

**Team Name:** Test Bayes

**Team Members:**

Adavya Bhalla (bhalla)

Aditya Jhamb (aditya97)

Avidant Bhagat (avidant)

Ethan Mayer (emayer4)

Steven Austin (saaustin)

## **Project Report**

### **Motivation :**

More often than not, testing takes a lot more time and effort than the actual solution to the problem<sup>8</sup>. A tool that reduces the amount of time it takes to test a given solution to a problem would be a major boost to the current development process. If a test suite is very large and is ordered in a way such that the first half or more passes but the others fail, then the developer loses out on precious time, as the tests could have been ordered to indicate errors first, failing fast. In such a scenario, the developer is likely to save a lot of time and come up with a solution quicker since they were just working on the code. Another great feature is that tests that fail together i.e. they fail due to the same bug are not tested again and again but, if one such test fails we push the other ones back. This implies every initial fail points to new bugs. This makes the developer's job even easier since the tests that failed initially can help them discover a lot of bugs in different areas of the project faster.

### **Our approach:**

The re-ordering of tests is treated as a probabilistic problem where we take advantage of not only the independent probability of each test failing, but also the correlated probabilities of that test conditioned on other tests. The probabilistic treatment ensures that tests that usually fail after an update are tested first. This happens because the our tool reorders the tests at run-time. Tests are re-ordered based on the outcome of the previous tests, ensuring that the order is interesting (i.e. it is more likely to fail fast and indicate diverse errors). To achieve this goal, we can have the following high level approach:

1. We maintain a list of the independent probabilities of all tests failing on different versions of the same project. This list is constructed as we run the test suite on different versions of the implementation. For example, to construct the list of probabilities, we go through 100 different versions of the same large project and find out the probabilities of individual tests failing.

2. Next, we have to construct multiple data structures that store all the conditional probabilities of these tests. If there are  $n$  tests then we can implement this using two  $n \times n$  matrices. The first matrix stores the probability of a test failing given that another test succeeded, while the second matrix stores the probability of a test failing, given that another test fails.
3. Finally, we first run a test that has the highest probability of failing regardless of other tests. If this test fails, then we iterate through the second matrix and find out which test has the highest probability of failing given that the first test failed. If it passes, then we do the same thing with the first matrix.

