

## APPROACH TO AUTOMATION TEST FRAMEWORK FOR BIG DATA PLATFORM USING OPEN SOURCE TECHNOLOGIES

<sup>1</sup>Dr. Vasanthi Kumari P, <sup>2</sup>Ashok GV

<sup>1,2</sup>Department of Computer Science and Engineering, Dayananda Sagar University, Bangalore

### Abstract-

The key to a successful automation effort lies in designing a good and robust automation framework. Automation can improve quality activities and lead to higher testing efficiency. As we are dealing with huge data and executing on multiple nodes there are high chances of having bad data and data quality issues at each stage of process. Data functional testing is performed to identify these data issues because of coding errors or node configuration errors. Testing should be performed at each of the three phases of big data processing to ensure that data is getting processed without any errors. Functional testing includes, i) Loading source data files into HDFS, ii) Perform map reduce operations and iii) Extract the output results from HDFS. Issue occurs in i) data moving from source system to Hadoop are incorrect data captured from source systems, incorrect storage of data in complete or incorrect replication. ii) Issues faced in data processing or coding issues in map reduce jobs, jobs working correctly when run in stand alone node, but working incorrectly when run on multiple nodes, incorrect aggregations, node configurations, and incorrect output format. iii) issues faced while generating reports are report definitions not set as per the requirement, report data issues, layout and formatting issues. Automation promises to find more defects with less effort. It's a must for every new project. The mystique of automation is so strong that it even seems that you no longer need to write test cases, run manual testing, or do impact analyses. Open-source technologies are helping organizations across industries gain strategic insights from the torrents of data that now flow through IT systems. A robust big-data validation framework can significantly improve high-volume, big-data testing helping to fortify quality assurance, refine analytical thinking and enhance the overall customer experience. This paper gives a brief insight of general test automation frameworks and describes a test automation framework which helped in increasing productivity by simplifying the process of test script development.

**Keywords:** Test Automation Framework; Open Source tools; Validations; Reporting; Visualization; Analysis

### Introduction

The ever-increasing complexity and scale of IT applications has made quality and reliability of paramount importance. However, delivering the same through testing is becoming a challenge. Software test automation has the capability to increase the overall test coverage, while also improving software quality. In today's environment of software industry is faced with challenges of developing increasingly complex applications and products over shorter time periods and at lower costs, while maintaining consistently high quality. This need for shorter time to market tasks for significantly reduced development and testing time. This makes test automation increasingly necessary and even critical. Test Automation is the process of writing programs to do testing that would otherwise be done manually. Once tests have been automated they can be run repeatedly and quickly making it cost effective, especially for products with long maintenance life. It consists of adaptable, object-oriented components that ensure quick and easy implementation of customized solutions.

The primary objectives of test automation are:

- Reduce the overall test cycle time and accelerate the time to market.
- Reduce cost of repetitive manual testing.
- Increase the test coverage and reduce risk.

- Ensuring consistent and thorough testing leading to improved quality and reliability
- Optimize the time and cost by repeating the tests and Higher test coverage levels
- Maintainable & scalable scripts and Minimize lead time

The paper presents the literature review in section 2, Problem statement in section 3, Proposed work in section 4, Results and discussions in section 5 followed by Conclusion in the final section.

### II Literature Review

V. Vaidhehi, Arun M. B [1], The outcome of this research is on data processing and storage approaches are facing many challenges in meeting the continuously increasing demand of big data. Mainly focused on Map reduce, one of the key approaches for meeting the big data demand through highly parallel processing on a large amount of commodity nodes. Challenges and solutions on four dimensions like data storage, analytics, online processing and privacy and security.

Rabiul Islam Jony, Ahsan Habib, Nabeel Mohammed, Rakibull Islam Rony [2], This paper identifies potential Big data use case

domains for telecom operators. It also presents the corresponding examples, required data types and typical challenges to implement these use cases. However, newer use cases are emerging as operators are more focusing on their big data strategy. For example, potential use cases based on location data are emerging currently. As a result, location based services can be introduced as a separate use case domain. As the field of big data is still in a great state of change and development, the possibilities for research in this area are extensive and the interest on the findings will certainly rise if the adoption of big data rises as expected.

The Computing Big Data Review 2015 [3] summarizes the results of a comprehensive research program undertaken by Computing during the first quarter of 2015. The findings are compared with those published in the Computing Big Data Review 2014 in order to establish the degree to which the market has developed in this short space of time and the degree to which Big Data and the democratization of data analytics is changing business organizations. The review contains some unique insights from high-ranking IT decision makers into how organizations are approaching and deploying Big Data solutions and some of the specific challenges they face. Methodology of the research project was conducted using a combination of qualitative and quantitative methods.

