



Hochschule
Bonn-Rhein-Sieg
University of Applied Sciences



R&D Project Proposal

Testing and Verification of Deep Neural Networks using Synthetic Dataset

Rashid Ahamed Meeran Tajdeen

Supervised by

Prof. Dr. Nico Hochgeschwender (h-brs)

Deebul Sivarajan Nair (h-brs)

May 2022

1 Introduction

- An introduction to the general topic you are covering.

Though Deep Neural Networks (DNNs) are being used extensively in many fields including safety critical systems such as autonomous driving and medical diagnostics, there are chances for the DNN to exhibit erroneous behaviours.

(Should i talk a bit about the differences between verification and testing?)

Testing a DNN is not that feasible as classical software testing. In a classical software, we know whats happening at every point of time. When you have a problem, you exactly know where it occurs and how to rectify them. Hence we would use white-box testing which touches the internals of the software. But in a AI software, testing the software internals is of no use. After a model is trained, we have no clue whats happening inside all those networks, it just happens. Hence the best option is to follow the black-box testing approach in which we test the functionality of the software from the outside.

A domain specific language (DSL) is a computer language that is customised for a specific application. A DSL with necessary functions definitions can be also used to map user defined requirements into usable code for testing. The advantages of using DSL are:

- We would be creating usable functions to map the DSL which reduces the complexity when it comes to testing those functions from an software testing perspective.
 - It is user-friendly as the user need not understand all the interior technical parts to use the DSL.
 - A DSL is entirely scalable to use a new software. Only thing to add is the functions from the new software similar to the old ones.
- What is the goal of this RnD?

The goal of this RnD is to test deep neural networks. The focus is not on testing the performance of a deep neural network, instead to test its capability. To aid this testing, we use behaviour driven development methods which

takes in both requirements from test engineers and the capabilities of blender thereby generating synthetic test dataset in blender for testing the learned DNN models.

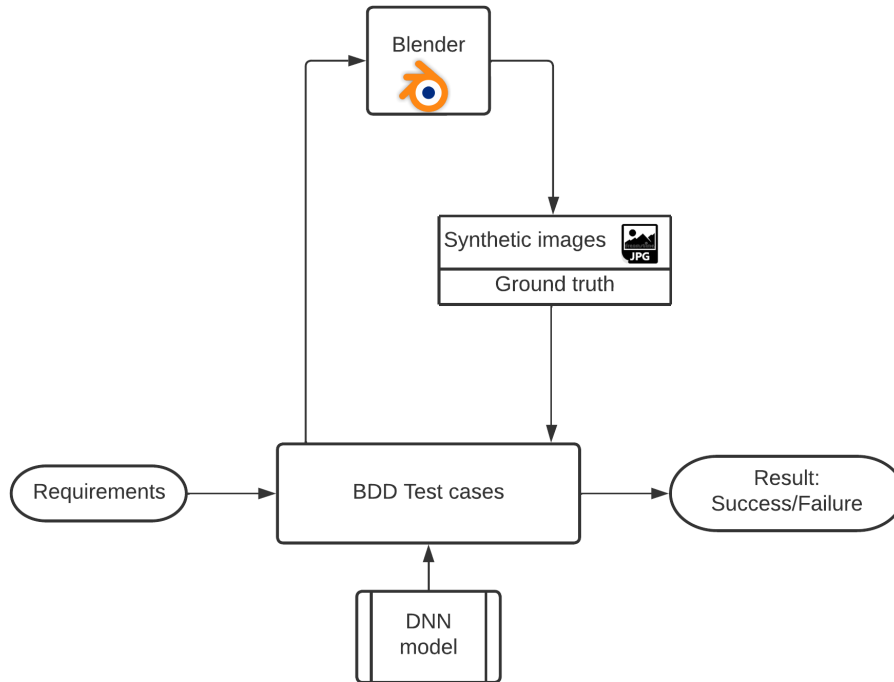


Figure 1: The basic block diagram of the testing process

- Why is it essential to test a DNN?
 - The most verified process of testing deep neural networks is to use test datasets. This aspect of testing is more viable and trusted than the other aspects of testing.
 - In most DNN image classifiers, the trained models misunderstands one class with the other, or even sometimes show one-sided biases. Hence the properties of the class are violated and misinterpreted.
 - Some DNNs perform well on simple tasks, but when it comes to the point to expand its use case, the user might begin to face capability issues from the model.