**Architecture:**

*Figure 1:*

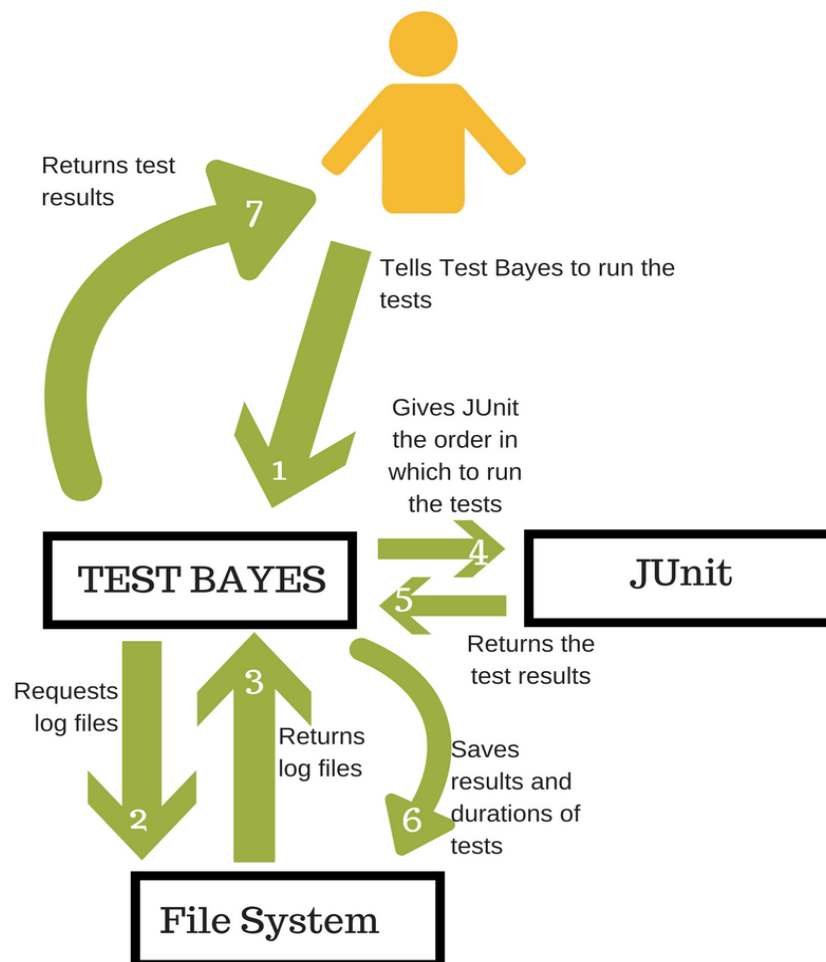
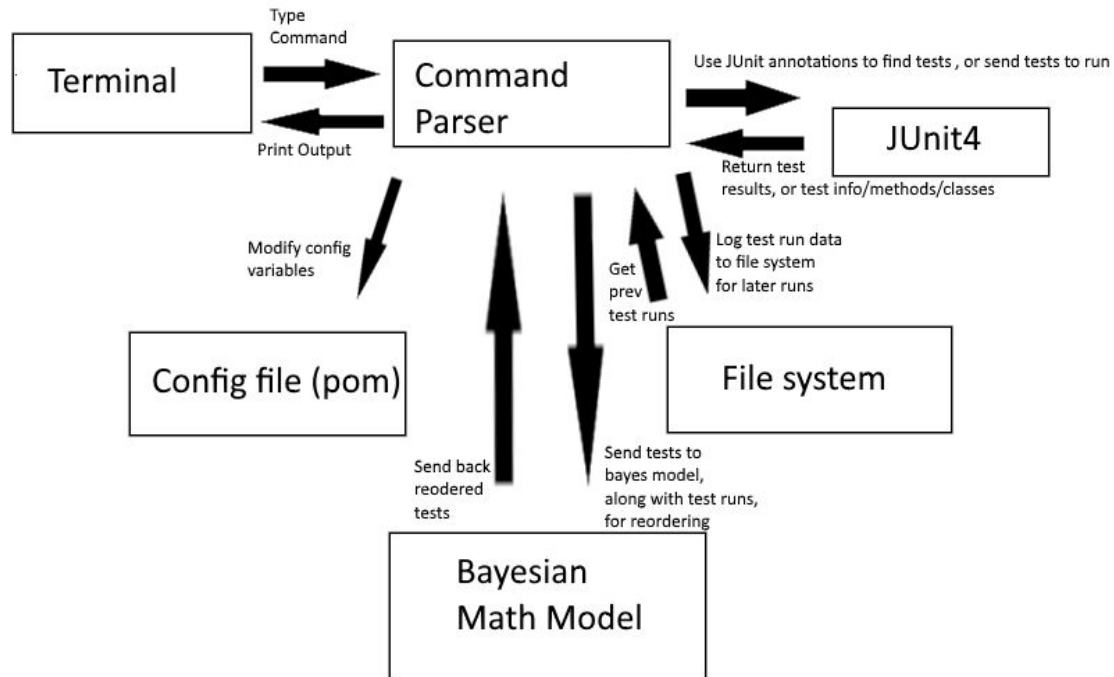


Figure 2:



These designs represent the way our modules interact with both the user, and themselves. Figure 1 depicts how the user and JUnit interact with our model, and the second model gives a more in depth understanding of how the “Test Bayes” module in Figure 1 works.

