affecting the receptive fields of the conv layer.

Pooling Layer

After some ReLU layers, you may choose to apply a pooling layer. It is also referred to as a down sampling layer. In this category, there are also several layer options, with MaxPooling being the most popular. This basically takes a filter (normally of size 2×2) and a stride of the same length. It then applies a filter to the input volume and outputs the maximum number in every sub region that the filter convolves around.

The intuitive reason behind this layer is that once it is known that a specific feature is in the original input volume (there is a high activation value), its exact location is not as important as its relative location to the other features. As you can imagine, this layer drastically reduces the spatial dimension (the length and the width change but not the depth) of the input volume. This serves two main purposes. The first is that the number of parameters or weights is reduced by 75%, thus lessening the computation cost. The second is that it controls over fitting. This term is used when a model is so tuned to the training examples that it is not able to generalize well for the validation and test sets. A symptom of over fitting is having a model that gets 100% or 99% on the training set, but only 50% on the test data.

Batch Normalization Layer

Batch normalization layer reduces the internal covariance shift. To train a neural network, some preprocessing to the input data are performed. For example, you can normalize all data so that it resembles a normal distribution (which means zero mean and a unitary variance). This prevents the early saturation of non-linear activation functions, such as sigmoid, and assures that all input data are in the same range of values.



An issue, however, appears in the intermediate layers because the distribution of the activations is constantly changing during training. This slows down the training process because each layer must learn to adapt them to a new distribution in every training step. This is known as internal covariate shift.

Batch normalization layer forces the input of every layer to have approximately the same distribution in every training step by following the process below during training:

- a. Calculate the mean and variance of the layers input.
- b. Normalize the layer inputs using the previously calculated batch statistics.
- c. Scale and shift to obtain the output of the layer.

This makes the learning of layers in the network more independent of each other and allows you to be carefree about weight initialization, works as regularization in place of dropout, and other regularization techniques.

Depthwise Convolution and 1 × 1 Convolution Layer

Depthwise convolutions are used to apply a single filter per each input channel (input depth). Pointwise convolution, a simple 1 × 1 convolution, is then used to create a linear combination of the output of the depthwise layer.

Depthwise convolution is extremely efficient relative to standard convolution. However, it only filters input channels. It does not combine them to create new features. An additional layer that computes a linear combination of the output of depthwise convolution through the 1×1 convolution is needed in order to generate these new features.

A 1 × 1 convolutional layer that compresses an input tensor with large channel size to one with the same batch and spatial dimension, but smaller channel size. Given a 4D input tensor and a filter tensorshape [filter height, filter_width, in_channels, channel_multiplier] containing in_channels convolutional filters of depth 1, depthwise_conv2d applies a different filter to each input channel, then concatenates the results together. The output has in channels × channel multiplier channels.

The architecture above provides nonlinearities and preservation of dimension that help to improve the robustness of the network and control over fitting.

4.2.2. Object Detection Network Output

From the input image model, it extracts the feature maps first and overlays them with a W × H grid. Each cell then computes K pre-computed bounding boxes called anchors. Each bounding box has the following:

- Four scalars (x, y, w, h)
 - A confidence score (Pr(Object)*IOU)
 - C° conditional class probability
- The current model architecture has a fixed output of WxHxK(4+1+C). where:
 - W, H = Grid Size
 - K = Number of Anchor boxes
 - C = Number of classes for which you want detection
- The model has a total of 13720 output values, which are derived from the following:
 - 14×14 grid
 - Seven anchor boxes per grid
 - Six values per anchor box. It consists of:
 - Four bounding box coordinates (x, y, w, h)
 - One class probability
 - One confidence score

As a result, there is a total of $14 \times 14 \times 7 \times 10 = 13720$ output values.

If your images are smaller, it is recommended to stretch the image to default size. You can also up-sample them beforehand.

If your images are bigger and you are not satisfied with the results of the default image size, you can try using a denser grid, as details might get lost during the downscaling.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.



4.2.3. Training Code Overview

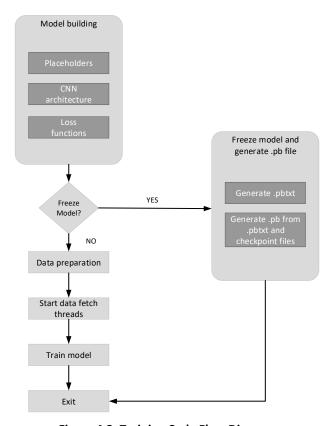


Figure 4.2. Training Code Flow Diagram

Training Code is divided into the following parts:

- Model Configuration
- Model Building
- Training for Overall Execution Flow

The details of each part can be found in subsequent sections.

4.2.3.1. Model Configuration

The design uses Kitti dataset and SqueezeDet model. *kitti_squeezeDet_config()* maintains all the configurable parameters for the model. Below is the summary of configurable parameters.

- Image size
 - Change mc.IMAGE_WIDTH and mc.IMAGE_HEIGHT to configure image size (width and height) in src/config/kitti_squeezeDet_config.py.

mc.IMAGE_WIDTH = 224 mc.IMAGE_HEIGHT = 224

Figure 4.3. Code Snippet – Input Image Size Config

• Since there are four pooling layers, grid dimension is H = 14 and W = 14. anchor_shapes variable of set_anchors() in src/config/kitti_squeezeDet_config.py indicates anchors width and heights. Update it based on anchors per gird size changes.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.
All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



```
def set_anchors(mc):
    H, W, B = 14, 14, 7
    div_scale = 2.0
```

Figure 4.4. Code Snippet – Anchors Per Grid Config #1 (Grid Sizes)

- Batch size
 - Change mc.BATCH SIZE in src/config/kitti squeezeDet config.py to configure batch size.
- Anchors per grid
 - Change mc.ANCHOR PER GRID in src/config/kitti squeezeDet config.py to configure anchors per grid.

Figure 4.5. Code Snippet - Anchors Per Grid Config #2

- Change hard coded anchors per grid in *set_anchors()* in *src/config/kitti_squeezeDet_config.py*. Here, B (value 7) indicates anchors per grid.
- To run the network on your own dataset, adjust the anchor sizes. Anchors are prior distribution over what shapes your boxes should have. The better this fits to the true distribution of boxes, the faster and easier your training is going to be.
- To determine anchor shapes, first load all ground truth boxes and pictures, and if your images are not of the same size, normalize their height and width by the images' height and width. All images are normalized before being fed to the network, so you need to do the same to the bounding boxes and consequently, the anchors.
- Second, perform a clustering on these normalized boxes (that is, you can use k-means without feature whitening and determine the number of clusters either by eyeballing or by using the elbow method).
- Check for boxes that extend beyond the image or have a zero to negative width or height.

Figure 4.6. Code Snippet - Anchors Per Grid Config #3

- Training Parameters
 - Other training related parameters such as learning rate, loss parameters, and different thresholds can be configured from src/config/kitti_squeezeDet_config.py.

FPGA-RD-02242-1 0



```
mc.WEIGHT DECAY
                       = 0.0001
mc.LEARNING RATE
                      = 0.01
mc.DECAY STEPS
                       = 10000
mc.MAX GRAD NORM
                       = 1.0
mc.MOMENTUM
                       = 0.9
mc.LR DECAY FACTOR
                       = 0.5
mc.LOSS COEF BBOX
                       = 5.0
mc.LOSS COEF CONF POS = 75.0
mc.LOSS COEF CONF NEG = 100.0
mc.LOSS COEF CLASS
                       = 1.0
mc.PLOT PROB THRESH
                       = 0.6
                       = 0.4
mc.NMS THRESH
mc.PROB THRESH
                       = 0.005
mc.TOP N DETECTION
                       = 10
                       = True
mc.DATA AUGMENTATION
                       = 150
mc.DRIFT X
mc.DRIFT Y
                       = 100
mc.EXCLUDE HARD EXAMPLES = False
```

Figure 4.7. Code Snippet – Training Parameters

4.2.3.2. Model Building

SqueezeDet class can be configured from *src/nets/squeezeDet.py*. SqueezeDet class constructor builds the model, which is divided into the following sections:

- Forward Graph
- Interpretation Graph
- Loss Graph
- Train Graph
- Visualization Graph
- Forward Graph

Forward Graph

- The CNN architecture consists of Convolution, Batch Normalization, ReLU, and Maxpool.
- Forward graph consists of seven fire layers as described in Table 4.1.

Figure 4.8. Code Snippet – Filter Values

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.
All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



Filter sizes of each convolutional block are mentioned in Table 4.1, which can be configured by changing the values
of depth, as shown in Figure 4.9.

```
if False: # Stride 2
    fire1 = self._fire_layer_v1 ('fire1', self.image_input, oc=depth[0], freeze=freeze_layers[0], w_bin=fl_w_bin, a_bin=fl_a_bin, pool_en=False, min_rng=min_rng, max_rng=max_rng, stride=2)
    else: # Max pool
        fire1 = self._fire_layer_v1 ('fire1', self.image_input, oc=depth[0], freeze=freeze_layers[0], w_bin=fl_w_bin, a_bin=fl_a_bin, min_rng=min_rng, max_rng=max_rng, stride=1,bias_on=bias_on)
        fire2 = self._mobile_layer('fire2', fire1, oc=depth[1], freeze=freeze_layers[1], w_bin=ml_w_bin, a_bin=ml_a_bin, min_rng=min_rng, max_rng=max_rng, bias_on=bias_on)
    if False: # stride 2
        fire3 = self._mobile_layer('fire3', fire2, oc=depth[2], freeze=freeze_layers[2], w_bin=ml_w_bin, a_bin=ml_a_bin, pool_en=False, min_rng=min_rng, max_rng=max_rng, stride=2,bias_on=bias_on)
    else: # Max_pool
        fire3 = self._mobile_layer('fire3', fire2, oc=depth[2], freeze=freeze_layers[2], w_bin=ml_w_bin, a_bin=ml_a_bin, pool_en=False, min_rng=min_rng, max_rng=max_rng, stride=1,bias_on=bias_on)
    ifire4 = self._mobile_layer('fire4', fire3, oc=depth[3], freeze=freeze_layers[3], w_bin=ml_w_bin, a_bin=ml_a_bin, min_rng=min_rng, max_rng=max_rng, bias_on=bias_on)
    fire6 = self._mobile_layer('fire6', fire4, oc=depth[4], freeze=freeze_layers[4], w_bin=ml_w_bin, a_bin=ml_a_bin, min_rng=min_rng, max_rng=max_rng, bias_on=bias_on)
    fire6 = self._mobile_layer('fire6', fire5, oc=depth[5], freeze=freeze_layers[6], w_bin=ml_w_bin, a_bin=ml_a_bin, min_rng=min_rng, max_rng=max_rng, bias_on=bias_on)
    fire7 = self._mobile_layer('fire6', fire6, oc=depth[6], freeze=freeze_layers[6], w_bin=ml_w_bin, a_bin=ml_a_bin, pool_en=False, min_rng=min_rng, max_rng=max_rng, bias_on=bias_on)
    fire7 = self._mobile_layer('fire6', fire6, oc=depth[6], freeze=freeze_layers[6], w_bin=ml_w_bin, a_bin=ml_a_bin, pool_en=False, min_rng=min_rng, max_rng=max_rng_bias_on=bias_on)
    fire7 = self._mobile_layer('fire6', fire6, oc=depth[6], freeze=freeze_layers[6], w_bin=ml_w_bin, a_bin=ml_a_bin, pool_en=False, min_
```

Figure 4.9. Code Snippet – Forward Graph Fire Layers for MobileNet v1

Figure 4.10. Code Snippet - Forward Graph Last Convolution Layer

Interpretation Graph

- The Interpretation Graph consists of the following sub-blocks:
 - interpret output

This block interprets output from network and extracts predicted class probability, predicated confidence scores, and bounding box values.

