

Applying Regression Test Selection on the example of Business Central

Diego Joshua Martinez Pineda (diem@itu.dk)

Supervisor: Mahsa Varshosaz (mahv@itu.dk)

December 13, 2021

Abstract

In this research project, we address the problem of implementing a test selection algorithm on the Business Central project. For the scope of this project, a single test selection algorithm was implemented, which is based on test coverage information. The findings of this work should serve as a foundation to explore the different selection and prioritizations techniques in this system. Much of the work presented here went into understanding the context and setting of the system, sufficiently to execute the algorithm in actual revision examples. The insights acquired provide a basis for future evaluation settings to increase the validity of the results here presented and of future algorithms to be evaluated.

1 Introduction

When working with large software systems, regression testing is a valuable technique that increases the confidence of the product stakeholders that future revisions, new features, and evolution of the system will not break previous functionality.

For this reason, it is common for products to have the execution of the regression test suite as part of the main pipeline of the product. However, when systems are sufficiently large, it becomes very expensive to run the full suite of tests.

A test selection algorithm aims at choosing a subset of a full test suite to run, that is not as expensive as running all tests cases. However, this selection must preserve the capabilities of fault detection that the full test suite has. This means that careful selection must be made so that the confidence in fault detection is not lost because of not selecting some test cases.

Business Central is an Enterprise Resource Planning software (ERP) software system from Microsoft, it is a large software system, with different sets of tests for the diverse aspects of the product.

For the test selection explored in this project, the test suite of the application code (business logic) is explored. The business logic of the system is written on

AL, a DSL similar to Pascal that has primitives for interacting with the runtime of Business Central.

We would like to explore different test selection techniques for this system, with the intention of reducing the time and computing power required to detect faults in the current CI pipeline.

1.1 Overview

The idea for the method explored in this research project is to use the line coverage information that each test case produces, and “read it backward” when given a code change in those lines. That means that by looking at which lines were modified on a given change, given that we know which tests were affecting those particular lines then we know that we should select those tests.

To further illustrate, see for example figure 1. Here we can see a **diff** of an AL code file on a given revision that added lines in between what previously were lines 30-31, and 39-40. Also what was previously line 29 was removed. Therefore, provided that we know which tests were hitting those lines, then we could propose them as tests to be selected.

```

28 | TestMethodLine.SetRange("Line Type", TestMethodLine);
29 | OnRunTestSuite(TestMethodLine);
30 |
31 | if TestMethodLine.FindSet() then
32 |     repeat
33 |         OnBeforeCodeunitRun(TestMethodLine);
34 |         CODEUNIT.Run(TestMethodLine."Test Codeu
35 |         TestMethodLine.Find();
36 |         OnAfterCodeunitRun(TestMethodLine);
37 |     until TestMethodLine.Next() = 0;
38 |
39 | OnAfterRunTestSuite(TestMethodLine);
40 | end;
41 |
42 |
30 | TestMethodLine.SetRange("Line Type", TestMethodLine);
31 | OnRunTestSuite(TestMethodLine);
32+ | CodeCoverageEnabled := ALCodeCoverage.Initializ
33 |
34 | if TestMethodLine.FindSet() then
35 |     repeat
36 |         OnBeforeCodeunitRun(TestMethodLine);
37 |         CODEUNIT.Run(TestMethodLine."Test Codeu
38 |         TestMethodLine.Find();
39 |         OnAfterCodeunitRun(TestMethodLine);
40 |     until TestMethodLine.Next() = 0;
41 |
42+ | if CodeCoverageEnabled then
43+ |     ALCodeCoverage.StopCodeCoverageForTestSuite
44+ |
45 | OnAfterRunTestSuite(TestMethodLine);
46 | end;
47 |

```

Figure 1: A revision change on AL code.

Intuitively, this seems like a good heuristic to select tests. However, we do not claim it is sufficient, further examples can show where this heuristic fails. For example, when a part of the code is dependent on the side effects it produces; then changing the environment through the side effects of a revision, can cause a seemingly unrelated set of test cases to fail.

