

Handling `#ifdef` Expressions in CPPSTATS

Claus Hunsen
`hunsen@fim.uni-passau.de`

September 2014

^{test} CPPSTATS was initially developed by Jörg Liebig at University of Passau for a set of studies [LAL⁺10, LKA11]. In 2013, Claus Hunsen from the same university has taken over development. Later, CPPSTATS was used in another study [HZS⁺14]. The most current version of CPPSTATS is available at <https://github.com/clhunsen/cppstats/>. Further information can be found at <http://fosd.net/cppstats>.

1 Introduction

CPPSTATS is a tool for analyzing software systems regarding their variability. Therefore, we focus on software systems written in C using the capabilities of the CPP (the C pre-processor) to express variability. CPPSTATS handles the expressions of the CPP inclusion-guards (`#if`, `#elif`, `#endif`, etc.) in a special way, which is why we provide the used procedure in this document.¹ The handling directly affects the scattering and tangling analyses of CPPSTATS.² Processing steps that affect these analyses are marked with `*` in the headline.

There is handling of `#ifdef` expressions within CPPSTATS during three different parts of the whole program: 1) light adaption of the expressions during *preparation* part of CPPSTATS, before generating SRCML files from the source code; 2) collecting the expressions from the SRCML files and rewriting them by making implicit tangling explicit; and 3) building a global expression pool for the analyzed software-project.

2 Example

As `#ifdefs` are explained best by means of an example, the three common pattern of CPP usage are shown in Figure 1. Each part of the example is present in a different file.

¹We mainly describe the mechanism that are performed in the GENERAL analysis (`analyses/general.py`) as the main analysis of CPPSTATS, but the mechanisms relate to the other analyses, too.

²Only the scattering and tangling measurements over `#ifdefs`, not the ones that measure over files!

Part (a) shows nesting of `#ifdefs` and duplicate expressions, while Part (b) and (c) show the use of `#else` and `#elif` branches in an `#ifdef` cascade, while Part (c) also contains an include guard.

During this document, the reader is referenced to these small examples to illustrate all matters. For reproducibility, the example files are included in the folder `example`, right next to this document. Run an analysis of CPPSTATS directly on this folder to reproduce the data.

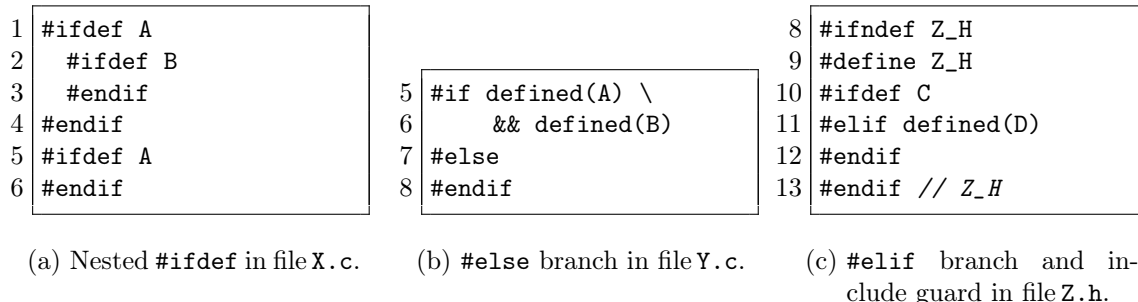


Figure 1: Short examples of the patterns that occur while using CPP and that are treated by CPPSTATS. Each example and their rewriting rules are explained in this very document.

3 Source-Code Preparation Before Generating SRCML files

The first part of a CPPSTATS run is the source-code preparation that transforms the plain-text source code to SRCML³ files. SRCML is an XML format for source code that preserves line numbers and preprocessor information.

Before transforming the source code to SRCML, there are several code-normalization steps, which heavily depend on the code analysis to be performed. For the different preparation types, please refer to the file `preparation.py` of CPPSTATS.

There are up to three possible steps in the preparation that can have effect on the `#ifdef` expressions, though all have only cosmetic effect regarding the actual variability implementation: 1) the handling of multi-line expressions; 2) the rewriting of the shortcut expressions `#ifdef` and `#ifndef`; and 3) the removal of include guards.

After all steps of the preparation, CPPSTATS has generated proper SRCML files for each source file. Based on these files, CPPSTATS performs the actual measurements.

3.1 Multi-Line `#ifdef` Expressions (`preparations/rewriteMultilineMacros.py`)

A multi-line `#ifdef` expression is shown in Fig. 1b on Lines 5 and 6. This expressions is rewritten as a single-line expression, so that CPPSTATS does not have to handle line breaks in `#ifdef` expressions later. The line numbers are preserved during this step.

³<http://www.srcml.org/>

The output, after applying this change to Fig. 1b, is shown in Fig. 2b. The other files are not changed.

<pre> 5 #if defined(A) \ 6 && defined(B) 7 #else 8 #endif </pre>	<pre> 5 #if defined(A) && defined(B) 6 7 #else 8 #endif </pre>
(a) Repetition of Fig. 1b, containing a multi-line expression on Lines 5 and 6.	(b) The same code with single-line expressions.

Figure 2: An `#ifdef` expression in file `Y.c`, (a) before and (b) after rewriting multi-line expressions into single-line fashion.