Output of the convnet is a $14 \times 14 \times 70$ tensor – there are 70 (7 × (4 + 1 + 5)) channels of data for each of the cells in the grid that is overlaid on the image and contains the bounding boxes, confidence score, class predictions. This means the 70 channels are not stored consecutively but are scattered all over and need to be sorted. Figure 4.11 and Figure 4.12 show the details.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice



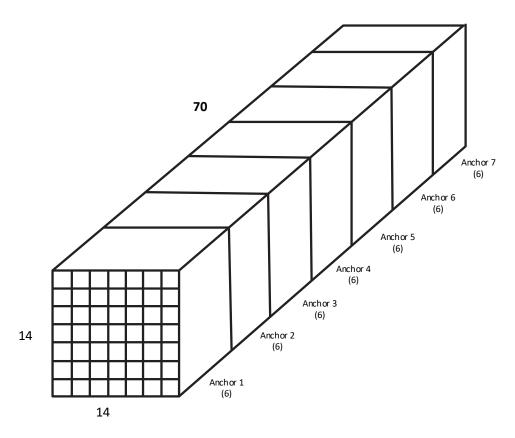


Figure 4.11. Grid Output Visualization #1

For each grid, cell values are aligned as shown in Figure 4.11.

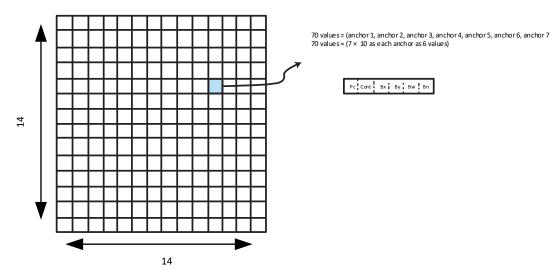


Figure 4.12. Grid Output Visualization #2

Figure 4.12 shows the output from the conv12 layer (4D array of batch size \times 14 \times 14 \times 70) that needs to be sliced with the proper index to get all values of probability, confidence, and coordinates.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



```
# probability
num_class_probs = mc.ANCHOR PER GRID*mc.CLASSES
print ('ANCHOR_PER_GRID:', mc.ANCHOR_PER_GRID)
print ('CLASSES:', mc.CLASSES)
print ('preds2:', preds[:, :, :, :num_class_probs])
print ('ANCHORS:', mc.ANCHORS)
self.pred class probs = tf.reshape(
    tf.nn.softmax(
        tf.reshape(
            preds[:, :, :, :num class probs],
             [-1, mc.CLASSES]
    [mc.BATCH_SIZE, mc.ANCHORS, mc.CLASSES],
    name='pred_class_probs'
# confidence
num confidence scores = mc.ANCHOR PER GRID+num class probs
self.pred_conf = tf.sigmoid(
    tf.reshape(
        preds[:, :, :, num_class_probs:num_confidence_scores],
        [mc.BATCH SIZE, mc.ANCHORS]
    ).
    name='pred confidence score'
# bbox delta
self.pred box delta = tf.reshape(
    preds[:, :, :, num_confidence_scores:],
    [mc.BATCH SIZE, mc.ANCHORS, 4],
    name='bbox delta
)
```

Figure 4.13. Code Snippet - Interpret Output Graph

For confidence score, this must be a number between 0 and 1, as such, sigmoid is used.

For predicting the class probabilities, there is a vector of NUM_CLASS values at each bounding box. Apply a softmax to make it probability distribution.

bbox

This block calculates bounding boxes based on the anchor box and the predicated bounding boxes.

This block calculates Intersection over Union for detected bounding boxes and actual bounding boxes.

Probability

This block calculates detection probability and object class.

- This block calculates different types of losses, which needs to be minimized. To learn detection, localization, and classification, model defines a multi-task loss function. There are three types of losses which are considered for calculation:
 - **Bounding Box**

This loss is regression of the scalars for the anchors.

```
with tf.variable_scope('bounding_box_regression') as scope:
  self.bbox loss = tf.truediv(
          mc.LOSS COEF BBOX * tf.square(
              self.input_mask*(self.pred_box_delta-self.box_delta_input))),
      self.num objects,
      name='bbox_loss
  tf.add_to_collection('losses', self.bbox_loss)
```

Figure 4.14. Code Snippet – Bbox Loss

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal. FPGA-RD-02242-1 0



FPGA-RD-02242-1 0

- Confidence Score
 - To obtain meaningful confidence score, the predicted value of each box is regressed against the real and predicted box. During training, compare the ground truth bounding boxes with all anchors and assign them to the anchors with the largest overlap (IOU).
 - Select the *closest* anchor to match the ground truth box such that the transformation needed is reduced to minimum. Equation evaluates to 1 if the k-th anchor at position-(i, j) has the largest overlap with a ground truth box, and to 0 if no ground truth is assigned to it. This way, you only include the loss generated by the *responsible* anchors.
 - As there can be multiple objects per image, normalize the loss by dividing it by the number of objects (self.num objects).

Figure 4.15. Code Snippet – Confidence Loss

- Class
 - The last part of the loss function is cross-entropy loss for each box to do classification, as you would for image classification.

```
with tf.variable_scope('class_regression') as scope:
    # cross-entropy: q * -log(p) + (1-q) * -log(1-p)|
    # add a small value into log to prevent blowing up
    self.class_loss = tf.truediv(
        tf.reduce_sum(
            (self.labels*(-tf.log(self.pred_class_probs+mc.EPSILON)))
            + (1-self.labels)*(-tf.log(1-self.pred_class_probs+mc.EPSILON)))
            * self.input_mask * mc.LOSS_COEF_CLASS),
            self.num_objects,
            name='class_loss'
)
tf.add_to_collection('losses', self.class_loss)
```

Figure 4.16. Code Snippet – Class Loss

In one model architecture, you obtain the bounding box prediction, the classification, as well as the confidence score.

Train Graph

This block is responsible for training the model with momentum optimizer to reduce all losses.

Visualization Graph

• This block provides visitations of detected results.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



4.2.3.3. Training

Figure 4.17. Code Snippet - Training

sess.run feeds the data, labels batches to network, and optimizes the weights and biases. The code above handles the input data method in case of multiple threads preparing batches, or data preparation in the main thread.

4.3. Training from Scratch and/or Transfer Learning

To train the machine:

1. Go to the top/root directory of the Lattice training code from the command prompt. The model works on 224×224 images.

Current Object detection training code uses mean = 0 and scale = 1/128 (0.0078125) in pre-processing step. Mean and scale can be changed in training code @src/dataset/imdb.py as shown in Figure 4.18.

```
images, scates = [], []
for i in batch_idx:
    im = cv2.imread(self._image_path_at(i))
    im = im.astype(np.float32, copy=False)
    im -= mc.BGR_MEANS
    im /= 128.0 # to make input in the range of [0, 2)
    orig_h, orig_w, _ = [float(v) for v in im.shape]
    im = cv2.resize(im, (mc.IMAGE_WIDTH, mc.IMAGE_HEIGHT
    if mc.GRAY:
        im = cv2.cvtColor(im, cv2.COLOR_BGR2GRAY)
```

Figure 4.18. Training Code Snippet for Mean and Scale

The dataset path can be set in the training code @src/dataset/kitti.py and can be used in combination with the -data_path option while triggering training using train.py to get the desired path. For example, you can have <data_path>/train/crop_images and <data_path>/train/crop_labels.

```
def __init__(self, image_set, data_path, mc,mode='train'):
    imdb.__init__(self, 'kitti_'+image_set, mc)|
    self._image_set = image_set
    self._data_root_path = data_path
    self.mode = mode
    if self.mode == 'train':
        self._image_path = os.path.join(self._data_root_path, 'train', 'crop_images')
        self._label_path = os.path.join(self._data_root_path, 'train', 'crop_labels')
    else:
        self._image_path = os.path.join(self._data_root_path, 'val', 'crop_images')
        self._label_path = os.path.join(self._data_root_path, 'val', 'crop_labels')
    self._label_path = os.path.join(self._data_root_path, 'val', 'crop_labels')
    self._classes = self.mc.CLASS_NAMES
```

Figure 4.19. Training Code Snippet for Dataset Path

FPGA-RD-02242-1 0



Notes:

- train.txt file name of dataset images (Gets generated when crop data.py is executed)
- image_set train (ImageSets/train.txt)
- data path \$ROOT/DB/Object Detector
 - Images \$ROOT/DB/Object_Detector /train/crop_images
 - Annotations \$ROOT/DB/Object_Detector/train/crop labels
- 2. Modify the training script. @scripts/train.sh is used to trigger training.