This method is meant to provide a foundation for further explorations on selection methods using different criteria besides this straightforward heuristic. We describe the different ways we can extend this method on section 5.

1.2 Background

Before explaining details of the work done, we present for reference and clarity some of the minimal background knowledge required to understand how code and tests in AL are defined.

As mentioned in the introduction, AL is a DSL for this specific system and extensions for this system. The code is organized in several *AL objects*, which do not resemble objects in the traditional OOP sense, but objects from the system runtime. *AL objects* can be from different types, for example:

- *Tables*. Corresponding to underlying SQL Server tables.
- *Pages*. Corresponding to interactive pages the user can see.
- *Codeunits*. Modules with *procedures*, that is code that can be reused throughout the codebase

Each AL file in the codebase has as part of its definition, its type of object, and a numeric ID. Such pair must be unique for a given codebase.

To make things more concrete, figure 2 shows an example of AL objects of types *page* and *table*. They are identified by both the type of object and ID, notice in this example, they share the same id, but they differ on the type of object.

```

1  page 5200 "Employee Card"
2  {
3      Caption = 'Employee Card';
4      PageType = Card;
5      PromotedActionCategories = 'New,Process,Report,Employee,Navigate';
6      SourceTable = Employee;
7
8      layout
9      {
10         area(content)
11         {
12             group(General)
13             {

```

(a) AL Page

```

1  table 5200 Employee
2  {
3      Caption = 'Employee';
4      DataCaptionFields = "No.", "First Name", "Middle Name", "Last Name";
5      DrillDownPageID = "Employee list";
6      LookupPageID = "Employee list";
7
8      fields
9      {
10         field(1; "No."; Code[20])
11         {
12             Caption = 'No.';
13             trigger OnValidate()
14             begin

```

(b) AL Table

Figure 2: Examples of *AL Objects*.

A mechanism for code reuse in this language is, as listed above, *codeunits*, which is similar to the concept of “modules” in other languages. They consist of several *procedures*. See an example in figure 3

```

1  codeunit 5200 "Employee/Resource Update"
2  {
3      Permissions = TableData Resource = rimd;
4
5      trigger OnRun()
6      begin
7      end;
8
9      var
10         Res: Record Resource;
11
12         procedure HumanResToRes(OldEmployee: Record Employee; Employee: Record Employee)
13         begin
14             if (Employee."Resource No." <> '') and IsResourceUpdated()
15             then ResUpdate(Employee)
16             else
17                 exit;
18         end;
19
20         procedure ResUpdate(Employee: Record Employee)
21         begin
22             Res.Get(Employee."Resource No.");

```

Figure 3: Example of an *AL Codeunit* with procedures `HumanResToRes` and `ResUpdate`.

Finally, we explain how test scenarios for AL code are defined. Unit and scenario tests for application code are written in AL themselves, they correspond to *codeunits* that are marked as tests. Procedures within such a *test codeunit* may correspond to scenarios to be tested. See figure 4 for an example on how *codeunits* are marked as tests and a scenario example.

```

1  codeunit 139522 "VAT Group Retr. From Sub Test"
2  {
3      Subtype = Test;
4      TestPermissions = Disabled;
5
6      var
7          Assert: Codeunit Assert;
8          LibraryVATGroup: Codeunit "Library - VAT Group";

```

(a) AL Test Codeunit Definition

```

[Test]
0 references
procedure TestVATStatementReportLineMissing()
var
    VATStatementReportLine: Record "VAT Statement Report Line";
begin
    // [SCENARIO 374187] VAT Statement Report Line Missing

    // [WHEN] The table "VAT Statement Report Line" is empty
    VATStatementReportLine.DeleteAll();

    // [WHEN] Running the codeunit "VAT Group Retrieve From Sub."
    // [THEN] An error is expected
    asserterror Codeunit.Run(Codeunit::"VAT Group Retrieve From Sub.");
    Assert.ExpectedError(SuggestLinesBeforeErr);
end;

```

(b) AL Test Procedure - Scenario

Figure 4: Examples of definition of AL Test Scenarios.

