



EÖTVÖS LORÁND UNIVERSITY
FACULTY OF INFORMATICS
DEPARTMENT OF PROGRAMMING LANGUAGES
AND COMPILERS

Enhancing pattern matching-based static analysis of C-family software projects with project-level knowledge

Supervisor:

Richárd Szalay
Ph.D. Student

Author:

Benedek Attila Bahrami
Computer Science B.Sc.

Budapest, 2022

EÖTVÖS LORÁND TUDOMÁNYEGYETEM

INFORMATIKAI KAR

SZAKDOLGOZAT TÉMABEJELENTŐ

Hallgató adatai:

Név: Bahrami Benedek Attila

Neptun kód: C66EPC

Képzési adatak:

Szak: programtervező informatikus, alapképzés (BA/BSc/BProf)

Tagozat : Nappali

Belső témavezetővel rendelkezem

Témavezető neve: Richárd Szalay

munkahelyének neve, tanintézménye: ELTE-IK, Programozási nyelvek és Fordítóprogramok Tanszék

munkahelyének címe: 1117, Budapest, Pázmány Péter sétány 1/C.

beosztás és iskolai végzettsége: PhD student; Computer science MSc

A szakdolgozat címe: Enhancing pattern matching-based static analysis of C-family software projects with project-level knowledge

A szakdolgozat téma:

(A témavezetővel konzultálva adja meg 1/2 - 1 oldal terjedelemben szakdolgozat témájának leírását)

Static analysis is a method to analyse the source code of software projects without performing a real execution of the application. It is widely used in industry to find bugs and code smells during development, to aid in the prevention of bad code that misbehaves in production.

Among various methods, the most important techniques are the ones that are based on pattern matching on a syntactic representation of the software project.

Unfortunately, for programming languages in the C family, such as C++, the concept of "separate translation" causes issues for static analysis.

As most static analysers are built upon compilers, and in C++, each compiler only sees the local information in the source file it is to compile or analyse (as opposed to a more project-level knowledge), crucial details might be hidden, which lowers the accuracy of the analysis.

Clang-Tidy is a static analysis rule collection that is built upon the LLVM Compiler Infrastructure's C-family compiler, Clang.

It performs pattern matching on Clang's "Abstract Syntax Tree" (AST) representation, and generating diagnostics based on which analysis modules -- called "checks" -- the user turns on.

Unfortunately, the current versions of Clang-Tidy checks only access what is visible to the compiler, which is a local information.

Several classes of security issues and bad coding patterns might be diagnosed if the implemented checks would be capable of creating per-compilation knowledge, and reusing the full knowledge about the project during diagnosis.

The work of the thesis is to enhance Clang-Tidy on the infrastructure level to support multi-pass analyses in a generic manner, by utilizing the ideas similar to that of MapReduce.

This is achieved by allowing individual checks to store check-specific data on a thread-safe location.

A subsequent execution of the analysis will be able to do the pattern matching fine-tuned with the data stored in the previous step also available.

To prove the usability of the solution, a new safety and security related check, currently not provided by Clang-Tidy, will be developed utilizing the new infrastructure created in this work.

In the end, the results of the thesis will allow the international community behind LLVM to develop and make available a wider potential of checks.

Contents

1	Introduction	3
1.1	Motivation	3
1.2	Problem Statement	3
1.3	Static Analysis	4
1.4	Infrastructure Limitations	5
1.5	Related Works	7
1.6	Thesis layout	8
2	User documentation	9
2.1	Install guide	9
2.1.1	System Requirements	9
2.1.2	Building from source	10
2.2	Running by Translation Units	12
2.2.1	Configuring the threshold	12
2.2.2	Supressing Warnings	13
2.3	Multiple Phase Version	14
2.3.1	Collect	15
2.3.2	Compact	16
2.3.3	Diagnose	17
3	Developer documentation	20
3.1	Infrastructure	20
3.1.1	Collect Function	24
3.1.2	PostCollect Function	24
3.1.3	Compact Function	24
3.1.4	Additional Functions	25
3.2	Discarded Return Value Check	25
3.2.1	Abstract Syntax Trees	25

CONTENTS

3.2.2	Pattern Matching	26
3.2.3	Checker Logic	31
3.2.4	Multipass phases	34
3.3	Evaluation	39
4	Conclusion	46
4.1	Future Work	46
A	Project results	48
	Bibliography	85
	List of Figures	88
	List of Tables	89
	List of Codes	90

Chapter 1

Introduction

1.1 Motivation

There are programming languages, like Pascal, Ada or BASIC, that distinguish functions that do and do not have return values, the latter are known as *procedures*. C family languages do not. Pre-standardisation C language (also known as K&R C) [1] did not have `void` functions, instead an unspecified return value defaulted to `int`. This resulted in functions declared with `int` return value not returning anything and the supposed return value was unused on purpose. One of its consequences was for example, that for a while if you wrote a function, that was declared to return `int` without the return value, the compiler did not act on it. When writing code in C++ we often use functions from C, and a lot of standard library POSIX functions are C functions, taking `const char *`, `void *`, returning `int`.

1.2 Problem Statement

There are quite a few functions whose return value is often ignored, which could lead to potential bugs. Some examples:

- POSIX `read`: returns the number of bytes read; this return value can also indicate errors with it being -1.
- POSIX `scanf`: returns the number of items in the argument list successfully filled; also indicates errors with EOF return.
- `<cstdio>`'s `std::remove`: return indicates success or error of removal of a file.

- <algorithm>'s `std::remove`: Does not remove. It returns an iterator, and we still need to use container specific `erase` for all elements after this iterator.
- `std::remove_if`: Same as `remove`.
- container-specific `erase`: Returns an iterator to the next element after the removed.
- container-specific `insert`: Returns an iterator to the first of the new elements inserted.

Later, the `[[nodiscard]]` attribute was introduced [2] to notify and give warnings to the user if the return value was unused in case of a function with this attribute. However, in order to ask the compiler to give warnings on unchecked values, we would need permission to modify the library code. In case of external source code such as POSIX, STL or any third-party project, we will not have permission to do so. However, we still need to notify the user on the cases where they do not check non-void return value. This brings us to static analysis.

1.3 Static Analysis

Static analysis is a method to analyse the source code of software projects without performing a real execution of the application. It is widely used in industry [3] to find bugs and code smells during development, to aid in the prevention of bad code that misbehaves in production [4]. Among various methods, the most important techniques are the ones that are based on pattern matching on a syntactic representation of the software project.

Clang-Tidy is a declarative, object oriented, strong typed static analysis rule collection that is built upon the LLVM Compiler Infrastructure's C-family compiler, Clang. It performs pattern matching on Clang's “*Abstract Syntax Tree*” (AST) representation, and generating diagnostics based on which analysis modules, called “*checks*”, the user turns on. Pattern matching is an important tool in the industry as well. It can for example be used for type migration [5]. An example of an AST can be seen at fig. 1.1. We will address both the LLVM library and AST matchers in later chapters of the thesis.

Clang-Tidy's infrastructure is a powerful tool that can even be used to generate FixIts [6], but since a valid automatic fix does not exist for my checker's warnings, (as it is up to the user to decide whether or not the diagnostics is valid, or how they should acknowledge the result) this feature will not be used in this thesis.

```
TranslationUnitDecl
|-FunctionDecl <line:1:1, col:23> col:5 used foo 'int ()'
| `-CompoundStmt <col:11, col:23>
|   `-ReturnStmt <col:13, col:20>
|     `-IntegerLiteral <col:20> 'int' 1
`-FunctionDecl <line:3:1, line:5:1> line:3:5 main 'int ()'
  `-CompoundStmt <col:12, line:5:1>
    `-ReturnStmt <line:4:5, col:16>
      `-CallExpr <col:12, col:16> 'int'
        `-ImplicitCastExpr <col:12> 'int (*)()' <FunctionToPointerDecay>
          `-DeclRefExpr <col:12> 'int ()' lvalue Function 0x5629742c0820 'foo' 'int ()'
```

Figure 1.1: An example of an Abstract Syntax Tree of `return foo();`.

Let us imagine a checker, that keeps statistics on how many times a return value of a non-`void` function is used, or otherwise known as checked. This property is, as we previously stated, important, because there exist a great amount of functions, whose return value should be checked in most situations but remain unchecked in quite a lot.

1.4 Infrastructure Limitations

Unfortunately, for programming languages in the C family, such as C++, the concept of “separate translation” causes issues for static analysis. As most static analysers are built upon compilers, and in C++, each compiler only sees the local information in the source file (also known as the *Translation Unit*) it is to compile or analyse (as opposed to project-level knowledge), crucial details might be hidden, which lowers, or in most cases, completely distorts the accuracy of the analysis. This means that the way our infrastructure works, a checker like this would be of very limited use.

A function can obviously exist outside of the translation unit of its declaration. If such analysis with our imagined checker is done separately on each translation unit, it is easy to see how that can affect the outcome. A function might be called 100 times and checked 95 times through. We would want to give warnings for the unchecked 5% of those calls, but if, for example, these are in a separate translation

```
1 // First translation unit
2
3 std::vector<int> c{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 11, 13, 15};
4
5 auto iter = c.erase(c.begin());
6 iter = c.erase(iter + 2, iter + 5);
7 for (std::vector<int>::iterator it = c.begin(); it != c.end();) {
8     if (*it % 2 == 0)
9         it = c.erase(it);
10    else
11        ++it;
12 }
```

```
1 // Second translation unit
2
3 std::vector<int> c{0, 1, 2, 3, 4};
4 for (std::vector<int>::iterator it = c.begin();
5      it != c.end(); ++it) {
6     if (*it % 2 == 0)
7         c.erase(it);
8 }
```

Code 1.1: An example of the infrastructure’s limitations.

unit, then the analysis would return with 0 checks out of 5 calls. We would not want to give warnings to a function that is unchecked in all 100% of its calls. The detail of the statistics, that it was only unchecked in 5% is lost, unless we use project-level knowledge during our analysis. Consider another example, with code:

Here, with separate analysis, the ratio of checked and unchecked calls for `erase` are 100% in the first TU (2 checked out of 2 calls), and 0% in the second TU (0 checked out of 1call). With project-level knowledge, however we would get 66% (2 checked out of 3 calls). We used `erase` without a care towards its return value, and this could lead to potential bugs, but without the project level knowledge, however we can not diagnose it.

The separate analysis could also create false positive results. Imagine we have a function whose return value we could use but it is mostly optional. We could give our checker a threshold of percentage, to only diagnose unused values if we usually use them in most cases. Again, this means that with different translation units, we do not know how many times we actually ignored the value and can not use our

threshold properly. This leads to false positive diagnostics.

Unfortunately, the current versions of Clang-Tidy checks can only access what is visible to the compiler, which is a local information. Several classes of security issues and bad coding patterns might be diagnosed [7] if the implemented checks would be capable of creating per-compilation knowledge, and reusing the full knowledge about the project during diagnosis. The work of the thesis is to enhance Clang-Tidy on the infrastructure level to support multi-pass analyses in a generic manner, by utilizing the ideas similar to that of MapReduce [8].

This is achieved by allowing individual checks to store check-specific data in a thread-safe location. A subsequent execution of the analysis will be able to do the pattern matching fine-tuned, with the data stored in the previous step also available. To prove the usability of the solution, a new safety and security related check, currently not provided by Clang-Tidy, will be developed utilizing the new infrastructure created in this work. In the end, the results of the thesis will allow the international community behind LLVM to develop and make available a wider potential of checks, as we are planning to upstream the contribution into LLVM.

- D124446 [9]: Add the misc-discarded-return-value check
- D124447 [10]: Add infrastructure support for running on project-level information
- D124448 [11]: Add project-level analysis support to misc-discarded-return-value

1.5 Related Works

The topic of cross translation unit analysis shares other different solutions and implementations as well, one great example is Clang CTU [12]. As previously stated, the infrastructure enhancement follows the philosophy of MapReduce [8]. There exists a static analysis rule in Coverity Scan static analysis tool, that is similar to our checker [13].

1.6 Thesis layout

After the Introduction, in Chapter 2, the User Guide will have instructions on how to download, compile and set up the static analysis tool, and how to run it on a C++ project. In Chapter 3, the Developer Guide will explain how the checker, and the infrastructure itself works, in detail.

Chapter 2

User documentation

Both the changes in the Clang-Tidy infrastructure and our new checker focuses heavily on LLVM’s Clang-Tidy. Clang-Tidy is a Clang-based C++ “linter” tool. Its purpose is to provide an extensible framework for diagnosing and fixing typical programming errors, like style violations, interface misuse, or bugs that can be deduced via static analysis. Clang-Tidy is modular and provides a convenient interface for writing new checks.

The tool, after upstreaming, will be found in the LLVM project repository at <http://github.com/llvm/llvm-project>.

2.1 Install guide

2.1.1 System Requirements

Table 2.1 shows the system requirements and supported compilers for building LLVM. The checkers were developed with Ubuntu 20.04 and tested on Ubuntu 18.04, and WSL Ubuntu 20.04.

Building and using LLVM’s Clang-Tidy takes a lot of time on weaker computers. The minimum recommended memory size for building is 16 GiB, and the optimal amount is 64 GiB of memory.

Operating system	Processor architecture	Compiler
Linux	x86_64	gcc, clang
Linux	amd64	gcc, clang
Linux	arm	gcc, clang
Linux	Mips	gcc, clang
Linux	PowerPC	gcc, clang
Solaris	V9	gcc
FreeBSD	x86_64	gcc, clang
FreeBSD	amd64	gcc, clang
NetBSD	x86_64	gcc, clang
NetBSD	amd64	gcc, clang
macOS2	PowerPC	gcc
macOS	x86	gcc, clang
Cygwin	x86	gcc
Windows	x86	Visual Studio
Windows64	x86_64	Visual Studio

Table 2.1: System requirements and supported compilers for building LLVM.

Software requirements include (at least) GCC version *7.1.0*, CMake version *3.13.4*, Python version *3.6* and GNU Make version *3.79*.

2.1.2 Building from source

These commands will compile LLVM from source. The building process with parameters can be found on the `README.md` of LLVM project Github repository,¹ or the Getting Started² page of Clang documentation. The commands I used for the compilation are shown in listing 2.1.

Explanation for some flags: at `LLVM_USE_LINKER` we can change the linker we are using, either Gold or LLD (Linker for LLVM). The latter needs to be installed. `LLVM_PARALLEL_LINK_JOBS` and `-j` at Ninja sets the CPU multiprocessing job capacity. The recommended amount for `LLVM_PARALLEL_LINK_JOBS` is one quarter of the amount of cores available, and full Ninja job amount should be *cores* – 2. I used these commands on a server with 32 GiB memory and 14 (virtual) CPU cores.

¹<http://github.com/llvm/llvm-project#readme> accessed 2022. 05. 03.

²http://clang.llvm.org/get_started.html, accessed 2022. 05. 03.

```
1 git clone http://github.com/llvm/llvm-project.git
2 cd llvm-project
3 mkdir Build
4
5 # cmake -S llvm -B build -G <generator> [options]
6 cd Build/
7 cmake \
8   -DCMAKE_EXPORT_COMPILE_COMMANDS=ON \
9   -DLLVM_ENABLE_PROJECTS="llvm;clang;clang-tools-extra" \
10  -DLLVM_TARGETS_TO_BUILD="X86" \
11  -DLLVM_APPEND_VC_REV=OFF \
12  -DLLVM_ENABLE_BINDINGS=OFF \
13  -DLLVM_USE_RELATIVE_PATHS_IN_FILES=OFF \
14  -DBUILD_SHARED_LIBS=ON \
15  -DLLVM_USE_LINKER="lld" \
16  -DLLVM_PARALLEL_LINK_JOBS=3 \
17  -DCMAKE_BUILD_TYPE=Release \
18  -DLLVM_ENABLE_DUMP=ON \
19  -DLLVM_ENABLE_ASSERTIONS=ON \
20  -G Ninja \
21  .. llvm
22
23 # cmake --build build [-- [options] <target>] or your build system
24   ↳ specified above directly.
24 ninja -j12 clang-tidy llvm-symbolizer
```

Code 2.1: Compiling LLVM from source on a weaker machine.

2.2 Running by Translation Units

The user can give Clang-Tidy multiple translation units to run on, and it will emit diagnostics separately for each one. You run it by using `clang-tidy -checks='-* ,misc-discarded-return-value' -p ./Build a/main.cpp b/main.cpp`, where the “checks” first disables all checkers with `-*`, then enables our checker, the `-p` flag gets the build path that reads a compile command database, and finally we give the paths to the paths of our source files [14]. Here we are getting two separate diagnostics for our two separate files or translation units.

```

1 bahramib@cc:~/MyFolder/TestFolder$ clang-tidy
  ↵ -checks='-* ,misc-discarded-return-value' -p ./Build a/main.cpp
  ↵ b/main.cpp

2
3 /home/bahramib/MyFolder/TestFolder/a/main.cpp:12:9:
4   warning: return value of 'maybe_check_this' is used
5   in most calls, but not in this one [misc-discarded-return-value]

6
7   MyClass::maybe_check_this();
8   ^
9

10 /home/bahramib/MyFolder/TestFolder/a/main.cpp:12:9: note:
11   value consumed or checked in 75% (3 out of 4) of cases

```

Code 2.2: Diagnostic output without project-level knowledge.

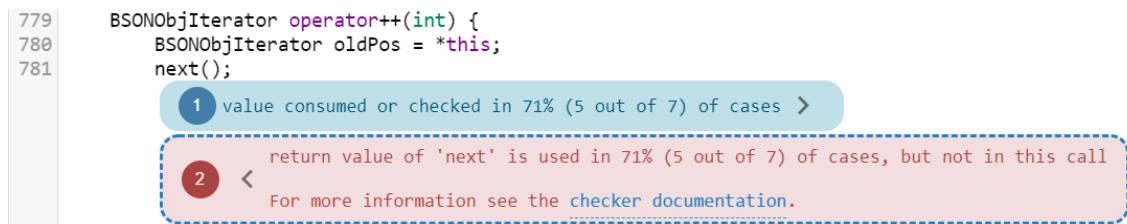


Figure 2.1: A report of the diagnostics of MongoDB’s source on CodeChecker with 50% threshold on a single TU.

Clang-Tidy is supported by CodeChecker. [15]

2.2.1 Configuring the threshold

As previously discussed, the ratio of the amount of checked function calls and all the function calls is very important for the diagnosis, since we do not necessarily

want to be notified of all the unused calls, but let us say if a function's return value is used in at least 65% of all calls, we might want to be warned about the missing unchecked 35%.

The checker has an option for this called the `ConsumeThreshold`, that is the percentage above which we would like to be warned, but not under. You can configure `ConsumeThreshold`, and any Clang-Tidy checker option with Tidy's `-config` flag. Here you will have to type the name of the checker and its option member, with the desired value as a *key-value pair*. This is how it looks with 50%: `-config='{CheckOptions:[{key: misc-discarded-return-value.ConsumeThreshold, value: 50}]}'`. With this option the full command looks like this: `clang-tidy -checks='-* ,misc-discarded-return-value' -config='{CheckOptions:[{key: misc-discarded-return-value.ConsumeThreshold, value: 50}]}' -p ./Build a/main.cpp b/main.cpp`.

2.2.2 Supressing Warnings

Let us review what is wrong with our code in listing 2.3. In function `foo`, we forgot to use our parameter `p`, and we do not check the return value when calling `foo`. These triggers one warning for each mistake as seen in fig. 2.2.

```

1 [[nodiscard]] int foo(int p) {
2     return 1;
3 }
4 int main(int argc, char**) {
5     foo(argc);
6 }
```

Code 2.3: An example of both unused parameter and ignored return value with `nodiscard`.

```

<source>:1:27: warning: unused parameter 'p' [-Wunused-parameter]
[[nodiscard]] int foo(int p) {
^
<source>:5:3: warning: ignoring return value of function declared with 'nodiscard' attribute [-Wunused-result]
foo(argc);
^~~~ ^~~~
```

Figure 2.2: The compiler warnings.

If we do not want to fix these issues, but rather silence our warnings, we can simply write `(void)` before the call to cast it to `void` and trick the compiler into thinking that we did check the return value. Similarly, we could cast our parameter `p` to `void` as a distinct statement to make it seem like we used it. On listing 2.4 we can see the suppressed version that does not yield warnings.

```
1 [[nodiscard]] int foo(int p) {
2     (void)p;
3     return 1;
4 }
5 int main(int argc, char**) {
6     (void)foo(argc);
7 }
```

Code 2.4: The same example, now with suppressed warnings.

The suppression of our new checker warnings follows the exact same convention. Casting to `void` acts as usage of the return value and thus will not give warnings, despite not actually checking it.

Both the philosophy and practice of silencing a single case of an unchecked function call is essentially the same as silencing the warning of `[[nodiscard]]` or an unused parameter. By simply writing `(void)` before the call, to cast it to void, tricking the compiler and the checker into thinking it is used, but not affecting anything in reality. If we got a warning on the line `foo();`, from the checker or `[[nodiscard]]`, we will get neither if we change it to `(void)foo();`.

2.3 Multiple Phase Version

The updated infrastructure contains two new flags for running Clang-Tidy, `multipass-phase` and `multipass-dir`. `multipass-phase` is an `enum` flag, that has three values, `"collect"`, `"compact"` and `"diagnose"` with the latter as default. `multipass-dir` needs a path to a directory where the checkers that support the collect feature can dump their collection data that they are going to compact and use later.

2.3.1 Collect

The Collect phase, as the name suggest, will have the checkers collect data on each translation unit and write them into unique YAML files to later reuse this data. This is how you normally run collect phase on the desired files:

```
1 clang-tidy \
2 -checks='-* ,misc-discarded-return-value' \
3 --multipass-phase=collect \
4 --multipass-dir='MyCollectionDirectory' \
5 -p ./Build \
6 a/main.cpp b/main.cpp
```

What the Discarded Return Value checker does in this phase is count the amount the declared non-void functions were called, and count the amount that these function's return values were checked. After finishing counting in one translation unit, it writes the collected numbers and function names into a YAML file.

After collecting, you do not get any diagnostics or output text, but the desired amount of (in this case 2) YAML files are generated.

```
1 bahramib@cc:~/MyFolder/TestFolder/MyCollectionDirectory$ ls -l
2 -rw-r--r-- 1 bahramib 196 Apr 30 13:00
   ↳ misc-discarded-return-value.main.cpp.12949585208029997868.yaml
3 -rw-r--r-- 1 bahramib 196 Apr 30 13:00
   ↳ misc-discarded-return-value.main.cpp.4924802982073527590.yaml
```

Code 2.5: The YAML files containing the collection data.

```
1 # First TU, a/main.cpp
2 ---
3 - ID:          'c:@S@MyClass@F@check_that#S'
4   Consumed:    3
5   Total:       3
6 - ID:          'c:@S@MyClass@F@maybe_check_this#S'
7   Consumed:    0
8   Total:       3
9 ...
```

```
1 # Second TU, b/main.cpp
2 ---
3 - ID:          'c:@S@MyClass@F@check_that#S'
4   Consumed:    1
5   Total:       3
6 - ID:          'c:@S@MyClass@F@maybe_check_this#S'
7   Consumed:    3
8   Total:       4
9 ...
```

Code 2.6: Contents of the results of collecting data.

2.3.2 Compact

The Compact phase will iterate through the collected data per checker and have the checkers read, use and compact all the data collected into one YAML file. Flags aside from `checks`, `multipass-phase` and `multipass-dir` have no effect.

```
1 clang-tidy \
2 -checks='-* ,misc-discarded-return-value' \
3 --multipass-phase=compact \
4 --multipass-dir='MyCollectionDirectory' \
```

This checker reads the data on each function back and constructs new data similar to the previous ones. If one function is called in multiple translation units, then it adds the numbers from the TU's YAML file to the new structure. After it is finished, the data is written into a single YAML file.

This phase does not write anything on standard output either, but will construct the compacted YAMLs per checker.

```
1 bahramib@cc:~/MyFolder/TestFolder/MyCollectionDirectory$ ls -l
2 -rw-r--r-- 1 bahramib 196 Apr 30 13:00
   ↳ misc-discarded-return-value.main.cpp.12949585208029997868.yaml
3 -rw-r--r-- 1 bahramib 196 Apr 30 13:00
   ↳ misc-discarded-return-value.main.cpp.4924802982073527590.yaml
4 -rw-r--r-- 1 bahramib 196 Apr 30 13:06
   ↳ misc-discarded-return-value.yaml
```

Code 2.7: The new file containing the collected data.

```
1 ---
2 - ID:           'c:@S@MyClass@F@check_that#S'
3 Consumed:      4
4 Total:         6
5 - ID:           'c:@S@MyClass@F@maybe_check_this#S'
6 Consumed:      3
7 Total:         7
8 ...
```

Code 2.8: Contents of the compacted file: the grouped sum of the input.

2.3.3 Diagnose

For backwards compatibility, the default diagnose phase will do exactly what the non-multiple phase Clang-Tidy did. It gives diagnostics for each translation unit separately, as demonstrated in section 2.2, if compact has not been performed for a checker. Otherwise that checker will read and use the data compacted in the respective YAML file, and give its diagnostics calculated with project-level knowledge.

```
1 clang-tidy \
2 -checks='-* ,misc-discarded-return-value' \
3 --multipass-phase=diagnose \
4 --multipass-dir='MyCollectionDirectory' \
5 -p ./Build \
6 a/main.cpp b/main.cpp
```

My checker reads in the compacted data with the call and check amounts for each function and simply decides if diagnostics are needed or not for each unchecked return value in the current translation unit. The output is the diagnostics.

```

1 /home/bahramib/MyFolder/TestFolder/a/main.cpp:12:9:
2   warning: return value of 'check_that' is used
3   in most calls, but not in this one [misc-discarded-return-value]
4
5     MyClass::check_that();
6     ^
7
8 /home/bahramib/MyFolder/TestFolder/a/main.cpp:12:9: note:
9   value consumed or checked in 66% (4 out of 6) of cases
10
11
12 /home/bahramib/MyFolder/TestFolder/a/main.cpp:15:5:
13   warning: return value of 'check_that' is used
14   in most calls, but not in this one [misc-discarded-return-value]
15
16     MyClass::check_that();
17     ^
18
19 /home/bahramib/MyFolder/TestFolder/a/main.cpp:15:5: note:
20   value consumed or checked in 66% (4 out of 6) of cases

```

Code 2.9: Diagnostic output with project-level knowledge.

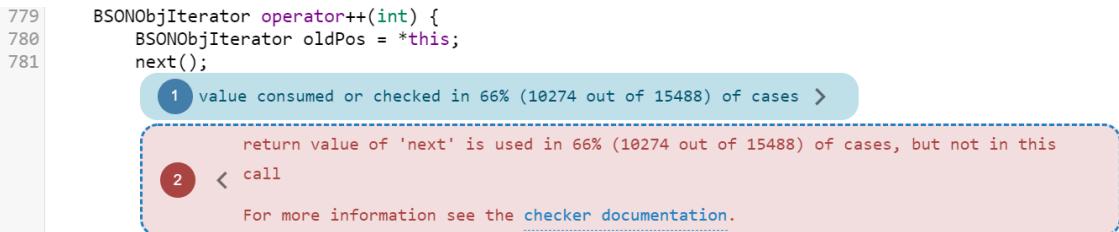


Figure 2.3: A report of the diagnostics of MongoDB’s source on CodeChecker with 50% threshold and project-level knowledge.