Figure 4.20 shows the input parameters, which can be configured.

```
python ./src/train.py \
    --dataset=KITTI \
    --pretrained_model_path=$PRETRAINED_MODEL_PATH \
    --data_path=$TRAIN_DATA_DIR \
    --train_set=$IMAGE_SET \
    --val_set=$VAL_SET \
    --train_dir="$TRAIN_DIR/train_weights" \|
    --net=$NET \
    --max_steps=1000 \
    --summary_step=50 \
    --checkpoint_step=100 \
    --gpu=$GPUID
```

Figure 4.20. Training Input Parameter

- \$TRAIN_DATA_DIR dataset directory path. /DB/Object_Detector is an example.
- \$TRAIN_DIR log directory where checkpoint files are generated while model is training.
- \$GPUID gpu id. If the system has more than one gpu, it indicates the one to use.
- --summary step indicates at which interval loss summary should be dumped.
- --checkpoint_step indicates at which interval checkpoints are created.
- --max steps indicates the maximum number of steps for which the model is trained.
- 3. Execute the train.sh command script which starts training.

```
nn.so.7
2021-09-12 07:58:54.539300: I tensorflow/stream_executor/platform/default/dso_loader.cc:42] Successfully opened dynamic library libcu las.so.10.0
train_conf_loss: 26.8348388671875, train_bbox_loss: 9.123150825500488, train_class_loss: 2.502012252807617
val_conf_loss: 27.00383949279785, val_bbox_loss: 8.954672813415527, val_class_loss: 2.502012252807617
val_conf_loss: 27.00383949279785, val_bbox_loss: 8.954672813415527, val_class_loss: 12.502012014389038
2021-09-12 07:58:56.332957: step 0, train_loss = 38.46,val_loss = 38.46 (3.3 images/sec; 6.036 sec/batch)
val_bbox_loss improved from inf to 8.954672813415527
val_class_loss improved from inf to 2.502012014389038
2021-09-12 07:59:09.482434: step 10, train_loss = 4.49,val_loss = 4.48 (109.6 images/sec; 0.183 sec/batch)
2021-09-12 07:59:09.488097: step 30, train_loss = 2.39,val_loss = 2.11 (45.5 images/sec; 0.440 sec/batch)
2021-09-12 07:59:13.107948: step 30, train_loss = 2.81,val_loss = 2.45 (46.7 images/sec; 0.488 sec/batch)
2021-09-12 07:59:13.107948: step 40, train_loss = 2.61,val_loss = 2.45 (46.7 images/sec; 0.488 sec/batch)
train_conf_loss: 0.6213857531547546, train_bbox_loss: 0.082319475710392, val_class_loss: 1.660585880279541
val_conf_loss: 0.6213857531547546, train_bbox_loss: 0.082319475710392, val_class_loss: 1.839078426361084
2021-09-12 07:59:12.430031: step 60, train_loss = 2.37,val_loss = 2.240 (6.3 images/sec; 0.178 sec/batch)
2021-09-12 07:59:21.430031: step 60, train_loss = 1.96,val_loss = 1.29 (44.9 images/sec; 0.178 sec/batch)
2021-09-12 07:59:21.310546: step 90, train_loss = 1.58,val_loss = 1.72 (41.1 images/sec; 0.486 sec/batch)
2021-09-12 07:59:28.367665: step 80, train_loss = 1.58,val_loss = 1.72 (41.1 images/sec; 0.447 sec/batch)
2021-09-12 07:59:28.367665: step 80, train_loss = 1.58,val_loss = 1.79 (48.0 images/sec; 0.4486 sec/batch)
2021-09-12 07:59:28.367665: step 80, train_loss = 1.58,val_loss = 1.79 (48.0 images/sec; 0.4486 sec/batch)
2021-09-12 07:59:28.367665: step 80, train_loss = 1.58,val_loss = 1.96 (48.0 images/sec; 0.4
```

Figure 4.21. Execute train.sh Script

4. Start Tensorbaord.

```
$ tensorboard --logdir <log directory of training>
```

For example: tensorboard --logdir ./weights/Object_Detector/train_weights

- 5. Open the local host port on your web browser.
- 6. Check the training status on tensorboard.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.
All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



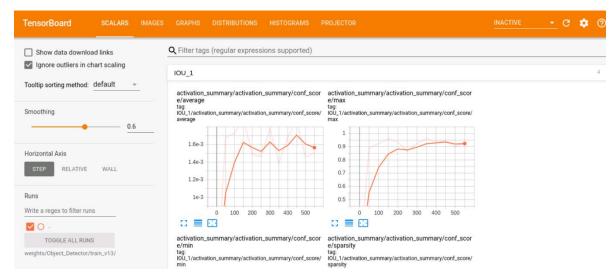


Figure 4.22. TensorBoard

7. Check if the checkpoint, data, meta, index and events (if using TensorBoard) files are created at the weights directory. These files are used for creating the frozen file(*.pb)



Figure 4.23. Example of Checkpoint Data Files at weights/Object_Detector/train_weights/folder



FPGA-RD-02242-1 0

5. Testing the Trained Model

This section describes how to test the trained model on the PC.

5.1. Running the Test Shell Script

The script below shows how to run the test shell. The results of detections are stored in the DB/Object Detector/val/images out v1 directory.

\$./test.sh

```
keep_idx=[0, 1]
final_boxes=[array([105.37656, 112.21976, 70.8886, 62.71424], dtype=float32), array([51.338802, type=float32)]
final_probs=[0.70821524, 0.7969094]
# of final boxes= 2
Image detection output saved to ./DB/Object_Detector/val/images_out_v1/out_Combined2_2057.jpg
```

Figure 5.1. Test Script Last Log

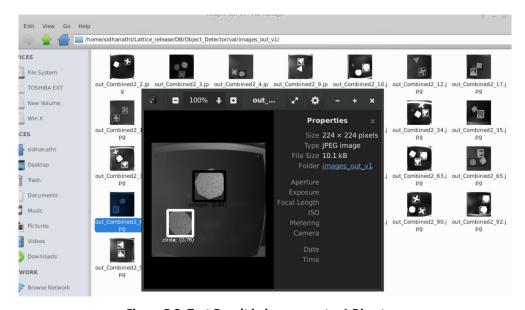


Figure 5.2. Test Result in images_out_v1 Directory

5.2. Checking the Mean Average Precision (mAP) of the Model

To check the mAP of the model, perform the following steps:

- 1. First the ground truth is generated
 - \$ python generate groundtruth.py
- 2. Run the map calculation script
 - \$ python MAP/map.py



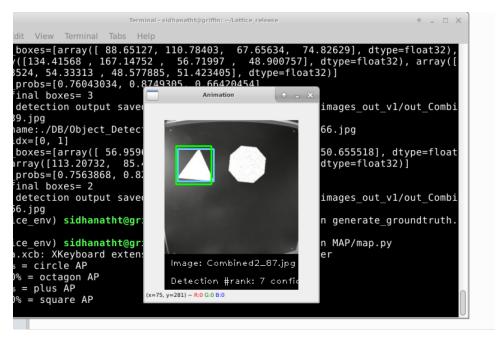


Figure 5.3. Running map.py

3. The results are stored in the MAP/results directory.

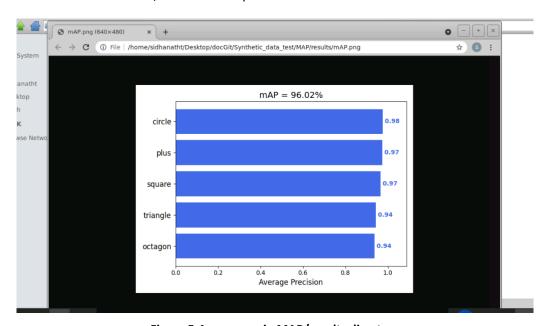


Figure 5.4. map.png in MAP/results directory



5.3. Testing the Centroid Tracking Algorithm

To test the centroid tracking algorithm, flash with_448x448p_window.bit into the board. This bit file makes the board a normal USB webcam.

5.3.1. Testing the Images Captured using the Board Camera

Open the desktop application and start the video. Click the *Capture images* button in the application to record images from the board. To stop capturing the images, click *Capture images* again to stop capturing. Images are stored in the *captured_images* folder in the application root directory.

To test the centroid tracking on the generated video:

1. Create a directory named *board_test* and copy the captured_images directory to it. Since the captured images are 640×480p, you need to crop it to 448×448p by running the script below.

```
$ python crop.py ./board test/captured images
```

```
(lattice_env) sidhanatht@griffin:~/Lattice_release$ python crop.py ./board_test/
captured_images/
./board_test/captured_images/
./board_test/crop_frames
(lattice_env) sidhanatht@griffin:~/Lattice_release$
```

Figure 5.5. crop.py Terminal Output

2. Edit the test_fixedpoint.sh to set the path for the weights, input, and output.

```
test fixedpoint.sh
Open ▼
                                                                                            Sav
1 #rm ./DB/test/img out/*
3 # cap 32 96 r n200.png cap_32_96_r_n76.png cap_32_96_r_p176.png
4 # cap 0608 2.png cap 0611 face 1.png
6 #rm ./DB/test/img out/*
7 python ./src/inference_fixedpoint.py --mode=tracker \
                       --checkpoint=./weights/Object Detector/train weights/model.ckpt \
               --input_path=./board_test/crop_frames/* \
0
               --out dir=./board test/tracker out/crop centriod/ \
.1
               --out_label_dir=./board_test/tracker_out/crop_centriodlabel/ \
.2
               --demo net=squeezeDet
13
```

Figure 5.6. Input Parameter for the Script

3. After saving, run ./test_fixedpoint.sh and click **Enter** to load the next frame. Press and hold **Enter** to load the images rapidly.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



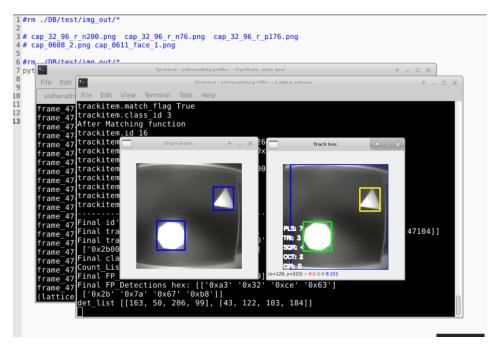


Figure 5.7. Output Window with Overlay on Images

5.3.2. Live Testing of the Centroid Tracking using Real-time Data from the Board

In this procedure, the board should be connected to the PC and two monitors are recommended. The Ubuntu virtual machine is used to run the following ML scripts.

- 1. Connect the board to the PC using a micro-USB 3.0 cable. Make sure no other camera is connected to the PC.
- 2. Open the desktop app in one monitor. Do not click open camera.
- 3. Run the script.
 - \$./test_fixedpoint live.sh
- 4. Click Play in the desktop app. The result should look similar to the one shown in Figure 5.8.

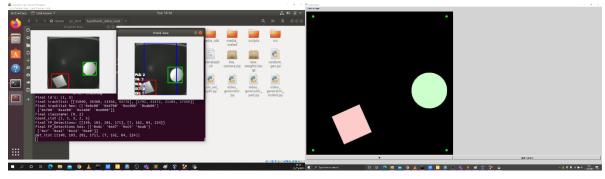


Figure 5.8. Object Detector on PC



6. Creating Frozen File

This section describes the procedure for freezing the model, which is aligned with the Lattice sensAl[™] tool. Perform the steps below to generate the frozen protobuf file.

6.1. Generating the Frozen .pb File

Generate .pb file from the best mAP accuracy model using the command below from the root directory of the training code.

\$./generatepb.sh.

(lattice_env) sidhanatht@gamora:~/projects/lattice/Package_test/synthetic_data_o
bject_detector-dev1\$./generatepb.sh

Figure 6.1. pb File Generation from Checkpoint

Figure 6.2 shows the generated .pb file in the weights/Object_Detector/train_weights folder.