2 Test Selection Method

2.1 Coverage information of test runs on AL

AL provides primitives to record line coverage information when executing arbitrary sections of code. This line coverage information has the following components, relevant for the selection method:

- Object ID
- Object type
- Line number
- Number of hits

As explained in section 1.2, *Object ID* and *Object Type* are properties related to the way the code is organized in this language. Both the *Object ID* and *Object Type* identify a file in the codebase, for which the related line number refers to.

It is relevant to stress that this information does not provide a file path on the codebase directly, but only a mean for identification.

Currently on the system, this primitive is used throughout the execution of the test suite. These measures test coverage information, providing the set of objects and lines covered by the suite, and it is used to impose coverage requirements for new additions to the codebase.

This means, that the information obtained by recording coverage throughout the complete execution of the test suite, does not distinguish between different test cases. The coverage information currently available on each run does not have which lines were hit by a particular test, but it only determines that it was hit by some test on the test suite.

However, previously in section 1.1, we described the method as selecting those test cases which affected a particular set of lines. With the current test coverage information, it is not possible to make that distinction.

For this reason, changes on the test runner were needed, to allow for individual tracking of coverage information per test case. However, these changes are currently not on the main pipeline of the product ¹.

Another relevant technical consideration of how this information was dealt with, is how the output of the coverage information of each test case was stored. This was part of the work for this project. An SQLite database was proposed to allow for querying the information in different ways efficiently. The current usage of this information within the product is in form of CSV files, so the corresponding parsing-loading scripts were therefore in the scope of the project as well.

Figure 6 shows the proposed ER diagram in crow's foot notation of the DB schema proposed. AL Objects are modeled with their lines in a one-to-many

¹These changes were not work of the author, but from Nikola Kukrika nikolak@microsoft.com, as they are on the backlog to be added to the pipeline

relationship. These in turn are joined through the **coverage data** pivot table with the test procedures.

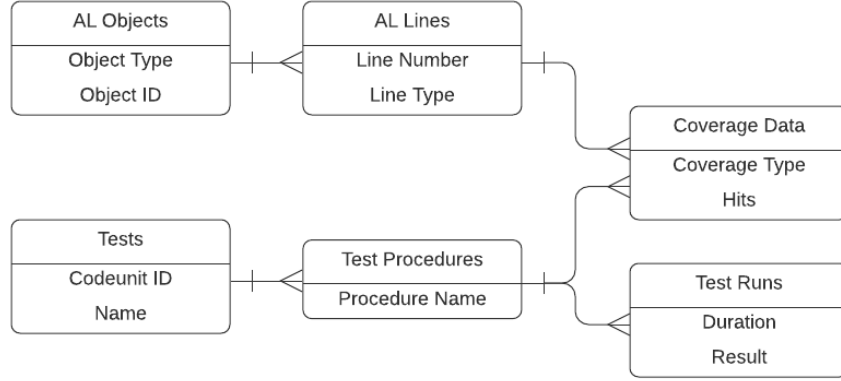


Figure 5: Proposed ER diagram for the coverage database schema

2.2 Representing revision differences

The codebase with the business logic is composed of several AL applications.

When working on a new revision, some of the files on these applications get changed, removed or added. The canonical way of treating with differences between revisions is with line **diffs** of the files, see for example figure 1.

The output information of such **diff** is the set of file paths added, removed, or changed, together with the lines where these changes occurred.

However, in the previous section we explained that the coverage information we are able to extract has as component the Object Id and Object Type where line coverage was recorded. While these properties correspond uniquely to a file in the codebase, it is not a direct inspection of the file path where code was changed that will tell us the information of which Object Id and Object Type that file has.

Even more, we are considering changes on the whole codebase and that may imply changes on tests as well, for example adding new test procedures. For these kinds of changes, we would not like to have the same interpretation on tests that should be selected, we would like for example that whenever test procedures get added, that they are at least selected. This would not be possible by only restricting to coverage information since these tests did not exist previously.

These two situations give way to thinking that more information than just the file path and lines changed is required from a given revision change. For this selection method, a lightweight parsing was performed on the given changed files to extract which Object Id is being affected, and whether or not such object corresponds to a test.

