

~~CSE 503: Project Proposal: Semantical Merge~~

Ji Luo (Student No. 1940053; UW NetID luoji)

23 April 2019

Motivation

Version control systems (VCSes) are ubiquitously adopted in modern software development for **codebase tracking**. In distributed VCSes like Git, the history of a project is modeled as a directed acyclic graph, where each revision (commit) is a vertex pointing to its parents. Coding by a programmer corresponds to creating vertices with one parent — the parent represents the version on which the change is based, and the new vertex represents the updated version. When the software is developed by a team, team members contribute to the codebase concurrently, resulting in a graph where one vertex might have multiple children. Consolidation of team effort is achieved by merging multiple revisions. Automatic merging of revisions makes collaboration more effective, and finding a good merging algorithm is of great interest.

Computing the difference (**diff**) between revisions is closely related to merging. A widely used strategy (if not the only one) for merging is to compute the difference between the two vertices to be merged **with their nearest common ancestor**, then combine **the differences** to produce the result. Therefore, one can improve the merging algorithm by improving the **diff** algorithm, the way differences are combined, or both. In addition to being a stage of the merging algorithm, **diff** is often displayed so that a programmer can read out the changes between revisions, making it useful on its own.


Git computes **diff** of text files by line **and binary files as a whole**, which naturally leads to line-based merging strategy. The method takes advantage of the fact that code files are often formatted so that each line roughly corresponds to a part in the semantical hierarchy (both for semantics of the programming language and semantics to human). However, line-based methodology has its problems in both computing **diff** and merging. The reason is that lines **are a bad proxy of** the semantical hierarchy.


Many programming languages ignore the difference between a whitespace character and consecutive whitespace characters, so line separators and indentation can be inserted or removed without changing the semantics of the program. Insertion and removal of line separators are often done for aesthetic reasons (e.g., to keep each line shorter than 80 characters). Insertion and removal of indentation take place with refactoring or change of logic external to the chunk being indented/unindented. When a line-based **diff** is produced, such less important changes are mixed with other semantically significant changes, making the **diff** less easy to understand for human (e.g., in Figure 1, the indentation is mixed with the real change — the introduction of a new condition). Another problem of line-based **diff** is that it sometimes matches irrelevant lines that are identical by accident (Figure 2).

Bad **diffs** make merging difficult. The previous example continued, if a block of code is indented in one revision and has an expression changed in the other, Git merging algorithm gives up and reports a merge conflict. For another example, even if there is no merge conflict, line-based (or more generally, text-based) merging might yield an undesirable result. Suppose a function is renamed in one revision, and in the other revision, another call to the function is added (using its old name). Merging the two revisions result in a state where the function is renamed, yet there is one call to that method using its old name. This subtle bug can be hard to notice if the old name happens to be reused in the first revision.


Summary of Comments on semantic-merge-commented.pdf

Page: 1


 Number: 1 Author: mernst Subject: Cross-Out Date: 4/24/2019 2:20:23 PM

 Number: 2 Author: mernst Subject: Highlight Date: 4/24/2019 2:20:18 PM

What is "codebase tracking"?


 Number: 3 Author: mernst Subject: Highlight Date: 4/24/2019 2:24:33 PM

This is not clear. "Difference between the two vertices" makes it sound like this is computing the difference between those two vertices. Exactly how many differences are being computed, and exactly what are the differences between? Can two vertices have multiple nearest common ancestors?


 Number: 4 Author: mernst Subject: Highlight Date: 4/24/2019 2:31:30 PM

What is the type of "the differences"? The next sentence ("improving the diff algorithm, the way differences are combined, or both") implies that the type of this is "a diff".


This seems to me to be too narrow a conception of the problem. I can imagine algorithms that do merging without ever running the diff program -- for example, algorithms that compute a difference over the AST rather than the program's textual representation. If you want to restrict yourself to textual algorithms, that's OK, but it is not the full space of possible solutions. You need to acknowledge, somewhere in the paper, other possible approaches, and here you should not say they are impossible.

 Number: 5 Author: mernst Subject: Highlight Date: 4/24/2019 2:29:37 PM

I think it might make no attempt to diff binary files. This text suggests that it does diff binary files.

 Number: 6 Author: mernst Subject: Highlight Date: 4/24/2019 2:30:38 PM

You just said that "each line roughly corresponds to a part in the semantical hierarchy", so this text seems to contradict your previous statement.

 Number: 7 Author: mernst Subject: Highlight Date: 4/24/2019 2:33:53 PM

Please write a topic sentence that clarifies the topic of this paragraph.