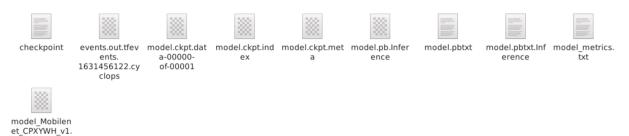


Figure 6.2. Frozen .pb File



7. Creating Binary File with Lattice sensAl

This chapter describes how to generate the binary file using the Lattice Neural Network Compiler Software 4.0 program.



Figure 7.1. sensAl Home Screen

To create the project in sensAl tool:

- Click File > New.
- 2. Enter the following settings:
 - Project Name
 - Framework TensorFlow
 - Class CNN
 - Device CrossLink-NX
 - 'Compact Mode' should be unchecked.
- 3. Click **Network File** and select the network (.pb) file.



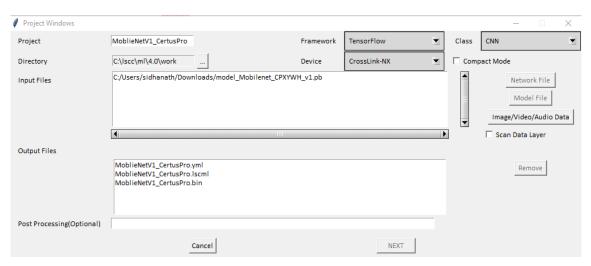


Figure 7.2. sensAI - Network File Selection

4. Click Image/Video/Audio Data and select the image input file.

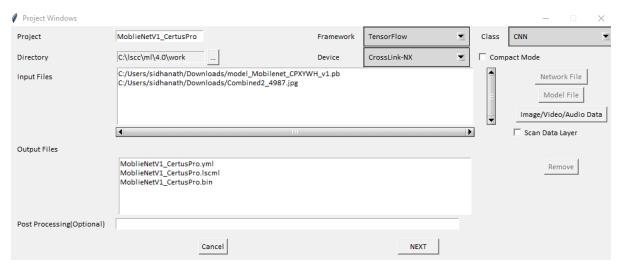


Figure 7.3. sensAI - Input Image File Selection

- 5. Click NEXT.
- 6. Configure your project settings as shown in Figure 7.4.



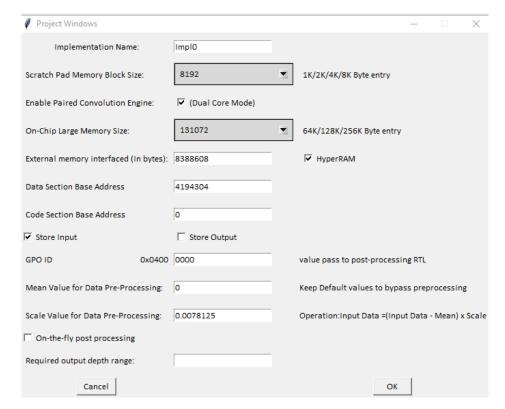


Figure 7.4. sensAI - Project Settings

Scratch Pad Memory Block Size and Data Section Base Address should match with FPGA RTL code.

- 7. Click **OK** to create the project.
- 8. Double-click analyze.

FPGA-RD-02242-1.0

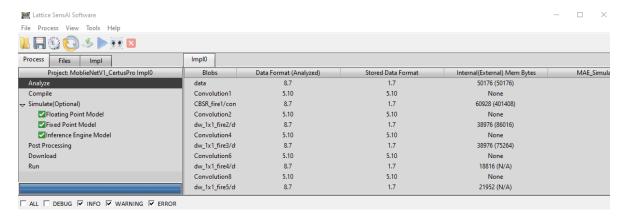


Figure 7.5. sensAI – Analyze Project

- 9. After analyzing the project, the tool generates the correct format for each layer since you perfromed the quantization in training.
- 10. Confirm the Q format of each layer as shown in Figure 7.6. Edit the fractional bit for each layer by double-clicking on the values against each layer one by one.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



Blobs	Data Format (Analyzed)	Stored Data Format	Internal(External) Mem Bytes
data	8.7	1.7	50176 (50176)
Convolution1	5.10	5.10	None
CBSR_fire1/con	8.7	1.7	60928 (401408)
Convolution2	5.10	5.10	None
dw_1x1_fire2/d	8.7	1.7	38976 (86016)
Convolution4	5.10	5.10	None
dw_1x1_fire3/d	8.7	1.7	38976 (75264)
Convolution6	5.10	5.10	None
dw_1x1_fire4/d	8.7	1.7	18816 (N/A)
Convolution8	5.10	5.10	None
dw_1x1_fire5/d [,]	8.7	1.7	21952 (N/A)

Figure 7.6. Q Format Settings for Each Layer

- 11. After changing the fractional bit, double-click on Analyze again.
- 12. Double-click **Compile** to generate the Firmware file.

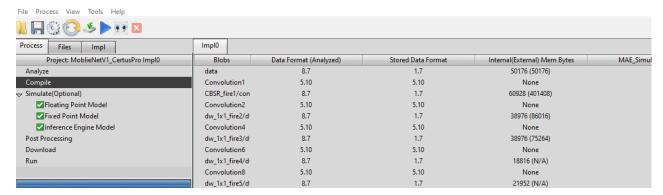


Figure 7.7. Compile Project



8. Hardware Implementation

8.1. Top Level Information

8.1.1. Block Diagram

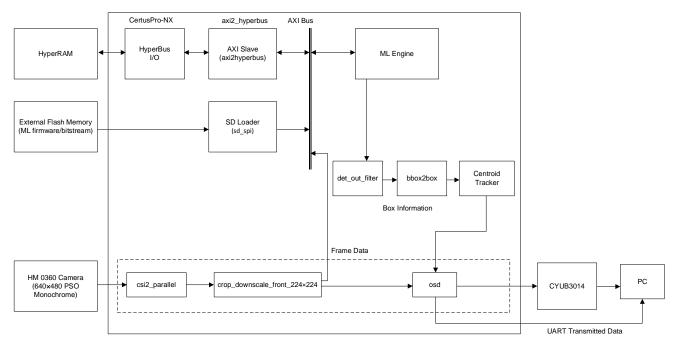


Figure 8.1. RTL Top Level Block Diagram

8.1.2. Operational Flow

This section provides a brief idea about the data flow across the CertusPro-NX board.

- The CNN module is configured with the help of a binary (.bin) file stored in a SD card. The .bin file is a command sequence code, which is generated by the Lattice Machine Learning software tool.
- The command code is written in hyperRAM through AXI before the execution of CNN Accelerator IP Core starts. CNN reads command code from hyperRAM during its execution and performs calculation with it per command code. Intermediate data may be transferred from/to hyperRAM per command code.
- The RAW8 data from the csi2_to_parallel module is downscaled to 224 × 224 image resolution by the crop_downscale_front_224x224 module to match CNN input resolution. This data is written into hyperRAM memory through axi2_hyperbus through the axi_ws2m AXI interface module.
- After the command code and input data are available, the CNN Accelerator IP Core starts calculation at the rising edge of start signal.
- The output data of CNN is passed to *det_out_filter* for post processing. *det_out_filter* generates bounding box X, Y, W, and H coordinates associated with top 5 confidence value indexes for 224 × 224 image resolution.
- These coordinates are passed to osd_back_128x128_human_count for resizing to fit the actual image resolution on the PC.



8.1.3. Core Customization

Table 8.1. Core Parameter

Constant	Default	Description
	(Decimal)	
OVLP_TH_2X	5	Intersection Over Union Threshold (NMS)
NUM_FRAC	10	Fraction Part Width in Q-Format representation.
EN_INF_TIME	0	Enable Timing measurement logic
		By default, it is zero and the memory file used is human_count.mem.
		If assigned 1, timing measurement is enabled and the memory file used is human_count_INF.mem.
		In order to configure the respective memory file, follow the steps below:
		1. Open dpram8192x8_human_count.ipx from the File List in Radiant.
		2. Click Browse Memory File from Initialization section.
		3. Update the mem file path:
		For 0 – /src/jedi_common/human_count.mem
		For 1 – /src/jedi_common/human_count_INF.mem
INF_MULT_FAC	15907	Inference time multiplying factor calculated as per CNN clock frequency and using Q-Format (Q1.31).
		CNN Clock Frequency = 135 MHz
		Hence, CNN clock period
		$= 1/(135 \times 10^{-6}) \mu s$
		= 0.000007407 ms
		Now, Q1.31 = 0.000007407 × 2 ³¹ = ~15907
FLASH_START_ADDR	24'h300000	SPI Flash Read Start Address (keep the same address in programmer while loading the firmware file)
		For example, for the current start address, programmer address should be 0x00300000.
FLASH_END_ADDR	24'h400000	SPI Flash Read End Address (keep the same address in programmer while loading the firmware file)
		The address must be in multiple of 512 bytes.
		For example, for the current end address, programmer address should be: 0x00400000.
	(Constant Parameters (Not to be modified)
NUM_ANCHOR	1372	Number of reference bounding boxes for all grids
NUM_GRID	196	Total number of Grids (X * Y)
NUM_X_GRID	14	Number of X Grids
NUM_Y_GRID	14	Number of Y Grids
PIC_WIDTH	224	Picture Pixel Width (CNN Input)
PIC_HEIGHT	224	Picture Pixel Height (CNN Input)
TOP_N_DET	10	Number of Top confidence bounding boxes detection
HYPERRAM_BASEADDR	0x400000	Indicates hyperRAM starting Base address location value. This should match with the sensAl compiler while generating the firmware.



8.2. Architecture Details

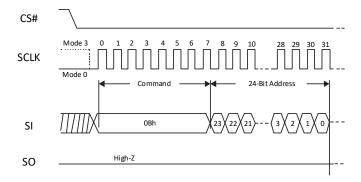
8.2.1. SPI Flash Operation

The RTL module spi_loader_spram provides SPI Flash read operation and writes that data into HyperRAM through the AXI interface. It reads from SPI Flash and as soon as the board gets powered up, the .bit and .bin files are loaded in the expected addresses.

- Expected Address for BIT File (Programmer) 0x0000000 0x00100000
- Expected Address for Firmware File (Programmer) FLASH START ADDR FLASH END ADDR

Typical sequence of the SPI Read commands for SPI Flash MX25L12833F is implemented using FSM in RTL as per the flow of the operation below.

- After FPGA Reset, RELEASE FROM DEEP POWER DOWN command (0xAB) is passed to SPI Flash memory. Then RTL
 waits for 500 clock cycles for SPI flash to come into Standby mode, if it is in Deep Power Down mode.
- RTL sends FAST READ command code (0x0B) on SPI MOSI signal for indication of Read Operation to SPI Flash.
- RTL sends three bytes of Address on SPI MOSI channel which determines the location in SPI flash from the position the data needs to be read.
- This SPI Flash has eight Dummy cycles as wait duration before read data appears on MISO channel. After waiting for eight dummy cycles, the RTL code starts reading the data.
- This read sequence is shown in Figure 8.2. The SPI Interface Signal Mapping with RTL signals are as follow:
 - CS (Chip Select) => SPI CSS
 - SCLK (Clock) => SPI CLK
 - SI (Slave In) => SPI_MOSI
 - SI (Slave Out) => SPI_MISO
- The Read Data on the MISO signal is stored in a FIFO in RTL, which then reads the data in multiples of 512 bytes. After 512 bytes chip select is de-asserted, the AXI FSM state is activated.