Notice that the choice of what information to extract from a given revision change was just to address these two situations. We could however extract more information, enhancing our knowledge of what kind of change the revision affected; this could be even more granular for example to detect changes on specific procedures, addition of procedures, definition of new columns on a table; and in general, semantical information related to changes. This could be a direction for future work, as input for more complex selection criteria. This is further explained in section 5.

For completeness, table 1 shows the information extracted for identifying a set of revision changes.

Type of change	Parameters
AddObject	None
DeleteObject	removedObjectId
AddLinesObject	objectId, afterLine
RemoveLinesObject	objectId, firstLineRemoved, nLinesRemoved
AddTestObject	testCodeunitId
DeleteTestObject	testCodeunitId
AddLinesTestObject	testCodeunitId
RemoveLinesTestObject	testCodeunitId

Table 1: Information extracted from revision changes

The code for obtaining such representation was as well part of this project. For reference, this consisted in taking the `diff` from the change, visiting each file on the appropriate revision, and parsing to obtain the required information.

2.3 Test selection

Given that coverage information per test procedure is available, and revision changes can be understood in terms of what the coverage information provides; we can specify the heuristic explained in section 1.1 by describing what tests should the algorithm select for each kind of revision change we observe. This is the main core procedure proposed for selecting tests on this project.

2.3.1 Adding lines to an AL object

Whenever lines are added to an AL object, they are added after a given line. The selected tests then are those that were hitting that line and the next.

For further illustration see figure 6, where we can see such a change, a new line was added between what previously was line 1107 and 1108. From the coverage information we can obtain the set of test procedures that on the previous revision were hitting those lines; this method selects those tests.

1104	FormatAddressFields(Header);	1104	FormatAddressFields(Header);
1105	FormatDocumentFields(Header);	1105	FormatDocumentFields(Header);
1106	if SellToContact.Get("Sell-to Contact No.") then;	1106	if SellToContact.Get("Sell-to Contact No.") then;
1107	if BillToContact.Get("Bill-to Contact No.") then;	1107	if BillToContact.Get("Bill-to Contact No.") then;
		1108 +	CompanyInfo.CopyBankAccountInfoForCodeOrDefault(Head
1108		1109	FillLeftHeader;
1109	FillLeftHeader;	1110	FillRightHeader;
1110	FillRightHeader;	1111	
1111		1112	

Figure 6: Example of a change that added a line between what previously was line 1107 and 1108.

2.3.2 Removing lines from an AL object

For removed lines, any test that was hitting what were previously those lines are selected.

2.3.3 Adding an AL object

We can not infer any test run from a new object, since coverage information would not include it.

2.3.4 Removing an AL object

We can understand this case as being the same as removing all the lines of it. Therefore any test hitting that object should be selected. It is worth noticing that this is not a traditional revision change, but it is taken into consideration for completeness.

2.3.5 Adding lines to an AL test codeunit

The current representation for the revision change does not distinguish if the lines added are because of adding new test cases on a given test codeunit, or because of changing a test procedure. Both cases could be treated differently if the model for the revision change allowed for that.

Instead, we only keep the information about which test codeunit lines were added. As a conservative approach, in this method we select all test procedures inside this test codeunit.

2.3.6 Removing lines from an AL test codeunit

Similar to adding lines to an AL test codeunit, as a conservative measure, all tests in this test procedure are selected.

2.3.7 Adding an AL test codeunit

For this case, select all the test procedures in the new test codeunit.

2.3.8 Removing an AL test codeunit

When the change is removing a test codeunit, there is we can not infer running any test. However the Object Id of the removed test codeunit is kept, to filter our selection by removing any references to such tests.

2.4 The full algorithm

We give an overview of how different parts of the selection method explained on the previous section compose together.

A codebase at a given point in time can be enriched when running the test suite with the information of the lines covered by each test case. This is achieved through the modifications on the *Test Runner* application explained on section 2.1. This information is stored in a coverage database.

When a revision is given, we know the set of file paths and lines changed for such. We can leverage from parsing such files on each revision to gather more information about what this revision change implies. The information extracted is to allow for different test selection criteria for each type of change, as it is listed in section 2.2, table 1. More generally, we can think of this step as extracting a helpful representation from the revision change.