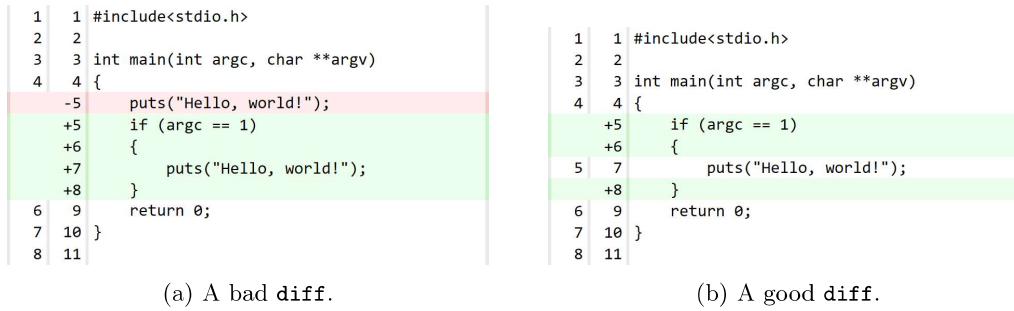


Figure 1: Example exhibiting undesired result from line-based diff. Created with <https://strcmp.cc/>.

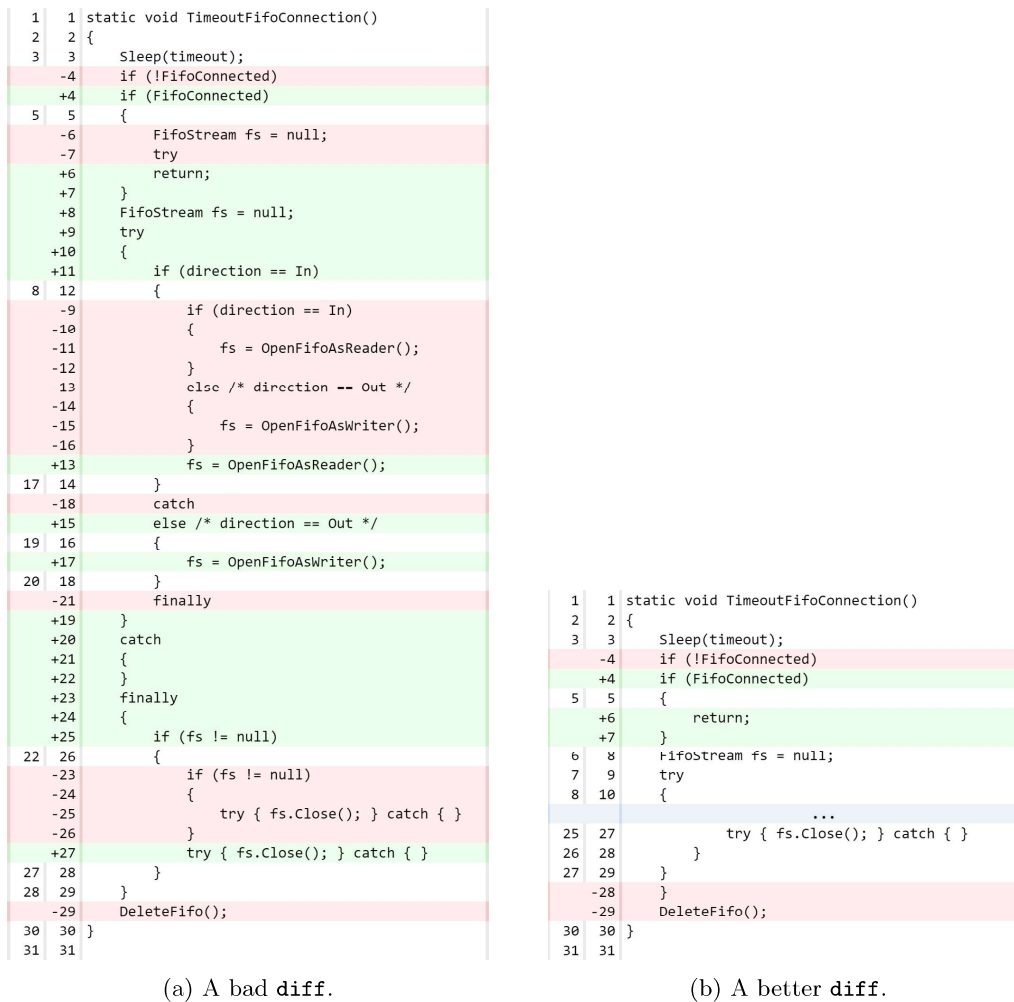


Figure 2: Example exhibiting result that forgets language structure from line-based diff. The change is from deep nesting to early return. Note how the block delimiters for **if-else** and **try-catch-finally** are mingled.