### III Problem Statement

As we are dealing with huge data and executing on multiple nodes there are high chances of having bad data and data quality issues at each stage of process. Data functional testing is performed to identify these data issues because of coding errors or node configuration errors. Testing should be performed at each of the three phases of big data processing to ensure that data is getting processed without any errors. Functional testing includes,

#### 1. Verification of Pre-Hadoop processing

Data from various sources like weblogs, social network sites, call logs, transactional data etc is extracted based on requirements and loaded into HDFS before processing it further. **Issues:** Some of the issues which we faced during this phase of the data moving from source system to Hadoop are incorrect data captured from source systems, incorrect storage of data in complete or incorrect replication.

#### 2. Verification of Hadoop map reduce process data output

Once data is loaded into HDFS Hadoop map reduce process is run to process the data coming from different process.

**Issues:** Some issues that we faced during this phase of the data processing or coding issues in map reduce jobs, jobs working correctly when run in stand alone node, but working incorrectly when run on multiple nodes, incorrect aggregations, node configurations, and incorrect output format.

#### 3. Verification of data extract and load into EDW.

Once Map-Reduce process is completed and data output files are generated, this processed data is moved to enterprise data warehouse or any other transactional systems depending on the requirement.

**Issues:** Some issues that we faced during this phase include incorrectly applied transformation rules, incorrect load of HDFS files into EDW and incomplete data extract from Hadoop HDFS.

#### 4. Verification of Reports

Analytical reports are generated using reporting tools by fetching the data from EDW or running queries on Hive.

**Issues:** Some of the issues faced while generating reports are report definitions not set as per the requirement, report data issues, layout and formatting issues. Apart from this functional verification non-functional testing including performance testing and failover testing need to be performed.

**Failover Testing:** Hadoop architecture consists of a name node and hundred of data nodes hosted on several server machines and each of them are connected. There are chances of node failure and some of the HDFS components become non-functional. Some of the failures can be node failure, data node failure and network failure. HDFS architecture is designed to delete these failures automatically and recover to process with the processing. Failover testing is an important focus area in big data implementations with the objective of validating the recovery process and to ensure the data processing happens seamlessly when switched to other data nodes. Some validation that needs to be performed during failover testing are validating that checkpoints of edit logs and FsImage of name node are happening at defined intervals, recovery of edit logs and FsImage files of name node, no data corruption because of the name node failure, data recovery when a data node fails and validating that replication is initiated when one of the data nodes fails or data becomes corrupted. Recovery Time Objective (RTO) and Recovery Point Objective (RPO) metrics are captured during failover testing.

### IV Proposed Methodology

Developing the right architectural framework for test automation will result in automation code that can be used for longer period of time with less maintenance than a simple record/playback solution. This translates to a significant saving over the course of long projects, and it provides ability to more thoroughly test an application with less manpower. The result of this architecture is reliable automation code with scripts that can last the entire life of the product (not just the project) and that can be used and enhanced by business analysts who have little to no knowledge of automation testing. This approach should be started at the planning of the automation work. The requirements should be correctly defined and finalized. The automation requirements can be as formal and specific as the software requirements, or they can be as informal as a greening what modules or portions of the applications should and should not be automated and which paths, particular test cases will take through the application.

Test Automation Framework is a generic framework based on hybrid model. It is a reusable framework that can be easily invoked by an automated test suite for an application under test. It optimizes the effort required to automate the testing of the application.

The framework hides the underlying challenges like communicating with the environment, database, configurations etc. from the application under test. The framework can work with multiple automation tools and applications under test.

The framework is divided into four logical sections and is a collection of various Resources, Driver Scripts, Function Libraries, Test Automation Architecture and Result Libraries. Figure 1 gives proposed Test Automation Framework.

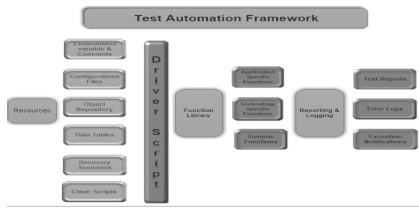


Figure 1 Proposed Test Automation Framework

**Resources:** These are various inputs that are required for the driver script to run. They may consist of test data files, object repositories, environment variables etc.

**Function Libraries:** These form the core part of the framework and consist of reusable functions. These basically represent business flows, generic functionalities etc. and are recalled from the driver scripts.

**Driver Scripts:** These consist of calls to various reusable functions and other resources forming a complete end-to-end business flow.

**Reports:** These consist of results of test execution.

### V Results And Discussions

Big data is still emerging and there is a lot of focus on testers to identify innovative ideas to test the implementation. Designed a test automation framework in the BDT (Behavior Driven Testing) model using Cucumber & Gherkins and Java. The framework helps in performing count and data validation during the data processing stage by comparing the record count between the Hive and SQL tables and confirmed that the data is properly loaded without any truncation by verifying the data between Hive and SQL tables.