In Figure 2, we have a command parser (main portion of the tool) that determines if the user wants to modify the test reordering config, or if the user wants to run tests. JUnit is used to parse annotations and run tests. The file system is either sent new data to log, or is queried for existing data. The main brain of this tool is the Bayesian Math Model. This model takes in the previous run information, and returns the mathematically best ordering of the tests, in the order they should be run. The file structure returns a JSON model of the previous test runs. The mathematical model takes in a list of method names, along with the JSON data needed for bayes calculations, and returns a list of methods to be run. The command going into the parser could be a list of test directories to run, or a flag that states all tests should be run. JUnit requires you to query it with a test suite class, requiring us to either have the user put all test classes in a dummy test suite using the `@SuiteClass` annotation, or to have our command parser generate the suite on the fly.

## Implementation:

We have decided to make two slightly different options for the actual test ordering; each is ordered in a different way. The first option involves a greedy heuristic where we first run the test with the most probability of failing and then based on the outcome of the first test, greedily decide which test to run next based on conditional probability. We would break ties on the greedy decisions by factoring in time taken to run the test, conditional probability, total probability, and a running average. As a note, the term running average, and moving average are used interchangeably, and mean the same thing; this moving/running average is simply the last  $n$  test runs that the user specified they want in the bayesian calculations our tool performs. The second option would be to just use everything from the start instead of waiting for ties. Evaluation will be critical in deciding which option will work better. Here are a few details about the data structures we plan on using:

1. **HashMap for total probability** - We want to store the total probability for each test passing in a HashMap where the hashCode is the name for each test and the value is the probability for each test passing. The value will be in a fraction object created by us, that stores the numerator and denominator as integers..
2. **HashMap for runtime** - We want to store the runtime for each test in a HashMap where the hashCode is the name for each test and the value is the runtime for the given test. The runtime will be stored as a double with seconds as the unit of time.
3. **Matrix for Conditional probability** - This is where we still have certain grey areas. We want to store the conditional probability in a matrix so that if we access the value for a test  $x$  by test  $y$  index we get the probability for test  $x$  passing given test  $y$  either passes or fails. There will be two matrices to facilitate faster calculations and decrease computing time. One way of doing this is to have a HashMap that itself has a HashMap as a value and then the second HashMap contains a fraction object as value which is the probability we seek. A 2D array implementation might also suffice. These are the two options we are currently looking into; right now we like the HashMap implementation better.

We've created two runners that can actually run the JUnit tests. One runner is an extension of the BlockJUnit4Runner class, which allows a user to use the `@RunWith(OurRunner.class)`, and simply run their test suites as if they were using JUnit natively. However, this implementation only allows for reordering within one class, meaning JUnit would spin up an instance of our runner for each test class, instead of one for all test classes. We had a second implementation that scans the testing directory for all test suites, using JUnit annotations, and manually runs the tests,

ordering them irrelevant from the test class they reside in. The drawback of the second implementation causes us to go under the hood with JUnit and run the tests manually.

As far as our mathematical model goes, we must use a fraction class that allows us to represent numerators and denominators separately, as it is mathematically impossible to perform our bayesian analysis without these. This is because we need fractions to be not in the lowest form. Our bayes calculations requires us to only add to the numerator or denominator, meaning that we cannot represent current running probabilities of test failure as only numerators or denominators. As such, we implemented a class that can properly represent fractions.

We decided to implement our log files in a JSON style format, with one file per test run, and one line per file. This makes reads and writes very fast, as we do not have to constantly open, read lines, edit, and close files, we can just open, read one line, and close; file creation is a very fast operation. Our data will be saved in a simple JSON format, denoting the test name, success or failure, and the time it took to actually run. This is the only information necessary for the bayes math to be calculated, since we are generating joint probability tables of tests based on their success/failures, and doing further epsilon reordering based on tests of similar failure, to reorder based on time at that point.

On the very first run, when there is no test data at all, each test will be run as if it were being run by JUnit alone, meaning that each test will be deemed a 50/50 chance of passing or failing. After that first run, the bayes probabilities will come into play, along with our moving average, causing the tests to be reordered.