We can clearly conclude, that the multipass diagnosis resulted in different warnings. Figure 2.3 differs from fig. 2.1 in only the percentage, but with a different threshold that could lead to different results, to either false positives, or false negatives. Our own example, however shows both, since we had a call that used to give warning when diagnosed separately which with project-level knowledge does not, and one that did not give any warnings, but now it gives two, because it passed the threshold.

In listing 1.1 we talked about our code getting no warnings without project-level knowledge. After phases collect and compact, however, the following diagnostics will be emitted:

```
1 /home/bahramib/MyFolder/TestFolder/b/vec.cpp:7:15:  
2   warning: return value of 'erase' is used  
3   in most calls, but not in this one [misc-discarded-return-value]  
4  
5     c.erase(it);  
6     ^  
7  
8 /home/bahramib/MyFolder/TestFolder/b/vec.cpp:7:15: note:  
9   value consumed or checked in 66% (2 out of 3) of cases
```

We can clearly see, that our diagnosis gave us the real results this time.

Chapter 3

Developer documentation

This chapter will cover the process of developing a Clang-Tidy checker, and the infrastructure that it is built on. The developer should have an advanced knowledge of C++ and strong typing, minimal user experience with Linux and Clang-Tidy, and familiarity with declarative programming paradigm.

The work of the thesis and all the code I have written is a contribution and enhancement to the originally existing library of LLVM Compiler Infrastructure. I have attached the patch file that contains all the modifications to the original code with the thesis.

3.1 Infrastructure

The enhancement of the infrastructure focused on implementing the three distinct phases previously discussed in Chapter 2 for the checkers to distinguish and the developers to use. Figure 3.1 and fig. 3.2 demonstrates difference between the old infrastructure ("`single`") and the new one ("`multi`"). Note that if the aforementioned compact phase has not been run for a checker, the third phase (diagnose) will perform the "single" path, as it is the default mode of Clang-Tidy for both backwards compatibility, and intentional usage.

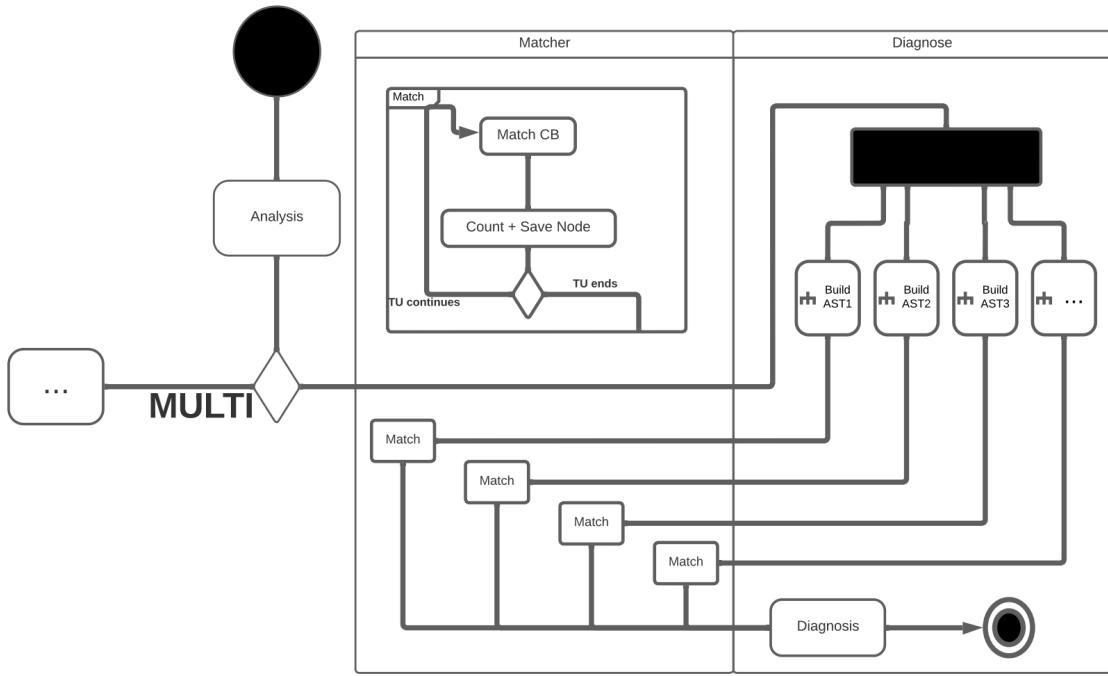


Figure 3.1: Activity diagram of Clang-Tidy infrastructure in “single” mode.

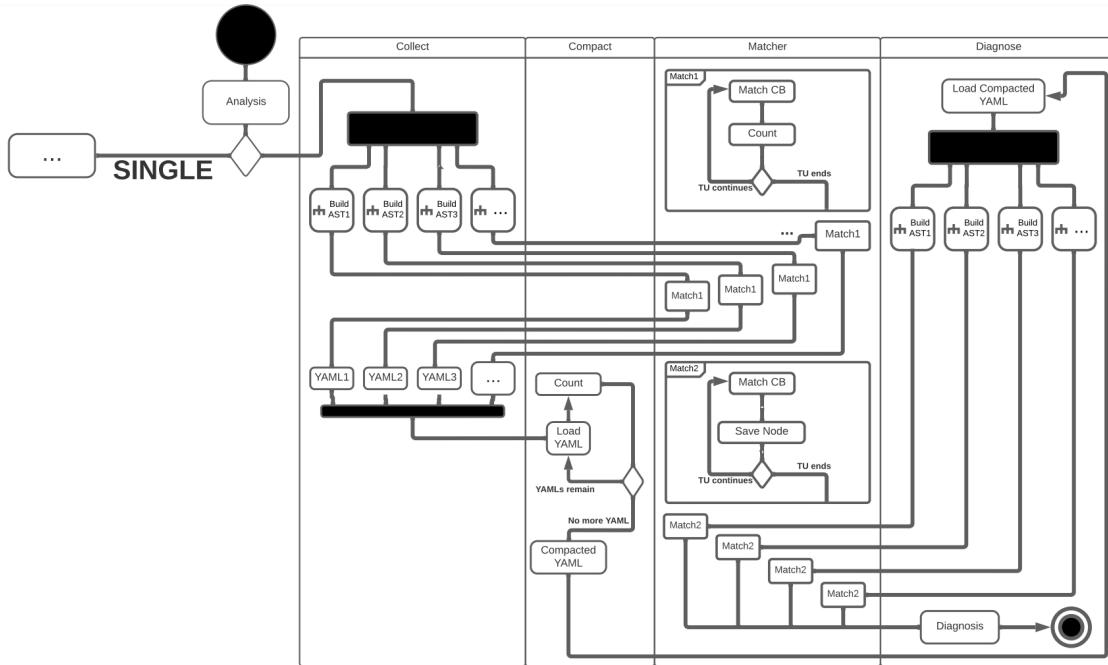


Figure 3.2: Activity diagram of Clang-Tidy infrastructure in “multi” mode.

The new flags I implemented are an `enum` flag for the current phase that we are in, and a regular `string` for the collection directory path. If the directory does not currently exist on the given path, the infrastructure will create one.

The base class of the checkers, `ClangTidyCheck` already has `virtual` functions for the checkers to override, with some being optional and some not. There are func-

tions for *language support*, registering *preprocessor callbacks*, registering *matchers*, checking matched *AST nodes*, and storing checker *options*. There are two additional functions from the base of `ClangTidyCheck`, `MatchCallback`, and these are `onStartOfTranslationUnit` and `onEndOfTranslationUnit`, to be able to do separate commands before and after the checker's actual work is done.

```
1 /// Override this to register ``PPCallbacks`` in the preprocessor.
2 virtual void registerPPCallbacks(const SourceManager &SM,
3     Preprocessor *PP,
4     Preprocessor *ModuleExpanderPP) {}
5
6 /// Override this to register AST matchers with Finder.
7 virtual void registerMatchers(ast_matchers::MatchFinder *Finder) {}
8
9 /// ``ClangTidyChecks`` that register ASTMatchers should do the
10 /// actual work in here.
11 virtual void check(
12     const ast_matchers::MatchFinder::MatchResult &Result) {}
13
14 /// Should store all options supported by this check with their
15 /// current values or default values for options that haven't been
16 /// overridden.
17 virtual void storeOptions(ClangTidyOptions::OptionMap &Options) {}
```

Code 3.1: Virtual functions from `ClangTidyCheck`'s header.

```
1 /// Called at the start of each translation unit.
2 /// Optionally override to do per translation unit tasks.
3 virtual void onStartOfTranslationUnit() {}
4
5 /// Called at the end of each translation unit.
6 /// Optionally override to do per translation unit tasks.
7 virtual void onEndOfTranslationUnit() {}
```

Code 3.2: Virtual functions from `MatchCallback`'s header.

Originally, when a *callback* happened (which we will address later in section 3.2) `ClangTidyChecker` simply called the derived `class`'s `check` function.

```

1 void ClangTidyCheck::run(const MatchFinder::MatchResult &Result) {
2     check(Result);
3 }
```

Code 3.3: The old infrastructure's way of calling `check`.

Now the checkers have three new `virtual void` member functions to override: `collect`, `postCollect` and `compact`. It is important to note that this override is optional, in a sense, that some checkers can not be improved with and do not need project-level knowledge. If one, however, decides to use the new multipass infrastructure with one's checker, then none of these overrides are optional.

`ClangTidyChecker` has a private `ClangTidyContext` member called `Context`. This contains information about the current phase, the directory path and the phase `enum` as well. This was used in the newly implemented member functions of the `ClangTidyChecker` parent `class`, and with protected getter functions, the derived checker can use some information from this context.

As we can see in listing 3.3, this function only called the checker's `check`, but now it distinguishes two out of our 3 separate phases, diagnose and collect. In diagnose phase, it still calls for `check`, but in collect phase it calls for the checker's separate `collect` function.

```

1 void ClangTidyCheck::run(const MatchFinder::MatchResult &Result) {
2     switch (getPhase()) {
3         case MultipassProjectPhase::Diagnose:
4             check(Result);
5             break;
6         case MultipassProjectPhase::Collect:
7             collect(Result);
8             break;
9         case MultipassProjectPhase::Compact:
10            llvm_unreachable("AST Matchers should not have run in "
11                           "compact mode.");
12    }
13 }
```

Code 3.4: Run function distinguishing Diagnose and Collect phase.

After the collection has been done on a translation unit, the destructor of the `class` that creates the checkers will call for the `runPostCollect` function of the base

`class`, which as the name suggests, runs the `postCollect` function of the derived checker.

During compact phase, the infrastructure does not need the translation units at all, so there are no checkers instantiated during the compilation, since this phase excludes compilation. Instead, the checker creation is now refactored into a function and that is called in a separate path of the code, where we only create the checkers and call `compact` on each of them once.

In the following subsections I will explain how the `virtual` functions should and how the member functions do work.

3.1.1 Collect Function

`collect` receives the same parameters as `check` does, a `MatchResult` reference. In this function the checker should do something similar to its original `check` function. Usually, it gets information out of the translation unit, and stores it in a data structure of choice that is fit for the checker's logic.

3.1.2 PostCollect Function

`postCollect` receives a `StringRef`, the name of the output YAML file. In this function a checker should take all the data collected from the translation unit, from the data structure, and write it in a YAML file whose filename is the function parameter.

3.1.3 Compact Function

`compact` has two parameters. One is a `StringRef` that, similarly to `postCollect`, contains the name of the output file. The other is a `vector<string>`s that contains all the YAML files in our multipass directory, that the current checker created in the collect phase and wrote the collection data into for each translation unit. In this phase, the checker should collect all the data from these YAML files and compact them into a single data structure. After this step, it should write it into the YAML file whose filename is the function parameter.

3.1.4 Additional Functions

There are newly implemented protected functions to use aside from these `virtual` ones:

- `getCollectPath`: returns a `string` to the path where the current check should write collected data to.
- `getCompactedDataPath`: returns a `StringRef`, the name of the file where the current check should write or read compacted data to/from.
- `getPhase`: returns the `enum` representing the current phase

3.2 Discarded Return Value Check

In section 1.2, we justified the existence of a checker which focuses on detecting potential bugs and code smells and giving warnings to function calls where the return value is ignored or unchecked, if we deem it necessary by analysing the ratio of the amount of unchecked calls and all calls for the same function.

3.2.1 Abstract Syntax Trees

Before making a checker, first we should know about how checkers work. Clang-Tidy checkers perform pattern matching on Clang's object oriented and strong typed representation of the abstract syntax tree of the source code. ASTs represent the syntactic structure of a code as seen in fig. 1.1.

In fig. 1.1 we have 2 function declarations, one of which is our main. We can see that our other function is called `foo`, that takes no parameters and returns an `int`. In the `body` we have a statement, more specifically a return statement. The subject of the return statement is a simple integer literal, with the value of 1. In `main`, we only have another return statement, that has a call expression of a function that has a declaration reference expression that refers to the declaration of our function `foo`.

Obviously, ASTs get increasingly more complex and harder to read as the source code grows. Another example can be seen on fig. 3.3.

```

TranslationUnitDecl
`-FunctionDecl <line:3:1, line:7:1> line:3:5 main 'int ()'
`-CompoundStmt <col:12, line:7:1>
`-ForStmt <line:4:5, line:6:15>
|-DeclStmt <line:4:10, col:19>
| `-VarDecl <col:10, col:18> col:14 used n 'int' cinit
| ` -IntegerLiteral <col:18> 'int' 0
|-<<<NULL>>>
|-BinaryOperator <col:21, col:26> 'bool' '<='
| |-ImplicitCastExpr <col:21> 'int' <LValueToRValue>
| | ` -DeclRefExpr <col:21> 'int' lvalue Var 0x5574beea9ad8 'n' 'int'
| ` -IntegerLiteral <col:26> 'int' 10
|-UnaryOperator <col:30, col:32> 'int' lvalue prefix '++'
| ` -DeclRefExpr <col:32> 'int' lvalue Var 0x5574beea9ad8 'n' 'int'
`-IfStmt <line:5:9, line:6:15>
|-BinaryOperator <line:5:13, col:22> 'bool' '==' 
| |-BinaryOperator <col:13, col:17> 'int' '%'
| | |-ImplicitCastExpr <col:13> 'int' <LValueToRValue>
| | | ` -DeclRefExpr <col:13> 'int' lvalue Var 0x5574beea9ad8 'n' 'int'
| | ` -IntegerLiteral <col:17> 'int' 2
| ` -IntegerLiteral <col:22> 'int' 0
`-UnaryOperator <line:6:13, col:15> 'int' lvalue prefix '++'
| ` -DeclRefExpr <col:15> 'int' lvalue Var 0x5574beea9ad8 'n' 'int'

```

Figure 3.3: The AST of listing 3.5

```

1 int main() {
2     for (int n = 0; n <= 10; ++n)
3         if (n % 2 == 0)
4             ++n;
5 }

```

Code 3.5: The code of fig. 3.3.

3.2.2 Pattern Matching

Clang-Tidy checkers have *matchers* that you can set up in the checker's overridden function `registerMatchers`. This has its matcher, a `MatchFinder` as the parameter, and in this function we have to use its `addMatcher` member function to add new patterns for it to match, and also to bind a string literal to the matches for later usage. If the matcher finds a match, a `callback` action is executed.

In order to utilize the matches and the callbacks, we have to know how to match nodes on an AST. Clang-Tidy has built in matchers, that we can see in the AST Matcher Reference [16]. This tells us which matcher accepts which matcher or matchers as parameter, and what their return type is. It also gives us examples of usage.

Theoretically, there are 3 distinct base types for AST nodes. `Decl` (declaration), `Stmt` (statement) and `Type`. In practice, Clang-Tidy has more distinct bases in its AST representation, but for now let us view a simple example of matching from *Compiler Explorer*:¹

```

1 int foo() { return 1; }

2

3 int main() {
4     foo();
5     int n = foo();
6     return foo();
7 }
```

Code 3.6: A very simple code for matching function calls.

```

TranslationUnitDecl
|-FunctionDecl <line:1:1, col:23> col:5 used foo 'int ()'
| `-CompoundStmt <col:11, col:23>
|   `-ReturnStmt <col:13, col:20>
|     `-IntegerLiteral <col:20> 'int' 1
`-FunctionDecl <line:3:1, line:7:1> line:3:5 main 'int ()'
  `-CompoundStmt <col:12, line:7:1>
    |-CallExpr <line:4:5, col:9> 'int'
    | `-ImplicitCastExpr <col:5> 'int (*)()' <FunctionToPointerDecay>
    |   `-DeclRefExpr <col:5> 'int ()' lvalue Function 0x564ea6f9a820 'foo' 'int ()'
    |-DeclStmt <line:5:5, col:18>
      |-VarDecl <col:5, col:17> col:9 n 'int' cinit
        |-CallExpr <col:13, col:17> 'int'
          |-ImplicitCastExpr <col:13> 'int (*)()' <FunctionToPointerDecay>
            `-DeclRefExpr <col:13> 'int ()' lvalue Function 0x564ea6f9a820 'foo' 'int ()'
      `-ReturnStmt <line:6:5, col:16>
        |-CallExpr <col:12, col:16> 'int'
          |-ImplicitCastExpr <col:12> 'int (*)()' <FunctionToPointerDecay>
            `-DeclRefExpr <col:12> 'int ()' lvalue Function 0x564ea6f9a820 'foo' 'int ()'
```

Figure 3.4: The AST of listing 3.6.

Here we can have 1 function declaration of the same function, and 3 call expressions, the first one in as a regular statement, the second one in a variable declaration, and the third one in a return statement. Let us build a matcher for calls that a checker like mine should be concerned about. First, we want to match for function calls, whose declaration reference refers to a function declaration. If we use `match callExpr(callee(declRefExpr(to(functionDecl().bind("definition")))).bind("call"))` we

¹<http://godbolt.org/z/ncToe315q>, accessed 2022. 05. 14.

get 3 matches and 6 binds. It matches all 3 calls of `foo` bound to the literal "call", and we also get 3 binds, all for the first line, the `FunctionDecl` of `foo`, bound to "definition".

We can also save the matcher bound to "definition" as a variable for easier usage and simplicity, if we wanted to. With `let Def functionDecl().bind("definition")` it will be saved as `Def`, and with `let Call callExpr(callee(declRefExpr(to(Def)))).bind("call")`, where we reuse our `Def`, we save our `Call` for later use. Now we only have to match for `Call` with `match Call` to get the same results.

We could improve it this matcher for the checker, with not matching functions declared with `void` as return value. In Clang-Tidy the "saved" matcher will be a `static const auto Call = callExpr(hasDeclaration(functionDecl(unless(returns(voidType())))))`. In reality, the checker will need more specification in order to work properly, but this distinction is important as well.

There are separate matchers in static variables for any kind of usage of a return value, and then these are added in a single `addMatcher` call of the matcher. Let us go through some of these matchers.

Call

Our actual specific call expression can be seen in listing 3.7. Since this is the base of all other matchers that we use, I will explain what we need to match for, in detail, starting with the bind. The `static constexpr` that `Call` is bound to is an `llvm::StringLiteral` type, that wraps "call". We want to match for function calls, but not all function calls. First, we will sort out all functions whose declaration is a function declaration to at least either a `void` function, because they do not return a value to be checked or ignored, and a function with `[[nodiscard]]`, since the compiler will already warn the user if needed. This is what we achieve with lines 2-3. Second, we want to sort out any operator calls except for the following: `operator()` or `operator[]` (this is done by lines 4-5), and `operator *` or `operator &`, but since those operators have multiple uses and we want to include a specific use only, we only need these if they take only one argument, also known as the *dereference* and the *address-of* operator (which is achieved by lines 6-7). All of this specification can be seen in fig. 3.5.

```

1 static const auto Call =
2 callExpr(hasDeclaration(functionDecl(unless(anyOf(
3         returns(voidType()), hasAttr(attr::WarnUnusedResult)))),
4         unless(cxxOperatorCallExpr(
5             unless(anyOf(hasAnyOverloadedOperatorName("()", "[]"),
6                 allOf(hasAnyOverloadedOperatorName("*", "&"),
7                     argumentCountIs(1)))))))
8 .bind(::Call);

```

Code 3.7: The matcher for our desired call expression.

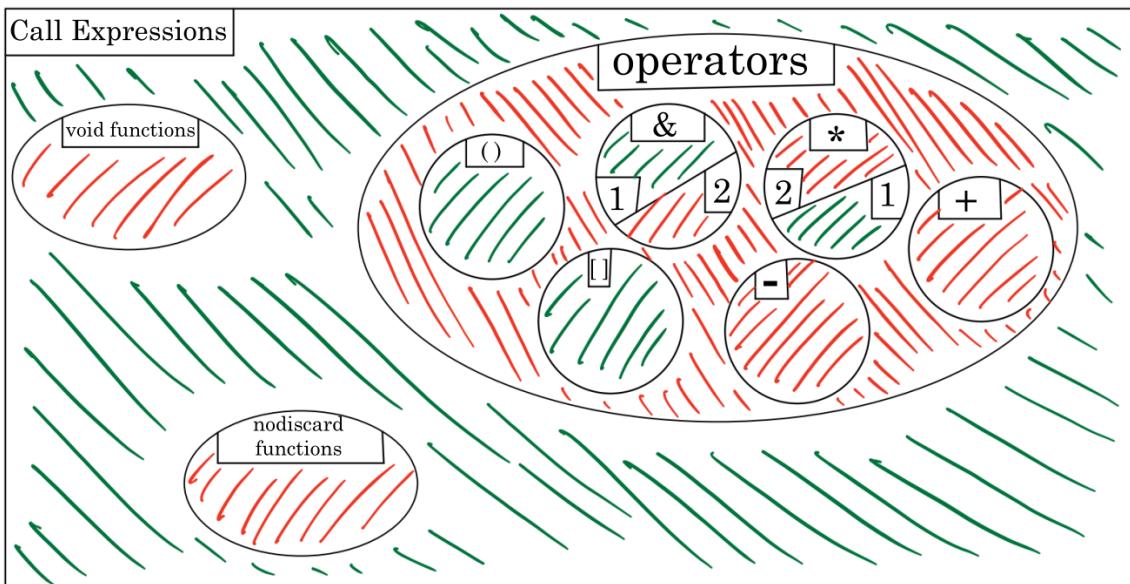


Figure 3.5: A Venn-diagram of what `Call` (see listing 3.7) matches, where 1 and 2 indicates the arity. The green subsets are to be matched.

While

In listing 3.8 we match for either a `while` statement or a `do` statement that uses our call in their condition, and then since the matcher `anyOf` does not automatically cast its return type to that of a statement, we do it manually with `stmt`.

```

1 static const auto While =
2     stmt(anyOf(whileStmt(hasCondition(Call)),
3                 doStmt(hasCondition(Call))));

```

Code 3.8: The matcher for usage in while expression.

Return

In listing 3.9 we simply match if a return statement uses our function as its return value.

```
1 static const auto Return = returnStmt(hasReturnValue(Call));
```

Code 3.9: The matcher for the usage in return statement.

For

In listing 3.10 we match for any and all of the following: if the `initialization`, the `condition`, or the `incrementation` of the for loop, or if a `range-based` for loop, or any descendants of either of these uses the return value of our call.

```
1 static const auto For = anyOf(
2     forStmt(eachOf(hasLoopInit(findAll(Call)),
3                         hasCondition(findAll(Call)),
4                         hasIncrement(findAll(Call)))),
5     cxxForRangeStmt(hasRangeInit(findAll(Call))));
```

Code 3.10: The matcher for the usage in for statement.

We can also create new matchers, should we need one that is not yet in the library. It provides *macros* that can be used to implement matchers with custom logic. For example in listing 3.11 the macro defines a single-parameter function called `hasAssertExpr` that takes a `Matcher<Expr>` called `InnerMatcher` and returns a `Matcher<StaticAssertDecl>` object. It passes the matcher into the assertion expression of assert, meaning if the inner matcher matches, the entire assertion will too.

```
1 AST_MATCHER_P(StaticAssertDecl, hasAssertExpr,
2                 ast_matchers::internal::Matcher<Expr>, InnerMatcher) {
3     return InnerMatcher.matches(*Node.getAssertExpr(), Finder,
4                                 Builder);
5 }
6
7 static const auto StaticAssert =
8     staticAssertDecl(hasAssertExpr(Call));
```

Code 3.11: Custom logic and usage of a new matcher called `hasAssertExpr`.

3.2.3 Checker Logic

Before going into how my checker works, we should go through the member functions of it. The checker obviously overrides all non-optional and some optional `virtual` functions as seen in listing 3.12. In listing 3.13 there are three additional member functions. These contain the logic of the checker and how it shapes its data structure according to the current phase.

```
1 void registerMatchers(ast_matchers::MatchFinder *Finder) override;
2 void storeOptions(ClangTidyOptions::OptionMap &Opts) override;
3 void onStartOfTranslationUnit() override;
4 void onEndOfTranslationUnit() override;
5 void collect(const ast_matchers::MatchFinder::MatchResult &Result)
6     override;
7 void postCollect(StringRef OutputFile) override;
8 void compact(const std::vector<std::string> &PerTuCollectedData,
9             StringRef OutputFile) override;
10 void check(const ast_matchers::MatchFinder::MatchResult &Result)
11    override;
```

Code 3.12: Overridden virtual functions in `DiscardedReturnValueCheck`'s header.

```
1 void matchResult(const MatchFinder::MatchResult &Result,
2                  bool ShouldCount);
3 void registerCall(const CallExpr *CE, const FunctionDecl *FD,
4                   bool IncrementCounters,
5                   const void *ConsumingContext);
6 void diagnose(const Function &F);
```

Code 3.13: Private member functions in `DiscardedReturnValueCheck`'s header.

For the data structure itself, The checker uses a `struct` representing the functions with their tracked amount of calls and checks (and its member function calculating the ratio of these in percentage), a `StringMap` of functions to store their name and their corresponding struct object, a `DenseMap` of `FunctionDecl - String` pairs, a `DenseMap` of `CallExpr - SmallVector` pairs to prevent multiple callbacks on a single call expression being processed twice. It also uses a `uint8_t` to save the threshold option, and an `Optional<bool>` to track whether or not compacting has been done on the translation units.

```
1 struct Function {
2     std::size_t ConsumedCalls;
3     std::size_t TotalCalls;
4     const FunctionDecl *FD;
5     llvm::SmallPtrSet<const CallExpr *, 32> DiscardedCEs;
6
7     std::uint8_t ratio() const;
8 };
9 using FunctionMapTy = llvm::StringMap<Function>;
10
11 const std::uint8_t ConsumeThreshold;
12 Optional<bool> CacheProjectDataLoadedSuccessfully;
13
14 llvm::DenseMap<const CallExpr *, SmallVector<const void *, 2>>
15     ConsumedCalls;
16 llvm::DenseMap<const FunctionDecl *, std::string> FunctionIDs;
17 FunctionMapTy CallMap;
```

Code 3.14: Members of the data structure.

