# Robotics Inference

A. Cámara

**Abstract**: Two image inference projects are studied. In first place, a given project with images of a conveyor belt where bottles and candy boxes are classified using deep neural networks. The second project classifies images of the hand game "stone, paper, scissors". Popular neural networks has been chosen and carefully trained to get good results in these classification problems.

## Introduction

Machine vision is the technology and algorithms used to recognize objects in images. It is an important field of study in robotics because vision can give machines a lot of information about their environment.
In first place, vision requires acquisition of images, usually using cameras and lightning that has been designed to provide the differentiation required by the following processing. Once good images are acquired, multiple stages of processing are applied to get the desired classification result. A typical image processing starts with filters that modify the image, followed by objects extraction, then objects data extraction to end with data analysis to get final results.

A modern approach to image classification is the use of deep neural networks (DNNs). They allow to learn all the necessary processing stages making faster and easier prototyping a machine vision system. DNNs require a lot of examples of classified images to learn and evaluate progress and a big quantity of computing power but they avoid the engineer the hard image processing work.
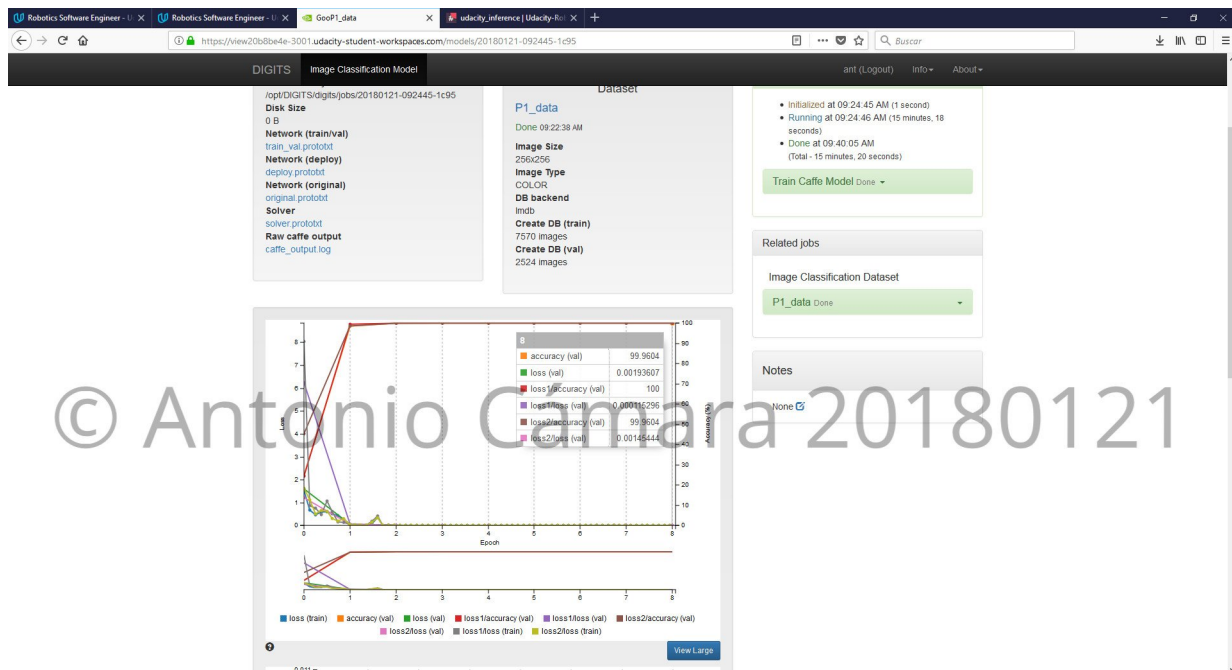
Recent success of DNNs has brought a good number of predesigned networks that can be used for fast prototyping like LeNet, AlexNet and GoogleNet. Moreover, there are several libraries (CNTK, Tensoflow, Caffe) and development environments like NVIDIA DIGITS that allow to make, train and deploy a DNN project in very little time.

