

Using Deep Neural Networks for autonomous UAV navigation in an apple orchard

BRESILLA, Trim

Abstract—With the increase of population in the world, the demand for quality food is increasing too. One of the biggest and the base of raw food production comes from Agriculture. In the recent years, due to this demand, and other environmental factors have heavily influenced the way agricultural production is done. Automation and robotics for fruit and vegetable production and monitoring has become the new standard. In this paper we discuss an autonomous Unmanned Aerial Vehicle (UAV) that would be able to navigate through the rows in an orchard environment. The UAV is comprised of a flight controller (PIXHawk), a microcontroller (Arduino) for analog reading from different sensors, and an On-Board Computer OBC (Raspberry Pi gen. 3). Pictures are taken through PiCamera and streamed through WiFi to a Ground Control Computer GCC running a convolutional neural network model. Based on prior trainings the model sends back to the UAV the direction to the drone using MAVLink protocol, thus performing autonomous navigation.

Index Terms—Robotics, Agriculture, Udacity, Orchard, Deep Learning.

1 INTRODUCTION

Automation in Agriculture is recent years a highly increasing. Even though Agriculture as one of the oldest occupation in the world, has seen many changes during centuries. Before Industrial Revolution was estimated that more than 80% of population were working as farmers, while now is estimated that number to be 2%. One of the dominant changes that characterizes a growing economy is the proportionate decline in Agriculture Sector. This phenomenon is commonly attributed to two facts: food is not as demanding as other goods and services, and the rapid development of new farming technologies lead to expanding food supplies per hectare and per worker.

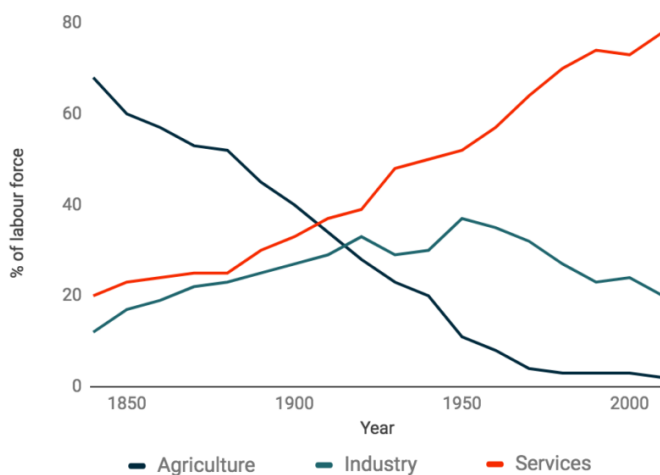


Fig. 1: Distribution of the labor force by sector, 1850-2000

The decrease shows significance every step of the industrial revolution. Industry 1.0 brought mechanisation, water management etc.; Industry 2.0 brought steam engines, electricity, protected environments etc.; Industry 3.0 brought

computers and smart appliances, GPS tractors and so on; while Industry 4.0 is foreseen to be the most significant one, bringing AI, Automation and Robotics.

The challenge of how we'll feed the evergrowing world population in the future - is sustainable, cost-effective and most importantly environmentally friendly. In order to feed 9.5 billion people that Food and Agriculture Organisation (FAO) predicts to inhabit the planet by 2050 while climate change is making more difficult to grow crops - is going to be done by Smart Farming, a high-tech and AI driven agricultural management system. Agriculture is highly repetitive, and such, many tasks can and are being automated. Individual agricultural activities on the farm takes effort, for example planting, maintaining, and harvesting crops need money, energy, labor and resources. What if we can use technology to replace some of the human activities and guarantee efficiency? That's where artificial intelligence comes in. Agriculture is slowly becoming digital and AI in agriculture is emerging in three major categories, (i) agricultural robotics, (ii) soil and crop monitoring, and (iii) predictive analytics [1]

For a farmer robot to be fully autonomous, it needs to navigate through quiet diverse and harsh environment without the human supervision, then perform a set of actions at specific location like: pickin a fruit, evaluate a site, spray pesticide, cut branches, plant a seed, image and scan a whole plant and take specified measurement. Controlled environments like greenhouses are more manageable because of controllable environment and better engineered, where the sensor measurements produce less noise. Whereas outdoor environment are much harsher and generally not controllable, thus making far more difficult than indoor environments. Most of outdoor robot are equipped with GPS for sensing the location, but due to accuracy, they are often companioned with other sensors like IMUs, 3DCameras, Rotary Encoders to create a sensor fusion for a much precise action taking process. Robots nowadays are wirelessly connected to a central operator to both receive updated in-

structions regarding the mission, and report status and data. However, making an autonomous farm robot requires clever controllers, localisation, communication and action taking systems. The technology is similar to that of autonomous cars applied to agtech. Where it differs is that farming robots often need to manipulate their environment, picking vegetables or fruits, applying pesticides in a localised manner, or planting seeds. All these tasks require sensing, manipulation, and processing of their own.