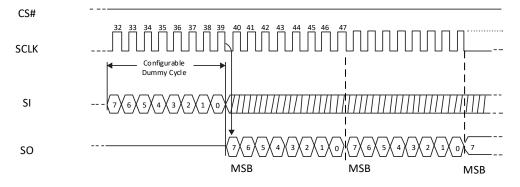


Figure 8.2. SPI Read Command Sequence



- AXI logic reads the data from FIFO in bursts of four on the AXI write channel, with each burst having 128 bytes.
- In accessing the HyperRAM, the axi_ws2m module is used as a Muxing module among the multiple input slave AXI interfaces as shown in Figure 8.6. The spi_loader_spram module is considered as SLAVE 0 and given priority to write into HyperRAM. The Master Interface connects to the axi2_hyperbus module, which provides output interface for accessing HyperRAM.
- After writing to HyperRAM is complete, the 512 bytes are fetched from the SPI Flash using the same command sequence as explained above until the FLASH_END_ADDR is reached.

8.2.2. Pre-processing CNN

The output from the *csi2_to_parallel* module is a stream of RGB data that reflects the camera image, which is given to the *crop_downscale_front_224x224* module.

The $crop_downscale_front_224x224$ module processes that image data and generates input of 224 × 224 image data interface for CNN IP.

8.2.2.1. Pre-processing Flow

- RAW8 data values for each pixel are fed serially line by line for an image frame.
- These RAW8 data values are considered as valid only when horizontal and vertical masks are inactive. The mask parameters set to mask out boundary area of resolution (1920x1080) to 448 × 448 are shown below.
 - Left masking = 96
 - Right masking = 544 (Obtained as 96 + 448)
 - Top masking = 16
 - Bottom masking = 464 (Obtained as 16 + 448)
- The image obtained after masking is shown in Figure 8.3.

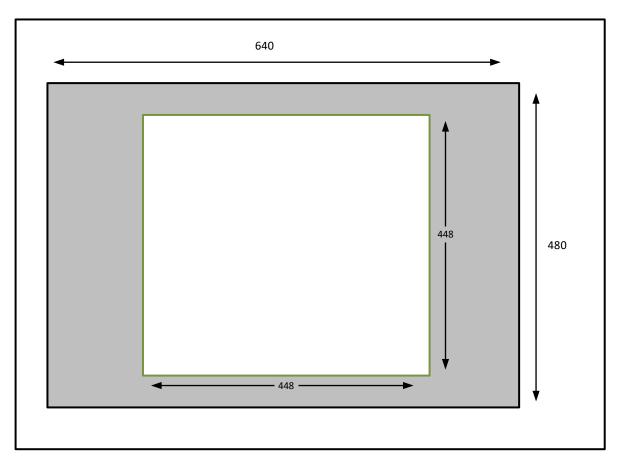


Figure 8.3. Masking



• The 448 × 448 frame block is downscaled into 224×224 resolution image as shown in Figure 8.4 by accumulating 2×2 pixels into single pixel (that is $448/2 \times 448/2 = 224 \times 224$).

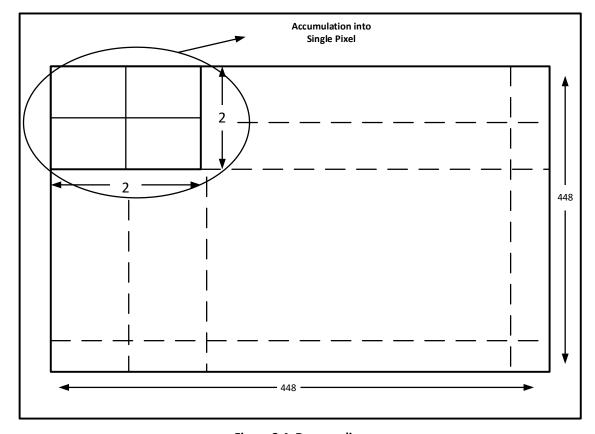


Figure 8.4. Downscaling

- This accumulated value is written into Line Buffer. Line Buffer is a True Dual-Port RAM. Accumulated RAW8 pixel values for 2 × 2 grids are stored in the same memory location.
- When data is read from memory, each RAW8 value is divided by 4 (that is the area of the 2 × 2 grid) to take the average of 2 × 2 grid matrix.
- The data from memory is read and stored in HyperRAM for CNN input through axi2_hyperbus, through the axi_w2sm module, which acts as an AXI interface to write data from slave (crop_downscale_front_224x224) to master (axi2_hyperbus). This process is described in the next section.

8.2.3. HyperRAM Operations

The CertusPro-NX board uses external HyperRAM for faster data transfer mechanism among the internal blocks and enhances the system performance. The *crop_downscale_front_224x224* module uses HyperRAM to store the downscaled image data.



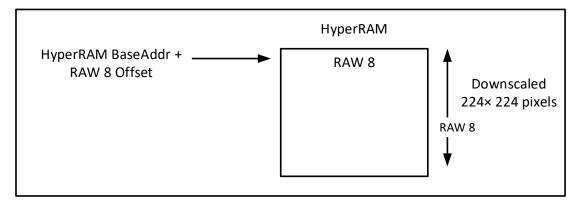


Figure 8.5. HyperRAM Memory Addressing

- The 448×448 image is distributed into 224 horizontal and 224 vertical lines, and each block consists of 2×2 pixels as shown in Figure 8.4. Thus, there is a total of 224×224 pixel values for the downscaled image.
- Primarily, the *crop_downscale_front_224x224* module stores 224 values each of RAW8 into a local FIFO for all 224 horizontal blocks. Later, this stored data is written to HyperRAM through the AXI write data channel.
- As shown in Figure 8.5, when final data is written out, 224 × 224 RAW8 pixels are initially stored into HyperRAM starting from HyperRam Base address location.
- The 224 × 224 pixel values stored in HyperRAM are serially obtained by the CNN engine after getting command sequence through the AXI interface.
- In order for the *crop_downscale_front_224x224* module to access HyperRAM for the operations explained above, the *axi_ws2m* module functions as a Muxing module for multiple input slave AXI interfaces as shown in Figure 8.6.
- For the internal blocks to access HyperRAM, the axi_ws2m module considers the sd_spi module as SLAVE 0, the cnn_opt module as SLAVE 1, the crop_downscale_front_224x224 module as SLAVE 2, and the MASTER connects these slaves to the axi2_hyperbus module.
- The priority to select write channel is given, respectively, to the pi_loader slave, cnn slave, and crop-downscale slave. Whenever valid address is available from the respective Slave on its write address channel, that slave is given access of master channel if other priority slaves are not accessing it. Thus, when valid write address is obtained from crop_downscale_front_224x224 module, access is given to Slave 2 to use HyperRAM.



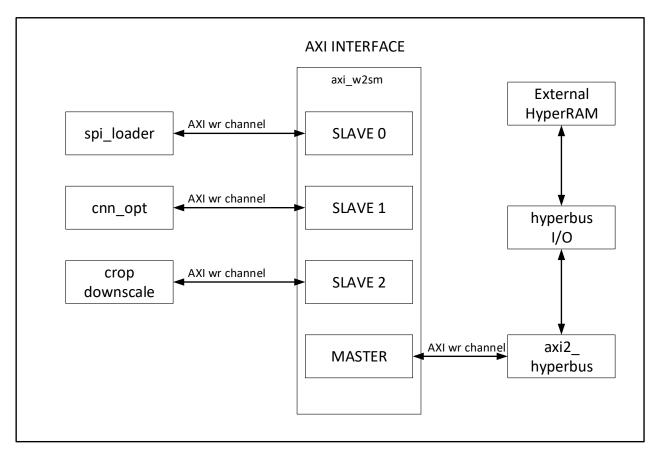


Figure 8.6. HyperRAM Access Block Diagram

8.2.4. Post-processing CNN

CNN provides a total of 13720 [1372 \times 10 (C,P, X, Y, W, H)] values, which are given to the *det_out_filter* module. The CNN output data consists of the following parameters.

Table 8.2. Data Parameters of CNN Output

Parameter	Description
С	This parameter indicates the confidence of detected object class. For each grid cell (14×14), one confidence value (16 -bit) for each anchor box (7) is provided making total values of confidence $14 \times 14 \times 7 = 1372$ from CNN Output.
Р	This parameter indicates the probability of detected object class. For each grid cell (14 \times 14), one probability value (16-bit) for each anchor box (7) is provided making total values of probability $14 \times 14 \times 7 = 1372$ from CNN Output.
Х	This parameter indicates the Relative X coordinate to transform the anchor box into a predicted bounding box for detected object. For each grid cell, one Relative X value (16-bit) for each anchor box is provided making total values of 14 × 14 × 7 = 1372 for X from CNN Output.
Y	This parameter indicates the Relative Y coordinate to transform the anchor box into a predicted bounding box for detected object. For each grid cell, one Relative Y value (16-bit) for each anchor box is provided making total values of 14 × 14 × 7 = 1372 for Y from CNN Output.
W	This parameter indicates the Relative W (Width) coordinate to transform the anchor box into a predicted bounding box for detected object. For each grid cell, one Relative W value (16-bit) for each anchor box is provided making total values of 14 × 14 × 7 = 1372 for W from CNN Output.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice



Parameter	Description
Н	This parameter indicates the Relative H (Height) coordinate to transform the anchor box into a predicted bounding box for detected object. For each grid cell, one Relative H value (16-bit) for each anchor box is provided making total values of 14 × 14 × 7 = 1372 for H from CNN Output.

Figure 8.7 shows the format of CNN output.