### **How is it currently done:**

Developers will usually run their tests sequentially, and wait for them to pass or fail or, manually reorder them such that the tests they think will fail are executed first. But, their general feeling of which test might fail is not mathematically backed up and most times will not yield good results. Testing frameworks, especially with JUnit, do not have support for dynamic ordering. This dynamic ordering would be the first of its kind. Read the “Similar Works” section of report to gain more insight into why other projects out there do not quite address the problems we are tackling in our design.

**Similar work:** Most work related to our project are frameworks that allow for didn't types of testing, rather than the actual act of reordering tests. One framework that is popular is the Spring Framework<sup>7</sup>, which is designed around testing a program by giving

control flow back to the framework (often called Inversion of Control, or IoC). However, this is (as previously stated), more concerned with the type of testing, rather than test reordering itself. After searching on github, and google scholar, the only thing similar to our project was a research paper conducted on the time benefits of reordering test based on test failure rates. The paper stated that “the effectiveness of prioritization techniques... on JUnit has not been investigated”<sup>6</sup>. However, the paper stated that “numerous empirical studies have shown that prioritization can improve a test suite's rate of fault detection”<sup>6</sup>, showing that our project has promise in being a useful tool for people to use.

### Evaluation:

To evaluate our code, we will be using two metrics to compare against two alternative ways of ordering tests. We will track both the time and number of tests run until the first failure, and compare our results for these figures to a test suite run in the default junit ordering, and a randomly ordered test suite. We will run these tests on multiple repositories (with commits that have a failed test build), along with general debugging and intermediary evaluation with the CSE 331 RatPoly code, and compare how our ordering fairs on each execution. We expect that test Bayes ordering will improve dramatically after a few runs of the test suite. An example data table we would generate is as follows:

| Repository Names | Time taken to first failure |                 |                    | Number of tests to first failure |                 |                    |
|------------------|-----------------------------|-----------------|--------------------|----------------------------------|-----------------|--------------------|
|                  | JUnit Ordering              | Random Ordering | TestBayes Ordering | JUnit Ordering                   | Random Ordering | TestBayes Ordering |
|                  |                             |                 |                    |                                  |                 |                    |
|                  |                             |                 |                    |                                  |                 |                    |
|                  |                             |                 |                    |                                  |                 |                    |
|                  |                             |                 |                    |                                  |                 |                    |

### Main Demographics:

Since this is not an existing method of testing in industry, there is a huge user base to use this resource. Companies have codebases that can take days to run through all tests even for a small package. Moreover, nearly every developer has had to run tests that take a significant amount of time, only for them to fail near the very end.

This cases development to be slowed, making it very cumbersome for a developer to actually debug large projects.

**Impact:**

As a group, we feel as if the impact of this project is obvious to every developer who has ever run a test suite. Our tool will save developers a huge amount of time in running test suites with its ordering feature. Individual coders will feel the benefits in their personal time being saved, and the frustration of waiting for tests to complete reduced. Larger organizations could literally measure their monetary savings if they compared the average wait time of a test suite before and after implementing our tool, correlating to higher productivity. Overall, this tool will dramatically cut down on idle time in the testing process.

Every developer has had to write test suites before. No matter the size, testing is a great benchmark that many developers use to gauge how correct their code is. Sometimes they even going so far as to writing the tests before implementing their code, as we have done occasionally in our studies. For this reason, it's not hard to see that every developer could benefit from this tool, since they would be able to rapidly determine the failure points of their code instead of waiting for long test suites to run. Our group has had several members work at internships where test suites could take anywhere from ten to thirty minutes, even overnight! Of course, failures near the end of the suite were common, and wasted a lot of time. Our code would be not only generally useful to the common developer, but could be crucial to many industry codebases, speeding up the development process significantly.