Project Plan

[AGMO17] uses tree-based `diff` and merging algorithms. We plan to pursue and push forward the methodology based on [1] `ASTs of the code`. We will use [2] `static semantics`, a step [3] `forward` from `ASTs`, to do `diff` and merging.

[4] `the diff algorithm` [5] `will have a subprocedure` that is similar to the first few stages of a compiler. The subprocedure takes in source code, parses it into `AST` and [6] `perform syntax-directed translation` on it to obtain static semantics of the code (which functions are defined, which the arguments are, which function this expression calls and what parameters are passed, etc.). The `diff` algorithm, given two sets of code files, compares the static semantics, and determines and outputs the difference. The difference is not a textual one, but [7] `symbolic one`. To merge two revisions, the merging algorithm compares both against their nearest common ancestor, [8] `combines the symbolic difference`, and applies the combined difference to the ancestor.

[10] `analyzing static semantics` [11] `'t an easy-to-program task`. The C programming language [9] `has a good trade-off` between language complexity and expresiveness, and it is also a popular programming language. So in this project, we plan to implement the algorithms that `diffs` and merges C source code.

We plan to evaluate the project by testing the algorithms on several publicly available repositories that mainly consist of C code files with the following metrics:

- [12] `the length of diffs` produced for neighboring revisions (implicitly compared to the necessary length).
- [13] `the frequency and size of merge conflicts`.
- The consistency of automatic merging results [14] `th manual merging results`.
- The frequency of cases where the algorithm raises a merge conflict [15] `st for which Git` silently produces a textually correct and semantically incorrect merge.

Lastly, we point out the novelties of our project:

- It explicitly decomposes merging into 2 stages: [17] `roducing diffs` and `combining diffs`, and [16] `poses the methodology of improving diff` (structurally, like from lines to trees; or within the same structure) `to improve merging algorithm`.
- It does not try to use a textual encoding for the internal structure to be `diffed` (as is done for `ASTs` in [AGMO17]). [18] `ther, it operates natively on the structured (in-memory) representation of the semantics`.

References

- [AGMO17] Dimitar Asenov, Balz Guenat, Peter Müller, and Martin Otth. Precise version control of trees with line-based version control systems. In Marieke Huisman and Julia Rubin, editors, *Fundamental Approaches to Software Engineering*, pages 152–169, Berlin, Heidelberg, 2017. Springer Berlin Heidelberg.