After all the possible usages of a return value have been refactored into these matchers, we put all of them into one call of `addMatcher` and bind them to "`consume`". In reality, the usages contain all 3 of the base node types, so we put these matchers into 3 `addMatcher` calls in respect of which matcher returns which node type. The matchers `ExplicitCast`, `New`, `Delete`, `Argument`, `Unary`, `Dereference`, `BinaryLHS`, `BinaryRHS`, `If`, `While`, `For`, `Switch`, and `Return` should return statements; `StaticAssert`, `VarDecl`, and `CtorInits` should return declarations; and `Decltype`, `TemplateArg`, and `VLA` should return types. Now all we have to do is match for a simple call.

```
1  Finder->addMatcher(
2      traverse(
3          TK_IgnoreUnlessSpelledInSource,
4          stmt(eachOf(ExplicitCast, New, Delete, Argument, Unary,
5              Dereference, BinaryLHS, BinaryRHS, If, While, For,
6              Switch, Return)))
7          .bind(Consume),
8          this);
9  Finder->addMatcher(traverse(TK_IgnoreUnlessSpelledInSource,
10      decl(eachOf(StaticAssert, VarDecl, CtorInits)))
11      .bind(Consume),
12      this);
13 Finder->addMatcher(traverse(TK_IgnoreUnlessSpelledInSource,
14      type(eachOf(Decltype, TemplateArg, VLA)))
15      .bind(Consume),
16      this);
17
18 Finder->addMatcher(traverse(TK_IgnoreUnlessSpelledInSource, Call),
19      this);
```

Code 3.15: Separating the matchers to statement, declaration, and type.

The checker will match for all function calls (that are not declared `void` et cetera) bound to "`call`", and for the ones that are checked it will match once more with the bind "`consume`". For every match the callback happens and since matching "`consume`" deterministically comes before matching "`call`", we can be sure that if we matched with "`consume`" it is a checked call, and if we did not, it is an unchecked call so we can count the checks of our matched call expressions accordingly. In listing 3.16 we try to get our "`consume`" bound node for either of the 3 base node types, and if either gets a node back, we count it as a checked return. Otherwise we count it as unchecked.

```

1 void DiscardedReturnValueCheck::matchResult(
2     const MatchFinder::MatchResult &Result, bool ShouldCount) {
3     const auto *CE = Result.Nodes.getNodeAs<CallExpr>(Call);
4     assert(CE && "Bad matcher");
5
6     const void *ConsumeNode = nullptr;
7     if (const auto *D = Result.Nodes.getNodeAs<Decl>(Consume))
8         ConsumeNode = D;
9     if (const auto *S = Result.Nodes.getNodeAs<Stmt>(Consume))
10        ConsumeNode = S;
11     if (const auto *T = Result.Nodes.getNodeAs<Type>(Consume))
12        ConsumeNode = T;
13     if (ConsumeNode)
14         return registerCall(CE, CE->getDirectCallee(), ShouldCount,
15                             ConsumeNode);
16     if (ConsumedCalls.find(CE) == ConsumedCalls.end())
17         return registerCall(CE, CE->getDirectCallee(), ShouldCount,
18                             nullptr);
19 }
```

Code 3.16: Catching nodes of used return values.

Then the call expression is emplaced in our data structure for counting and because we do not want to count the same call twice, and if the call belongs to a function declaration that have previously been saved in our data structure, then we also increment its data in respect of it being checked or not.

At the end of the translation unit, we give our diagnostics.

3.2.4 Multipass phases

It is very clear that a checker like this is almost completely useless for projects larger than a single translation unit, unless we use the enhanced infrastructure. This makes the checker a perfect fit to test it. First, we have to implement the new `virtual` functions.

`collect` does the same thing as the old `check` callback did. Because of this, the logic of this step was put into a member function and is called in both functions, `check` and `collect`.

We save declarations and adjust the data structure with each new call of the same function. At the end of a translation unit, `postCollect` writes this data into a new YAML file under the unique name that was generated by the infrastructure.

`compact` reads in all the data from the collection files, fills up the data structure and then it writes the compacted data into the output YAML. Because both writing and reading into our data structure is necessary more than once in the checker, this also have been refactored, and `compact`, `postCollect` and `check` all call these functions.

In listing 3.17, we can see that we first extract the data from the YAML file into an optional `vector` of the function representation, and if the extraction actually happened, we emplace and modify the data of the files into the current state of the data structure. At the end, it returns whether succesful extraction happened or not. In listing 3.18 we see that `writeYAML` puts our data into the output `stream`. These are both examples of possible writing and loading methods for YAML serialization that are easy to use and implement using LLVM's YAML library [17].

```
1 using FunctionVec = std::vector<SerializedFunction>;
2
3 static Optional<FunctionVec> loadYAML(StringRef File) {
4     using namespace llvm;
5
6     ErrorOr<std::unique_ptr<MemoryBuffer>> IStream =
7         MemoryBuffer::getFileAsStream(File);
8     if (!IStream)
9         return None;
10
11    FunctionVec R;
12    yaml::Input YIn{**IStream};
13    YIn >> R;
14
15    return R;
16 }
17
18 static bool loadYAML(StringRef FromFile,
19                       DiscardedReturnValueCheck::FunctionMapTy &ToMap) {
20     Optional<FunctionVec> OV = loadYAML(FromFile);
21     if (!OV)
22         return false;
23
24     for (const SerializedFunction &SF : *OV) {
25         DiscardedReturnValueCheck::Function &F =
26             ToMap
27             .try_emplace(SF.ID,
28                         DiscardedReturnValueCheck::Function{0, 0, nullptr, {}})
29             .first->second;
30
31         F.ConsumedCalls += SF.ConsumedCalls;
32         F.TotalCalls += SF.TotalCalls;
33     }
34
35     return true;
36 }
```

Code 3.17: Functions for loading.

```
1 static void writeYAML(StringRef Whence, FunctionVec Elements,
2                       StringRef ToFile) {
3     std::error_code EC;
4     llvm::raw_fd_ostream FS(ToFile, EC, llvm::sys::fs::OF_Text);
5     if (EC) {
6         llvm::errs() << "DiscardedReturnValueCheck: Failed to write "
7                     << Whence << " output file: " << EC.message();
8         llvm::report_fatal_error("", false);
9         return;
10    }
11
12    llvm::yaml::Output YAMLOut(FS);
13    YAMLOut << Elements;
14 }
```

Code 3.18: Function for writing.

Both `postCollect` and `compact` uses `writeYAML` with the YAML compatible representation of the function representation. For the purpose of writing, these functions both use a refactored function that converts the latter into the YAML compatible version. These are seen in listing 3.19 with the method of making a `struct` YAML compatible as well.

```
1 struct SerializedFunction {
2     std::string ID;
3     std::size_t ConsumedCalls, TotalCalls;
4 };
5
6 template <> struct MappingTraits<SerializedFunction> {
7     static void mapping(IO &IO, SerializedFunction &F) {
8         IO.mapRequired("ID", F.ID);
9         IO.mapRequired("Consumed", F.ConsumedCalls);
10        IO.mapRequired("Total", F.TotalCalls);
11    }
12 };
13
14 LLVM_YAML_IS_SEQUENCE_VECTOR(SerializedFunction)
15
16 static FunctionVec
17 yamlize(const DiscardedReturnValueCheck::FunctionMapTy &Map) {
18     FunctionVec SFs;
19     llvm::transform(Map, std::back_inserter(SFs), [] (auto &&E) {
20         return SerializedFunction{E.first().str(),
21                             E.second.ConsumedCalls,
22                             E.second.TotalCalls};
23     });
24     return SFs;
25 }
```

Code 3.19: Steps for making the data YAML writable.

Check works in a similar way. First, we check if compact has happened already. If not, check runs without project-level knowledge and emits diagnostics per translation unit, the same way as before. If compact phase has not been run for this checker, we fill up our data structure, but only if this step has not happened before. For this, we use the same refactored functions as in `collect`, `matchResult` and thus `registerCalls`, with the exception that adjustments in the data structure after each call does not need to happen, since all the data on the amount of checks and calls had been acquired before, from the compacted file. Then we take the unique name (as a `string`) of a function declaration, or generate one if one does not already exist. After this we check again, whether or not we stored this declaration in the data structure, and if not, we store it, but only if its call was unused, and then if compact has not happened, we increment its data accordingly. If it has, we have nothing left to do. Now we only need to emit diagnostics to the right calls, which

happens at the end of each translation unit.

3.3 Evaluation

This checker's upgrade to the multipass infrastructure differs from the original one only by a little. However, it gives us very different results. The checker has run on large projects. It ran with 50% and 80% threshold both with and without collection and compacting. The results are shown in table 3.1

Project name	Multi-pass?	Threshold	# reports
Bitcoin v0.20.1 [18]	✗	50%	140
	✓	50%	320
	✗	80%	25
	✓	80%	92
CodeChecker v6.17.0 [15]	✗	50%	2
	✓	50%	5
	✗	80%	0
	✓	80%	0
Contour v0.2.0.173 [19]	✗	50%	38
	✓	50%	107
	✗	80%	10
	✓	80%	8
cURL 7.66.0 [20]	✗	50%	55
	✓	50%	210
	✗	80%	6
	✓	80%	104
FFmpeg n4.3.1 [21]	✗	50%	904
	✓	50%	1 846
	✗	80%	148
	✓	80%	613
libWebM 1.0.0.27 [22]	✗	50%	9
	✓	50%	10
	✗	80%	0
	✓	80%	1
LLVM 12.0.0 [23]	✗	50%	4 615
	✓	50%	6 837
	✗	80%	1 077
	✓	80%	2 150
Memcached 1.6.8 [24]	✗	50%	32
	✓	50%	49
	✗	80%	5
	✓	80%	16

Project name	Multi-pass?	Threshold	# reports
MongoDB r4.4.6 [25]	✗	50%	1 671
	✓	50%	2 463
	✗	80%	370
	✓	80%	1 223
OpenSSL 3.0.0 [26]	✗	50%	427
	✓	50%	1 060
	✗	80%	56
	✓	80%	244
PostgreSQL r13.0 [27]	✗	50%	799
	✓	50%	644
	✗	80%	62
	✓	80%	215
Protocol Buffers v3.13.0 [28]	✗	50%	51
	✓	50%	122
	✗	80%	20
	✓	80%	75
Qt Base v6.2.0 [29]	✗	50%	1 195
	✓	50%	2 764
	✗	80%	200
	✓	80%	1 093
SQLite 3.33.0 [30]	✗	50%	284
	✓	50%	288
	✗	80%	42
	✓	80%	48
tmux 2.6 [31]	✗	50%	35
	✓	50%	48
	✗	80%	1
	✓	80%	1
twin v0.8.1 [32]	✗	50%	54
	✓	50%	69
	✗	80%	9
	✓	80%	22

Project name	Multi-pass?	Threshold	# reports
Vim v8.2.1920 [33]	✗	50%	227
	✓	50%	383
	✗	80%	32
	✓	80%	54
Xerces v3.2.3 [34]	✗	50%	125
	✓	50%	197
	✗	80%	9
	✓	80%	38

Table 3.1: Results of different runs on live projects.

In table 3.1, we can see that the numbers between both single and multiple phase, and 50% and 80% threshold runs largely differ in a considerable amount of test projects. There were results where the multiple phase run gave more warnings, indicating the aforementioned bugs hiding in separate translation units only, and some results where it gave less warnings, indicating false positives in single runs.

In fig. 3.6 and fig. 3.7, we can see diagrams of the results. Each line represents a function, with its height as the total amount of calls. The green and red parts represent the checked and unchecked calls, respectively. The descending graph of black dots display the ratio of checked to all calls, which is important regarding our threshold.

In these runs, the amount of separate unchecked functions were similar, but using the multi-pass infrastructure, the amount of calls for these functions was obviously a lot higher. In the old infrastructure run, the ratios are closer to the 80% threshold, and most functions here were called less than 20 times, and were only unchecked once or twice. However, if we look at the multi-pass run, the ratios are a lot higher. They were unchecked only a couple of times in this run as well, but the amount of total calls per function was bigger than in the previous run, resulting in the functions staying more consistently above our 80% threshold.

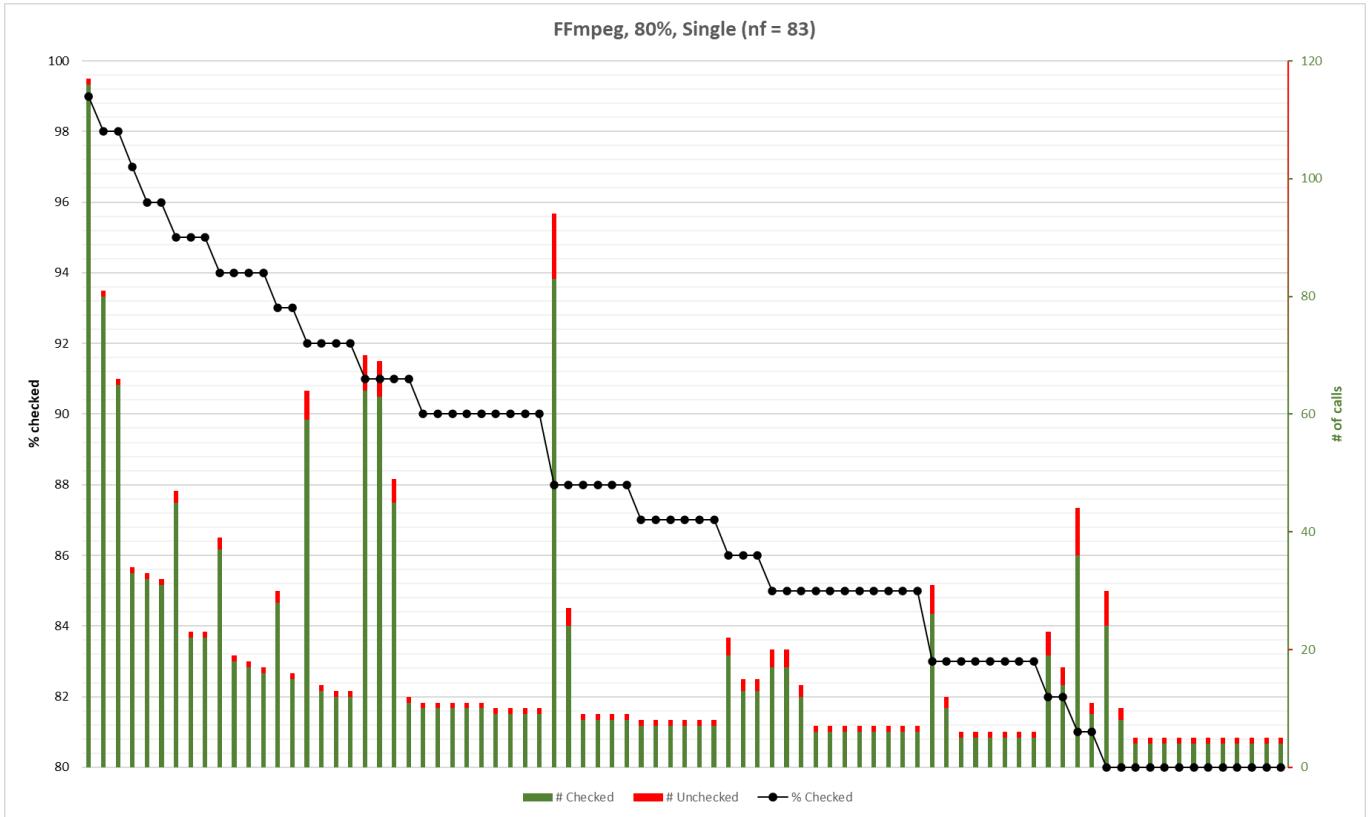


Figure 3.6: Diagram of 80% threshold single run on FFmpeg.

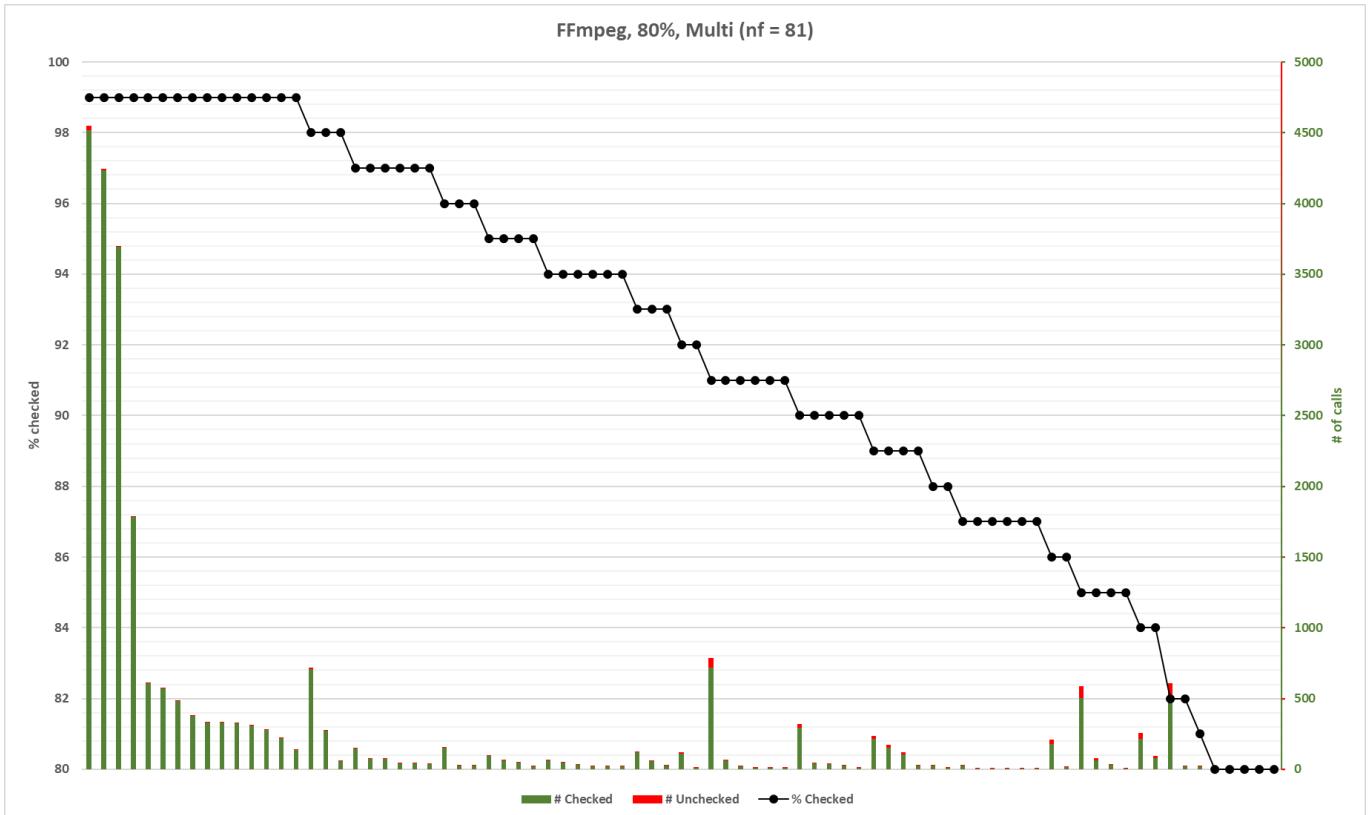


Figure 3.7: Diagram of 80% threshold multi run on FFmpeg.

Let us take a look at a 50% threshold run comparison as well. In fig. 3.8 and fig. 3.9 we get very different results again. Here we can still see the higher ratios in the multi run, as well as higher overall amount of function calls. This time we have more unchecked calls for functions that have been called considerably more than the rest, in both amount and ratio.

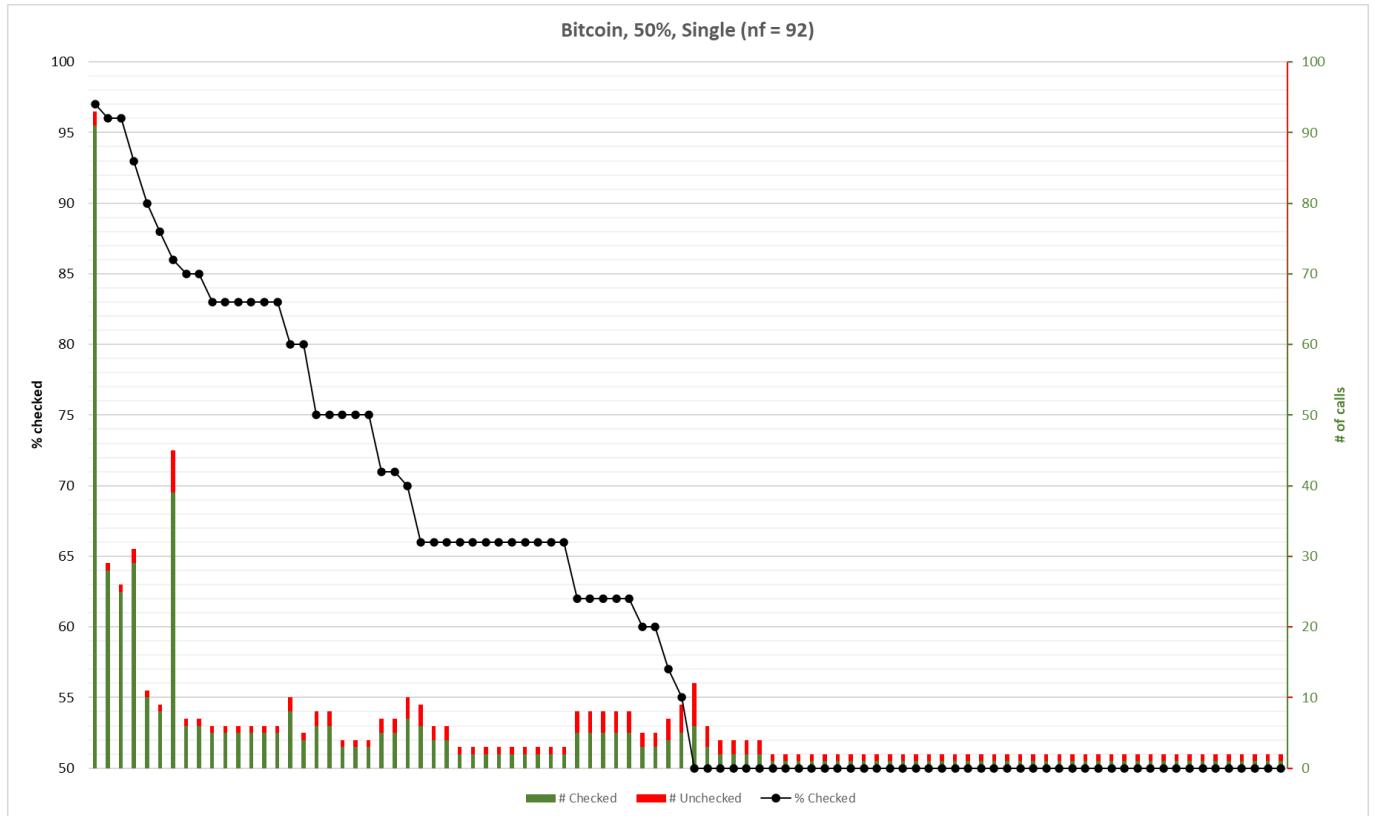


Figure 3.8: Diagram of 50% threshold single run on Bitcoin.

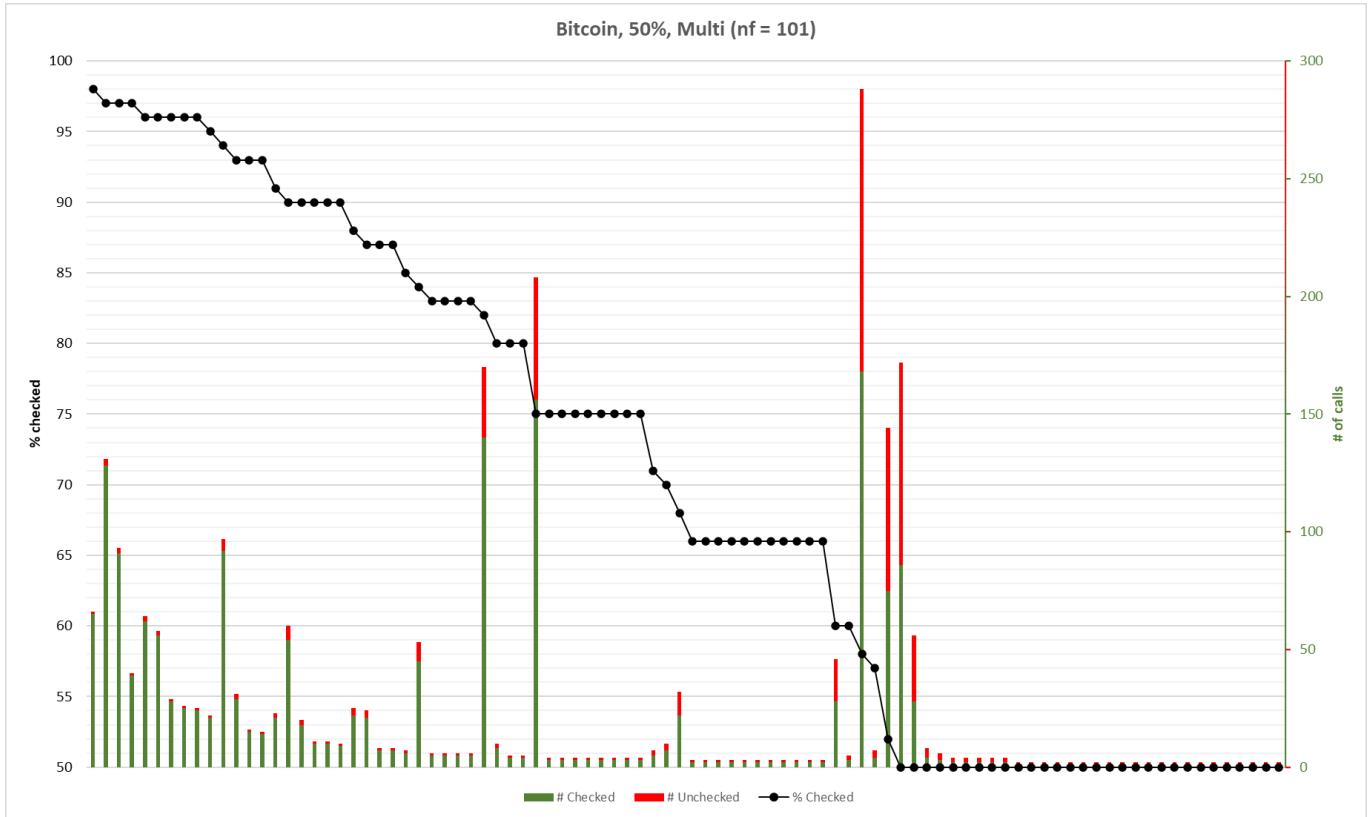


Figure 3.9: Diagram of 50% threshold multi run on Bitcoin.

The diagrams for all of the aforementioned project tests can be found in appendix A. What can be said about almost all of the larger results is that they have a great amount of functions barely passing the threshold whose total amount of calls do not exceed a considerable number compared to the size of the project and some other functions within the results. These could potentially be ignored in some cases in the future.

Chapter 4

Conclusion

This thesis demonstrated the problem of unchecked return values, to which problem it gave a solution in form of a statistical method without actually modifying the source code of said unchecked return values. This solution was made as an implementation in the Clang-Tidy pattern matching based static analysis library, however the statistical method would fail because of the way Clang-Tidy checkers work, separately for each translation unit. This leads to invalid statistics. In order to fix this, I have enhanced the infrastructure and created the option of multiphase analysis, that will collect the needed information per translation unit, accumulate these and make a project-level set of information, and finally emit diagnosis with this newly compacted information. The finished checker was accomplished using this new infrastructure. This checker method can obviously run without the use of multiphase analysis, because the new infrastructure is made to be backwards compatible with the single phase mode, but the new infrastructure was vital for the statistics of the checker. This proved the checker to be an excellent example for the usage and utility of the multiphase mode.

4.1 Future Work

Previously in section 3.3, I talked about ignoring a set of functions whose total amount of calls do not exceed a given value. This could potentially be implemented into the checker as an option, so that the user could give the checker the number, and the checker will not give warnings to a function call if the function's amount of calls is less than said number. It would make the results clearer in some cases.