In fruit production, as is with all other fields of agriculture, crop monitoring is extremely important as there can be estimation ahead of time thus making to the farmer very easy to plan logistics and distributions. In this paper is discussed an autonomous Unmanned Aerial Vehicle (UAV) flying under tree canopy, between two rows and under the anti-hail nets. In order for the robot to successfully follow the row, it has firstly to know the orchard and where is the starting row, then has to perceive the path between two rows while maintaining the altitude and avoiding any collision with lateral branches from the trees and avoiding any other obstacles. The approach taken in this paper is that it considers the whole navigation as a classification task. By analysing the front face camera images, by using a convolutional neural network to classify the video / image frames stream into direction with respected weight on a single shot.

2 BACKGROUND / FORMULATION

Orchards nowadays are very complex with multiple components and management procedures moving during vegetative period of the plant. There are many management decisions that often change structure and visuals of the orchard, in addition to the nature's randomness and complexity, thus making it as an ever changing organism. In this perspective, hard-coding algorithms for specific task where randomness is infinitely high is not a good approach. The path between two rows, is maintained by farmers in different ways, differently during the season, the same with the canopy, where plants start without leaves and then later growing them. The robot itself has to fly under anti-hail nets, 1.5m above the ground, making the whole path as a corridor. Using deep learning approach the model is able to accommodate for changes and progressively learn how to navigate even when new scenarios are being dealt with.

2.1 Materials

The UAV/robot uses a RaspberryPi gen.3 as a On-Board Computer OBC which in turn is connected to flight controller - a PixHawk board - through serial link. The OBC due to performance limitations can not run deep-learning algorithms on its own, however it is an excellent power efficient computer running full LINUX inside. It serves as intermediary layer between flight controller and other devices like more powerful computer, arduino boards with different sensors, camera, and other communication devices. The communication is done through MAVLink protocol and MAVProxy is used to multiply the device infos and status to other devices connected to it through wireless hotspot. As it runs full LINUX and Robotics Operating System ROS, with the help of MAVROS package, the whole UAV is controlled



(a) LEFT



(b) STRAIGHT

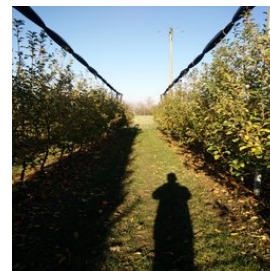


(c) RIGHT

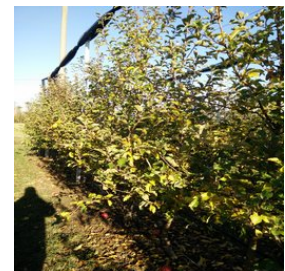
Fig. 2: Images taken from UAV - winter (trees without leaves)



(a) LEFT



(b) STRAIGHT



(c) RIGHT

Fig. 3: Images taken manually - autumn (trees with leaves)

as any robot inside ROS. OBC's camera is a RaspPi Cam gen.2 which makes the image processing and encoding in its own chips, thus making it very easy to directly stream in network.

The input images taken from front-facing camera, are sent to Ground Control Computer GCC (CUDA capable) through WiFi. The computer runs the picture stream through a trained model that has three outputs: right, left, and straight. The moving average of three outputs is sent as MAVLink RC_CHANNEL_OVERRIDE through OBC to PixHawk.

2.2 Data Acquisition

Data collection is done through the camera of the OBC in the UAV moving inside the orchard, controlled manually with radio controller. When OBC starts, it automatically launches few scripts (the process is managed with crontab):

- 1) Connect to known WiFi if any exist, else create hotspot
- 2) Use MAVproxy protocols to sent and receive flight plans and commands
- 3) Automatically output camera stream to network

In the GCC, the camera feed is piped from network with netcat to a python program. The program takes the camera input, divides in frames, labels it and saves in a proper dataset.