With this information, the test selection criteria described in section 2.3 is then applied, resulting on a set of selected test cases.

Figure 7 shows an informal diagram to further illustrate the description here given.

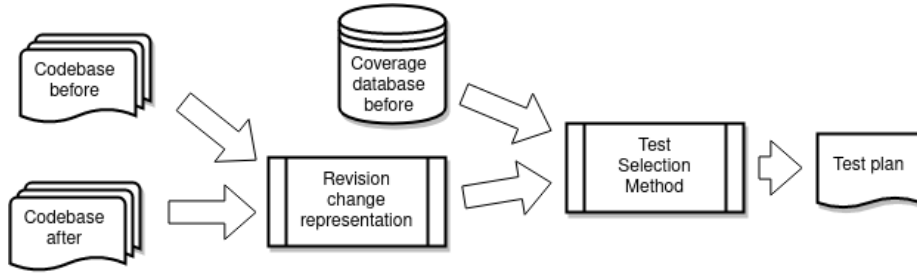


Figure 7: Informal diagram of the procedure for test selection.

3 Evaluation

3.1 Metrics

In the literature for the problem of test selection, several metrics have been proposed to quantify the different aspects addressed by a selection algorithm.

For example, measuring how many methods were selected, relates to the expected performance improvement on the execution of the test suite. After all, choosing all tests is a valid selection algorithm, although not ideal for our goal.

Another example is measuring how effective the algorithm was to select the faulty test cases. If our selection is very small it would be efficient, but if such selection is not capturing the faulty tests, then it is damaging to the product quality.

We select metrics based on a survey done by Pan, Bagherzadeh, et. al. in [Pan2021TestCS] on this topic. One of the research questions they dealt with, was the metrics used to evaluate the Test Selection and Prioritization (TSP) problem. Additionally other common metrics were used.

3.1.1 Test selection execution time

This metric measures the total time of running the selected subset of tests, against the total time that it takes to run all of them.

Ideally we would like this value to be as small as possible.

3.1.2 Test selection time to first failure

This metric measures the time taken for test execution until the first failure is detected, according to the proposed plan; with respect to the total execution time of all tests. This is a metric for evaluating a test prioritization technique, since a selection may have different orderings, and each ordering would change the value of this metric.

In such context, we take the selection to be the prioritized set, and take one of the orderings arbitrarily. This is a relevant metric to implement since in further work we would like to also evaluate prioritization techniques.

Ideally, we would like to detect failures as quickly as possible, that is, the time taken to run the test that fails the quickest.

3.1.3 Inclusiveness

This metric refers to the ratio of the failing test cases included out of the total number of failing test cases. Therefore, for a test selection algorithm that includes all the failing test cases, we would aim to have a value close to 1 for this metric.

3.1.4 Selection size

For this metric, the ratio of the selected tests against the total amount of tests is taken. Ideally, we aim at as low number as possible while detecting every failing test case.

3.1.5 Normalized APFD

Average of the Percentage of Faults Detected is another commonly used metric to measure the effectiveness of Test Prioritization algorithms. First introduced by Rothermel et. al in [962562].

This metric measures the area under the curve of the plot of tests run against faults detected. The normalized version takes into account possibly not detecting faults, which is the case for a selection algorithm. For completeness the expression for computing it is given in (1).

$$nAPFD = p - \frac{\sum_{i=1}^m TF_i}{nm} + \frac{p}{2n} \quad (1)$$

Where p is the ratio of detected over total faults; n is the total number of test cases; m is the total number of faults; and TF_i is the ranking of the first test case that revealed the fault of the i -th failure.

This value measures the proportion of area covered by the described plot, therefore we would like this area to be as close to 1 as possible.

3.2 Selecting real revision examples

Gathering real revision changes that made a test fail, would greatly increase the validity of the evaluation of any test selection method.

In this section we describe the process and difficulties of this, while also outlining the required changes for future work. We present the metrics obtained of running the evaluation in a single failing revision example; while aware of the limitations to argue with these results, the main outcome was the knowledge of the process to further elaborate on.