The first possible future work ignores functions based on their statistics, but ignoring based on function property could be useful too. The checker should categorically ignore some functions. There are also functions whose return value can not always be checked. A classic example would be the `<<` operator or `std::append`, since the first will always return the `stream`, and the latter will always return `*this` for continuous use of the same function. `std::cout << "1" << "2" << "3" << std::endl;` will always be unchecked in the end, since the last `<<` can not be checked with merit and without creating another unchecked return value.

The hash generation for the YAML filename does not respect multiple compilations of the same file with different compile flags. This is something that could be fixed in the future as well.

Acknowledgements

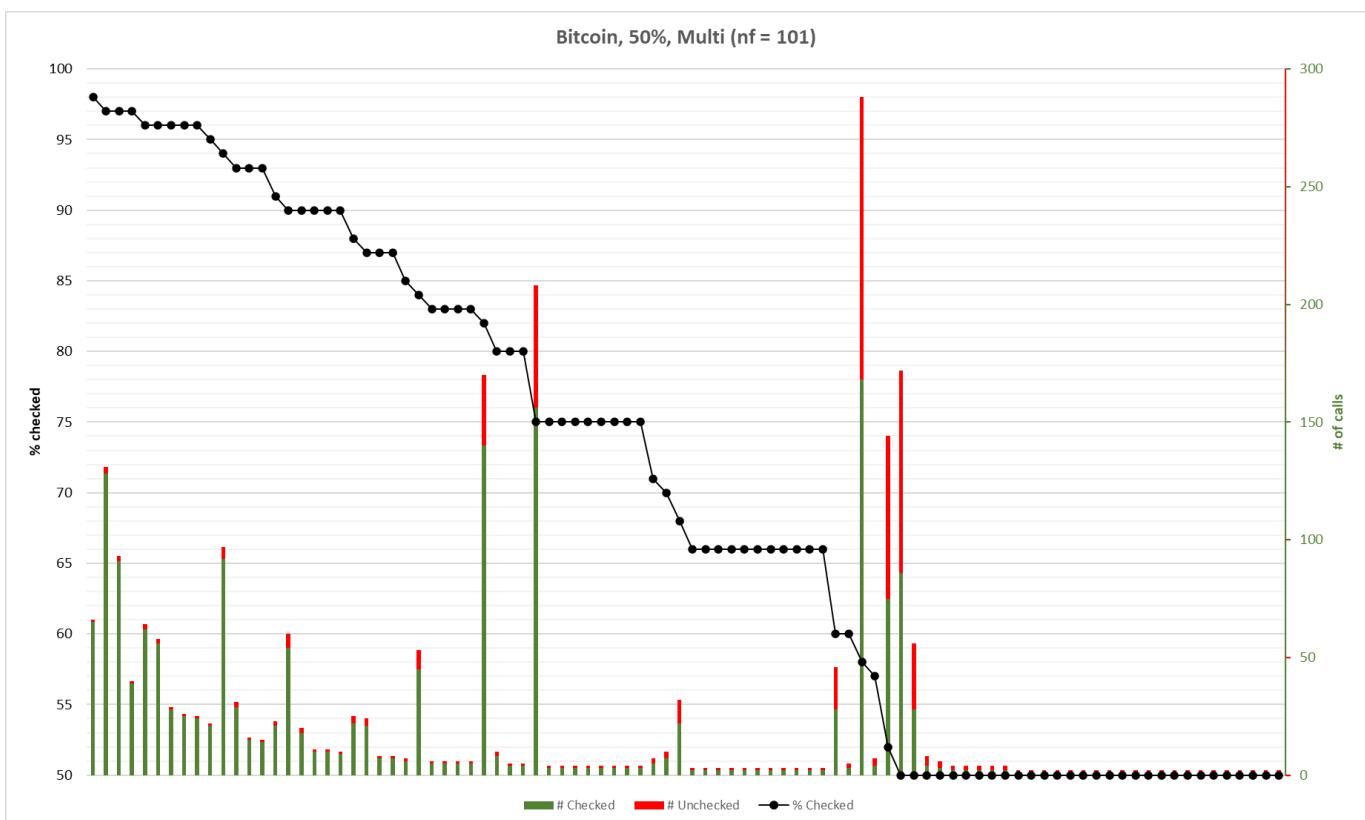
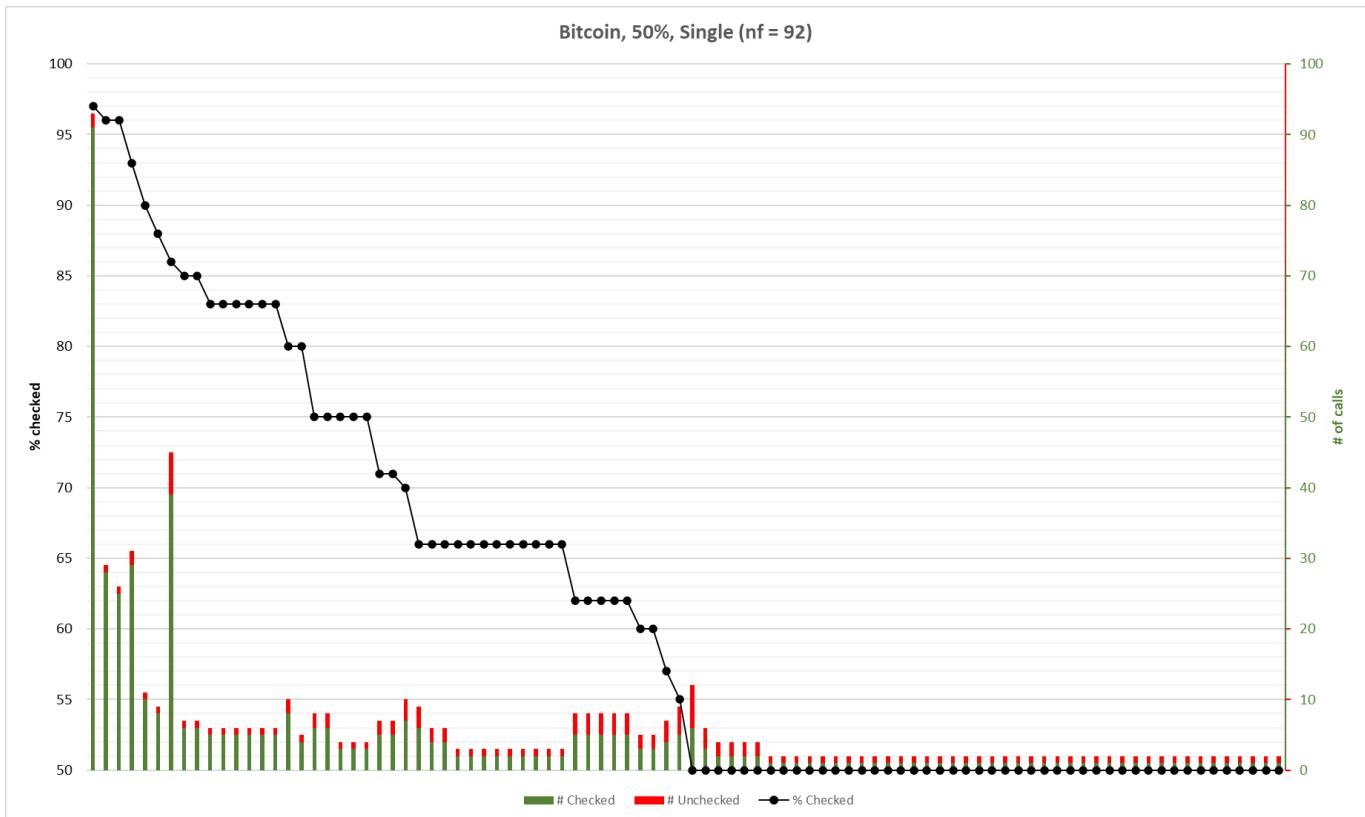
I would like to express my deepest appreciation to my supervisor, Richárd Szalay for his invaluable patience and feedback, and without whom I could not have completed this thesis in time. I am also grateful to my professor, Dr. Zoltán Porkoláb, who generously provided his knowledge, expertise and inspiration. Furthermore I would highlight Aaron Ballman's (Intel Corporation) and Gábor Márton's (Ericsson AB.) insights during the open-source code review process.

Appendix A

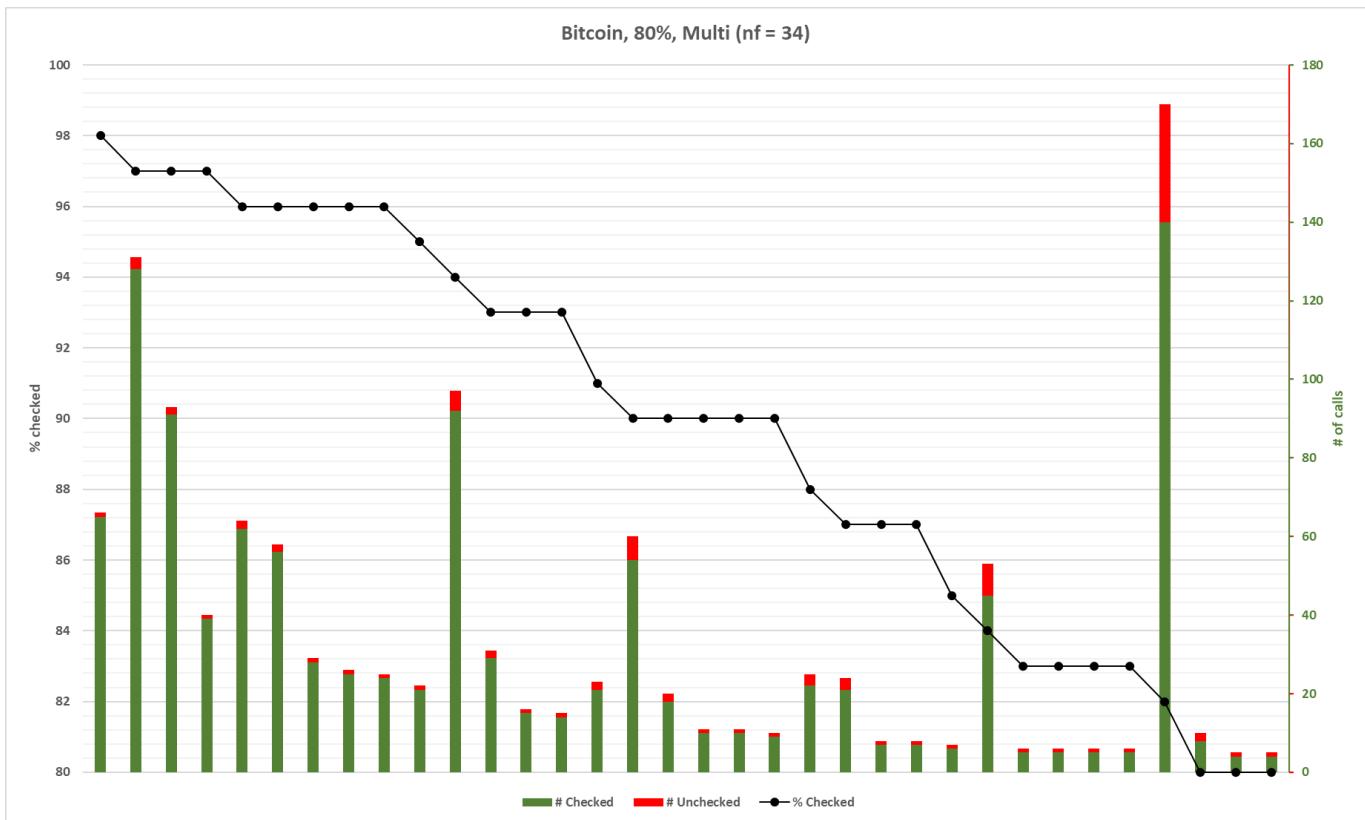
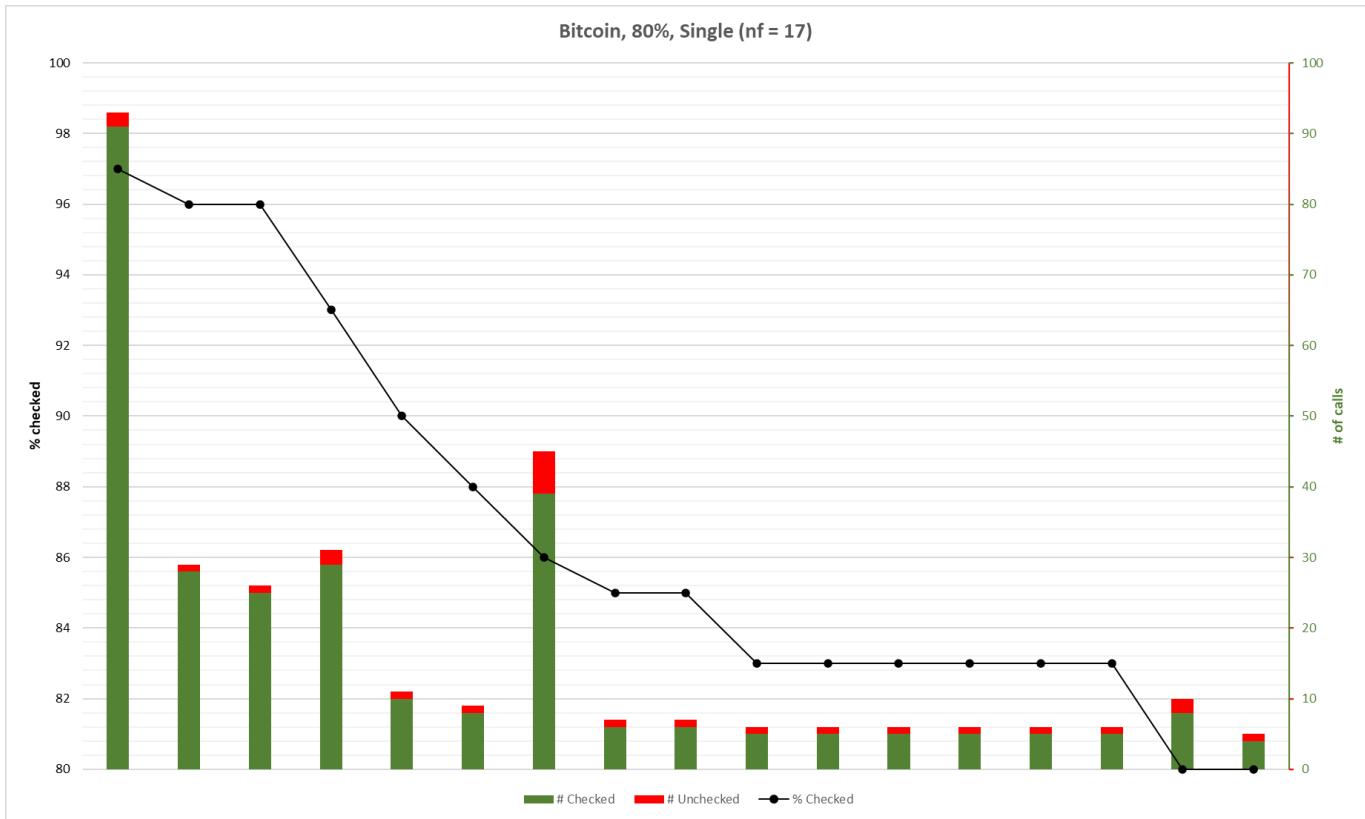
Project results

The Discarded Return Value Checker was run on 18 different open-source, actively developed projects with 4 different setting. Two using the new multiple phase infrastructure and two with the single phase one, at 50% threshold and at 80% threshold. These are the results of the runs regarding the amount of different functions found, their amount of calls both checked and unchecked, and the ratio of these two amounts, similarly to the depiction in section 3.3. In each of the graphs, nf refers to the number of functions in the reports.

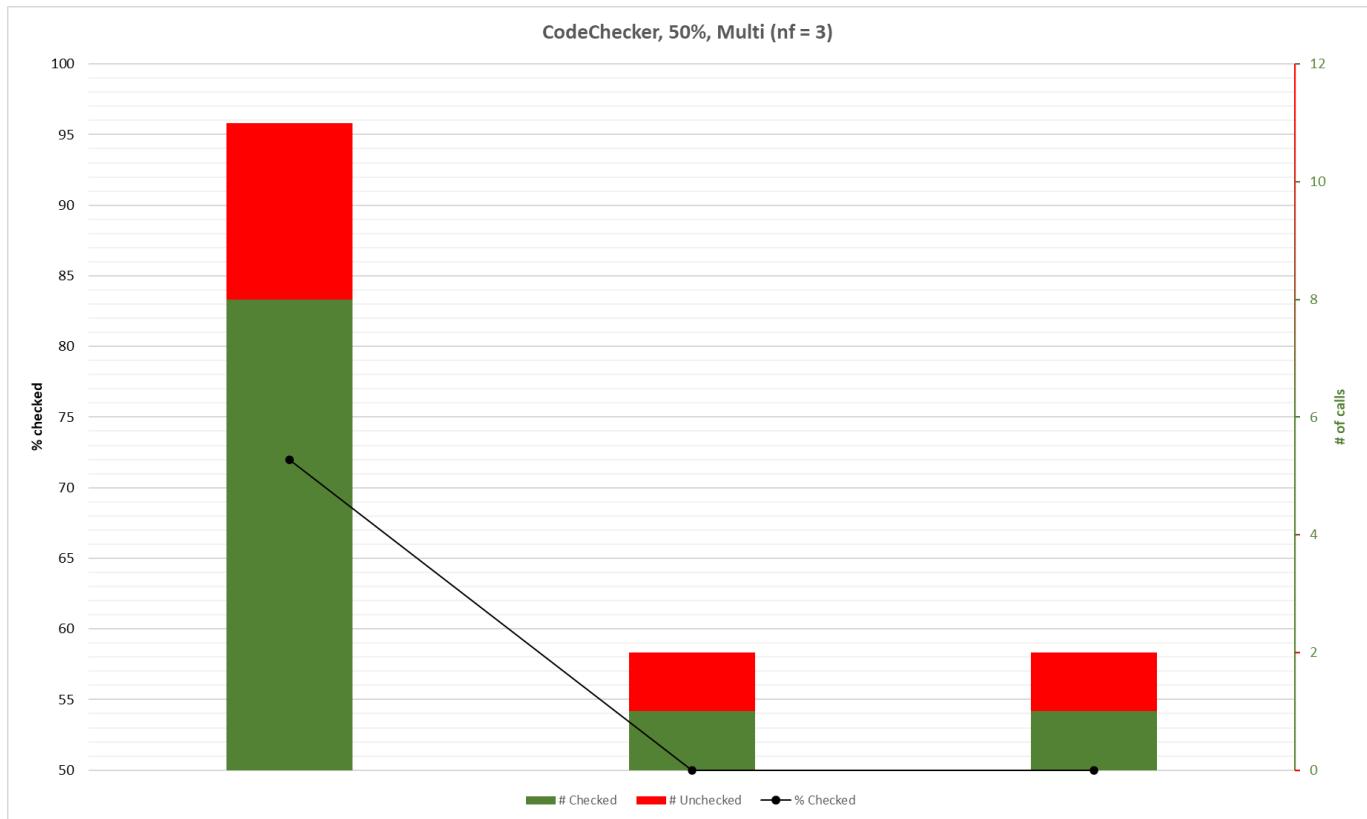
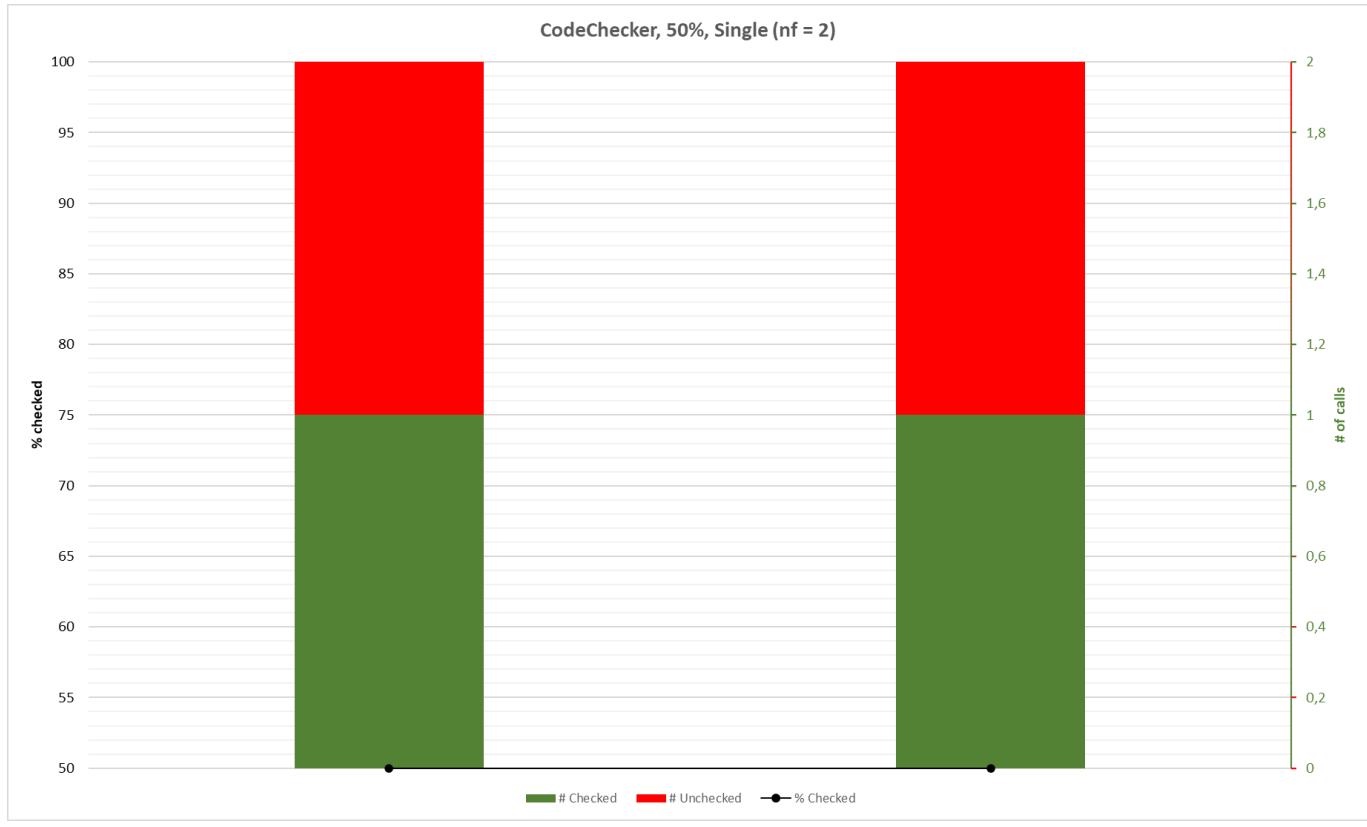
Bitcoin