Index No. 0-195 196 - 391 1176 - 1371				Output Data			С			
Grid No. 0-195 0-195 0-195 Output Data P1 P5 Index No. 1372 - 1567 1568 - 1763 2548 - 2743 6860 - 7055 7056 - 7251 8036 - 8231 Grid No. 0-195 0-195 0-195 0-195 0-195 Anchor No. 1 2 7 1 2 7 Output Data X Y W H X Y W H Index No. 8232 - 8427 8428 - 8623 8624 - 8819 8820 - 9015 12936 - 13131 13132 - 13327 13328 - 13523 13524 - 13719 Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195				-			C			
Output Data P1 2 7 Index No. 1372 - 1567 1568 - 1763 2548 - 2743 6860 - 7055 7056 - 7251 8036 - 8231 Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 Anchor No. 1 2 7 1 2 7 Output Data X Y W H X Y W H Index No. 8232 - 8427 8428 - 8623 8624 - 8819 8820 - 9015 12936 - 13131 13132 - 13327 13328 - 13523 13524 - 13715 Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195				Index No.	0-195	196 - 391		1176 - 1371		
Output Data P1 P5 Index No. 1372 - 1567 1568 - 1763 2548 - 2743 6860 - 7055 7056 - 7251 8036 - 8231 Grid No. 0-195 0-195 0-195 0-195 0-195 Anchor No. 1 2 7 1 2 7 Output Data X Y W H X Y W H Index No. 8232 - 8427 8428 - 8623 8624 - 8819 8820 - 9015 12936 - 13131 13132 - 13327 13328 - 13523 13524 - 13719 Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195				Grid No.	0-195	0-195		0-195		
Index No. 1372 - 1567 1568 - 1763 2548 - 2743 6860 - 7055 7056 - 7251 8036 - 8231 Grid No. 0-195 0-195 0-195 0-195 0-195 Anchor No. 1 2 7 1 2 7 Output Data X Y W H X Y W H Index No. 8232 - 8427 8428 - 8623 8624 - 8819 8820 - 9015 12936 - 13131 13132 - 13327 13328 - 13523 13524 - 13719 Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 0-195				Anchor No.	1	2		7		
Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195 7 1 2 7 1 2 7 1 2 7 X Y W H X Y W H X Y W H 12936 - 13131 13132 - 13327 13328 - 13523 13524 - 13719 9-195 0-195<	Output Data		P1						P5	
Anchor No. 1 2 7 1 2 7 Output Data X Y W H X Y W H Index No. 8232 - 8427 8428 - 8623 8624 - 8819 8820 - 9015 12936 - 13131 13132 - 13327 13328 - 13523 13524 - 13718 Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195	Index No.	1372 - 1567	1568 - 1763		2548 - 2743		6860 - 7055	7056 - 7251		8036 - 8231
Output Data X Y W H X Y W H Index No. 8232 - 8427 8428 - 8623 8624 - 8819 8820 - 9015 12936 - 13131 13132 - 13327 13328 - 13523 13524 - 13719 Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195	Grid No.	0-195	0-195		0-195		0-195	0-195		0-195
Index No. 8232 - 8427 8428 - 8623 8624 - 8819 8820 - 9015 12936 - 13131 13132 - 13327 13328 - 13523 13524 - 13719 Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 0-195 0-195	Anchor No.	1	2		7		1	2		7
Grid No. 0-195 0-195 0-195 0-195 0-195 0-195 0-195	Output Data	Х	Y	W	Н		Х	Y	W	Н
	Index No.	8232 - 8427	8428 - 8623	8624 - 8819	8820 - 9015		12936 - 13131	13132 - 13327	13328 - 13523	13524 - 13719
And the National Control of the Cont	Grid No.	0-195	0-195	0-195	0-195		0-195	0-195	0-195	0-195
Anchor No.	Anchor No.		1						7	

Figure 8.7. CNN Output Data Format

The primary functionality of the *det_out_filter* module is to capture the CNN valid output and modify it to work with the *osd_back_128x128_human_count* module.

The det out filter module contains two sub-modules: det sort conf and det st bbox.

- 1372 values of confidence are passed to the *det_sort_conf* module. It sorts out the top 10 highest confidence values and stores their indexes. Index values are passed to the *det_st_class* and *det_st_bbox* modules.
- 1372 values of probability are passed to the *det_st_class* module. It provides the valid class probability bitmap, which is passed to the det st bbox module.
- 1372 × 4 values of coordinates are passed to the *det_st_bbox* module. It calculates the bounding box coordinates for 448 × 448 image from 224 × 224 coordinates.

The *crop_downscale* module contains logic for post processing.

- The draw_box module calculates the box coordinates for 896 × 896 image from 224 × 224 coordinates.
- The *lsc_osd_text* module generates character bitmap for text display on PC.

8.2.4.1. Confidence Sorting

- All input confidence values (1372) are compared with threshold parameter CONF_THRESH(500) value. Confidence values that are greater than threshold are considered as valid for sorting.
- The *det_sort_conf* module implements an anchor counter (0-1371), which increments on each confidence value. It provides the index of confidence value given by the CNN output.
- Two memory arrays are generated in this module: (1) sorted top 10 (TOP_N_DET) confidence value array, and (2) sorted top 10 confidence index array.
- For sorting, a standard sorting algorithm is followed. As input confidence values start arriving, each value is compared with stored/initial value at each location of the confidence value array.
- If the input value is greater than stored/initial value on any array location and lesser than stored/initial value of previous array location, the input value is updated on current array location. The previously stored value of current location is shifted into the next array location.
- Refer to Figure 8.8 for sorting of new value of confidence into existing confidence value array. Calculated
 confidence index (anchor count value) is also updated in the confidence index array along with the confidence
 value array.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.
All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



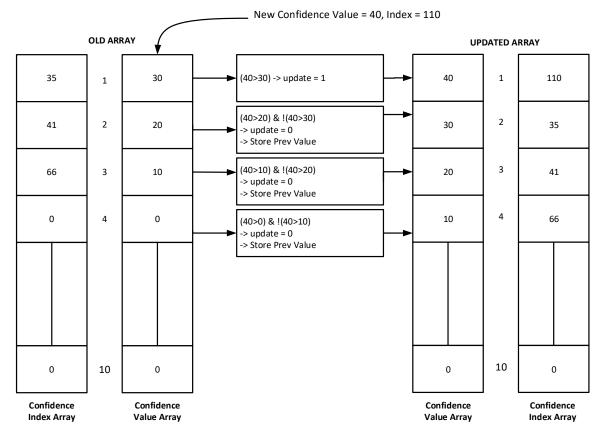


Figure 8.8. Confidence Sorting

This process is followed for all 1372 confidence values. This module provides 10 indexes (o idx 00 to o idx 09) as output along with the count of valid indexes (o_num_conf). o_idx_00 contains the highest confidence value index and o_idx_09 contains the lowest confidence value index.

8.2.4.2. Bounding Box Calculation

The Neural Network for Object Detection is trained with seven reference boxes of pre-selected shapes having constant W (Width) and H (Height). These reference boxes are typically referred as anchors.

Table 8.3. Pre-Selected Width and Height of Anchor Boxes

Anchor No.	1	2	3	4	5	6	7
W × H (pixel)	184 × 184	138 × 138	92 × 92	69 × 69	46 × 46	34 × 34	23 × 23

Anchors are centered around 14 × 14 grid cells of image. So each grid center has above seven anchors with preselected shape. 14 × 14 are the number of grid centers along horizontal and vertical directions. The grid center (X, Y) pixel values are shown in Table 8.4.

Table 8.4. Grid Center Values (X, Y) for Anchor Boxes

Grid No.	1	2	3	4	5	6	7	8	9	10	11	12	13	14
X (pixel)	15	30	45	60	75	90	105	119	134	149	164	179	194	209
Y (pixel)	15	30	45	60	75	90	105	119	134	149	164	179	194	209

CNN provides a total of 1372 (14 × 14 × 7) values of each relative coordinates X, Y, W, and H to transform the fixed size anchor into a predicted bounding box. Input X, Y, W, and H values associated with top 10 sorted confidence indexes are used for box calculation in det st bbox module.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.



Each anchor is transformed to its new position and shape using the relative coordinates as shown in logic 1.

```
LOGIC 1

X' = X coordinate of Predicted Box

X = Grid Center X according to Grid number

W = Width of Anchor according to Anchor number

DeltaX = Relative coordinate for X (CNN output)

X' = X + W × DeltaX

Y' = Y + H × DeltaY

W' = W × DeltaW

H' = H × DeltaH
```

The Box coordinates are passed to bbox2box module in jedi_human_count_top.v after NMS process.

8.2.4.3. NMS - Non Max Suppression

The NMS is implemented to make sure that in object detection, a particular object is identified only once. It filters out the overlapping boxes using OVLP_TH_2X value.

NMS process is started when the CNN output data is completely received.

- The process starts from the box having highest Confidence coordinates: 0th location in X, Y, W, H array.
- These coordinates are compared against the second highest Confidence coordinates: First location in X, Y, W, H array. From this comparison, Intersection and Union coordinates are found.
- From these coordinates, Intersection and Union area are calculated between the highest confidence box and the second highest confidence box as shown in Figure 8.9.

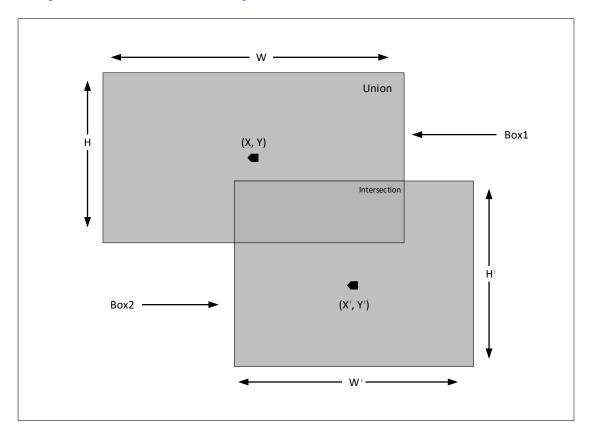


Figure 8.9. Intersection-Union Area NMS

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



- If Intersection Area × (OVLP_TH_2X/2) > Union Area, the box with the lower confidence value is blocked in final output.
- This NMS calculation is performed between all the combinations of two boxes.
- After all combinations are checked, output array o_bbox_bmap contains boxes, which are correctly overlapped or non-overlapped. o_out_en provides valid pulse for crop_downscale_human_count for further processing on these box coordinates.

8.2.4.4. Bounding Box Upscaling

The process of upscaling bounding boxes for 448 × 448 resolution is accomplished by two different modules *bboxbox* and *draw box simple*.

Initially, the bbox2box module in jedi_human_count_top.v obtains box coordinate outputs from det_out_filter.

Considering (X, Y) as center of the Box of Width W and Height H, it calculates extreme ends of the Box (X1, X2 and Y1, Y2) for 224 × 224 resolution. It also clamps the coordinate values so that the box remains out of masking area. This is shown in Logic 2.

```
LOGIC 2
X1 = If ((X' - W'/2) < 0)
                               =>
                                          else
                                                (X' - W'/2)
Y1 = If ((Y' - H'/2) < 0)
                                          else (Y' - H'/2)
                               => 0
X2 = If ((X' + W'/2) > 224)
                               =>
                                   223
                                          else
                                                 (X' + W'/2)
Y2 = If ((Y' + H'/2) > 224)
                               => 223
                                          else (Y' + H'/2)
```

The final calculated X1, X2, Y1, and Y2 values for all the boxes in bbox2box are then sent to *draw_box_simple* module through the *osd_back_128x128_human_count* module. The *draw_box_simple* module converts these input coordinates provided for 224 × 224 resolution into 448 × 448 resolution as shown in Logic 3.