This work is divided in two projects. The first one requires to build a DNN using DIGITS to classify a given set of images. There are 3 different types of images: candy boxes, bottles and nothing (empty conveyor belt). Images are 256x256 color pictures. Once the training set is loaded into DIGITS, a network is chosen to be trained in order to achieve an inference time of 10ms or less and an evaluation accuracy of more than 75 percent.
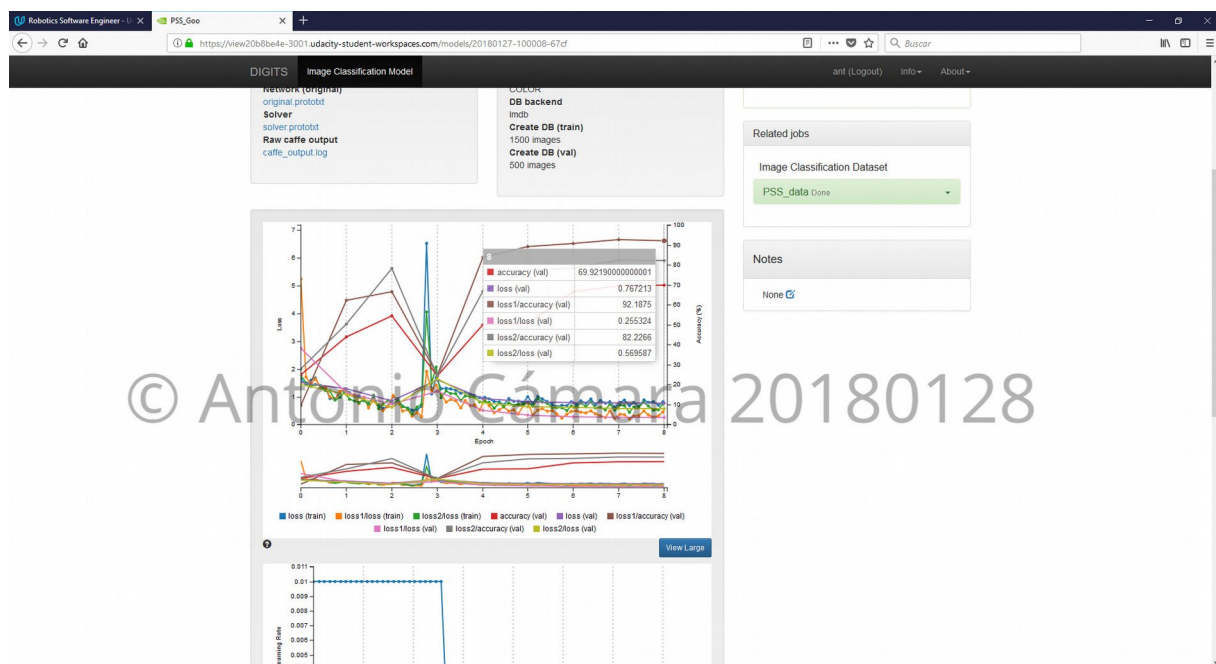
In the second project, a complete image inference project is implemented. Recognizing hand gestures of "stone, paper, scissors" game has been chosen as the goal of the project due to the importance of robots and human interaction. Gestures is a way of communication people use constantly and it is important for a social robots to understand and react to them. Our system should have into account 4 different classes: 3 gestures and nothing. After obtaining the necessary images, DIGITS is used again to train a network to do the classification. At the end, the network is deployed to a Jetson TX2 system and tested working in real time.
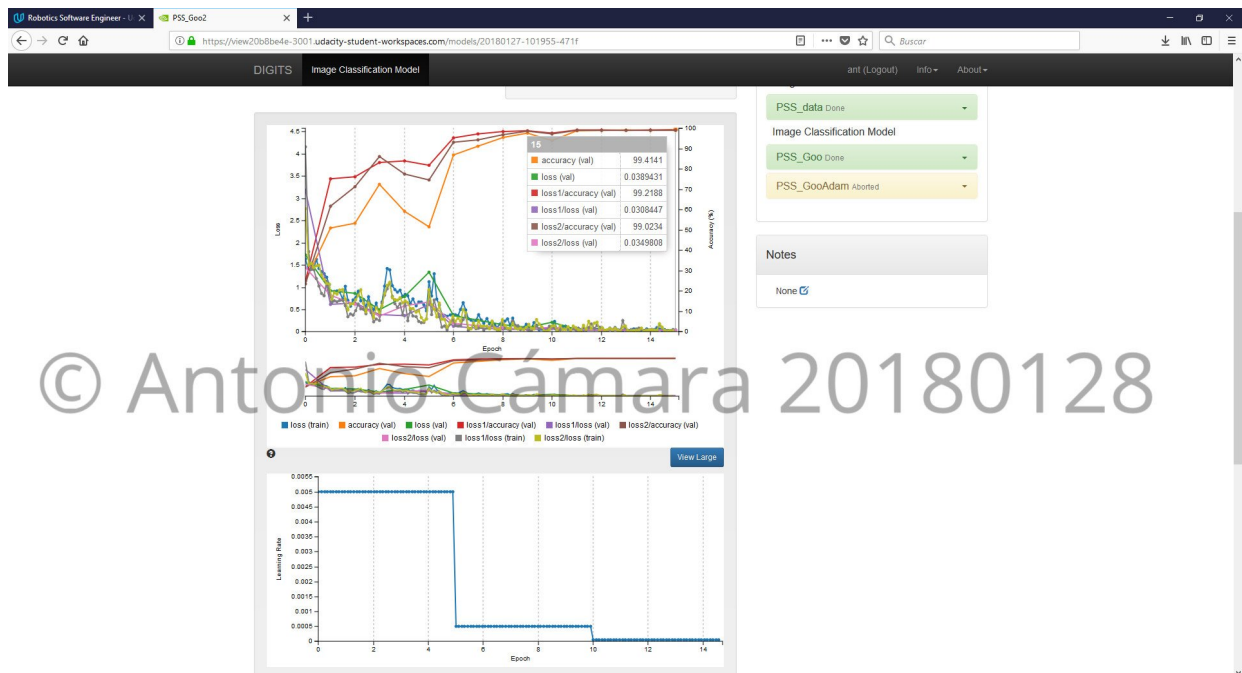
## Background / Formulation