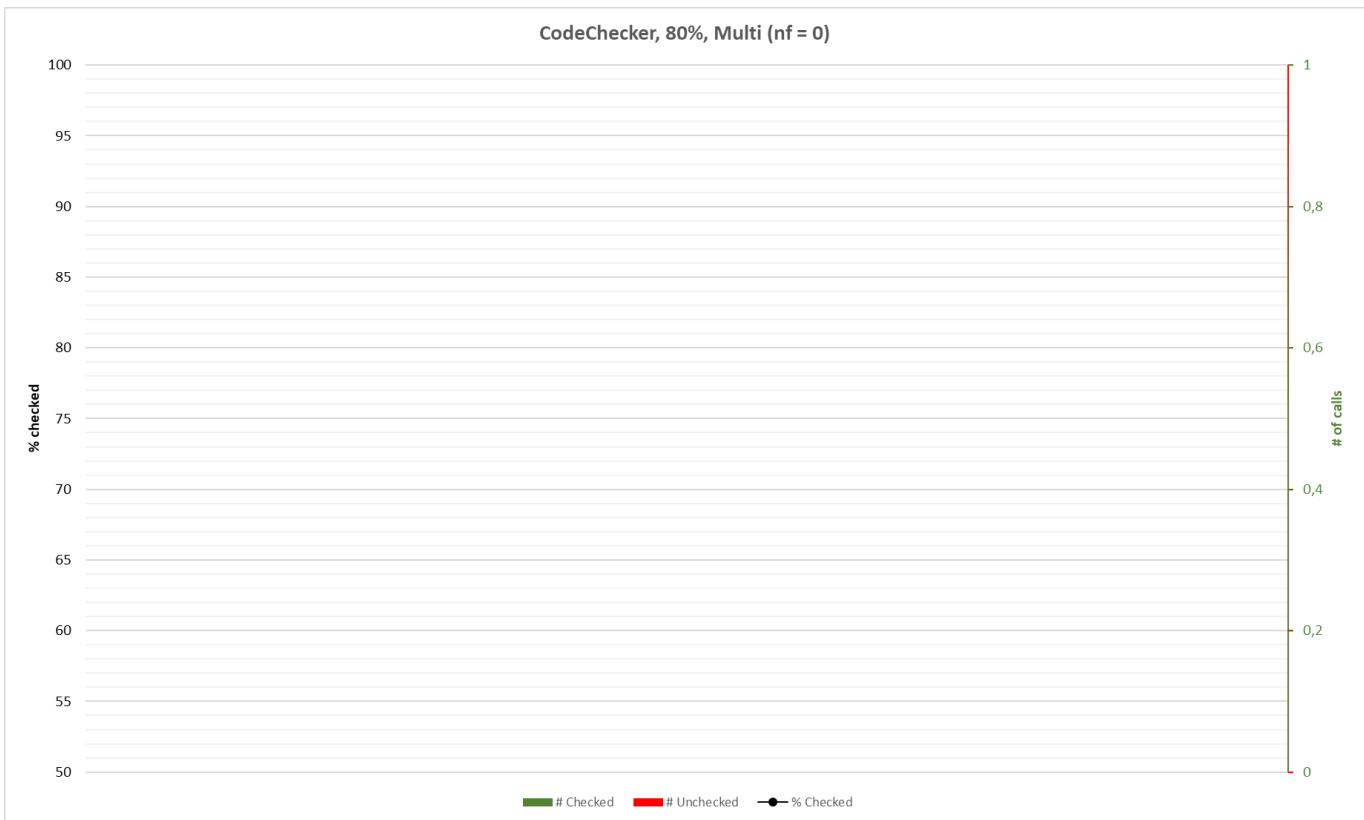
A. Project results



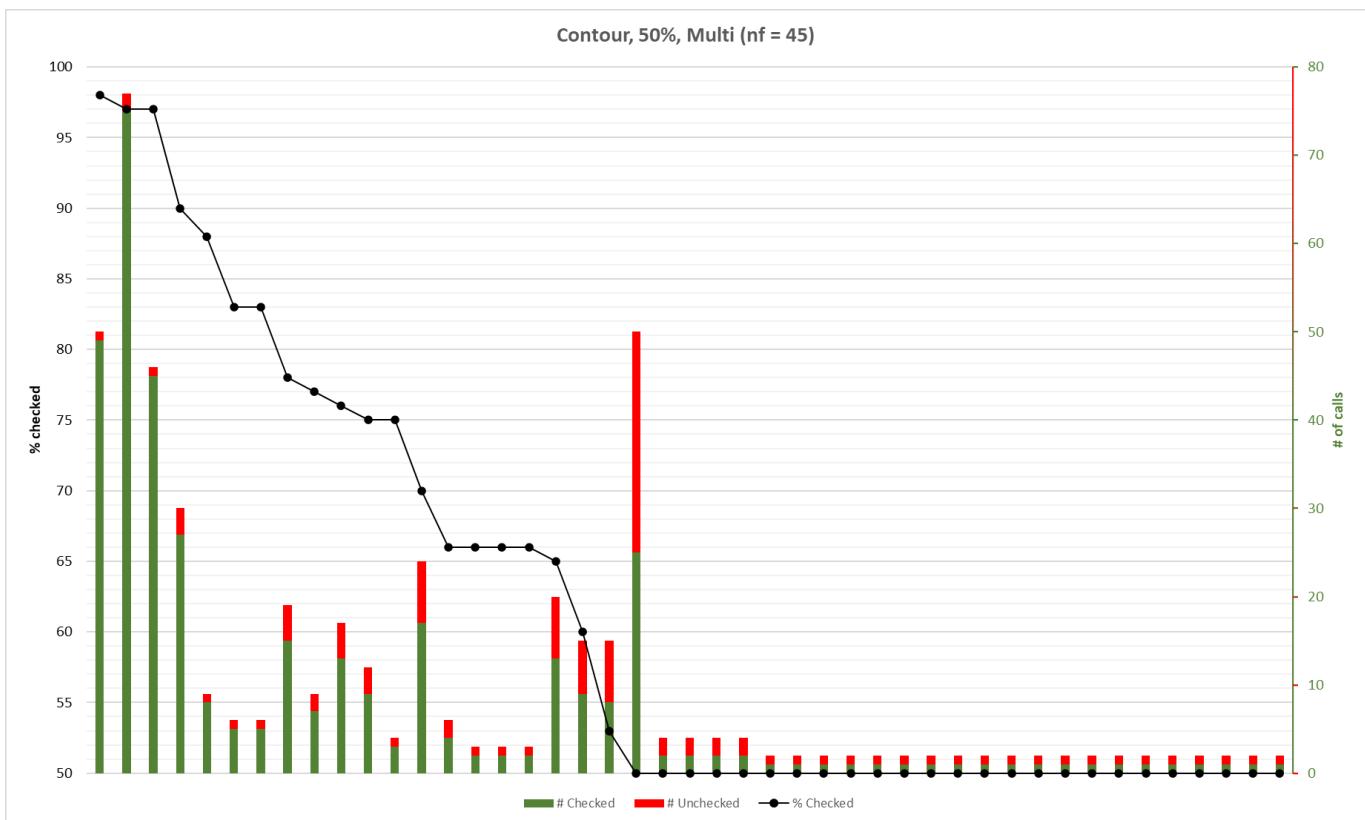
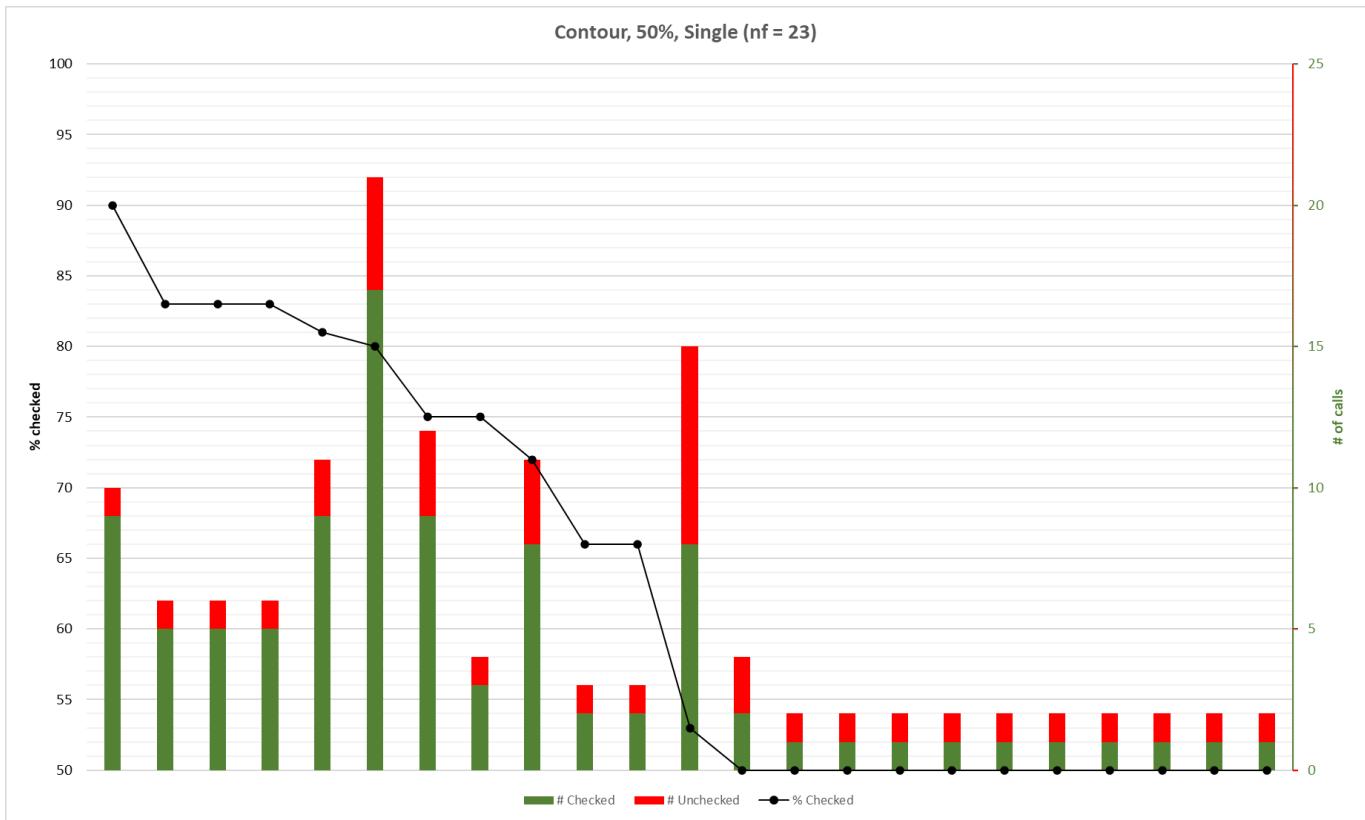
CodeChecker



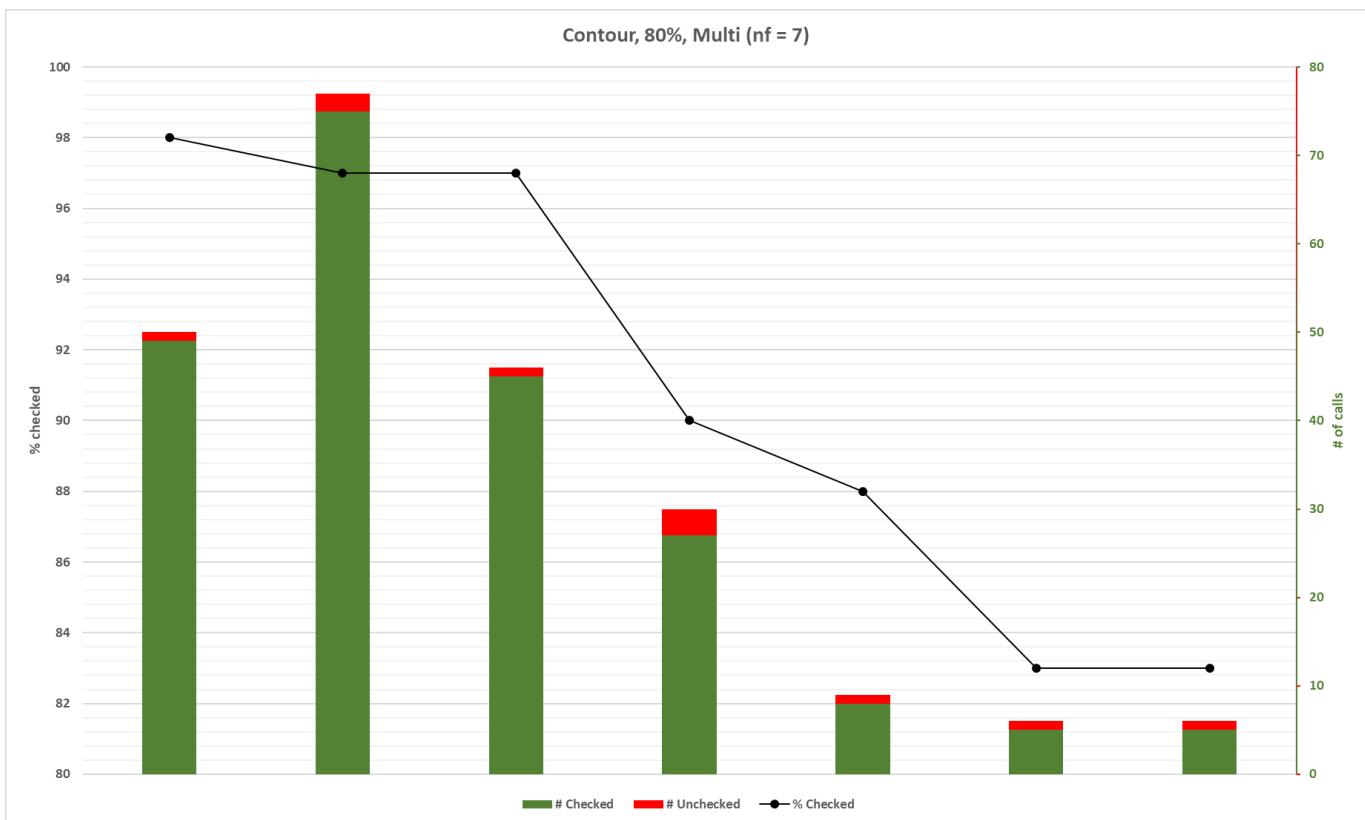
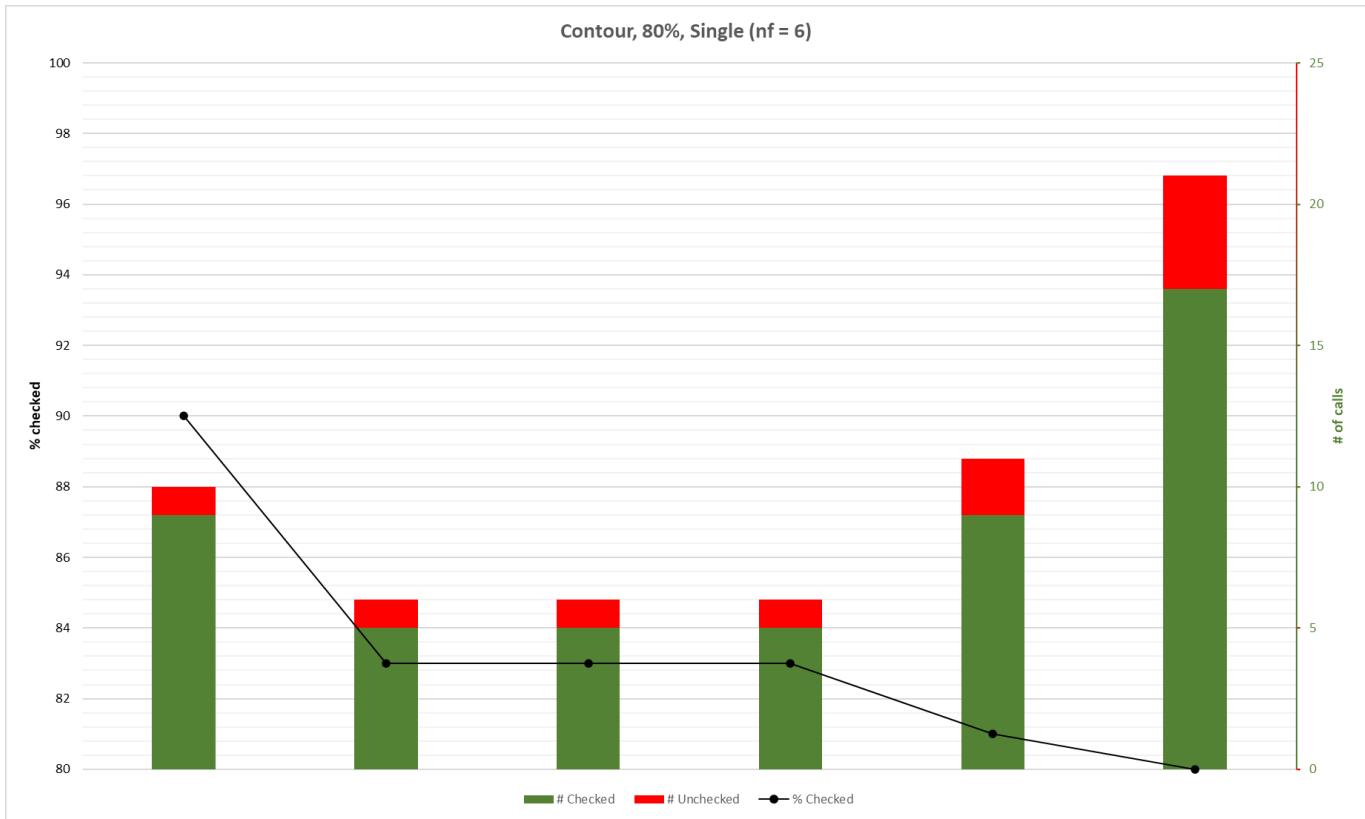
A. Project results



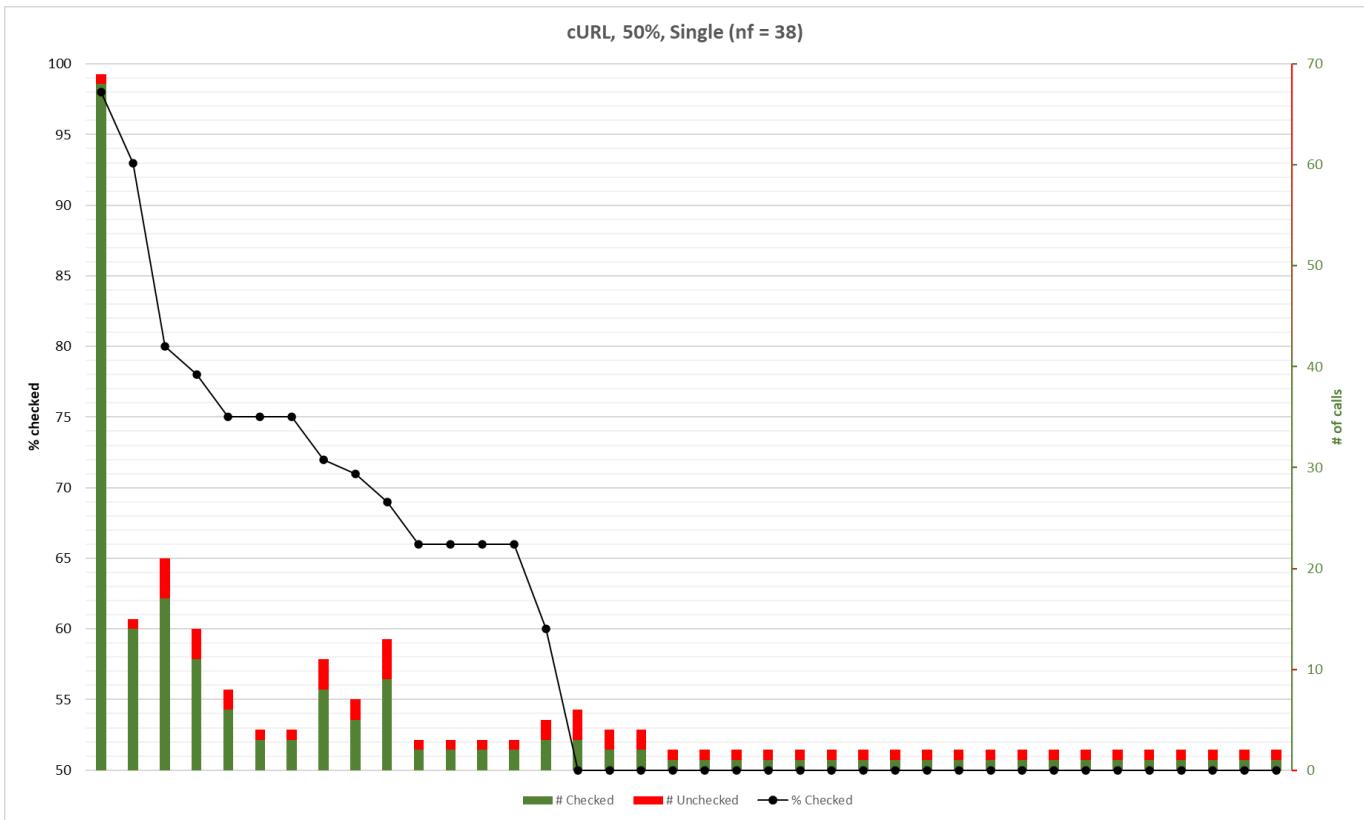
Contour



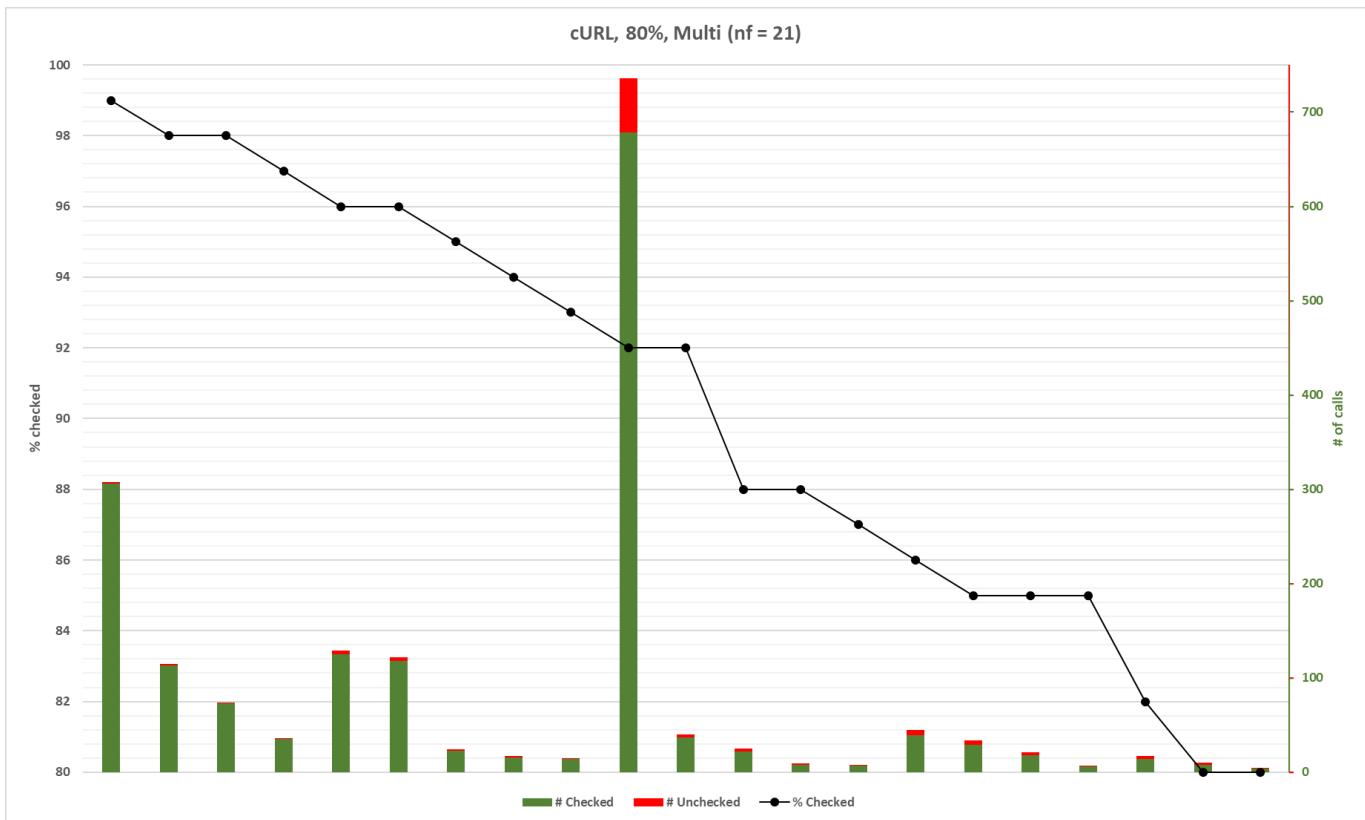
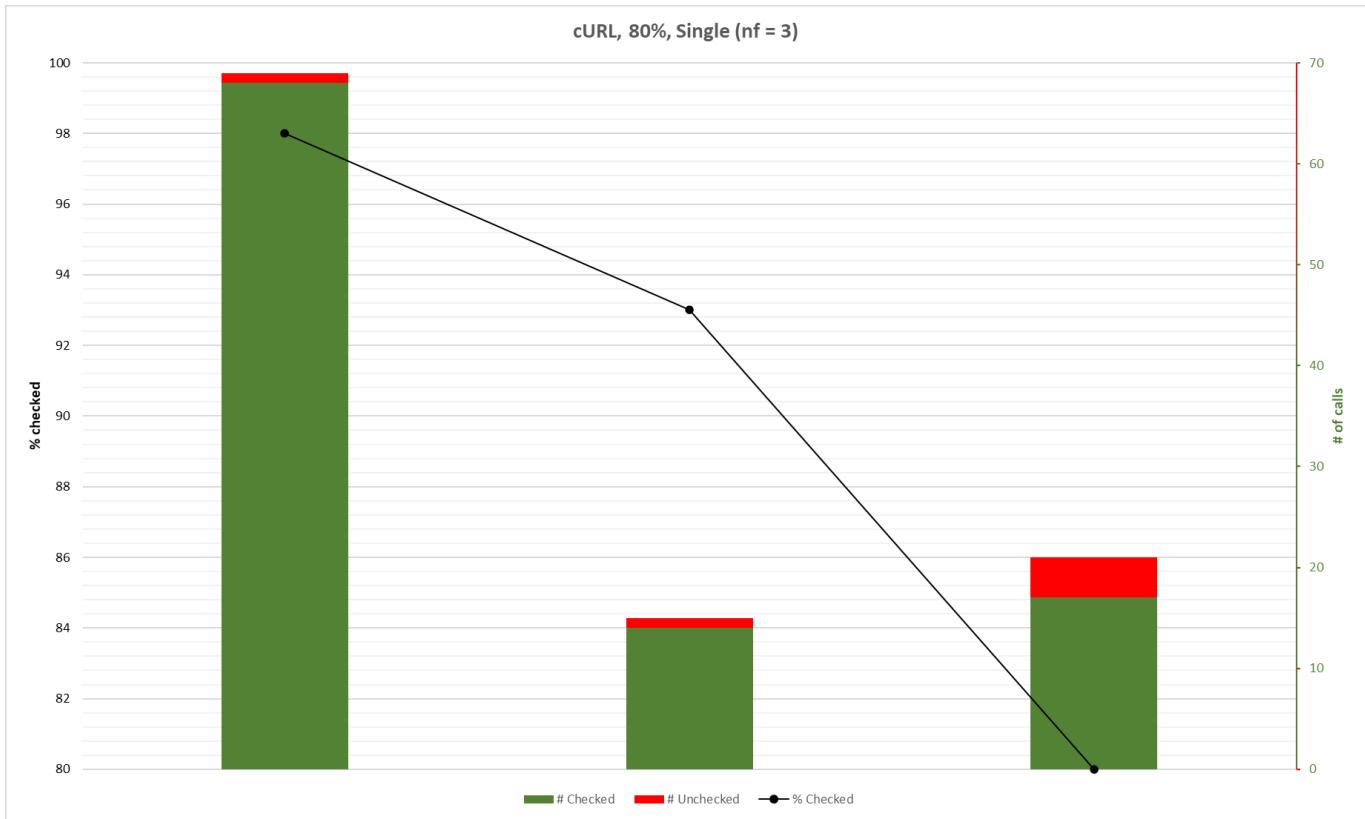
A. Project results