Similarly, when it comes to validation on the map-reduce process stage, it helps if the tester has good experience on programming languages. There is a reason because unlike SQL where queries can be constructed to work through the data MapReduce framework transforms a list of key-value pairs into a list of values. A good unit testing framework like JUnit can help validate the individual parts of the MapReduce job but they do not test them as a whole.

Building a test automation framework using a programming language like Java can help here. The automation framework can focus on the bigger picture pertaining to MapReduce jobs while encompassing the unit tests as well. Setting up the automation framework to a continuous integration server like Jenkins can be even more helpful. However, building the right framework for big data applications relies on how the test environment is set up as the processing happens in a distributed manner here. There could be a cluster of machines on the QA server where testing of MapReduce jobs should happen.

The open source community has also been hard at work on tools specific to analytics and big data. After the foundation is built using open source tools, big data scientists have the option of continuing to use open source to drill into data and produce reports.

DataTorrent, which was created by the team behind Apache Apex, is one such tool. It allows for real-time streaming data processing and batch processing. It's built on Hadoop and was named as one of these even leaders in big data streaming analytics solutions.

Genie is another open source big data tool, this one created by Netflix and also running on Hadoop. While Netflix created Genie for itself originally, Netflix also released Lipstick as open source, which allows for graphical representation of Hadoop Pig jobs.

LinkedIn developed Pinot, a distributed OLAP data store with real-time scalable analytics, as well as Taiga, which relies on the Kafka messaging system and Hadoop YARN to conduct distributed stream processing.

**Technology:** Java 1.8, Jenkins, Selenium WebDriver, TestNG, Maven, Skool, DataTorrent, Genie, Lipstick, Pinot, Taiga, AutoIT, Winium, Sikuli.

### V Conclusion And Future Work

Test automation framework could be difficult to implement in the big data platform. The framework underlying challenges like communicating with the environment, database, configurations etc. Proposed solution identifying the requirements and building a robust automation framework can help in doing comprehensive testing. In order to keep up with the pace of product development and delivery it is essential to implement effective, reusable test automation. Properly planned and implemented automated testing can significantly lower project risk and cost. If automation is well designed and implemented, it will not only reduce the risk of the current project, but can actually reduce the risk and cost of future projects. The ROI for test automation is easy to establish. Automation Framework provides a way to drive productivity and foster code reuse, ultimately enhancing the quality of resulting. However, a lot would depend on how the skills of the tester and how the big data environment is set up. And have identified some major challenges regarding big data supporting open source tools. But there are still some challenges with no mitigation strategies. Our future research will concentrate on developing a more complete understanding of challenges associated with big data and comp

lete Test Automation framework which includes maintenances of performance testing.

Challenges <https://www.hindawi.com/journals/tswj/2014/712826/>

### References

- [1] Author V. Vaidhehi, Arun M. B., Department of Computer Science, Christ University Bengaluru published the paper on "Managing Big data using Hadoop Map Reduce in Telecom Domain" at International Journal of Innovative Research in Advanced Engineering (IJIRAE) ISSN: 2349-2763 Issue 03, Volume 3 (March 2016) [www.ijirae.com](http://www.ijirae.com)
- [2] Authors Rabiul Islam Jony, Ahsan Habib, Nabeel Mohammed, Rakibul Islam Rony published research paper on 'Big Data Use Case Domains for Telecom Operators' from Primeasia University, Bangladesh in 2015 IEEE International Conference.
- [3] The Computing Big Data Review 2015 summarizes the results of a comprehensive research program undertaken by Computing during the first quarter of 2015. <http://www.informationweek.com/big-data/big-data-analytics/big-data-predictions-for-2016/d/d-id/1323671> [https://www.splunk.com/en\\_us/solutions/solution-areas/big-data.html](https://www.splunk.com/en_us/solutions/solution-areas/big-data.html) <https://www.bigdatareviews.org/?p=977>
- [4] In May 2012, Elsevier sponsored a 2-day conference in Canberra, Australia dedicated to the topics of Big Data, E-science and Science policy (see videos and links to the presentations here: <https://www.youtube.com/playlist?list=PL61DD522B24108837>).
- [5] <http://ieeexplore.ieee.org/>
- [6]. [https://en.wikipedia.org/wiki/Big\\_data](https://en.wikipedia.org/wiki/Big_data)
- [7]. International Journal of Innovative Research in Advanced Engineering (IJIRAE) ISSN: 2349-2763 Issue 03, Volume 3 (March 2016)
- [8] shwarappa, Anuradha J, "A Brief Introduction on Big Data 5V's Characteristics and Hadoop Technology", International Conference on Intelligent Computing, Communication & Convergence (ICCC-2015),
- [9] Bharti Thakur, Manish Mann, "Data mining for big data: A Review", International Journal of Advanced Research in Computer Science and Software Engineering, ISSN: 2277128x, Volume 4, Issue 5, May 2014.
- [10]. Big Data: Survey, Technologies, Opportunities, and