The build system used for Business Central is referred to as DME. Roughly described, when a developer submits a test job, its branch gets built and tests get run against it. After some time, the developer receives information of whether or not the changes can be integrated into the main branch for the product release being worked on.

A test job is composed of several interdependent tasks, and building the target system for which the tests are run is not a straightforward process. As an example of this, the codebase at a given point in time is a representation of several AL applications, each of which can be localized (country-specific logic organized within *layers*). Meaning that running tests, involves localizing the codebase for each country, building each application, building the appropriate databases, and executing different test solutions on it.

Trying to match this with the simple context we had so far analyzed increased the complexity for properly evaluating. However, the DME system is capable of being queried for further integrations, and in principle most of the process to get the evaluation point is possible, and to certain degree automated.

For this project a program was made to query and extract information of this system, to gather failing test jobs on real revisions.

The main complication faced was getting a reproducible revision example, together with the long time that these heavy operations take, which they are usually performed on clusters of machines for this purpose.

Nevertheless we present the case of job 2412636, a case where we were able to extract the line coverage information and get a prediction with the test selection method proposed. In this test job, changes to an extension `INTaxEngine`

were made, the corresponding test extension were failures were found were **TaxEngine-JsonExchange**. After obtaining the codebase for such revisions, the database and views were constructed for the failing version: the localized version corresponding to IN. To this setup, the modification to the test runner were published, and the test suite was run again on both revisions, collecting the coverage information.

Afterwards, the test selection algorithm was executed. Figure ?? shows the output of the selection algorithm, and figure ?? the failures that occurred in the build system. We can observe that for this case, the inclusiveness results in 0. No selected test was failing. When looking at the stack trace of the error, we can see that the failure was due to not having permissions on new tables after their modifications. This type of faults are hard for our current algorithm to propose, since they depend on the context that the changes introduce. This is not corresponding directly to statements or lines on the codebase, and therefore only using line coverage information does not suffice. In the next section, when evaluating against injected faults on the system, a class for the proposed changed made is context-faults that resemble faults like the ones detected for this job.

```
136802 TestInitializeTaxRates
136802 TestUpdateTaxRates
136802 TestUpdateFromAndToRangeAttributeOnTaxRates
136802 TestUpdateLinkdAttributeOnTaxRates
136802 TestUpdateFromAndToRangeAttributeForError
```

Figure 8: Output of the Selection algorithm for job 2412636.

```
137700 TestImportTaxTypes
137700 TestExportTaxTypes
137700 TestExportNodes
137700 TestUseCaseExecution
137700 TestUseCasePosting
137700 TestTaxInformationFactbox
137700 TestUseCasePostingExecution
```

Figure 9: Failing tests on job 2412636.

The program built for querying the build system provided different failing test jobs together with the log files associated to them. We could use them as starting point to collect the required information and attempt to reproduce them. However, we faced difficulties on streamlining this process. To list some of them, the failing job could be dismissed if the job was submitted as a private deployment for the developer, since developers can request jobs from the pipeline without publishing their changes to the remote repository, making it inaccessible for us to obtain. Or if the pull request associated to the test job was closed,

typically the branch gets deleted as well, making it impossible to get the commit where the failing job was executed. These difficulties were to just obtain the revisions, but also when obtaining the database and binaries associated to the build version, typically only the latest are provided to the developers since they are the ones required for their work, but to reproduce some of the errors, it required in some cases to change the version of the binaries used. Overall, this made the process prone to failure and it needed a large investment of time.

Future work should focus on integrating the collection of evaluation information on the pipeline, so that reproducibility errors are not an issue, and usage of the existing infrastructure to perform the different job tasks is leveraged from. In order to do this, changes would be required on the build system, which themselves live under the main application repository in the form of Infrastructure as Code. We expand more on this in section 5.

3.3 Injecting faults on the codebase

By limiting ourselves to fabricated faults, we reduce the validity of our evaluation. Since modifications were purposefully made to make certain tests fail. We did this however to state that we are capable of computing and gather some insight on the results of this selection, although we are aware that such might be skewed.