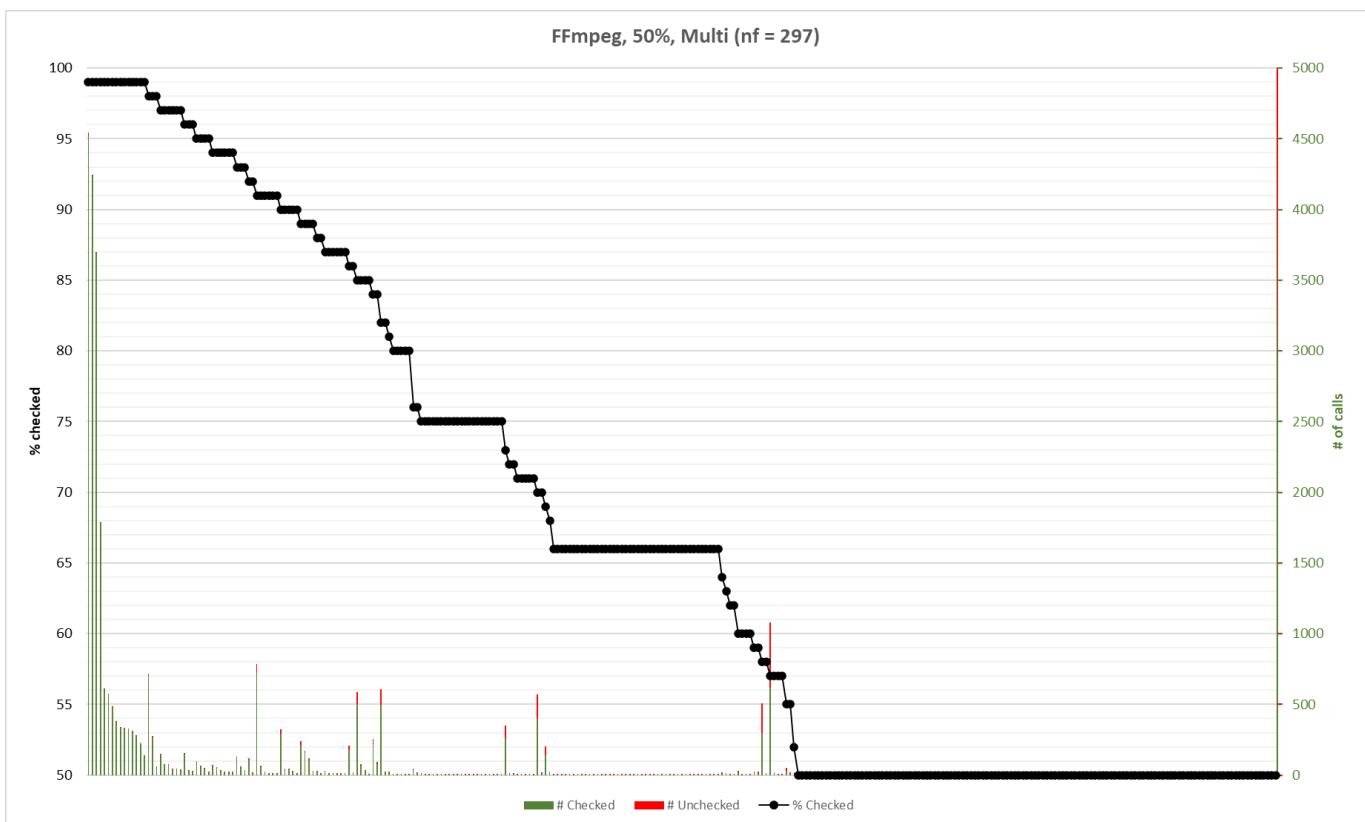
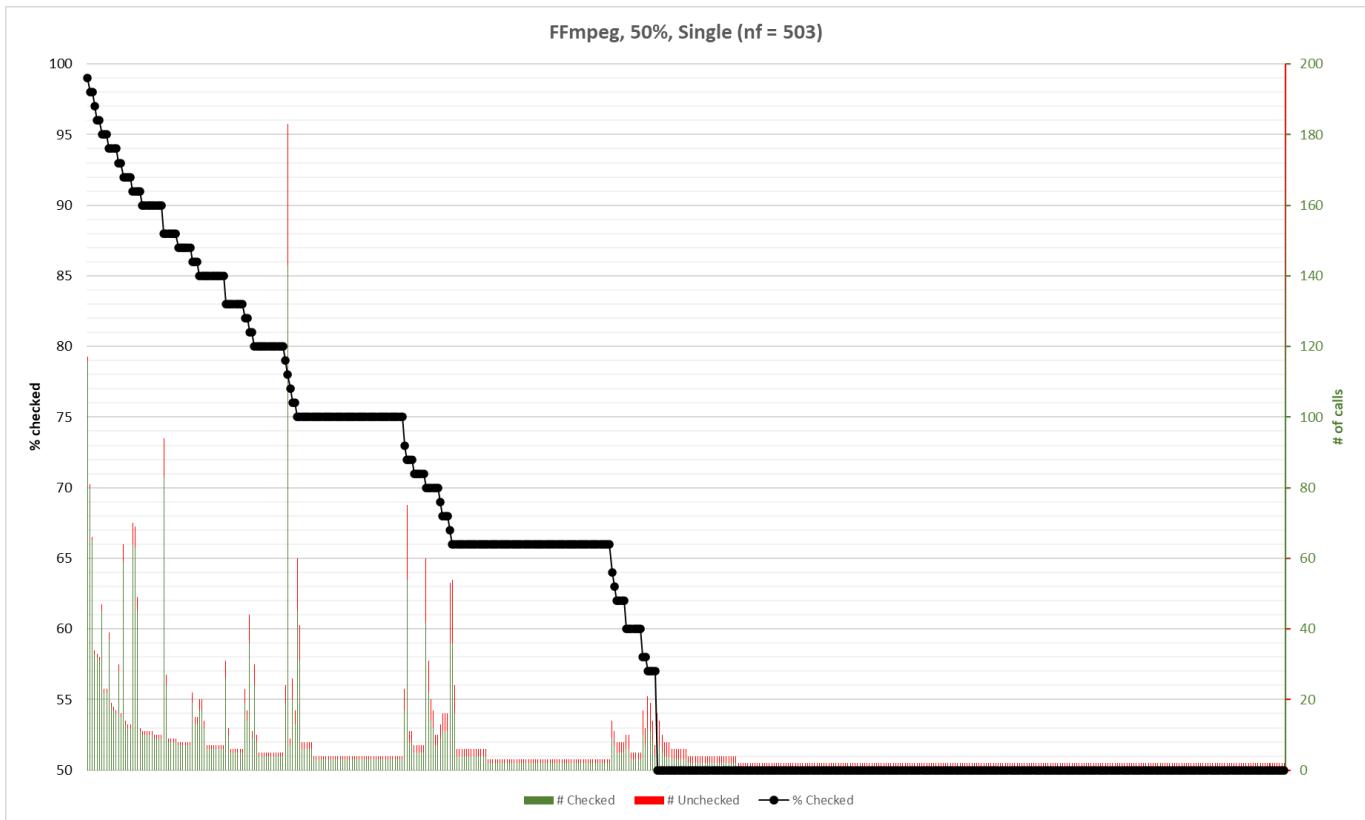
cURL



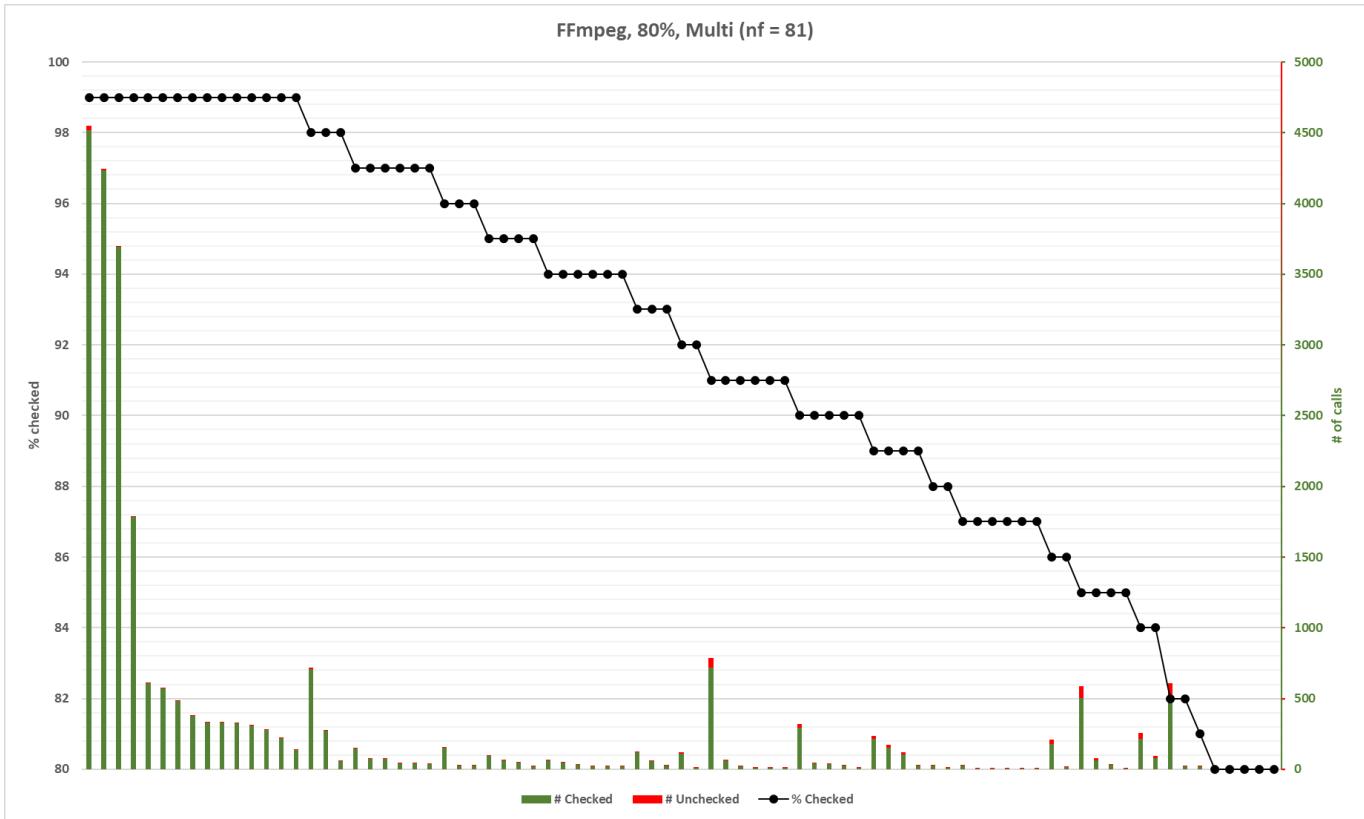
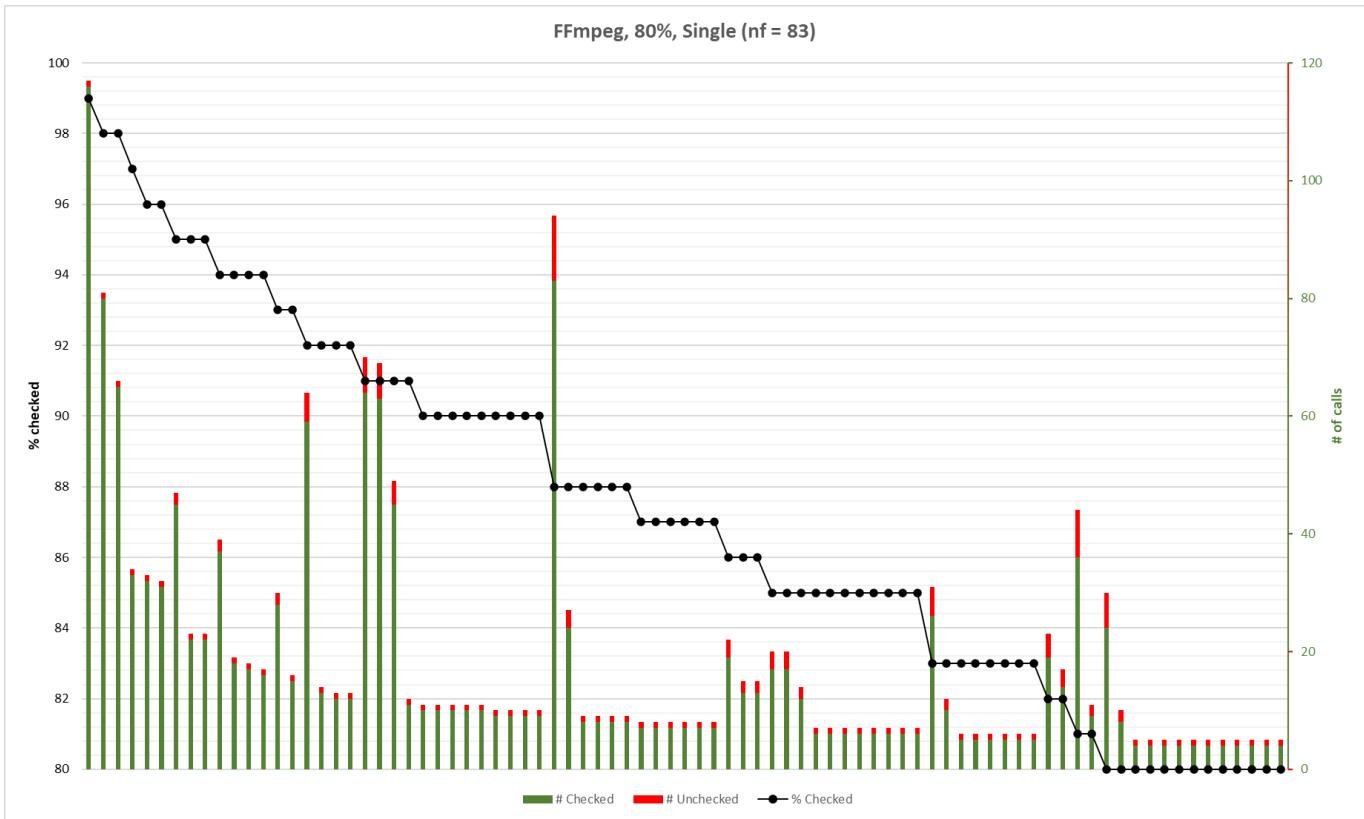
A. Project results



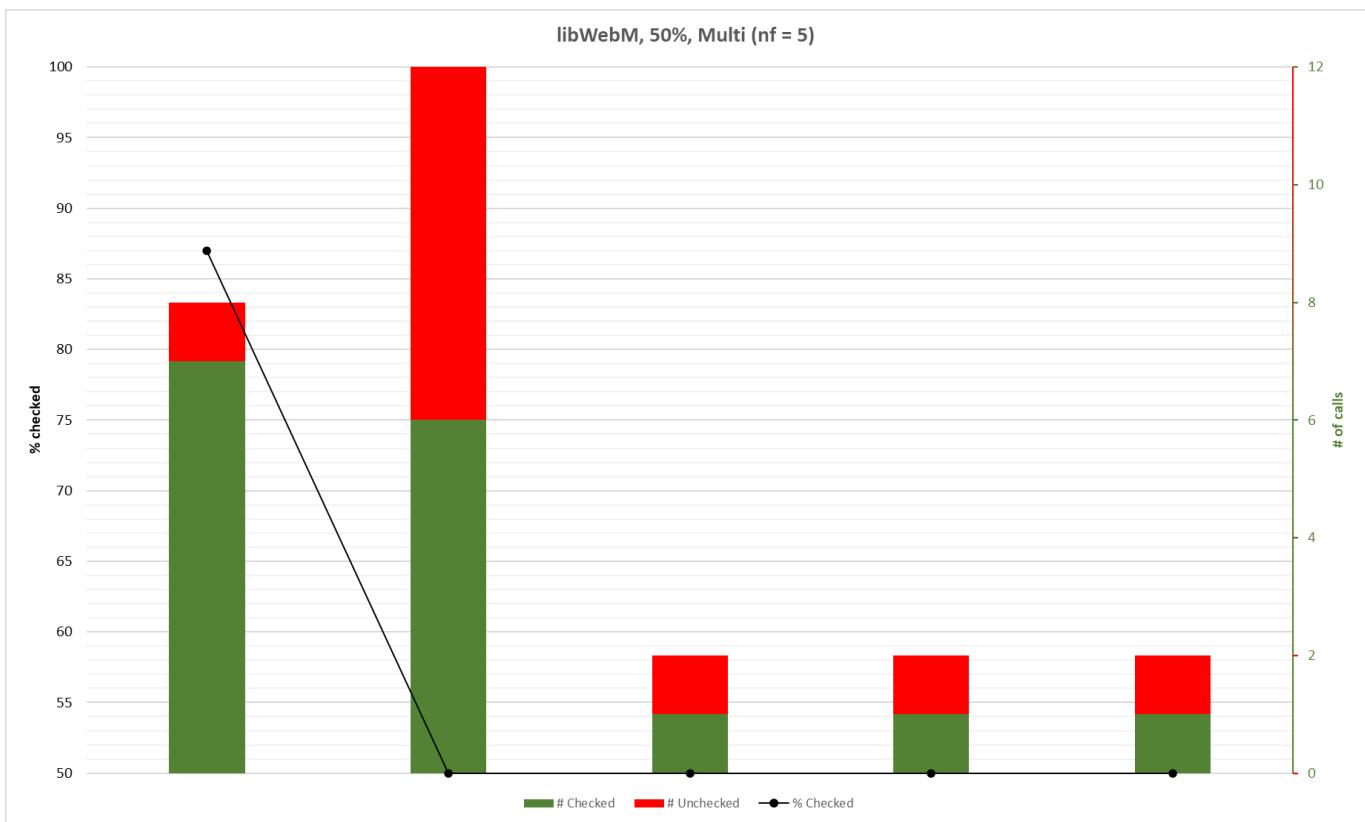
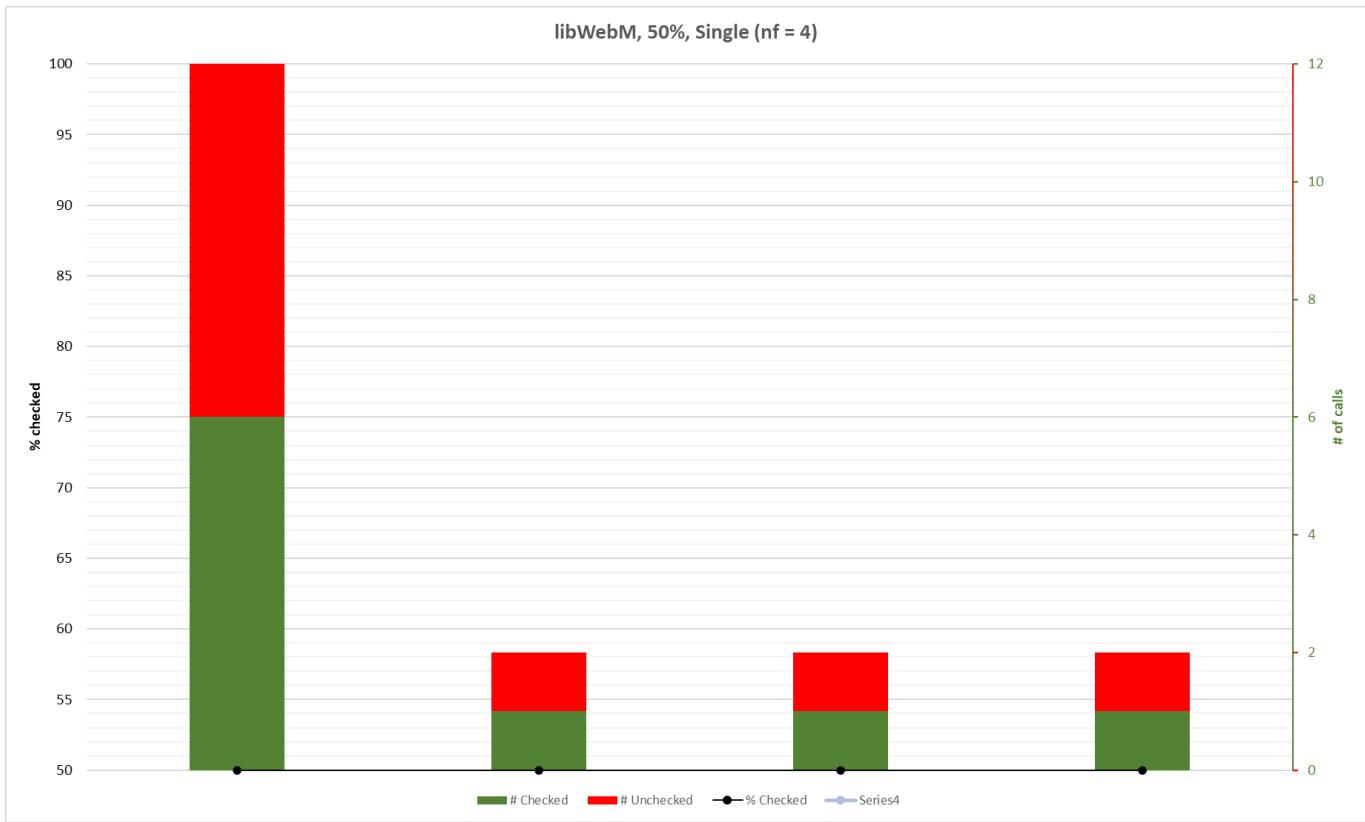
FFmpeg



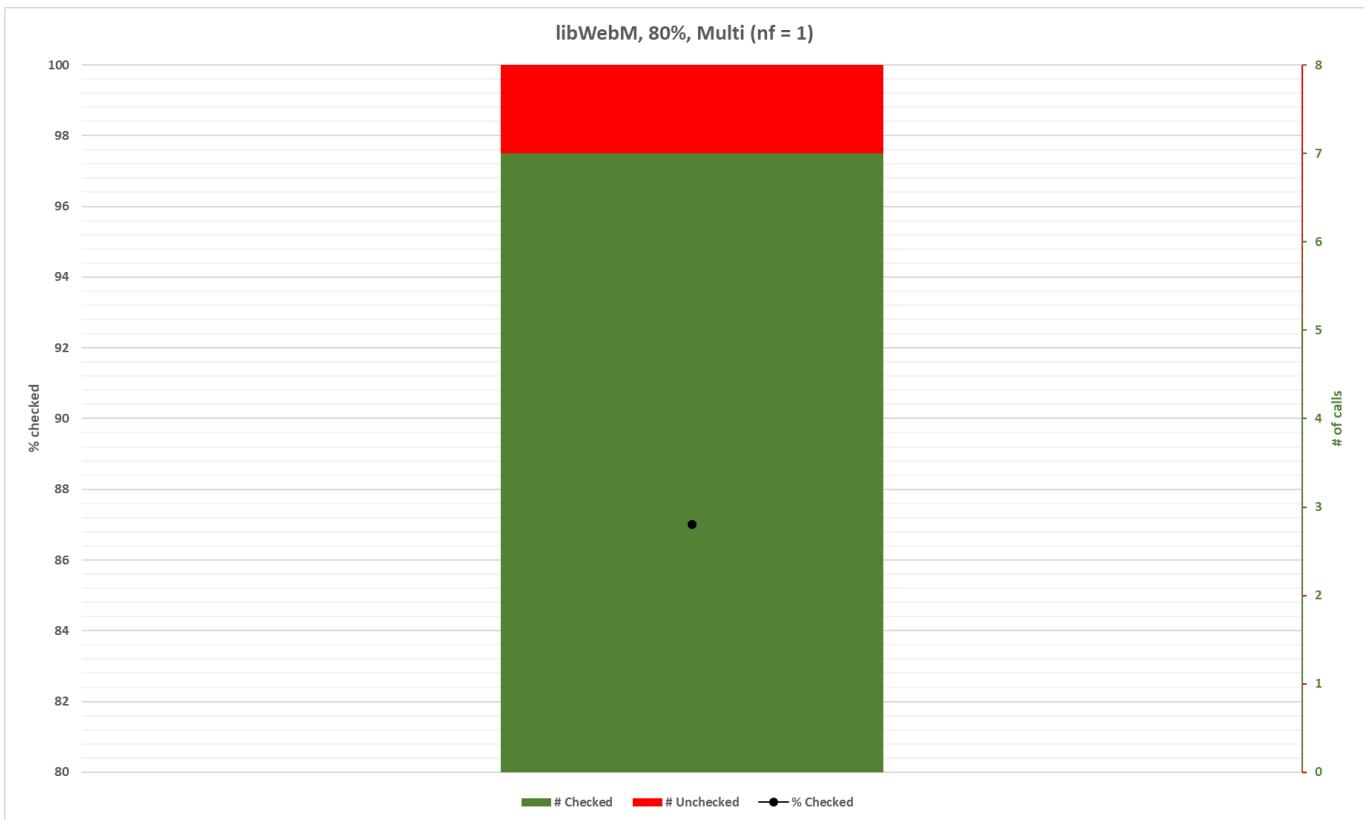
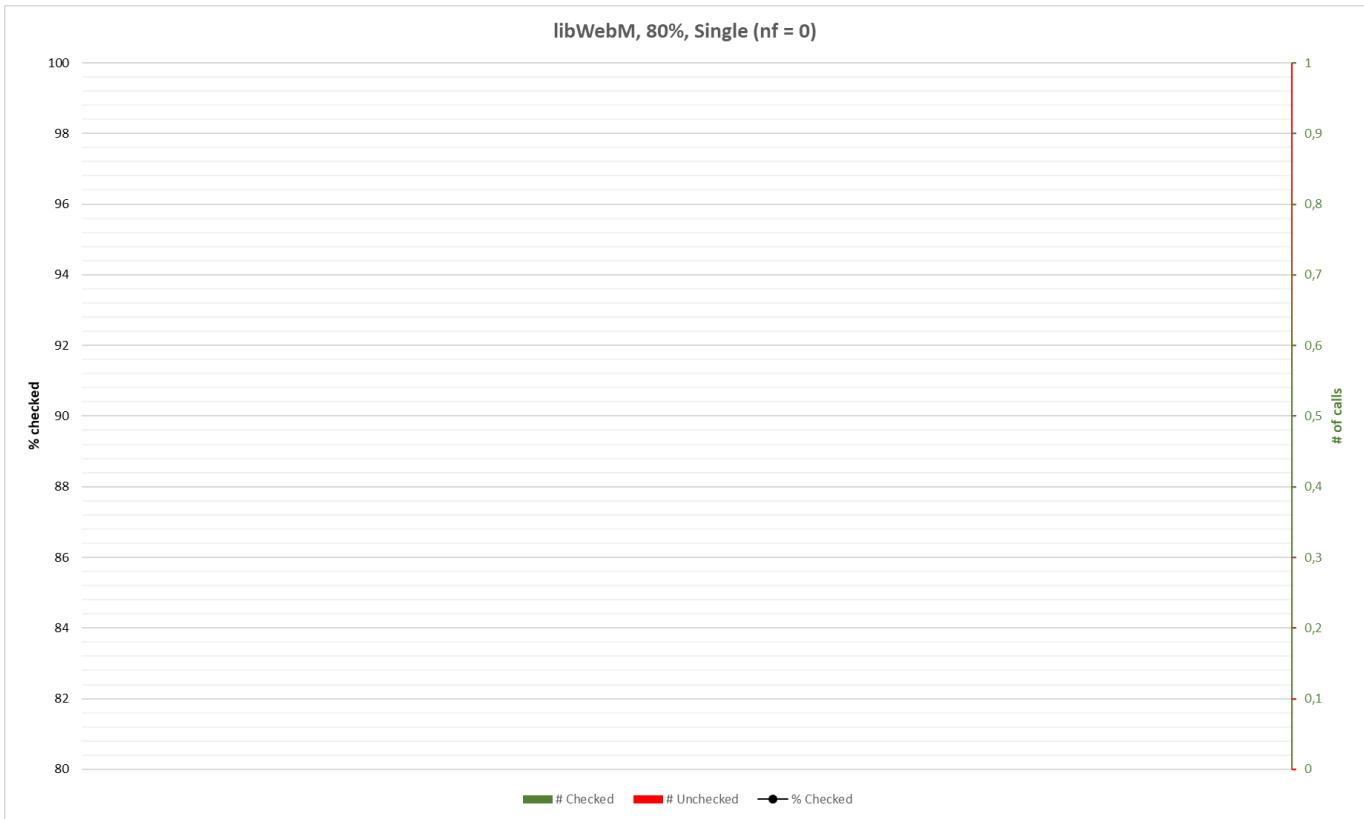
A. Project results