The UAV is flown very carefully inside the rows while streaming the camera feed to GCC. The flight is done many times in different rows and directions. However there have been three modes/categories, and for each mode a thousand pictures/frames were taken and labeled accordingly:

- 1) LEFT: The drone would fly closer to the left row, and/or yawed (facing) the left side. Data collected were labeled as LEFT, model would return LEFT and RIGHT CHANNEL OVERRIDE would be sent. Fig 3a.
- 2) STRAIGHT: The drone would fly the best position as much as it can, in the middle of the row, facing straight and having both rows symmetrical to each other. Data collected were labeled as STRAIGHT, model would return STRAIGHT and FORWARD CHANNEL OVERRIDE would be sent. Fig 3b.
- 3) RIGHT: The drone would fly closer to the right row, and/or yawed (facing) the right side. Data collected were labeled as RIGHT, model would return RIGHT and LEFT CHANNEL OVERRIDE would be sent. Fig 3c.

Pictures were captured during daytime in early late winter of 2018. Daytime is important as the RasPiCam is very sensitive to light quality and light exposure.

In addition to this dataset, another set of images is used. The later set is taken manually (using smartphone camera), but during autumn of 2017, while the tree had leaves and the chlorophyll was still green. The set has 100 images per mode/category. This number discussed in results in details, proved to be very small. However the same model is run through both sets separately, and then together too.

2.3 The Model

To better manage different datasets and models, Nvidia's Deep Learning GPU Training System (DIGITS) is used. DIGITS is not in itself a machine-learning framework, rather is a wrapper for most used frameworks available. It simplifies the commonly machine-learning tasks such as managing datasets including train/validation/test splitting, designing and training different neural networks (on CUDA capable GPUs), real-time monitoring of the training process and visualisation of the process.

GoogLeNet is used as a DNN classifier. Because of the use of Inception modules, GoogLeNet is more versatile and

computationally less expensive. A simple 3x3 kernel with 256 input channels and 256 output, would have an amount of $9 \times 256 \times 256$ calculations. Such a network where every output is connected to every input, is referred as dense connection. And while in most of CNNs activation layer for those connection is often either zero or not valuable, proving that not all those input channels are connected to output ones. Despite many techniques developed to cut off those unnecessary connections, the computation needed is huge. Inception module of the model we chose to work with approximates a sparse CNN with a normal dense construction, and since the effective number is low (because of zeros and unnecessary activations) the number of convolutional filters is kept small. In addition it uses convolutions of different sizes to capture details at different scales (5x5, 3x3) And it uses the so called bottleneck layer 1x1, for reduction of computation requirements. GoogLeNet is a 27 layer deep CNN, with 22 convolution and inception layers and 5 pooling layers. However the overall number of the independent blocks is well over 100.

type	patch size/ stride	output size	depth	#1x1	#3x3 reduce	#3x3	#5x5 reduce	#5x5	pool proj	params	ops
convolution	7x7/2	112x112x64	1							2.7K	34M
max pool	3x3/2	56x56x64	0								
convolution	3x3/1	56x56x192	2		64	192				112K	360M
max pool	3x3/2	28x28x192	0								
inception (3a)		28x28x256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28x28x480	2	128	128	192	32	96	64	380K	304M
max pool	3x3/2	14x14x480	0								
inception (4a)		14x14x512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14x14x512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14x14x512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14x14x528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14x14x832	2	256	160	320	32	128	128	840K	170M
max pool	3x3/2	7x7x832	0								
inception (5a)		7x7x832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7x7x1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7x7/1	1x1x1024	0								
dropout (40%)		1x1x1024	0								
linear		1x1x1000	1							1000K	1M
softmax		1x1x1000	0								

Fig. 4: GoogLeNet model structure

2.4 Training

Nvidia's DIGITS dealt with splitting training, validation and testing set. We decided to keep 15% of images for training and 5% for testing. Each images was of size 256x256. Later on we augmented (manually with a script), all LEFT images could be mirrored and produce RIGHT image and vice versa, while STRAIGHT images when mirrored created another STRAIGHT image. The model was trained for 60 epochs in Amazon's AWS S3 instance with NVIDIA's Tesla K80 with 12GB Memory. An Adaptive Movement Solver (ADAM Optimiser) with 0.002 base learning rate was chosen. A sigmoid decay of gamma 0.08 learning rate with 60 steps is used. After training for few hours, the model is downloaded as a pre-trained model, and transferred to a less powerful local computer (to the ground control computer GCC) this runs the transferred model through Nvidia's GTX 960 GPU.