3.2 Rewriting of `#ifdef` and `#ifndef` (preparations/rewriteIfdefs.py)

As another step, the `#ifdef` expressions are more streamlined by rewriting the conditionals `#ifdef` and `#ifndef` to their long forms. This yields `#if defined(E)` for the conditional `#ifdef E`, and `#if !defined(F)` for `#ifndef F`.

3.3 Removal of Include Guards (preparations/deleteIncludeGuards.py)

The last preparation step that affects `#ifdef` expressions is the removal of include guards. Include guards are used to prevent the multiple inclusion of header files and would bias the results of CPPSTATS, because they are not a mechanism to implement variability.

An inclusion guard, as shown in Fig. 1c consists of three parts: 1) a conditional-inclusion guard checking if the file was already included (Line 8); 2) the definition of the file guard for that is checked in Line 8 (Line 9); and 3) the closing `#endif` for the `#ifndef` (Line 13).

When a header file containing this construct is included again, the conditional in Line 8 will be false, because `Z_H` is already defined. The preprocessor will skip over the entire contents of the file, and the compiler will not see it twice.

As a result of this preparation step, the source-code lines corresponding to the include guard (Lines 8, 9, and 13) are removed from the file and exchanged with blank lines, so that line numbers are preserved.

4 Processing of Expressions During File Analysis *

CPPSTATS performs its analyses based on the SRCML files that are generated at the end of the preparation step. The `#ifdef` expressions are collected while traversing the XML tree of the input files; the main functionality is implemented in the function `_getFeatures` (`analyses/general.py`, Lines 565ff.). During this action, the nesting hierarchy of the conditionals is remembered and used to rewrite the `#ifdef` expressions accordingly in the

function `_getFeatureSignature` (Lines 479ff.). Implicit tangling of `#ifdef` expressions are made explicit in this step, while preserving all nesting-depth information.

For example, CPPSTATS rewrites all nested `#ifdefs`, each with a condition that conjoins their own conditional expression with the enclosing ones, because the inner expression depends on the evaluation of its own condition and, additionally, on the evaluation of the surrounding `#ifdef` expression. This rewriting is analogously done also for `#else` and `#elif` expressions.

The result of these mechanisms, after applying them on the files from Fig. 1, is shown in Fig. 3.

```
1 #if defined(A)
2   #if defined(A) && defined(B)
3   #endif
4 #endif
5 #if defined(A)
6 #endif
```

(a) Nested `#ifdef` in file `X.c`.

```
5 #if defined(A) && defined(B)
6
7 #elif !(defined(A) && defined(B))
8 #endif
```

(b) `#else` branch in file `Y.c`.

```
8
9
10 #if defined(C)
11 #elif (!(defined(C))) && (defined(D))
12 #endif
13
```

(c) `#elif` branch and include guard in file `Z.h`.

Figure 3: Rewritten `#ifdef` expressions from Fig. 1, also including all changes from Section 3. (Note that the files are formatted as SRCML actually, but, for better understanding, the source-code representations are shown here.)

5 Global `#ifdef` Expression Pool

There are two steps for building a global pool of expressions: 1) constructing a local pool of `#ifdef` expressions per file `*`; and 2) combining all local pools to a global pool.

5.1 Construction of Local Expression Pools *

The first step is already processed during the mechanisms of Section 4: Each distinct `#ifdef` expression is added to the list of expressions **only once per file!** The equality of expressions is checked via string equality (`wrapFeatureUp` in `analyses/general.py`, Lines 592ff.). Accordingly, the expression list for the file `X.c` from Fig. 3a is `[defined(A), defined(A) && defined(B)]`. Note that the second expression `defined(A)` from Line 5 is **not** added.

5.2 Construction of Global Expression Pool

After a file is processed entirely and its list of expressions is built, the list is added to the global pool of expressions in `_mergeFeatures` (`analyses/general.py`, Lines 1175ff.). Afterwards, the processed `#ifdef` expressions from Fig. 3 are added to a global pool of expressions, called `sigmap` (Line 1172). The resulting global pool for scattering and tangling analyses is shown in Fig. 4.

```
1 (defined(A)) && (defined(B))
2 defined(A) && defined(B)
3 (!(defined(C))) && (defined(D))
4 !(defined(A) && defined(B))
5 defined(A)
6 defined(C)
```

Figure 4: The resulting (unordered) global expression pool for the example files from Fig. 1 that is used for tangling and scattering analyses.

References

- [HZS⁺14] Claus Hunsen, Bo Zhang, Janet Siegmund, Christian Kästner, Olaf Leßenich, Martin Becker, and Sven Apel. Preprocessor-Based Variability in Open-Source and Industrial Software Systems: An Empirical Study. *Empirical Software Engineering*, 2014. Submitted.
- [LAL⁺10] Jörg Liebig, Sven Apel, Christian Lengauer, Christian Kästner, and Michael Schulze. An Analysis of the Variability in Forty Preprocessor-Based Software Product Lines. In *Proc. Int. Conf. Software Engineering (ICSE)*, pages 105–114. ACM, 2010.
- [LKA11] Jörg Liebig, Christian Kästner, and Sven Apel. Analyzing the Discipline of Preprocessor Annotations in 30 Million Lines of C Code. In *Proc. Int. Conf. Aspect-Oriented Software Development (AOSD)*, pages 191–202. ACM, 2011.