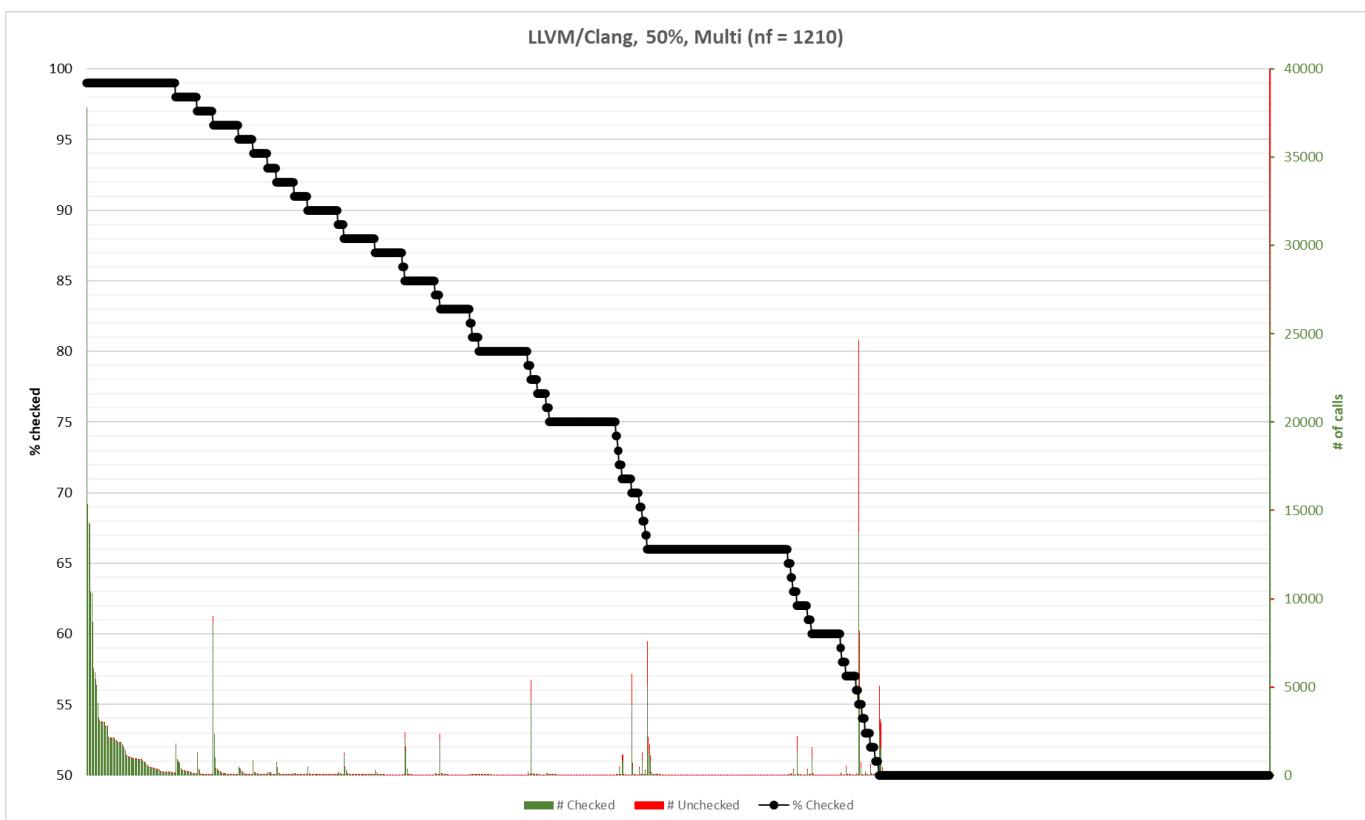
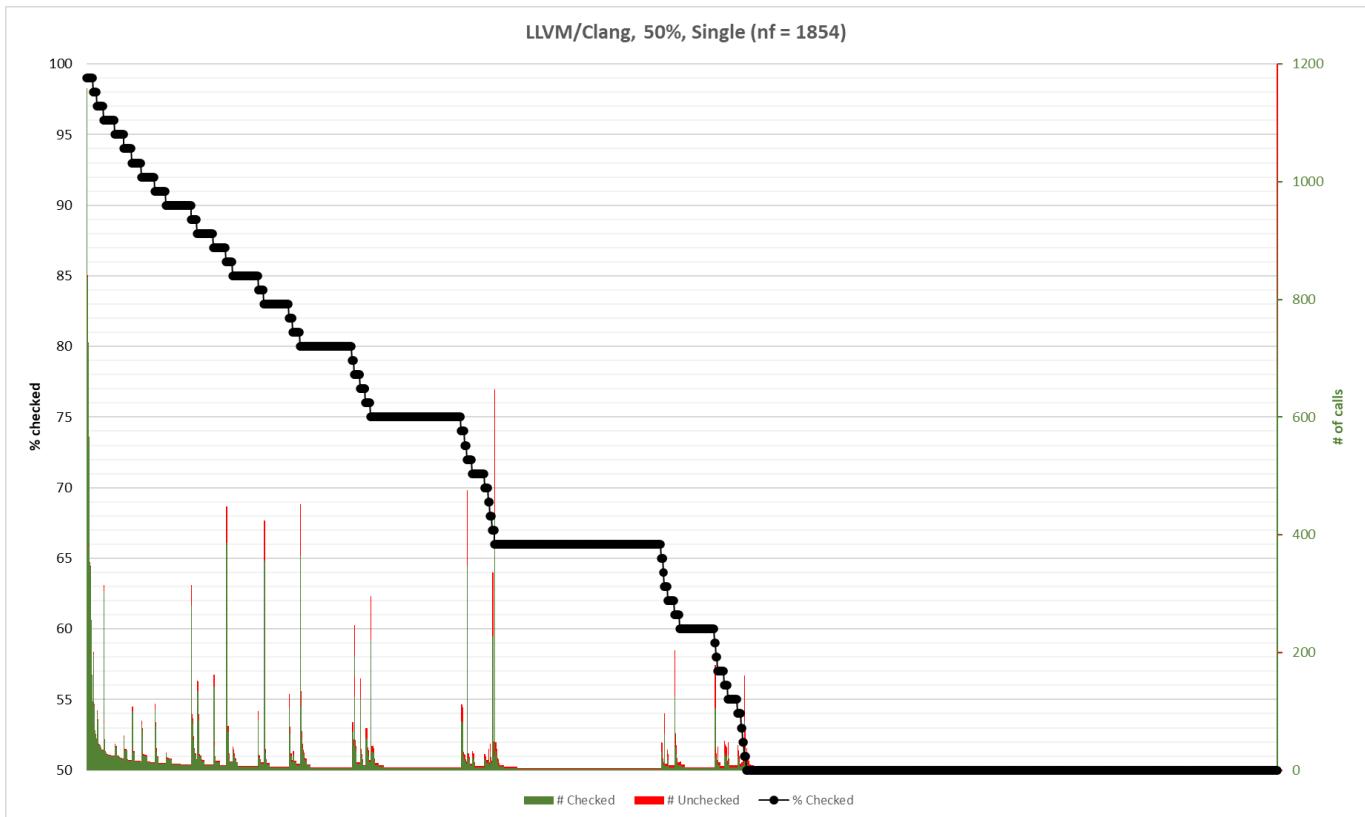
libWebM



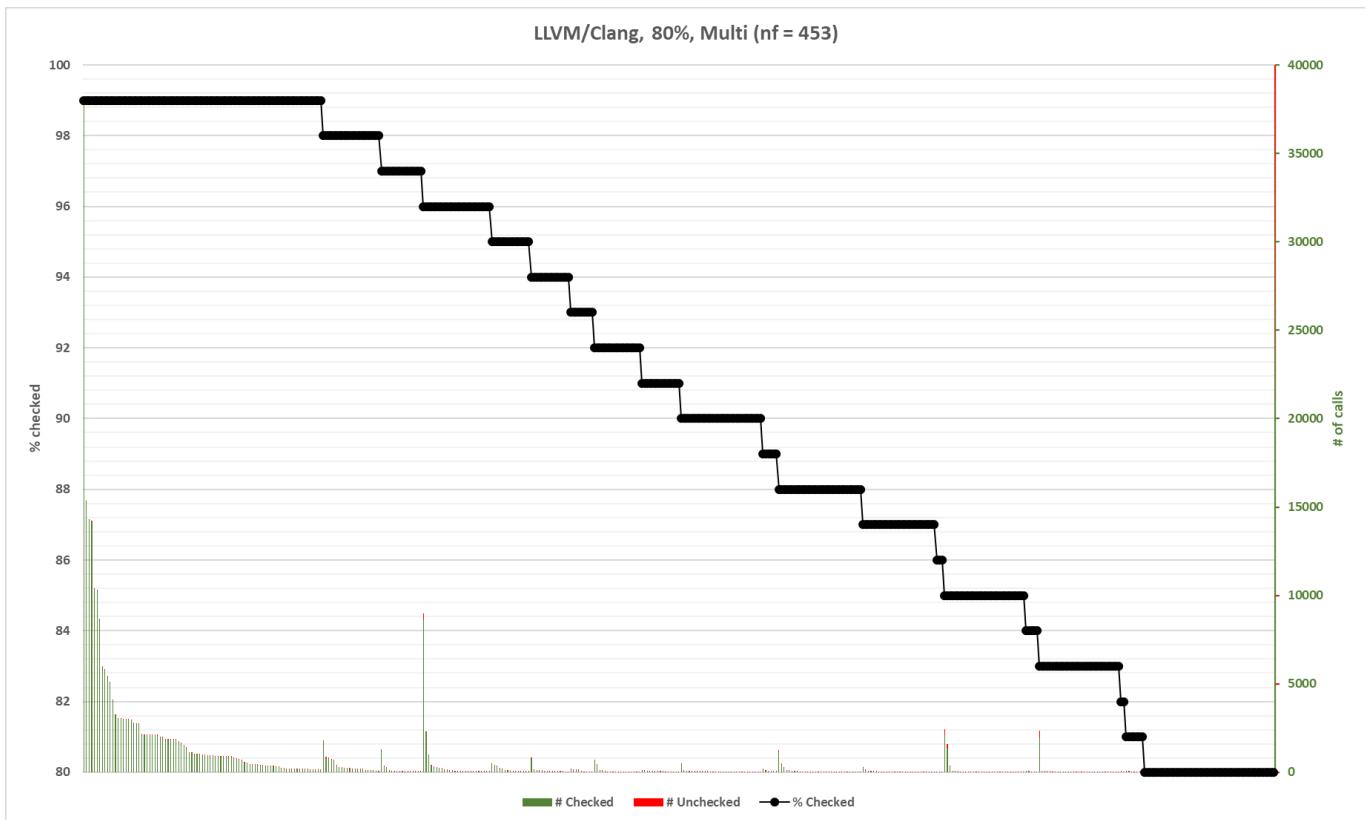
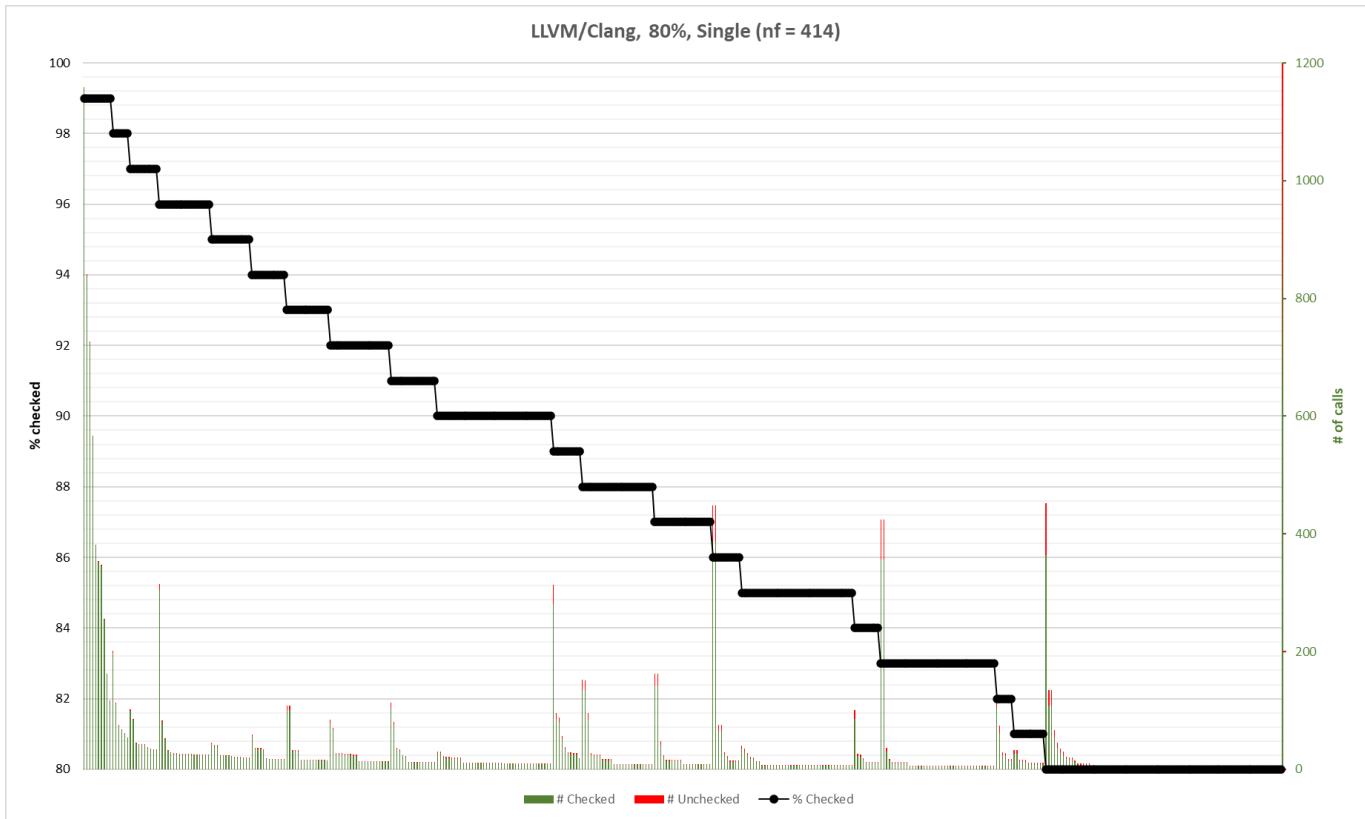
A. Project results



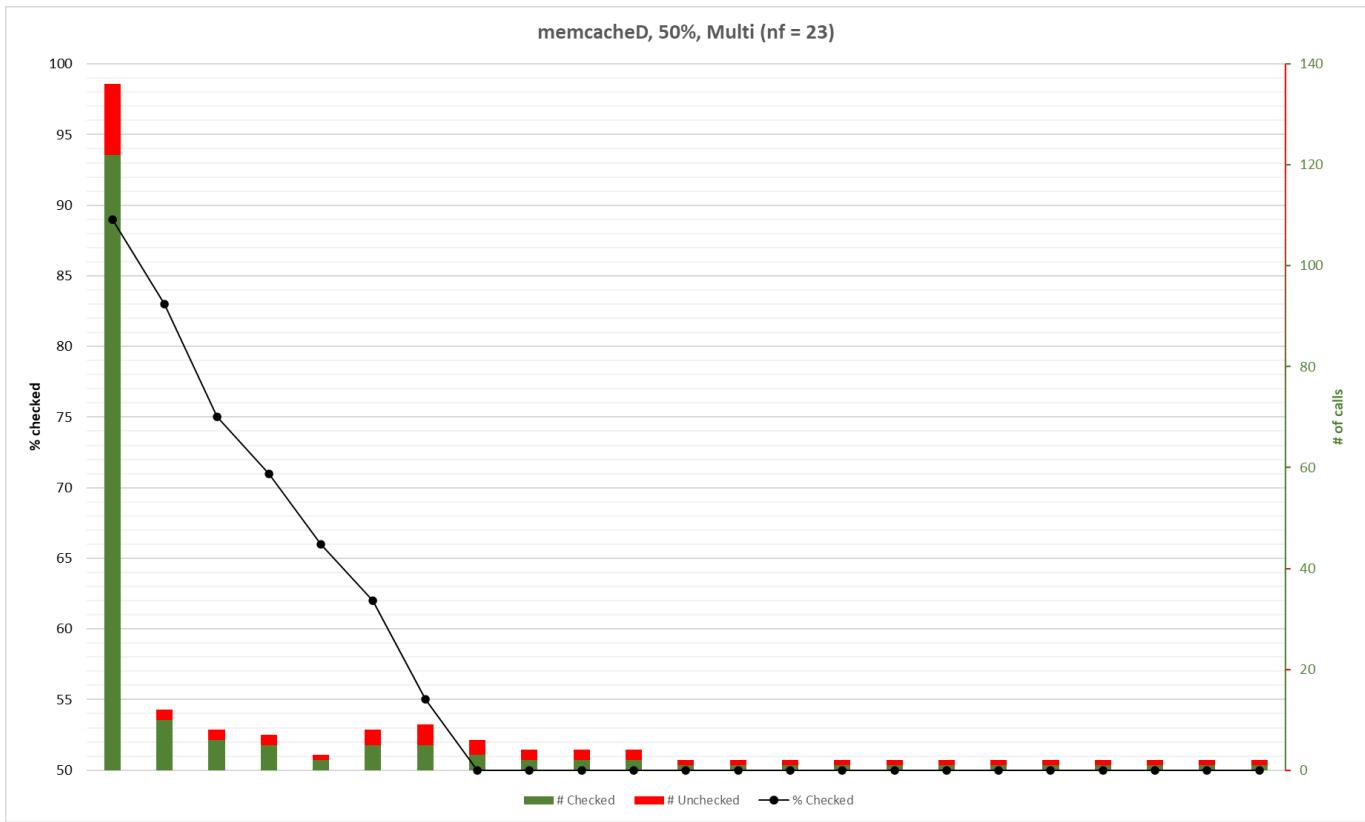
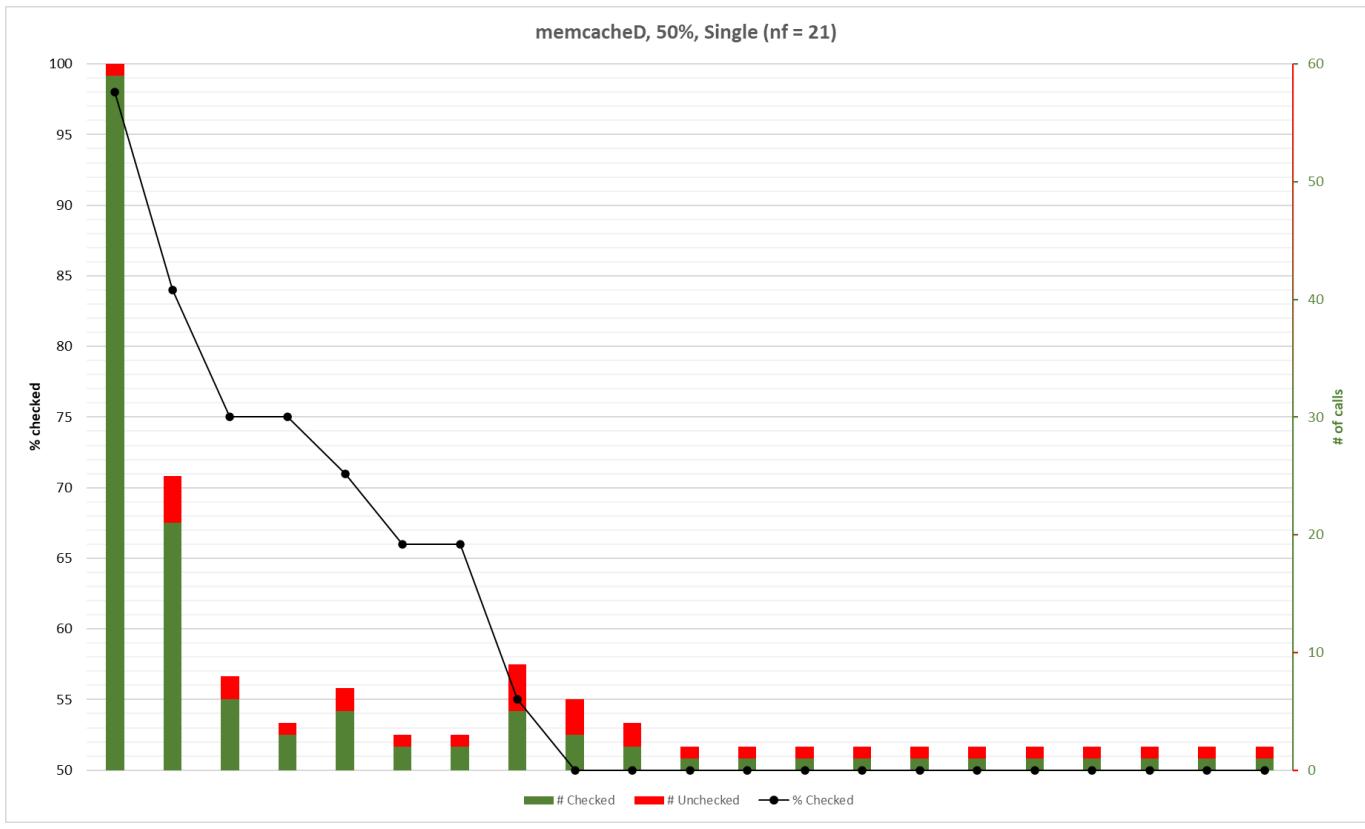
LLVM



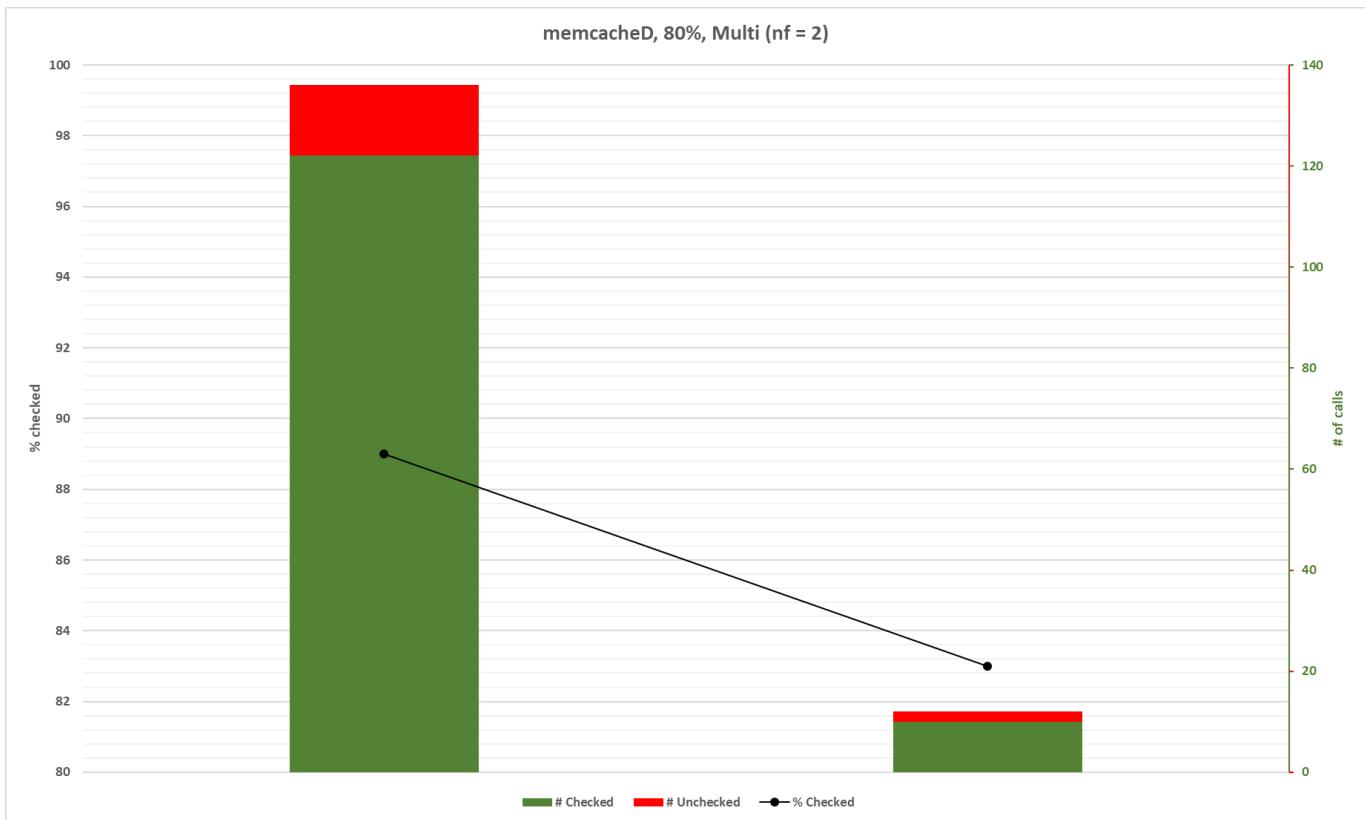
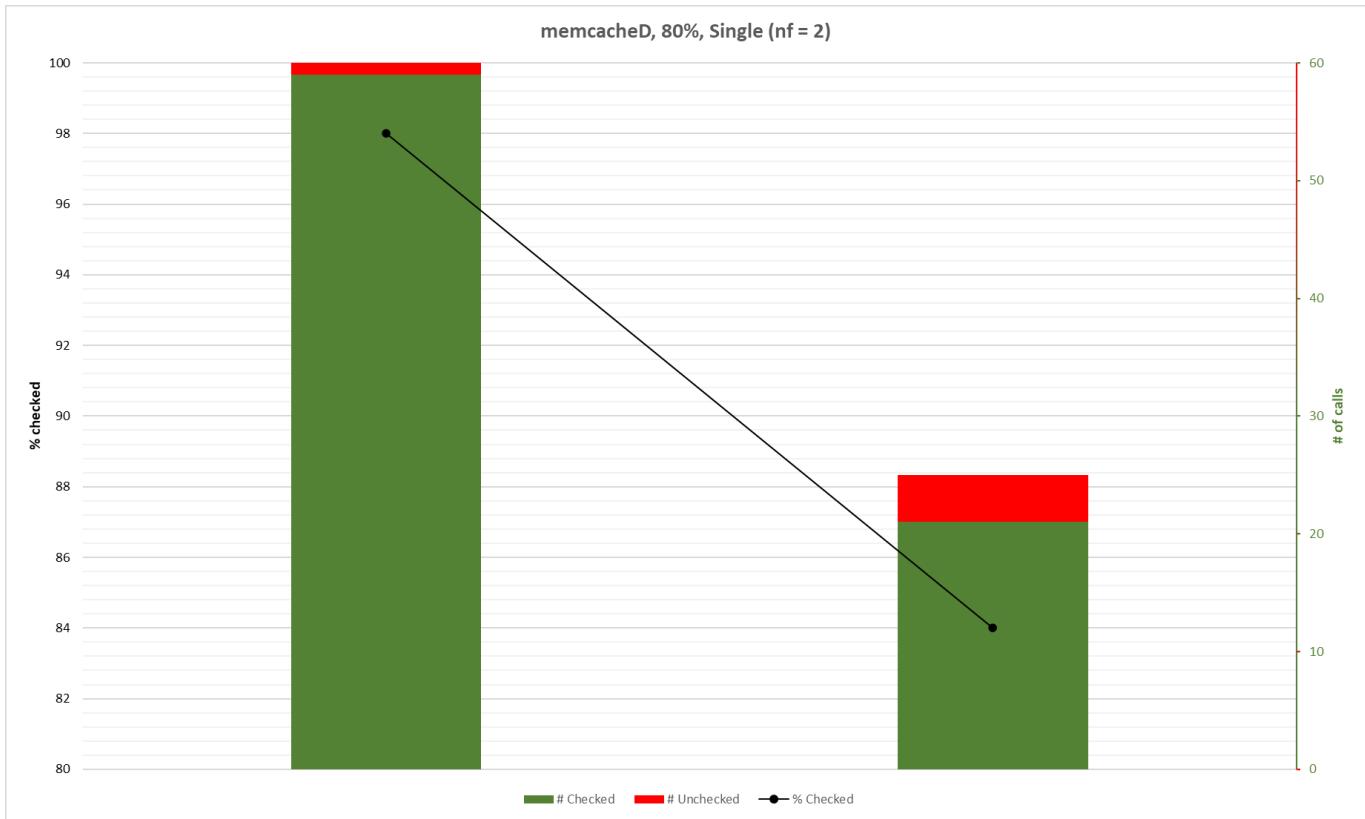
A. Project results