For converting from 224 to 448, the coordinates are multiplied with 2. Required offset value is added in coordinate calculations to keep the boxes out of mask area. X1, X2 and Y1, Y2 coordinates are calculated for each Box.

Pixel Counter and Line Counter keep track of the pixels of each line, and lines of each frame. The outer boundary of the box and inner boundary of the box are calculated when Pixel and Line counter reaches to coordinates (X1, X2) and (Y1, Y2) respectively. Calculations are done as per Logic 4.

```
LOGIC 4

Outer Box = (Pixel Count >= (X1 - 1)) and (Pixel Count <= (X2 + 1)) and (Line Count >= (Y1 - 1)) and (Line Count <= (Y2 + 1))

Inner Box = (Pixel Count > (X1 + 1)) and (Pixel Count < (X2 - 1)) and (Line Count > (Y1 + 1)) and (Line Count < (Y2 - 1))
```

Each bounding box is calculated by removing the intersecting area of outer and inner box. The box is only displayed if the

Box-Bitmap for that box is set to 1 (from the *det_st_bbox via bbox2box* module). Box on calculations are as done as Logic 5.

```
LOGIC 5

Box_on[1] = Outer Box[1] and ~Inner Box[1] and Box-Bitmap[1] Box_on[2] = Outer
Box[2] and ~Inner Box[2] and Box-Bitmap[2] . .
```

 $Box_on[20] = Outer Box[20]$ and $\sim Inner Box[20]$ and $Box_on[20]$

The o_box_obj signal is asserted when any of the above Box_on signal is set which is then connected to green_on signal and processed for Bounding Box display through the PC.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.

FPGA-RD-02242-1 0



8.2.4.5. Centroid Tracker

The Centroid Tracker module is used to track the object from beginning to end (left to right). It is also responsible for keeping the cumulative count of objects passing from left to right.

8.2.4.6. OSD Text Display

- The *lsc_osd_text* module provides bitmap of each ASCII character to be displayed with specified position on screen. It takes count of detected Humans.
- It sets an output signal (text_on) when text is to be displayed on the PC through the USB. When text_on is set, YCbCr value for that pixel location is assigned FF, 7F, and 7F respectively values (white color) and sent to USB output instead of original pixel value.

8.2.4.7. USB Wrapper

- The Wrapper_USB3 module is used to transmit 16-bit data to the output 16-bit interface every clock cycle.
- This module takes input data in YCbCr 24-bit format and gives output as 16-bit YCb and YCr format. This module does not change or regenerate input timing parameters.



9. Creating FPGA Bitstream File

This section describes the steps to compile RTL bitstream using Lattice Radiant tool.

9.1. Bitstream Generation Using Lattice Radiant Software

To create the FPGA bitstream file:

1. Open the Lattice Radiant software. Default screen is shown in Figure 9.1.

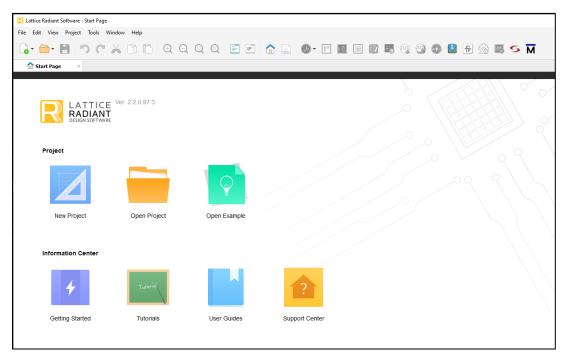


Figure 9.1. Radiant - Default Screen

- 2. Go to File > Open > Project.
- 3. Open the Radiant project file (.rdf) for CertusPro-NX Object Detection Demo RTL. As shown in Figure 9.2, you can also open the project by selecting the yellow folder shown in the user interface.



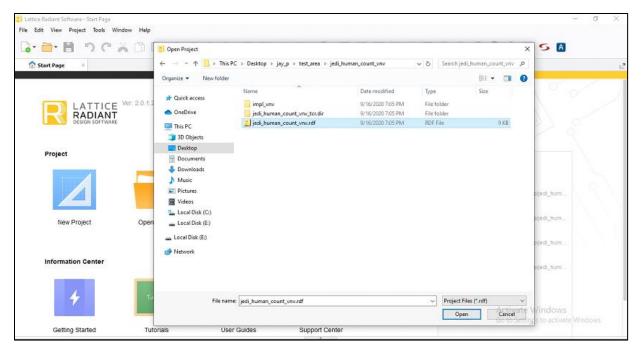


Figure 9.2. Radiant - Open CertusPro-NX Project File (.rdf)

- 4. After opening the project file, check the following points shown in Figure 9.3.
 - The design loaded with zero errors message shown in the Output window.
 - Check the following information in the Project Summary window.
 - Part Number LFCPNX-100-9BBG484I
 - Family LFCPNX
 - Device LFCPNX-100
 - Package BBG484

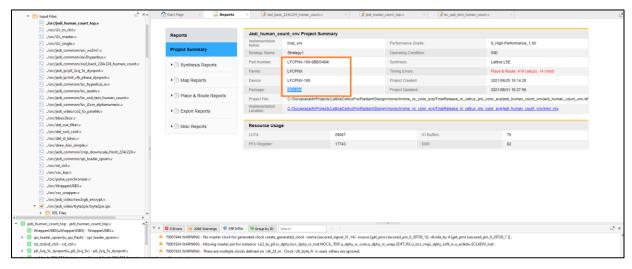


Figure 9.3. Radiant - Design Load Check After Opening the Project File

If the design is loaded without errors, click the Run button to trigger bitstream generation as shown in Figure 9.4.



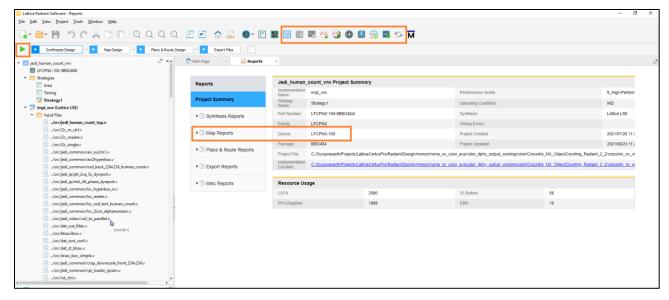


Figure 9.4. Radiant - Trigger Bitstream Generation

6. The Lattice Radiant tool displays *Saving bit stream in ...* message in the **Reports** window, as shown in Figure 9.5. The bitstream is generated at *Implementation Location*, as shown in Figure 9.5.

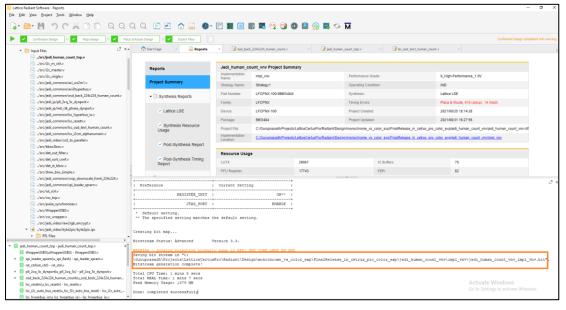


Figure 9.5. Radiant - Radiant Bit File Generation Report Window



9.2. Configuring the IP in Lattice Radiant Software

If you are going to uninstall an old version of the existing IP or if a latest version of the IP should be installed after loading the design without any errors, follow the procedure below and trigger the Bitstream generation *RUN* option. To configure the IP:

1. If the existing IP should be uninstalled, go to IP Catalog. In the IP tree, go to IP > DSP and select your IP. Click Delete as shown. Select Yes as shown in Figure 9.6 and the IP is successfully removed from the tree.

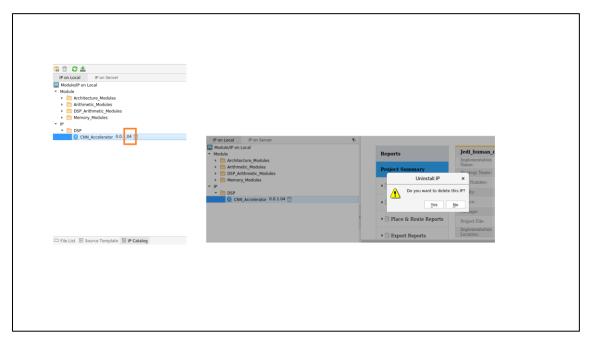


Figure 9.6. Radiant - Uninstall Old IP

2. To install a new IP in the tree, select **Install a User IP**, as shown in Figure 9.7.

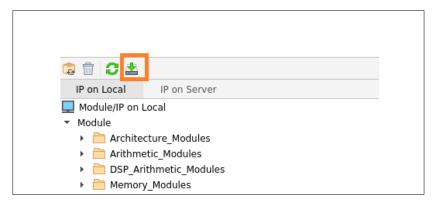


Figure 9.7. Radiant - Install New IP

3. Select your IP package (.ipk) and the IP License agreement window pops up. Click Accept and the user IP is installed in IP tree.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.
All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice.



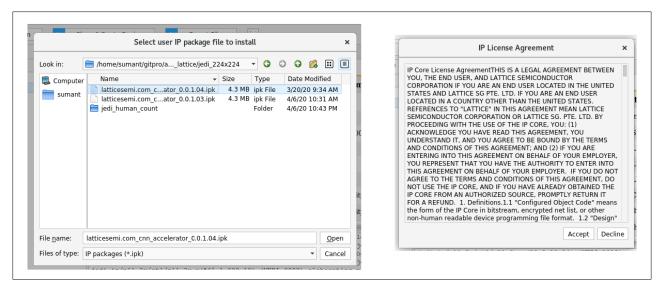


Figure 9.8. Radiant - Select User IP Package to Install

4. After the user IP is installed, the bitstream can be generated by triggering **RUN** button.



10. Programming the Demo

10.1. Programming the CertusPro-NX SPI Flash

10.1.1. Erasing the CertusPro-NX SRAM Prior to Reprogramming

If the CertusPro-NX device is already programmed (either directly, or loaded from SPI Flash), follow this procedure to first erase the CertusPro-NX SRAM memory before re-programming the CertusPro-NX's SPI Flash. If you are doing this, keep the board powered when re-programming the SPI Flash (so it does not reload on reboot).