- We could never say a DNN model is perfectly trained and can perform well on any situation, because the possibilities real world inputs to this model is infinite.

- An example of the testing process:

One sample requirement from the test engineer could be

Given there is a DNN classifier model trained using a particular dataset,

When the model is tested using synthetic images with varying illuminance of 100-200 lux

Then the model must correctly classify the images

During the when step of this requirement, the desired dataset along with its ground truth is generated and classified using the model, and in the then step the results are verified.

1.1 Problem Statement

Though Deep Neural Networks (DNNs) are being used extensively in many fields including safety critical systems such as autonomous driving and medical diagnostics, there are chances for the DNN to exhibit erroneous behaviours.

Unlike classical software testing, in AI testing after a model is trained, we have no clue whats happening inside all those networks, it just happens. Hence the best option is to follow the black-box testing approach in which we test the functionality of the software from the outside. During this process, this work aims to answer the following research questions.

Research Questions:

RQ1 What are all the parameters that can be varied in blender for dataset generation?

RQ2 What are the requirements for testing a vision based DNN?

RQ3 What are the requirements from the test engineers for verifying DNNs?

RQ4 How to design a DSL for DNN testing based on requirements from test engineers and capabilities of blender?

RQ5 How DSL generated synthetic dataset benefits AI testing?

2 Related Work

- In this paper [1], the authors investigate black-box input diversity metrics as an alternative to white-box coverage criteria. They first select and adapt three diversity metrics and study, in a controlled manner, their capacity to measure actual diversity in input sets. Then they analyse their statistical association with fault detection using two datasets and three DNN models. They further compare diversity with state-of-the-art white-box coverage criteria. From the experiments, it could be inferred that relying on the diversity of image features embedded in test input sets is a more reliable indicator than coverage criteria to effectively guide the testing of DNNs. One of the selected black-box diversity metrics far outperforms existing coverage criteria in terms of fault-revealing capability and computational time. Results also confirm that state-of-the-art coverage metrics are not adequate to guide the construction of test input sets to detect as many faults as possible with natural inputs.
- In this paper [6], inspired by the MC/DC coverage criterion, the authors propose a family of four novel test criteria that are tailored to structural features of DNNs and their semantics. They validate the criteria by demonstrating that the generated test inputs guided via the proposed coverage criteria are able to capture undesired behaviours in a DNN. Test cases are generated using a symbolic approach and a gradient-based heuristic search. By comparing them with existing methods, they prove that their criteria achieve a balance between their ability to find bugs (proxied using adversarial examples) and the computational cost of test case generation. Experiments are conducted on state-of-the-art DNNs obtained using popular open source datasets, including MNIST, CIFAR-10 and ImageNet.
- In this paper [7], the authors design, implement, and evaluate DeepTest, a systematic testing tool for automatically detecting erroneous behaviors of DNN-driven vehicles that can potentially lead to fatal crashes. A tool is designed to automatically generated test cases leveraging real-world changes

in driving conditions like rain, fog, lighting conditions, etc. DeepTest systematically explore different parts of the DNN logic by generating test inputs that maximize the numbers of activated neurons. DeepTest found thousands of erroneous behaviors under different realistic driving conditions (e.g., blurring, rain, fog, etc.) many of which lead to potentially fatal crashes in three top performing DNNs in the Udacity self-driving car challenge.