-
- 1** Number: 1 Author: mernst Subject: Highlight Date: 4/24/2019 2:34:46 PM
The entire previous section was about text-based merging and suggested it was the only possible approach. Therefore, this comes as quite a surprise.
-
- 1** Number: 2 Author: mernst Subject: Highlight Date: 4/24/2019 2:34:08 PM
Please define this term.
-
- 1** Number: 3 Author: mernst Subject: Highlight Date: 4/24/2019 2:38:47 PM
-
- 1** Number: 4 Author: mernst Subject: Highlight Date: 4/24/2019 2:36:01 PM
When written in that typeface, diff is a tool that works on text files. If you mean something different, then don't use the name of a specific Unix utility.
-
- 1** Number: 5 Author: mernst Subject: Highlight Date: 4/24/2019 2:37:04 PM
Don't introduce implementation details until you have defined the semantics and stated the inputs and outputs. Otherwise, the description is confusing.
-
- 1** Number: 6 Author: mernst Subject: Highlight Date: 4/24/2019 2:38:41 PM
I think this phrase means producing an executable. You don't mean that, though. Please clarify what you mean and use more precise terms.
I think that you have just described parsing, though perhaps also type resolution.
-
- 1** Number: 7 Author: mernst Subject: Highlight Date: 4/24/2019 2:41:37 PM
It is essential that you precisely define the output. From the description, I cannot understand or comment on your proposed approach. Showing some examples might be useful too -- certainly more useful than the page devoted to textual diffs.
-
- 1** Number: 8 Author: mernst Subject: Highlight Date: 4/24/2019 2:42:32 PM
How? This is the heart of the proposal. The most important part of the proposal is stated in three sentences. You must expand it.
-
- 1** Number: 9 Author: mernst Subject: Highlight Date: 4/24/2019 2:43:51 PM
What is the metric you are using to determine that it has a "good trade-off". Why is that even relevant? If you want to work on C programs, you can justify that by pointing out that C is widely used.
-
- 1** Number: 10 Author: mernst Subject: Highlight Date: 4/24/2019 2:43:01 PM
Do you mean analyzing it, or producing it? You don't need to produce it because you can use an existing parser.
-
- 1** Number: 11 Author: mernst Subject: Highlight Date: 4/24/2019 2:44:37 PM
Why is this relevant? It does not seem to be a topic sentence for the paragraph.
-
- 1** Number: 12 Author: mernst Subject: Highlight Date: 4/24/2019 2:45:26 PM
This is not the most important metric, in my view. Why do you put it first?
-
- 1** Number: 13 Author: mernst Subject: Highlight Date: 4/24/2019 2:46:43 PM
What is the metric? Is more better, or less? I don't see why a programmer cares about this, except insofar as the tool should not produce incorrect merges. So, you should measure that directly.
-
- 1** Number: 14 Author: mernst Subject: Highlight Date: 4/24/2019 2:47:06 PM
This seems like the only really important metric.
-
- 1** Number: 15 Author: mernst Subject: Highlight Date: 4/24/2019 2:49:44 PM
This is one small comparison with Git's built-in algorithm. It would be better to define metrics that are valid for any merging approach whatsoever, and then you can compare those metrics for a variety of algorithms, including Git's. Just comparing some subset of your tool's behavior with Git's does not present a full picture of the tradeoffs between them.
- Therefore, create a set of metrics that is relevant to any merge tool.
-
- 1** Number: 16 Author: mernst Subject: Highlight Date: 4/24/2019 2:52:03 PM
Do you make any change to the merging algorithm, or do you create textual diffs and then use existing merge algorithms?
Do existing algorithms work in the way you proposed, first creating a (textual) diff and then merging that diff? If so, in what respect is your approach novel?
-
- 1** Number: 17 Author: mernst Subject: Highlight Date: 4/24/2019 2:50:23 PM
I don't know what a diff is. Written this way, it seems to be a textual file.
-
- 1** Number: 18 Author: mernst Subject: Highlight Date: 4/24/2019 2:53:03 PM
This is not a representation of the semantics, but of the program or the parse tree. I remain deeply confused about what you propose to

Project Plan

[AGMO17] uses tree-based `diff` and merging algorithms. We plan to pursue and push forward the methodology based on `ASTs of the code`. We will use `static semantics`, a step `forward` from `ASTs`, to do `diff` and merging.

The `diff` algorithm will have a subprocedure that is similar to the first few stages of a compiler. The subprocedure takes in source code, parses it into `AST` and `perform syntax-directed translation` on it to obtain static semantics of the code (which functions are defined, which the arguments are, which function this expression calls and what parameters are passed, etc.). The `diff` algorithm, given two sets of code files, compares the static semantics, and determines and outputs the difference. The difference is not a textual one, but `a symbolic one`. To merge two revisions, the merging algorithm compares both against their nearest common ancestor, `combines the symbolic difference`, and applies the combined difference to the ancestor.

Analyzing static semantics `isn't an easy-to-program task`. The C programming language `has a good trade-off` between language complexity and expresiveness, and it is also a popular programming language. So in this project, we plan to implement the algorithms that `diffs` and merges C source code.

We plan to evaluate the project by testing the algorithms on several publicly available repositories that mainly consist of C code files with the following metrics:

- `The length of diffs` produced for neighboring revisions (implicitly compared to the necessary length).
- `The frequency and size of merge conflicts`.
- The consistency of automatic merging results `with manual merging results`.
- The frequency of cases where the algorithm raises a merge conflict, `yet for which Git` silently produces a textually correct and semantically incorrect merge.

Lastly, we point out the novelties of our project:

- It explicitly decomposes merging into 2 stages: `producing diffs` and `combining diffs`, and `proposes the methodology of improving diff` (structurally, like from lines to trees; or within the same structure) `to improve merging algorithm`.
- It does not try to use a textual encoding for the internal structure to be `diffed` (as is done for `ASTs` in [AGMO17]). `Rather, it operates natively on the structured (in-memory) representation of the semantics`.

References

- [AGMO17] Dimitar Asenov, Balz Guenat, Peter Müller, and Martin Otth. Precise version control of trees with line-based version control systems. In Marieke Huisman and Julia Rubin, editors, *Fundamental Approaches to Software Engineering*, pages 152–169, Berlin, Heidelberg, 2017. Springer Berlin Heidelberg.

do. You should already know what this is and you should have started to implement it already, so it should be possible to explain it in the document.