2.5 Navigation policy

As briefly mentioned above, the OBC (on-board computer) after capturing the data sends them through WiFi to the

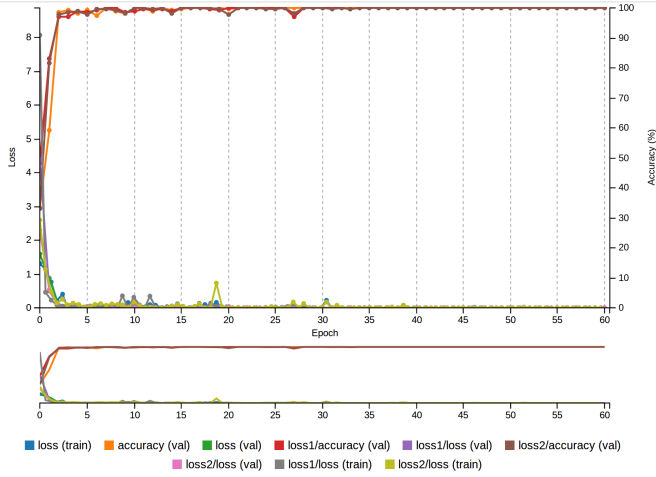


Fig. 5: Training on 60 epochs

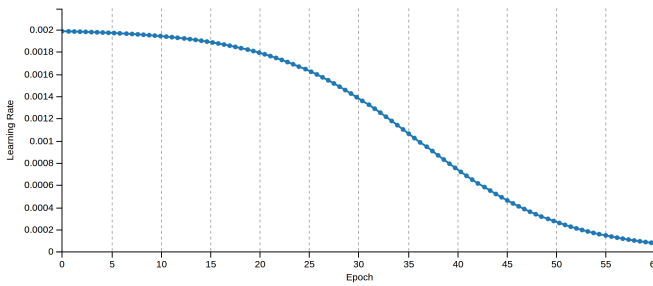


Fig. 6: Learning rate

GCC (ground control computer) to run through the model. Since the OBC is not a powerful computer, this step is necessary. Pictures received from the network are piped to the Coffe2 model through Netcat (a GNU/Linux utility for reading from and writing to network connections using TCP and UDP protocol). OpenCV captures the Netcat pipe and serves to the model as frames (a single image where the model performs the recognition).

The model that has three outputs (LEFT, RIGHT, STRAIGHT), and based on recognition/classification confidence to each frame the values from 0 (not confident) to 1 (very confident) per each class. The outputs are changed into RC_CHANNEL_OVERRIDE MAVLink protocol (with values from 1000 to 2000 per channel). Those values go through a moving average script, that takes the values for several frames (in our case all frames per one second) and extract the moving average or running mean values for each channel. Running mean of each channel is sent through ZeroMQ (an highly efficient, high-performance and reliable asynchronous messaging library) to the PID (Proportional Integral Derivative) controller that we build in the OBC.

2.6 PID and Sensor Fusion

(IN PROGRESS, when we finish more work on localisation and PID and add some references, the paper will be ready for publication, hopefully sooner than later)

(we are integrating the channel overrides, and other sensors (GPS, UltraSonic, LIDARv3 and Flow camera - facing down) all together and with MAVROS and ROS to a built



(a) LEFT



(b) STRAIGHT



(c) RIGHT

Predictions	
LEFT	100.0%
RIGHT	0.0%
STRAIGHT	0.0%

Predictions	
STRAIGHT	100.0%
RIGHT	0.0%
LEFT	0.0%

Predictions	
RIGHT	100.0%
STRAIGHT	0.0%
LEFT	0.0%

Fig. 7: Predictions per class

PID controller, for object avoidance, auto navigation and so on.)

(But as for now, we tested without them, just with simple CHANNEL-OVERRIDE. All this is very difficult inside a very small space - a 3x3m corridor).