In these projects DIGITS has been used to train DNNs that do image classification. In the robotics inference project, GoogleNet has been chosen due to the high accuracy required of 75 percent in the evaluation. Using that network with the default training parameters and 8 epochs (observed as enough in previous uses of DIGITS) was enough to meet the expectations.

Because GoogleNet is the best in accuracy of the pre-supplied networks in digits and given the good results of it with the first project, it has been used again to train the "stone, paper, scissors" project. Using the same parameters to train the second project gave a not very good training result of 69 percent of validation accuracy.

After lowering the base learning rate to 0.005 and increasing the number of epochs to 15 (to avoid problems of the first training doing it slower and training more to get to the desired validation accuracy), a good training result was obtained.

## Data acquisition

In the second project, four classes are learned from images: nothing, stone, paper, scissors. 500 RGB images were obtained per class. They are enough to make a fast DNN prototype.

Images were acquired using Jetson TX2 and its included camera. They were obtained with a 256x256 resolution and a fixed background. A controlled lightning environment was used with a lateral natural light source.



Color images with 256x256 resolution were chosen to train them easily with the selected GoogleNet DNN.

## Results

The robotics project model (P1) achieve the required 75 percent of accuracy and less than 10ms inference time as shown in the attached picture.

The second model obtained a good training accuracy and has been tested deploying it on the Jetson TX2. Using the inference-camera software, a good classification has been done. Nothing class is always well labeled. Paper and stone are almost always well classified. Scissors class is not always well identified been labeled as stone o paper in about 40% of the times.

## Discussion

The robotics inference project has achieved its desired results with the first attempt because of the good election of the network (GoogleNet). Better accuracy results could have been achieved using deeper models like ResNet or Inception.

Second project has good accuracy results except for scissors class that is sometimes mistakenly labeled. The training set of images is not very large, so with more images it can show better results. Another way to improve accuracy could be using better networks like ResNet or Inception.

As shown in the analysis (Canziani et al, 2016) accuracy and inference time are in a hyperbolic relationship. A little better accuracy costs a lot of computational time. Our robotics inference project requires real time response to the eyes of the person that interacts with it, therefore we need 0.2s of inference time or less. Having into account that the trained GoogleNet shows a maximum inference time of 0,05s there is margin to improve accuracy using Inception which needs 3 times more inference time.

After testing the network in the Jetson, it is recommended a larger testing set of images with different backgrounds and lights to achieve a high quality classifier that can perform well in a wide range of situations. Increasing the number of training images does not affect inference time but can produce better final results.

## Conclusion / Future work

Two image inference projects have been done. The first project has soon achieved desired results because of the good election of the network. The second project has been trained and deployed successfully. It has good results but has to improve scissors class results witch can be done with more training images.

The second project can be used in a social robot that will be able to play the popular stone, paper, scissors game like another human using only vision. That can be useful in leisure places like parks or in interactive museums. Another useful application. Another useful application of this project is to help working with children with social difficulties like autistic children

As a future work, a more complex gestures recognition can be done. Understanding the language of signs or gestures people do while they interact with others, could be a very important skill for robots that should work with humans. Therefore, the suggested future work is doing a more capable gestures classifier to allow robots understand humans better.