To erase the CertusPro-NX: device:

1. Launch Radiant Programmer. In the Getting Started dialog box, select Create a new blank project.

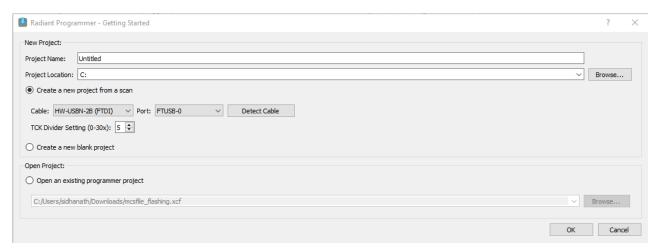


Figure 10.1. Radiant Programmer Default Screen

- 2. Click OK.
- 3. In the Radiant Programmer main interface, select **LIFMD** for Device Family, **LIFCL** for Device Vendor, and **LIFCL-40** for Device, as shown in Figure 10.2.

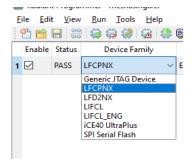


Figure 10.2. Radiant programmer- Device selection



- 4. Right-click and select **Device Properties**.
- 5. Select **JTAG** for Port Interface, **Direct Programming** for Access Mode, and **Erase Only** for Operation, as shown in Figure 10.3.

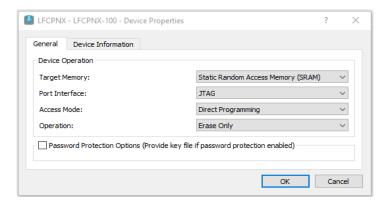


Figure 10.3. Radiant Programmer – Device Operation

- 6. Click **OK** to close the Device Properties dialog box.
- 7. Click the **Program** button <a> to start the erase operation.

10.1.2. Programming the CertusPro-NX Board

To program the CertusPro-NX Voice and Vision SPI Flash:

- 1. Ensure that the CertusPro-NX Voice and Vision device SRAM is erased by performing the steps in the Erasing the CertusPro-NX SRAM Prior to Reprogramming section.
- 2. In the Radiant Programmer main interface, right-click the CertusPro-NX Voice and Vision row and select **Device Properties**.
- 3. Apply the settings below:
 - a. Under Device Operation, select the options below:
 - Port Interface JTAG2SPI
 - Access Mode Direct Programming
 - Operation SPI Flash Erase, Program, Verify
 - b. Under Programming Options, select the bitstream file.
 - c. For SPI Flash Options, select the Macronix 25L12833F device, as shown in Figure 10.4.



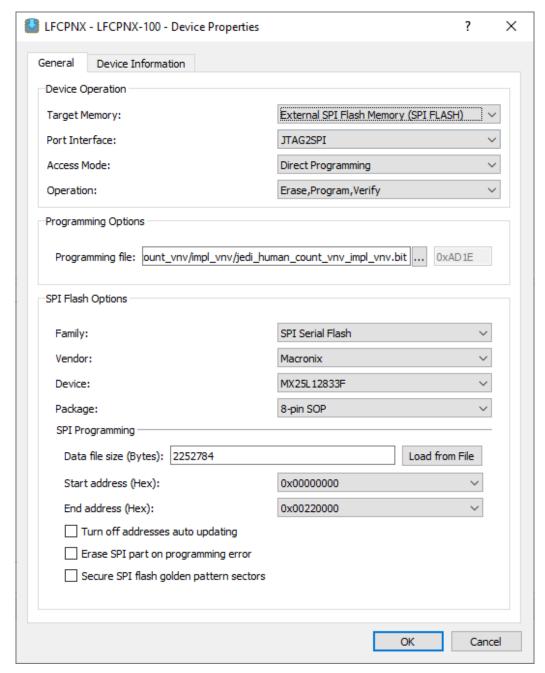


Figure 10.4. Radiant Programmer – Selecting Device Properties Options for CertusPro-NX Flashing

- d. Click **Load from File** to update the Data file size (bytes) value.
- e. Ensure that the following addresses are correct:
 - Start Address (Hex) 0x00000000
 - End Address (Hex) 0x00220000
- 4. Click OK.
- 5. Press the **SW5** push button switch before clicking the **Program** button, as shown in Figure 10.5. Hold it until you see the *Successful* message in the Radiant log window.



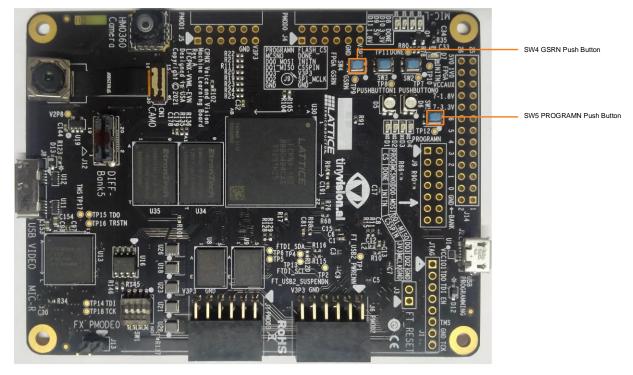


Figure 10.5. CertusPro-NX Flashing Switch - SW4 Push Button

- 6. Click the **Program** button to start the programming operation.
- 7. After successful programming, the **Output** console displays the result, as shown in Figure 10.6.



Figure 10.6. Radiant Programmer – Output Console

10.1.3. Programming sensAl Firmware Binary to the CertusPro-NX SPI Flash

10.1.3.1. Convert Flash sensAl Firmware Hex to CertusPro-NX SPI Flash

To program the CertusPro-NX SPI flash:

- 1. Ensure that the CertusPro-NX device SRAM is erased by performing the steps in the Erasing the CertusPro-NX SRAM Prior to Reprogramming section before flashing bitstream and sensAl firmware binary.
- 2. In the Radiant Programmer main interface, right-click the CertusPro-NX row. Select **Device Properties** to open the dialog box, as shown in Figure 10.7.
- 3. Select SPI FLASH for Target Memory, JTAG2SPI for Port Interface, and Direct Programming for Access Mode.
- 4. For Programming File, select the CertusPro-NX sensAl firmware binary file after converting it to hex (*.mcs).
- 5. For SPI Flash Options, follow the configurations in Figure 10.7.

© 2022 Lattice Semiconductor Corp. All Lattice trademarks, registered trademarks, patents, and disclaimers are as listed at www.latticesemi.com/legal.

All other brand or product names are trademarks or registered trademarks of their respective holders. The specifications and information herein are subject to change without notice



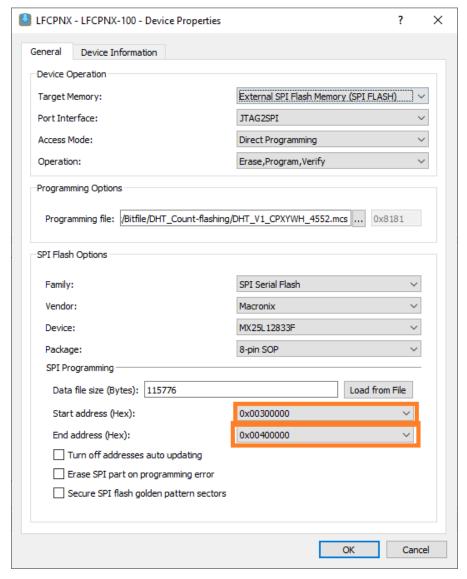


Figure 10.7. Radiant Programmer - Selecting Device Properties Options for CertusPro-NX Flashing

- 6. Click **Load from File** to update the data file size (bytes) value.
- 7. Ensure that the following addresses are correct:
 - Start Address (Hex) 0x00300000
 - End Address (Hex) 0x00400000
- 8. Click OK.
- 9. Press the **SW5** push button switch. Click the **PROGRAMN** push button and hold it until you see the *Successful* message in the Radiant log window.
- 10. Click the **Program** button with to start the programming operation.
- 11. After successful programming, the Output console displays the result, as shown in Figure 10.8.





Figure 10.8. Radiant Programmer – Output Console



11. Running the Demo

To run the demo:

- 1. Double-click **Executable** in the *Lattice_Object_Detector_App* folder to open the PC Application.
- 2. Connect the board using the USB3 cable. Once the USB device is connected, open the **Camera** button. The camera output from the board is displayed on the right side of the application window.
- 3. Press the **Play** button to start the video playback.
- 4. Hold the board in such a way that the camera sees the video playback. Objects detected will have bounding box drawn over them and cumulative count is displayed on the bottom left corner of the Camera output window.

The screen displays the video image with a bounding box when a DHT object is detected. The red LED D6 lights up when an object is detected. Cumulative count is displayed to the left of the Camera output window.

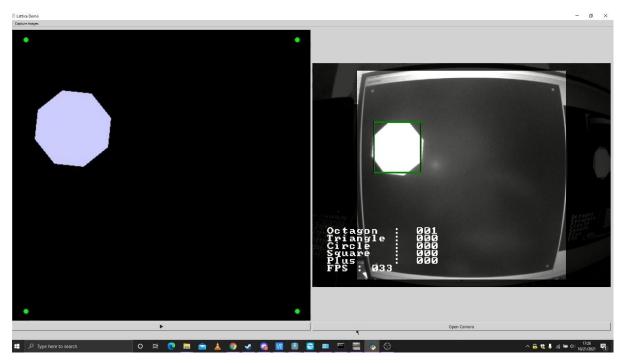


Figure 11.1. Running the Demo



Appendix A. Other Labelling Tools

Table A.1 provides information on other labelling tools.

Table A.1. Other Labelling Tools

Software	Platform	License	Reference	Converts To	Notes
annotate-to- KITTI	Ubuntu/Windows (Python based utility)	No License (Open source GitHub project)	https://github.com/SaiPrajwal95/annotate-to- KITTI	KITTI	Python based CLI utility that you can clone and launch.
LabelBox	JavaScript, HTML, CSS, Python	Cloud or On- premise, some interfaces are Apache-2.0	https://www.labelbox.com/	json, csv, coco, voc	Web application
LabelMe	Perl, JavaScript, HTML, CSS, On Web	MIT License	http://labelme.csail.mit.edu/Release3.0/	xml	Converts only jpeg images
Dataturks	On web	Apache License 2.0	https://dataturks.com/	json	Converts to json format but creates single json file for all annotated images
Labelimg	ubuntu	OSI Approved:: MIT License	https://mlnotesblog.wordpress.com/2017/12/ 16/how-to-install-labelimg-in-ubuntu-16-04/	xml	Need to install dependencies given in reference
Dataset_ annotator	Ubuntu	2018 George Mason University Permission is hereby granted, Free of charge	https://github.com/omenyayl/dataset- annotator	json	Need to install app_image and run it by changing permissions



References

- Google TensorFlow Object Detection GitHub
- Pretrained TensorFlow Model for Object Detection
- Python Sample Code for Custom Object Detection
- Train Model Using TensorFlow



Technical Support Assistance

Submit a technical support case through www.latticesemi.com/techsupport.



Revision History

Revision 1.0, January 2022

Section	Change Summary
All	Initial release.



www.latticesemi.com