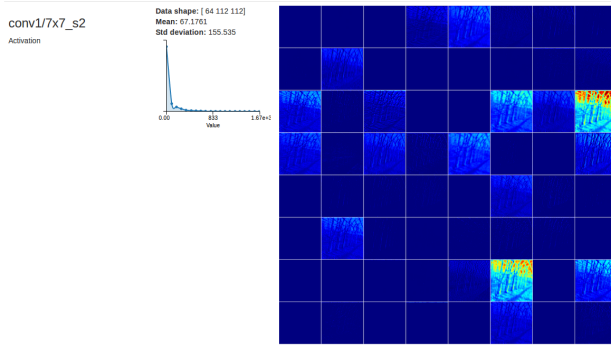
3 RESTRICTIONS

It does not go without mentioning that one of the major restrictions in this research is the uncontrolled environment. Since the navigation relies most on camera input (in addition to sonar and IR sensors), the light and sun exposure, the season of the year, and the vegetative growth of the trees makes the model to vary a lot. In this context, data on different weather conditions, different season and when tree leaves have different coloring and are in the tree itself, would make the model fit more scenarios and be more accurate.

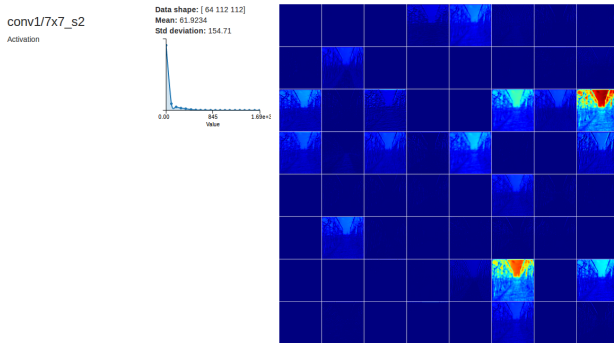
4 RESULTS

4.1 Deep-Learning model

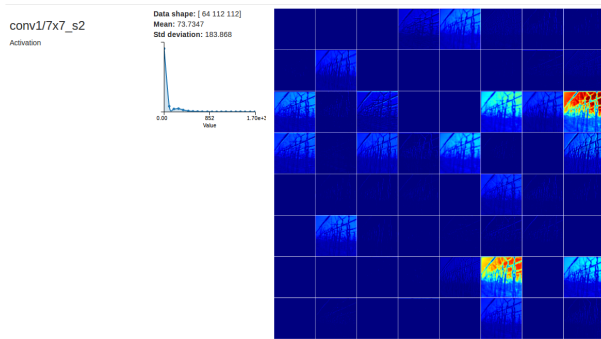
We run the model through test dataset and reached accuracy of more than 85%. In the first models, we trained the winter batch (set 1 of images) and autumn batch (set 2 of images) separately, and when tested images from second dataset in the first model, the accuracy is very low, under 20%. And the same vice versa. The third model, that runs on all images together, reaches an accuracy more than 76%. The testing set comprised of well over 200 images that were manually separated from DIGITS dataframe. In cases where the front path is not seen in the image, the model gives equal (or



(a) LEFT



(b) STRAIGHT



(c) RIGHT

Fig. 8: Activations of last convolutional layer

almost equal) confidence score to both RIGHT and LEFT, as when path is not in the frame, they are almost identical, later this would be handled from the PID on the OBC.

4.2 Autonomous navigation

OBC is connected through serial to flight controller and sends direct navigation messages as CHANNEL_OVERRIDE. Its astonishingly amazing to see how the OBS and flight controller handle the commands. The drone's facing down LIDAR and Flow camera keep it very steady in the position (in LOITER MODE - APM Flight stack) while the drone moves forward direction.

5 CONCLUSION / FUTURE WORK

Using deep learning and convolutional neural networks makes it very easy for UAV and other mobile robots to navigate paths and spaces not explicitly programmed. In

orchard environments, where the top of the orchard is usually covered with hail nets, its impossible to asses the trees and production for above. Our work shows that using a UAV inside the rows where, in the future this drone would handle specific sensors for precision measurement of fruit quality, fruit size estimation, tracking maturity of the fruits, controlling for diseases...

A better model, taking into consideration all the restrictions mentioned before, will contribute to a robustness of classifications (in turn batter navigation).

REFERENCES

- [1] A. Giusti, J. Guzzi, D. C. Cirean, F.-L. He, J. P. Rodriguez, F. Fontana, M. Faessler, C. Forster, J. Schmidhuber, G. D. Caro, D. Scaramuzza, and L. M. Gambardella, "A machine learning approach to visual perception of forest trails for mobile robots," *IEEE Robotics and Automation Letters*, vol. 1, pp. 661–667, jul 2016.