**Challenges and Risks:**

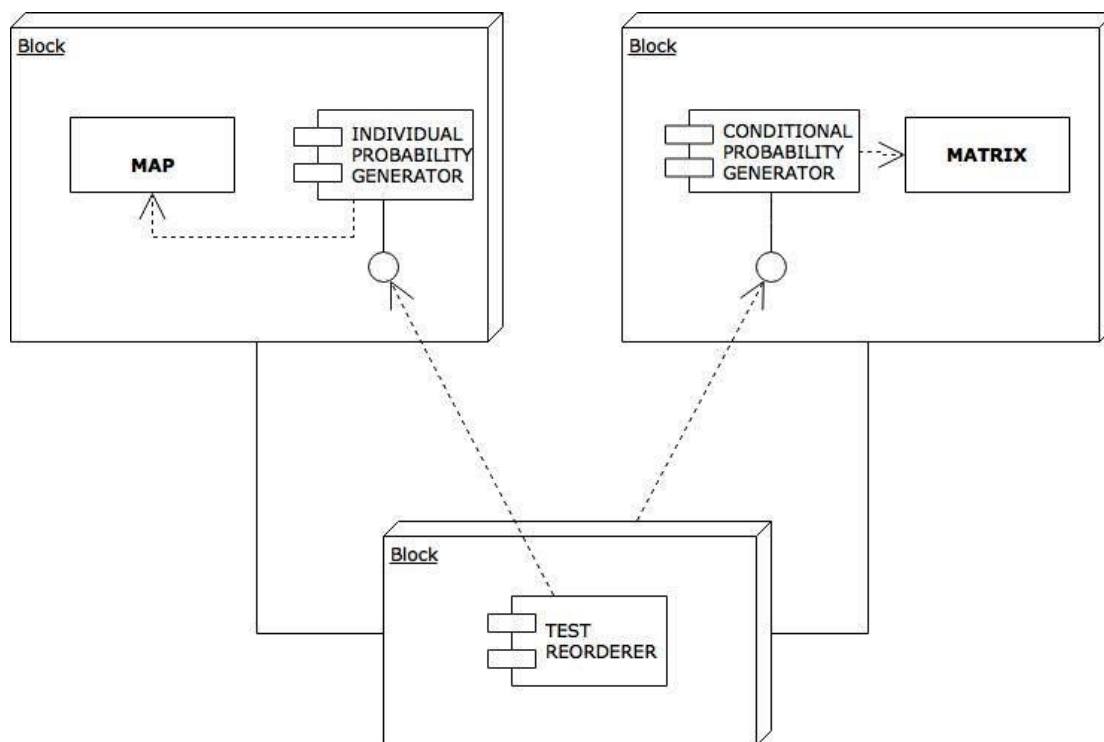
It is difficult to mathematically verify whether this approach will be successful or not. It sounds correct but may not be technically sound. The biggest risk is that the overhead from running this extension tool might take longer because of the initial training and construction of the two matrices and the list. Also, it might be challenging to find such a large open-source project that has so many previous versions.

Since there could potentially be a large amount of data from many test runs, we could run into the issue where the file/data processing required of the tool could outweigh the time the developer saves from test reordering. We will counteract this by incorporating a moving average with our Bayes model, so as to not consider extraordinarily large amounts of data.

Furthermore, if the data were too large, it could be a memory load on the user, meaning our approach will consist of using as little memory as possible in our file structure, only saving what is absolutely necessary.

Here is a small diagram that illustrates this procedure :

*Figure 3:*



### **Cost:**

In an ideal scenario, the costs of using this tool would be minimal. We are trying to make our solution take miniscule time to reorder, which should be much lower than the time saved by the re-ordered tests. To augment this, we would also be working with a moving average which will help reduce the amount of time we take to process the data.

### **Checks for Success:**

The most important factor to check for success would be the time to failure. The purpose of our project is to essentially reduce this as much as possible, and this makes



it the most important factor to test success. To do this check, we would want to compare the time taken by our test order to the time taken by the standard JUnit ordering, and also to the time taken by a random ordering of the tests. By doing so, we will not only be able to guarantee that our ordering saves the developer time as compared to the JUnit ordering, but also compared to a random ordering, which is often considered a better option.