We used two test extensions from the real application: **Bank** (Test cases 1-4) and **VATGroupManagement** (Test cases 5-8).

We can identify some classes on the mutations made, based on the type of modifications we made. However, we do not argue they are representatives from the actual test failures that the system usually experiences.

Depending on how lines are mapped to statements we classify the kinds of bugs introduced:

- Context-fault. Lines changed for this kind of bugs did not correspond to statements on the execution, but in the side effects they produced. Test cases #3, #4, #7 and #8 belong to this class.
- Statement-fault. Lines changed correspond to statements on a sequential procedure and therefore are easily captured by the heuristic we used.

The results in table ?? show some of the insight we had previously discussed. When changing lines from a sequential procedure, the coverage information is useful and results on the

3.4 Cost of the algorithm

For practical considerations it is relevant to know the complexity associated to running this procedure. It is important as well to highlight, that test selection is a process that involves several parts, and measuring the complexity of just one part skews the insight of the adoption of such technique in industrial settings. This was observed by Miranda, Cruciani, et. al. in [TODO]

Test case	$t_{Full}(\%)$	$t_{FF}(\%)$	i	$size$	$nAPFD$
#1	60.028	0.044	24.615	34.549	20.004
#2	0.058	0.027	100	0.796	99.867
#3	71.006	0	0	44.562	0
#4	84.198	0	0	49.337	0
#5	35.402	17.968	100	21.212	87.121
#6	33.216	33.216	100	21.212	79.545
#7	0	0	0	0	0
#8	37.058	0	0	24.242	0

Table 2: Evaluation of injected faults.

4 Related work

The heuristic presented as core procedure for selecting test cases is a very natural one, given line coverage information. Furthermore, literature on regression testing and test selection methods involving coverage information is quite extensive. So it is not a surprise that we can find similar approaches existing in literature.

In this section we expand on the comparison with methods found on literature, to also outline possible directions for future work. For this, two surveys were used: the work from Yoo and Harman in [Yoo2009RegressionTM] for general minimisation, selection and prioritisation techniques, to put in perspective the current approach with the academic context. And the work from Pan, Bagherzadeh, et. al in [Pan2021TestCS] where they focus on techniques using Machine Learning, to understand what inputs are required to understand how applicable these techniques can be in future work.

The implemented technique here explored follows quite precisely the formulation for identifying fault-revealing test cases given by Rothermel and Harrold in [TODO]. Here they define a modification-traversing test as such that results on different execution traces for two given revisions of a codebase. Such tests can happen when either they execute new or modified code in the new revision, or if they formerly executed code deleted in the new revision.

By using the *Controlled-Regression-Testing Assumption* they are able to prove that modification-traversing selected tests cases are a superset of the possible fault-revealing tests cases. Such is an strong assumption, and in our case it does not hold, as we gave insight for in section 1.1. There are faults that test cases reveal that are context-aware, they do not correspond to the execution traces of the test cases. For instance the real evaluation example described in section 3.

The algorithm they propose *SelectTests* compares Control Flow Graphs (CFG) generated from each test case, and when different, select all the tests that reached such nodes. We can understand our approach as a lose approximation of this, we do not rely on generating a CFG for the execution of the new revision of test cases, but we directly interpret line operations on files as

a difference and select the corresponding cases. We differ on how we decide whether or not execution traces of a test are equivalent.

Work and terminology introduced by Rothermel and Harrold was foundational for the area and it has been used as a framework to communicate results in surveys and other papers. That is the case of another similar approach by Volkos and Frankl, like in this project, they also take `diff` information as input for their selection; in the same fashion they characterize each line operation, and associate a set of tests depending on the execution traces from previous revisions. The main differences with what here was presented is the preprocessing of revisions into a normal form to make negligible stylistic choices, and to have a better correspondence between lines and execution traces. This is a direction that could be interesting to explore to increase the precision of selected test cases.

