

# Detection of failing tests in AI-based systems

Antonio Guerriero\*, Michael R. Lyu†, Roberto Pietrantuono\*, Stefano Russo\*

\* Università degli Studi di Napoli Federico II

{antonio.guerriero, roberto.pietrantuono, stefano.russo}@unina.it

† The Chinese University of Hong Kong

lyu@cse.cuhk.edu.hk

**Abstract**—With Artificial Intelligence (AI) being increasingly used in critical systems, reliability of AI-based systems is a great concern. Testing is a fundamental software reliability engineering technique, but its application to AI-based systems poses new challenges, as determining the correct output for an arbitrary input may be hard (*oracle problem*). Building on the notion of *failure detectors* in distributed systems, we propose a surrogate of a test oracle (*Failed Test Detector*, FTD) able to unveil when the output of an AI-based system is incorrect despite the correct output is uncertain. An FTD should conveniently trade off *completeness* and *accuracy*.

**Index Terms**—Testing, Artificial Intelligence, Reliability

## I. INTRODUCTION

Testing is fundamental in software reliability engineering [1], yet testers are facing new challenges with the increasing use of Artificial Intelligence (AI) in software systems [2]. Testing aims to expose failures of the system under test, and a *test oracle* is needed to establish if a test “passes” or “fails”. However an oracle may be hard to define for AI-based systems, due to possible uncertainty on their correct output [3].

This is attributable to the software development paradigm shift inherent to Machine Learning (ML) [4] and AI, from definition and coding of deterministic algorithms to training and learning from data with statistics algorithms. Most traditional testing techniques foresee the specification of input test data and context conditions and the *a priori* identification of the expected output. For an AI-based system, an objective and consistent specification, against which the system behavior can be tested, is rarely available, due to the uncertainty in system output for some input data, and for the same test data in varying context conditions [3]. As stated in [4], *generating reliable test oracle is sometimes infeasible for some ML systems*, and, according to Murphy *et al.*, more in general in AI-based systems *there is no reliable “test oracle” to indicate what the correct output should be for arbitrary input* [5].

We propose to address the *oracle problem* for AI-based systems by somehow “restricting” it to the definition of an oracle able to evaluate when a test *fails*, in situations where the exact output for that test is unknown or uncertain. Such an oracle is similar to a *failure detector* for processes in a distributed system [6]; we call it *Failed Tests Detector* (FTD).

## II. AI-BASED SYSTEMS

As stated in [3], *AI-based software and applications use machine learning models and techniques through large-scale data training to implement diverse artificial intelligent features*

*and capabilities*. We regard an AI-based system (Figure 1) as taking a feature vector ( $f_i$ ) as input; an internal component uses an AI to compute a prediction producing a response ( $r$ ), while other components ( $c_i$ ), not based on AI, produce additional outputs ( $a_i$ ) which, combined with  $r$ , yield the ultimate output  $o$ . For instance, an autonomous driving system exploits an AI to process the cameras’ images, while other components process on-board sensors’ data, ultimately determining a final output (speed or steering angle). The behavior of an AI-based system is strongly dependent on both the AI algorithms and the training data.

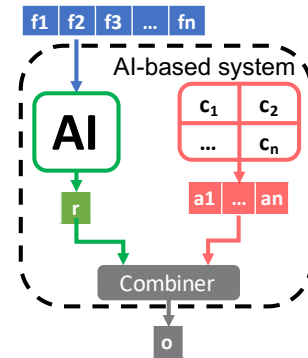


Fig. 1. AI-based system

## III. DETECTION OF FAILING TESTS

In partial analogy with distributed failure detectors [6], the output of an FTD for a given test is *fail* when the FTD is (“reasonably”) able to tell that the system has produced an incorrect output for the test case, otherwise it is *unknown*. The effectiveness of an FTD can be evaluated based on the following quality metrics:

- *Completeness*: all incorrect outputs of the system under test are correctly judged as *fail* from the FTD (the FTD is complete if it is able to avoid False Negatives);
- *Accuracy*: all correct outputs of the system under test are not judged as *fail* from the FTD (the FTD is accurate if it is able to avoid False Positives).

We envisage the architecture of a Failed Tests Detector as depicted in Figure 2. It is meant to exploit various ways to discover failed tests, using different sources of knowledge (domain knowledge, training set, system’s internal parameter values) to look for conditions, e.g. *invariants*, which hold when the system output can be judged incorrect, despite the