- In this paper [5], the authors summarize some of the available tools for creating a synthetic dataset for object segmentation and provide a detailed example using the Blender 3D creation suite. They discuss in detail about the vital points and considerations for automatic dataset generation for object segmentation.
- There are many parallels between meta-modelling—in the sense of model-driven engineering—and meta-learning. Both rely on abstractions, the meta data, to model a predefined class of problems and to define the variabilities of the models conforming to this definition. Both are used to define the output and input relationships and then fitting the right models to represent that behaviour. In this paper [3] , the authors envision how a meta-model for meta-learning can look like. They discuss possible variabilities, for what types of learning it could be appropriate for, how concrete learning models can be generated from it, and how models can be finally selected. They also discuss a possible integration into existing modelling tools.
- This survey paper [4] conducts a review of the current research effort into making DNNs safe and trustworthy, by focusing on four aspects: verification, testing, adversarial attack and defence, and interpretability. In total, they survey 202 papers, most of which were published recently after 2017.
- In this paper [2], The authors propose a technique that re-purposes software testing methods, specifically mutation-based fuzzing, to augment the training data of DNNs, with the objective of enhancing their robustness. This technique casts the DNN data augmentation problem as an optimization problem. It uses genetic search to generate the most suitable variant of an input data to use for training the DNN, while simultaneously identifying opportunities

to accelerate training by skipping augmentation in many instances. They instantiate this technique in two tools, Sensei and Sensei-SA, and evaluate them on 15 DNN models spanning 5 popular image data-sets. The evaluation shows that Sensei can improve the robust accuracy of the DNN, compared to the state of the art, on each of the 15 models, by upto 11.9% and 5.5% on average. Further, Sensei-SA reduces the average DNN training time by 25%, while still improving robust accuracy.

2.1 Subsection 1

2.2 Subsection 2

3 Project Plan

3.1 Work Packages

The bare minimum will include the following packages:

WP1 Literature review on testing and verification of Deep Neural Networks

T1.1 Understanding verification and testing of Deep Neural Networks and the concept of synthetic dataset generation.

T1.2 Literature search on the testing and verification methods introduced in the past.

T1.3 Identify and analyse the state-of-the-art methods.

WP2 Experimental setup

T2.1 Choose the best synthetic dataset generation tool and analyse its capabilities.

T2.2 Decide on the DNN models to be used for the verification.

T2.3 Gather requirements from a test engineer.

T2.4 Choose a suitable model to model transfer method to build a DSL that transfers requirements into dataset along with the ground truth.

WP3 Mid-term Report

T3.1 Submission of the mid-term report.

WP4 Experimental analysis

T4.1 Train the DNN model with the synthetic dataset.

T4.2 Follow a behaviour driven development method.

T4.3 The DSL creates dataset based on the requirements from the test engineer.

T4.4 The DNN model is tested using the created dataset.

T4.5 Verify the performance of the model.

WP5 Final Report

T5.1 Draft revision.

T5.2 Submission of the Final report.

3.2 Milestones

M1 Literature search

M2 Experimental setup

M3 Mid-term Report submission

M4 Experimental Analysis

M5 Final Report submission

3.3 Project Schedule

Include a gantt chart here. It doesn't have to be detailed, but it should include the milestones you mentioned above. Make sure to include the writing of your report throughout the whole project, not just at the end.