The similarity between the CFG traversing approach of Rothermel and Harrold with the textual difference approach from Volkos and Frankl was hinted when we compared how similar our approach was with the CFG approach. This was also observed by Yoo and Harman on their survey, noticing that the textual approach operated on a different representation of the system, but it's behaviour was in essence very similar.

Another related approach that could hint a direction on extension of the technique here proposed is the work by Chen et. al. [TODO] where the CFG is extended by changing the definition of how a program can be partitioned. Changing this representation allows for more entities to be considered. Extending the representation of AL files, and giving a way to relate coverage information with different representations could allow more information to be captured, increasing the inclusiveness of it.

Another variant to this same idea of testing only what is changed based on differences, is the given by Leung and White with their firewall approach, which decompose the system in modules and groups of tests to execute when changes arise. This leads to a more conservative approach, and it also would require the definition of what would be considered as a module.

A very similar approach was also found in [6405252] by Beszédes et. al., here they describe the implementation of a coverage based selection technique on the C++ components of the WebKit project. They follow a similar approach in the technical perspective, by implementing a procedure level coverage database instead of a line coverage database. They identify changed modules from a revision change and execute the related test cases according to their coverage database. This is an analogous process to the one we follow, only taking a coarser perspective by grouping the different changes on the procedures it affected.

On the other hand, a more recent survey focuses on the techniques using Machine Learning algorithms [Pan2021TestCS], as they have been subject of much interest in recent years. Rather than comparing them with the proposed selection algorithm, we present some of the applicable techniques that were found, that we would like to be focus of future work.

Conveniently, the survey exhibits a classification of the methods by artifacts and information provided by the papers. This can be very useful if we want to

compare different selection criteria

5 Future work

Throughout this report we have hinted possible directions for future work. In this section we summarize and expand on them.

5.1 Applicable techniques from literature

In section 4, we went over the different existing approaches to regression testing that followed similar ideas to the one proposed here. Some of the insight collected from such implementation can aid our algorithm and could be a way for extending the current algorithm, like preprocessing files into a normal form, or having different representation for AL files.

We also listed some applicable ML techniques, that use the information we have available, or that minor changes are required for it.

5.2 Different revision representations

Among the different techniques that were shown to be similar in section ??, a missing step on the current implemented technique was getting more information about the change made. For example, in [6405252] they identify the procedures where the changes were made by using different source code representation tools. Also in [] they depend on using a CFG representation of the tests, to match such information and identify in which node of the graph were changes made.

A future direction can be extracting more semantical information from a given change. AL compiles into C# code, which is used by the runtime server of the application to execute the required business logic. C# is a widely studied language with several source code representation tools that we can leverage from.

This is also relevant to explore Machine Learning techniques, since a common step for all of them is the feature extraction. This consists on extracting different kinds of information, both in the source code, coverage information and on historical information about the test cases.

5.3 Infrastructure for integrating evaluation on CI/CD pipeline

To increase our confidence on the validity of the results, we need a more robust evaluation method. For this project we obtained manually the revision information of existing failing jobs, and re-executed the test job failures gathering the coverage information to have the required inputs for our selection algorithm. This process was shown to be hard to reproduce, as explained in section 3.2.

A better approach to this problem would be not to reproduce the error as it happened, while collecting the required information. Instead we would like to

collect that information from the job itself. This would require changes in the infrastructure pipeline

They conceptualize their build systems in terms of metamodels and instances of such correspond to actual execution plans. Scope of future work here should focus on making the modifications to these models for evaluation data collection.

5.4 Efficiency of producing coverage database

A relevant engineering concern not addressed when introducing the proposed coverage database is the efficiency of producing the coverage database. While having a database allows for more efficient and flexible querying of information instead of the existing CSV schema of processing coverage; it comes currently with a toll in performance when inserting into the coverage database.

There are several optimizations possible regarding transaction management in SQLite, and . This work item becomes increasingly important if looking to integrate the evaluation on the CI/CD pipeline.

6 Conclusion

In this project, we have built some of the infrastructure required to run test selection techniques of real revision examples of the real codebase of Business Central. Namely: extracting a computer-friendly representation of the revision changes of AL code, getting

An important limitation and threat to validity is the correct evaluation of the methods we would like to explore.