### **Verifying Builds using Travis CI:**

When a pull request is made, Travis CI is run automatically, and upon viewing the pull requests status online, a green check mark will pop up next to the Travis CI section indicating that the build passed, an X denoting it failed, or a different symbol denoting the test partially passed. Past builds can be found in our repository at <https://travis-ci.org/avidant/test-bayes/builds/>. To view individual past pull requests instead of current states of branches, go to [https://travis-ci.org/avidant/test-bayes/pull\\_requests](https://travis-ci.org/avidant/test-bayes/pull_requests). You'll see that there are several passing and failing builds.

This can also be run from the command line with the “`mvn clean verify`” command from the terminal, as this is the command that Travis CI runs in the background.

### **Using Parameters to improve results:**

We have identified some places in our code where we are using made up numbers like the starting probability for each test, the importance we have to give to each part of our mathematical model. We would like these numbers to mean something more than just numbers that seem correct. Therefore, we plan on using optimization methods to optimize these parameters for maximum performance. Once we have the initial implementation complete, we will use the initial results to optimize these parameters for a better overall performance by our tool. We can achieve this by either using ML libraries or just small custom methods that minimize the time taken to failure based on different values of the parameter.

### **Schedule:**

- ✓ Week 3: Learn about extending JUnit, and hammer out details about file format, and the mathematical model.
- ✓ Week 4: Begin development of base functionality, meaning the ability to reorder tests randomly, or in some given order.
- ✓ Week 5: Incorporate our mathematical model into the code.
- Week 6: Test to see how the mathematical model performs on code.

Week 7: Measure benefit of our tool, as opposed to random ordering.

Week 8: Clean up code, add any small features or changes that are necessary.

Week 9: Test on larger code bases.

Week 10: Polish codebase, turn in

## **Resources Used:**

<sup>1</sup><http://www.baeldung.com/junit-4-custom-runners>

<sup>2</sup><http://junit.sourceforge.net/javadoc/org/junit/runners/BlockJUnit4ClassRunner.html>

<sup>3</sup><https://github.com/eugenp/tutorials/blob/776a01429eb6b7ec0d1153f88097e2cc9915e599/testing-modules/testing/src/test/java/com/baeldung/junit/BlockingTestRunner.java>

<sup>4</sup>[https://stackoverflow.com/questions/20976101/how-to-get-a-collection-of-tests-in-a-junit-4-test-suite?utm\\_medium=organic&utm\\_source=google\\_rich\\_qa&utm\\_campaign=google\\_rich\\_qa](https://stackoverflow.com/questions/20976101/how-to-get-a-collection-of-tests-in-a-junit-4-test-suite?utm_medium=organic&utm_source=google_rich_qa&utm_campaign=google_rich_qa)

<sup>5</sup><https://junit.org/junit4/javadoc/4.12/org/junit/runner/manipulation/Sortable.html>

<sup>6</sup><https://dl.acm.org/citation.cfm?id=1116157>

<sup>7</sup><http://spring.io/>

<sup>8</sup><http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.370.9611&rep=rep1&type=pdf>

<sup>9</sup><https://ieeexplore.ieee.org/abstract/document/1007961/>

<sup>10</sup><https://ieeexplore.ieee.org/abstract/document/5401169/>

<sup>11</sup><https://ieeexplore.ieee.org/abstract/document/1510150/>

<sup>12</sup><https://junit.org/junit4/javadoc/4.12/org/junit/experimental/max/MaxCore.html>

<sup>13</sup><http://automation-webdriver-testng.blogspot.com/2012/07/how-to-call-test-plans-dynamically.html>

<sup>14</sup><http://beust.com/weblog/2008/03/29/test-method-priorities-in-testng/>

<sup>15</sup><https://arxiv.org/pdf/1801.05917.pdf>

<sup>16</sup><https://lib.dr.iastate.edu/cgi/viewcontent.cgi?referer=https://www.google.com/&httpsredir=1&article=6493&context=etd>