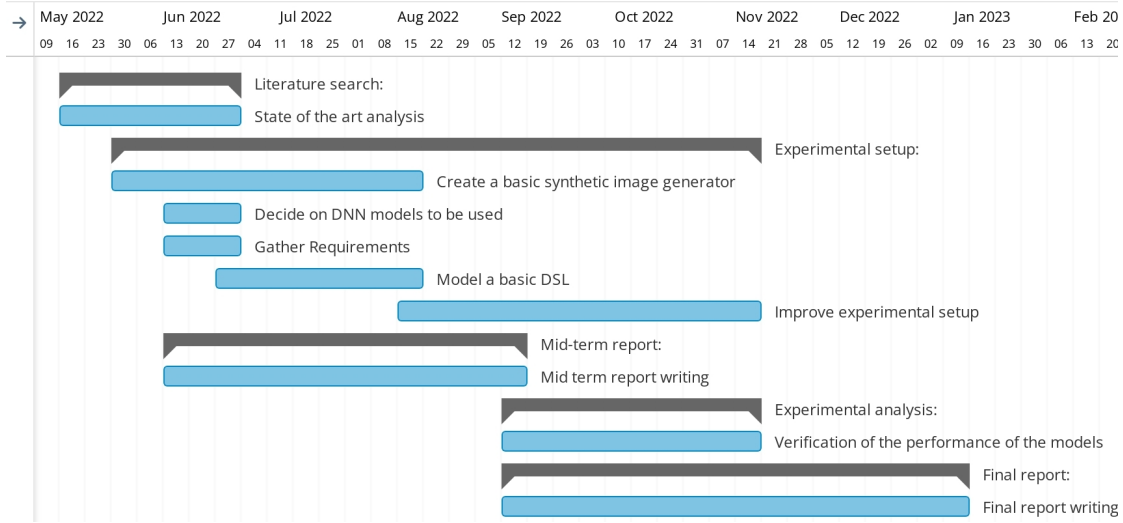


Figure 2: The gantt chart of the RnD project

3.4 Deliverables

Minimum Viable

- State of the art analysis.
- A simple and easy to use synthetic dataset generator with variabilies such as illuminance, sharpness, color, contrast, exposure, object poses, camera angles, FOV and perspective.
- Training the DNN with the synthetic data generated randomly irrespective of the requirements.
- A DSL that could map requirements to dataset generation functions and DNN testing functions.
- Verification of the DNN model for a classification or a regression application.

Expected

- A simple and easy to use synthetic dataset generator with additional variabilities such as material texture and reflection, crowded environment.

- Training the DNN with both synthetic and real world data on a ratio 3:1 respectively.
- Verification of the DNN model for a classification and a regression application.

Desired

- An easy to use synthetic dataset generator which could incorporate as much variabilities available in blender.
- Training the DNN with relatively more real world data than the synthetic data.

References

- [1] Zohreh Aghababaeyan, Manel Abdellatif, Lionel C. Briand, Ramesh S, and Mojtaba Bagherzadeh. Black-box testing of deep neural networks through test case diversity. *CoRR*, abs/2112.12591, 2021. URL <https://arxiv.org/abs/2112.12591>.
- [2] Xiang Gao, Ripon K. Saha, Mukul R. Prasad, and Abhik Roychoudhury. Fuzz testing based data augmentation to improve robustness of deep neural networks. In *2020 IEEE/ACM 42nd International Conference on Software Engineering (ICSE)*, pages 1147–1158, 2020.
- [3] Thomas Hartmann, Assaad Moawad, Cedric Schockaert, Francois Fouquet, and Yves Le Traon. Meta-modelling meta-learning. In *2019 ACM/IEEE 22nd International Conference on Model Driven Engineering Languages and Systems (MODELS)*, pages 300–305, 2019. doi: 10.1109/MODELS.2019.00014.
- [4] Xiaowei Huang, Daniel Kroening, Wenjie Ruan, James Sharp, Youcheng Sun, Emese Thamo, Min Wu, and Xinping Yi. A survey of safety and trustworthiness of deep neural networks: Verification, testing, adversarial attack and defence, and interpretability. *Computer Science Review*, 37:100270, 2020. ISSN 1574-0137. doi: <https://doi.org/10.1016/j.cosrev.2020.100270>. URL <https://www.sciencedirect.com/science/article/pii/S1574013719302527>.

- [5] Artx00FA;r I. Kx00E1;roly and Px00E9;ter Galambos. Automated dataset generation with blender for deep learning-based object segmentation. In *2022 IEEE 20th Jubilee World Symposium on Applied Machine Intelligence and Informatics (SAMI)*, pages 000329–000334, 2022. doi: 10.1109/SAMI54271.2022.9780790.
- [6] Youcheng Sun, Xiaowei Huang, and Daniel Kroening. Testing deep neural networks. *CoRR*, abs/1803.04792, 2018. URL <http://arxiv.org/abs/1803.04792>.
- [7] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. Deeptest: Automated testing of deep-neural-network-driven autonomous cars. In *Proceedings of the 40th International Conference on Software Engineering, ICSE '18*, page 303–314, New York, NY, USA, 2018. Association for Computing Machinery. ISBN 9781450356381. doi: 10.1145/3180155.3180220. URL <https://doi.org/10.1145/3180155.3180220>.