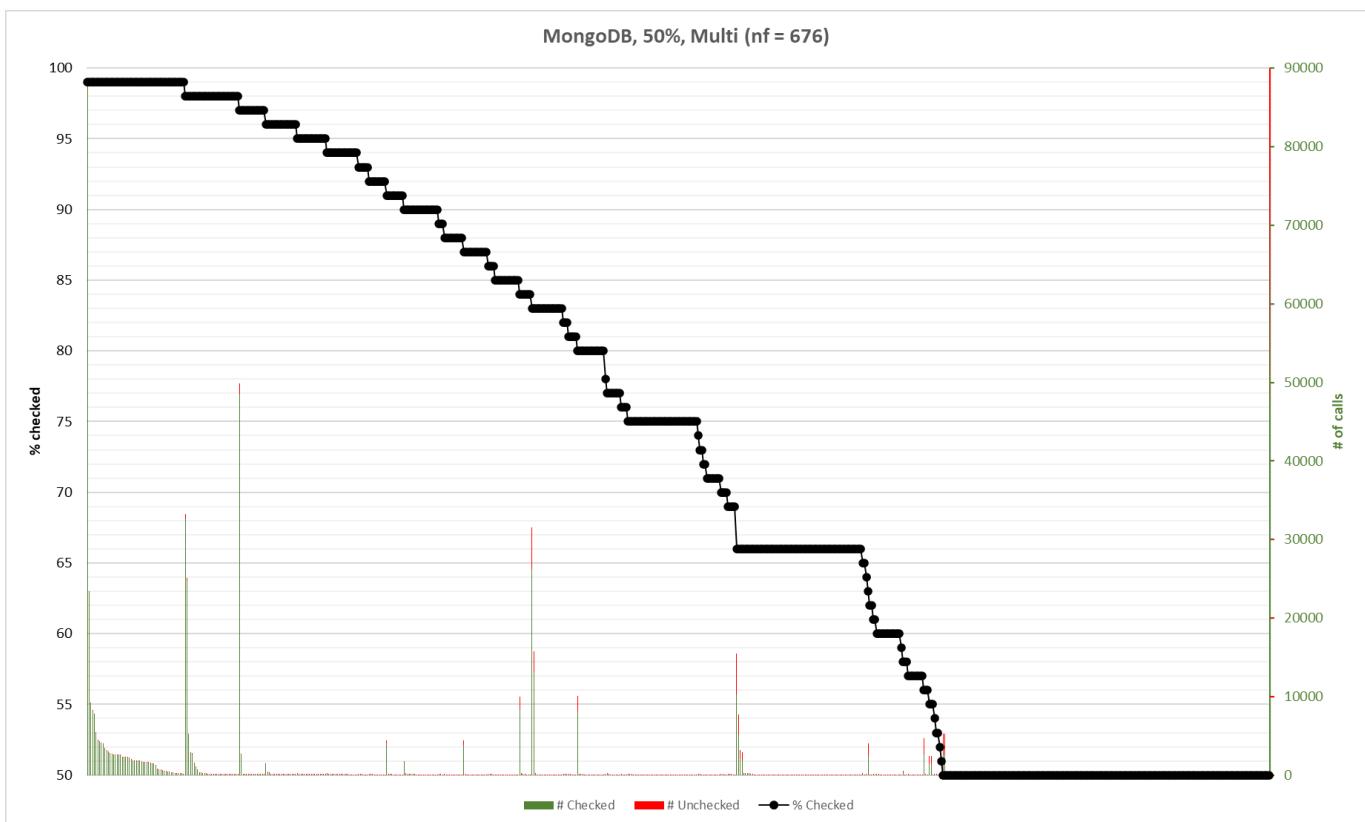
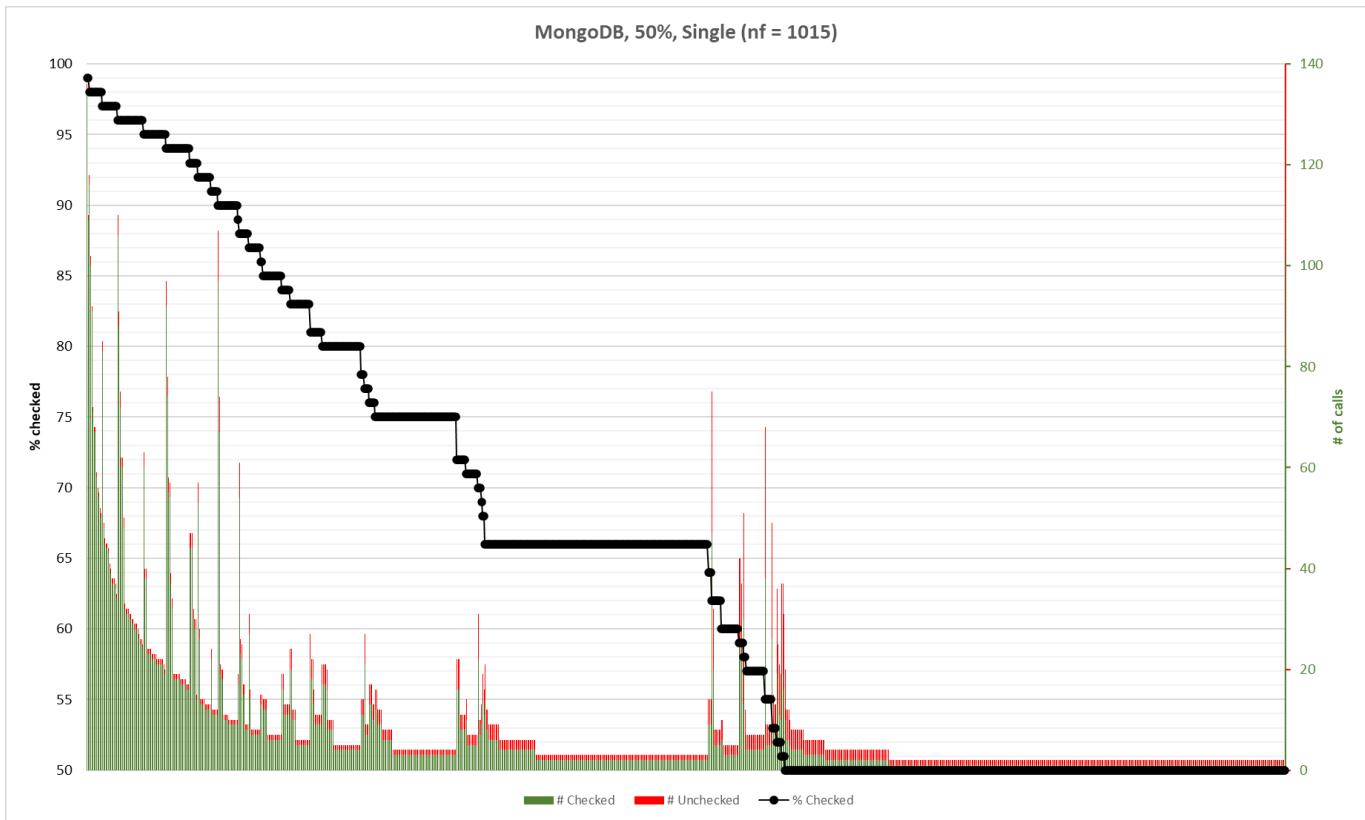
MemcacheD



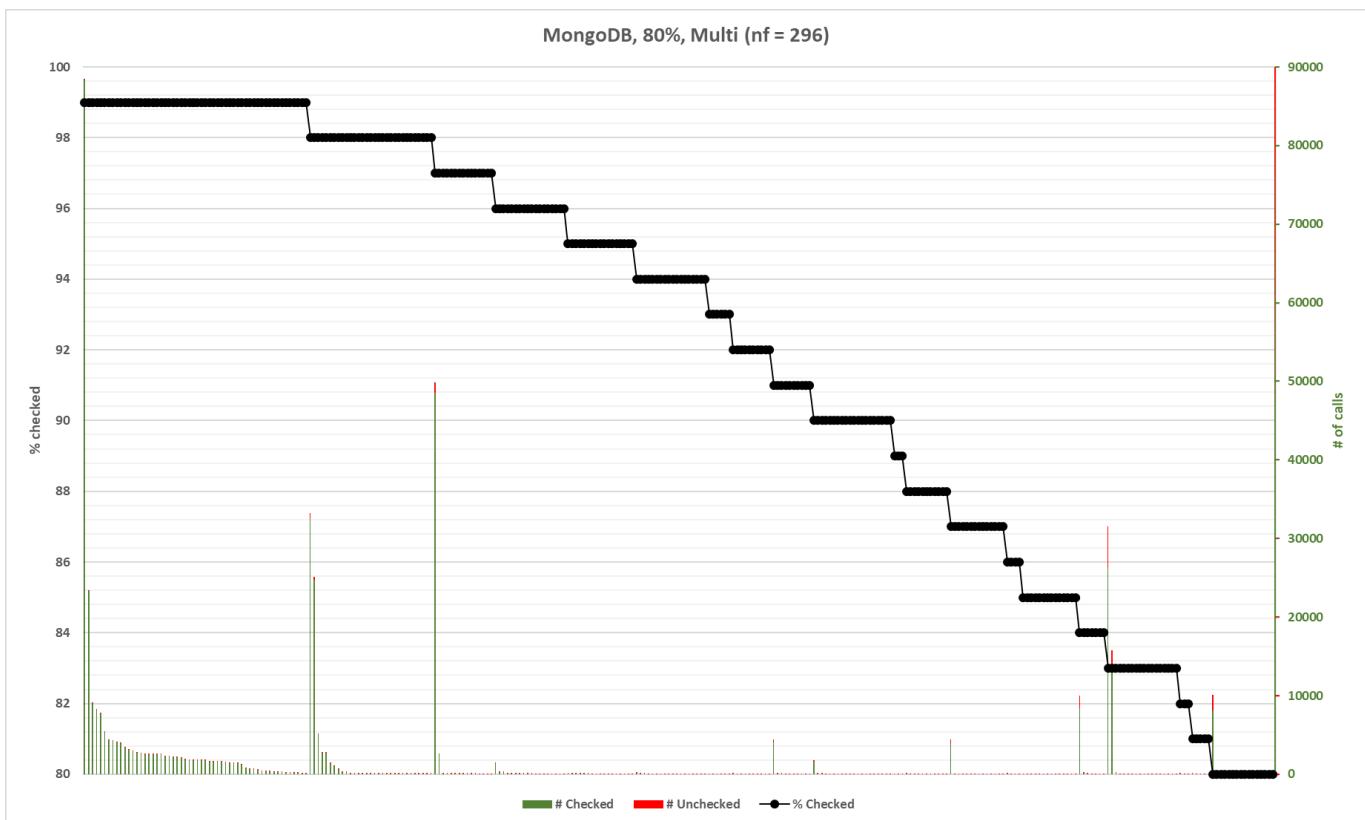
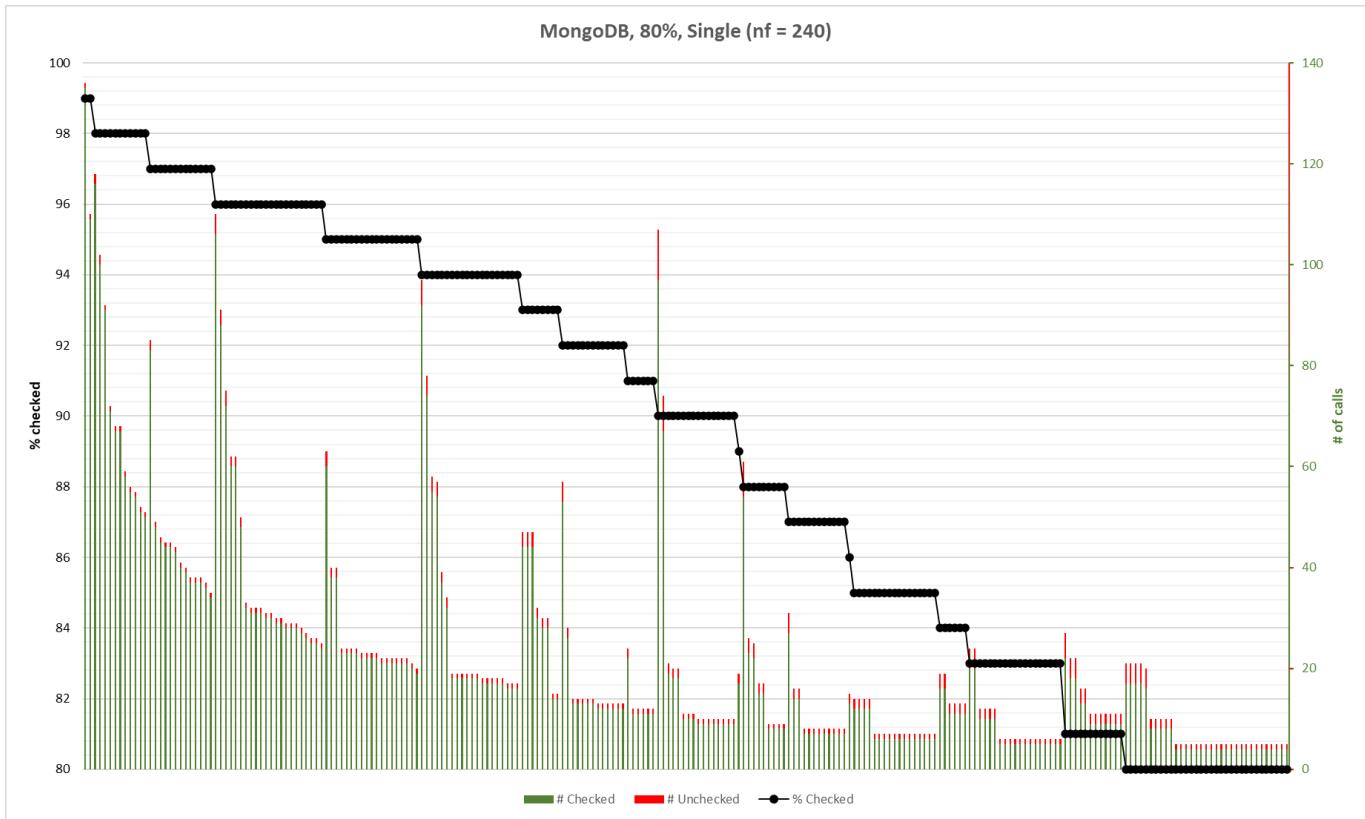
A. Project results



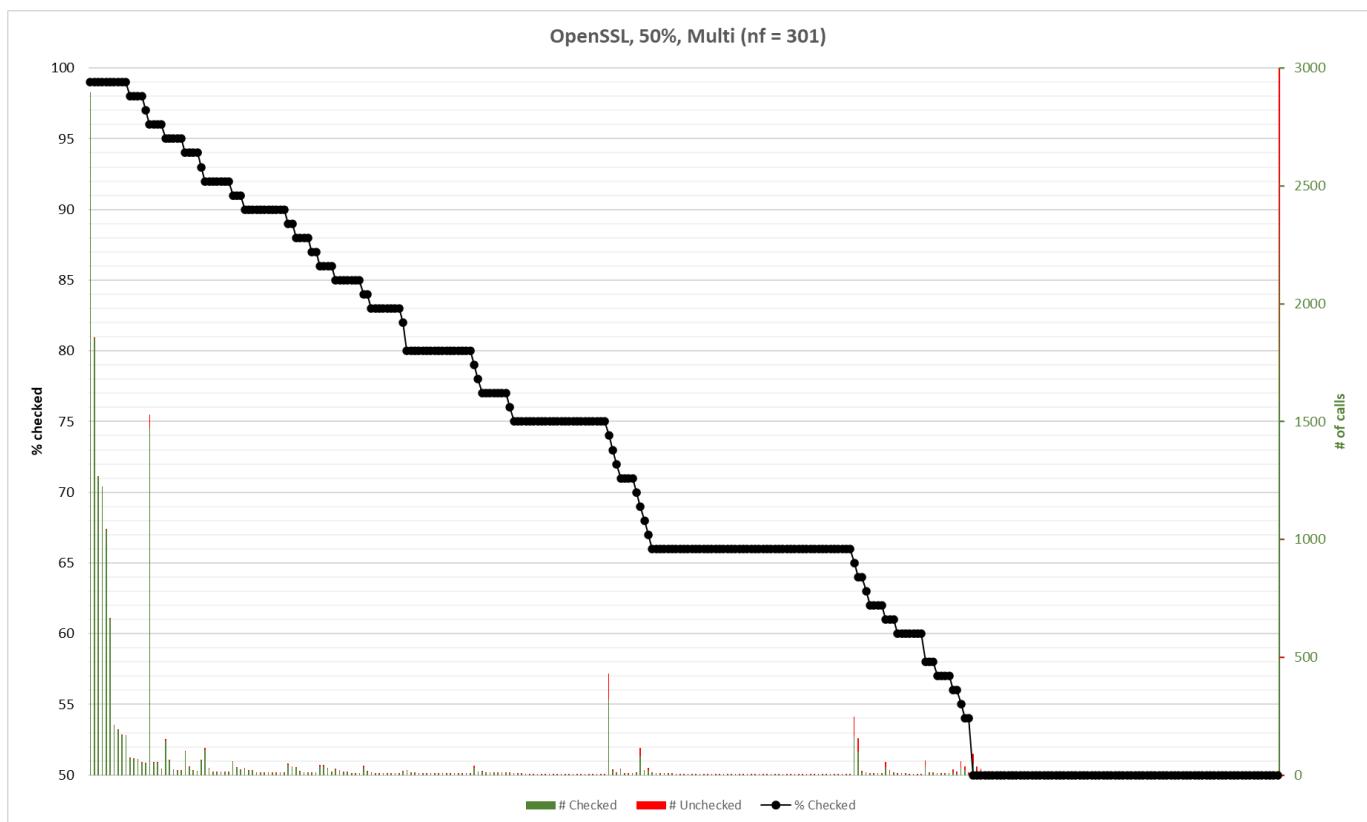
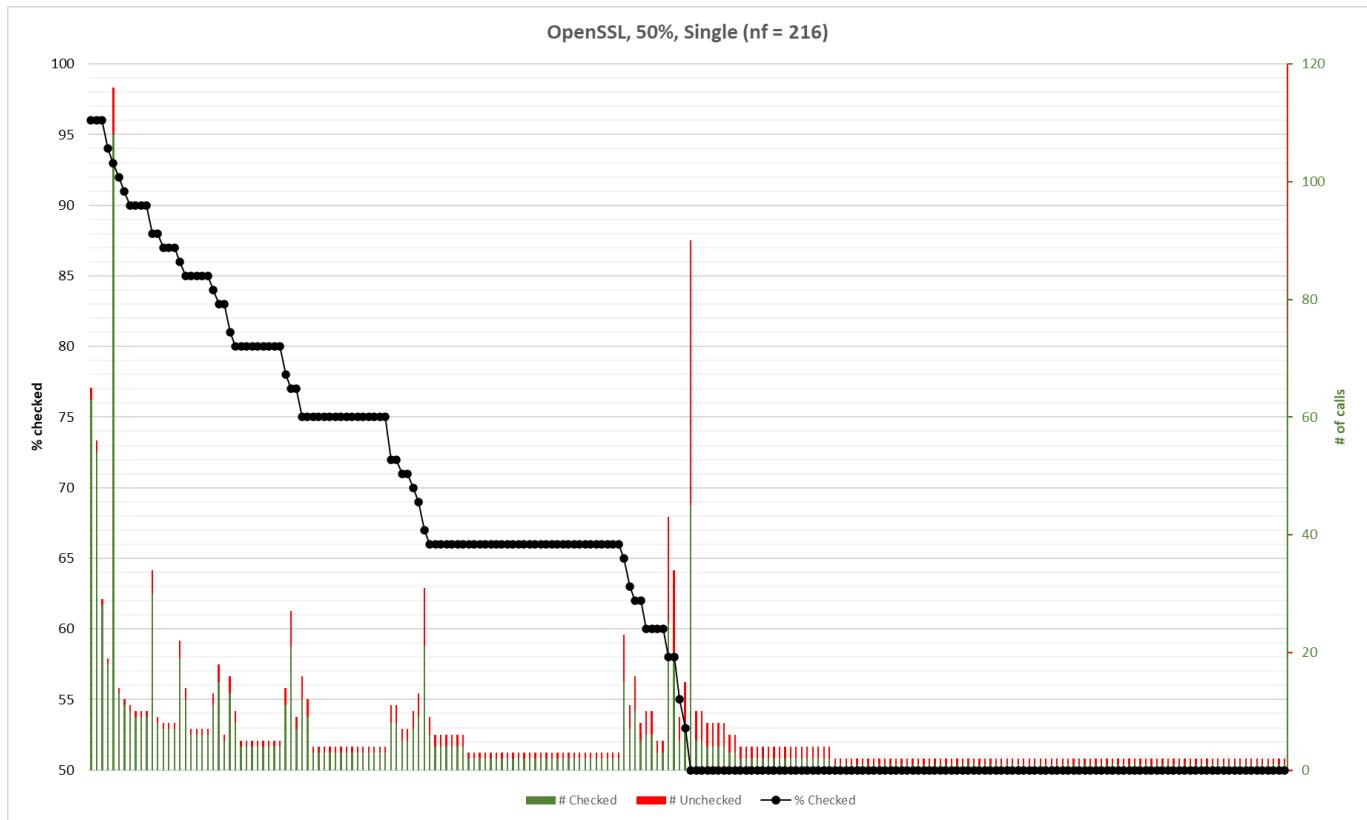
MongoDB



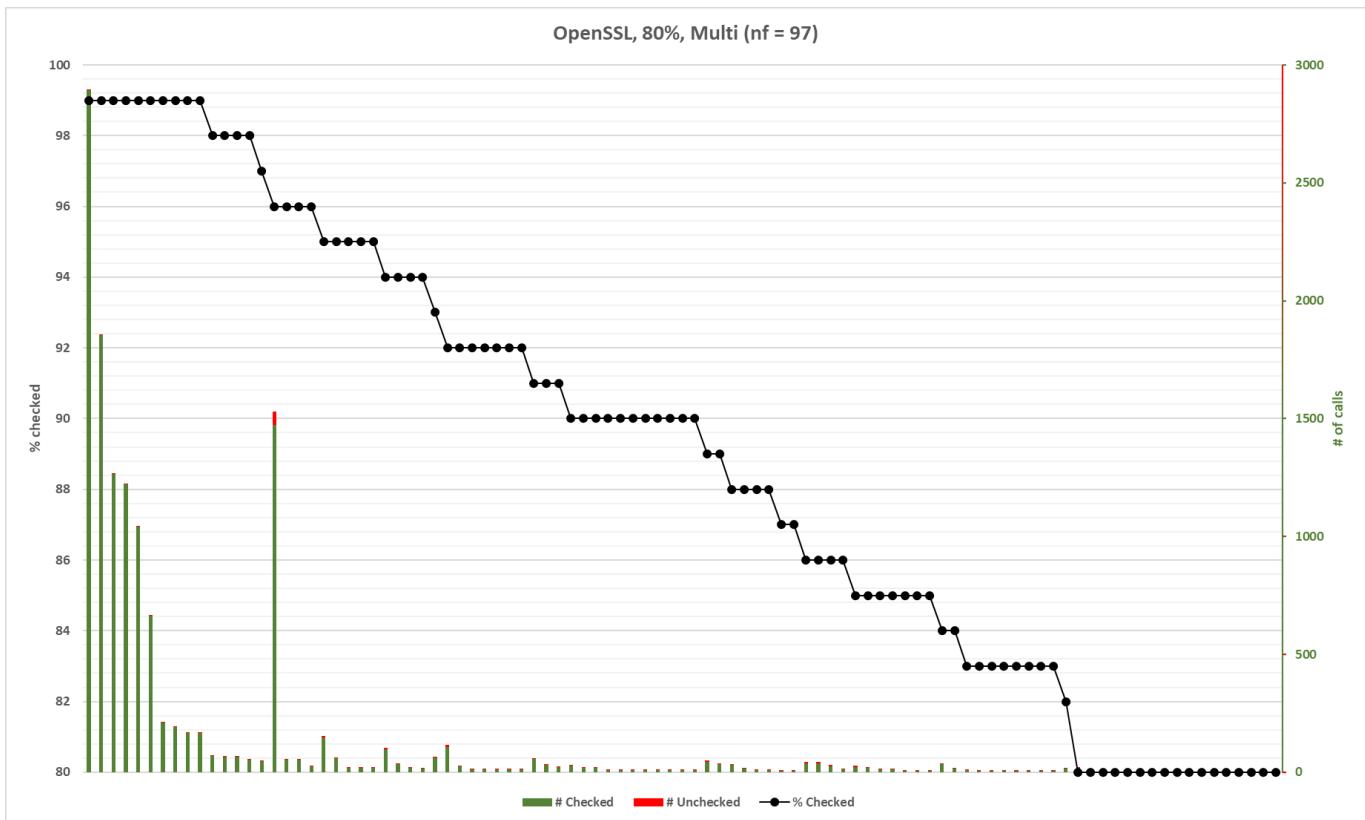
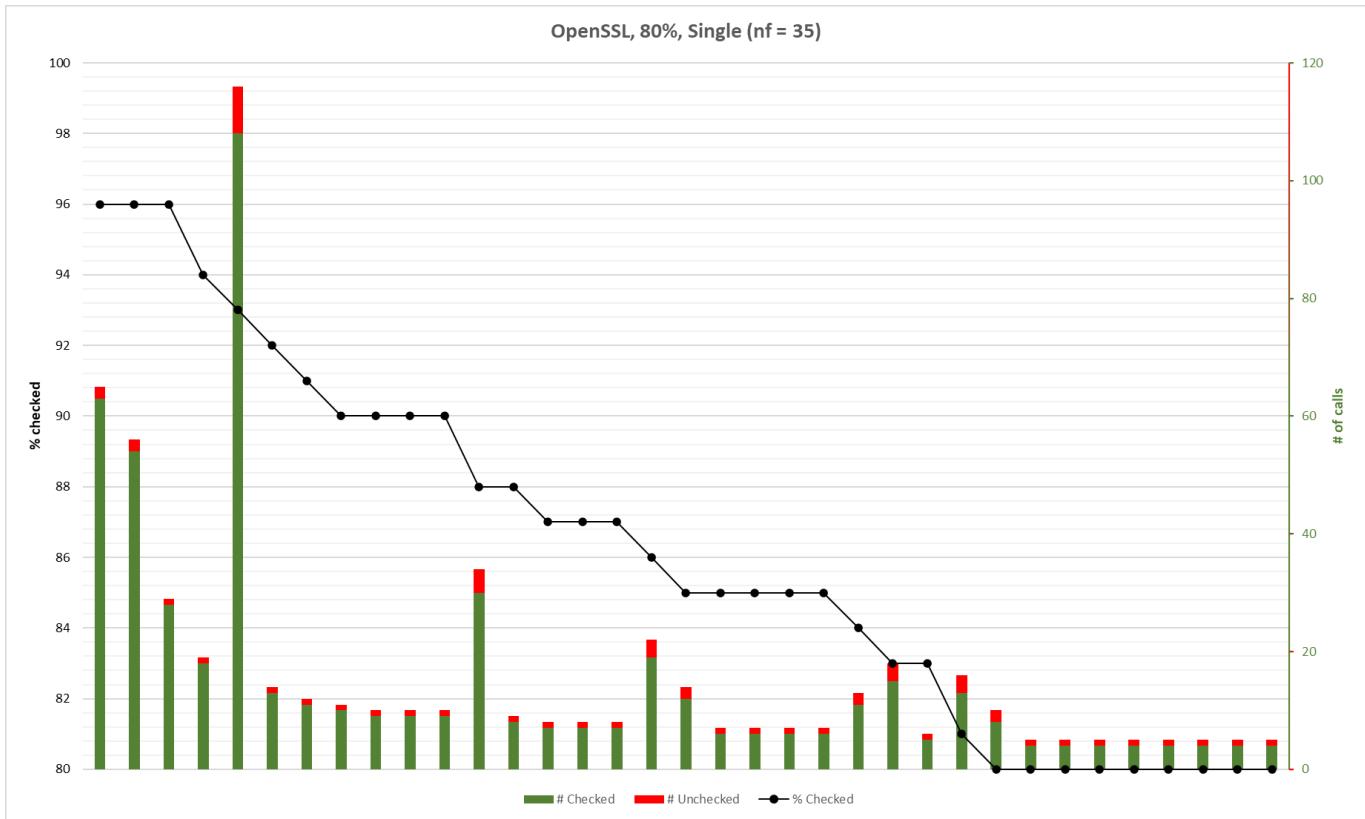
A. Project results



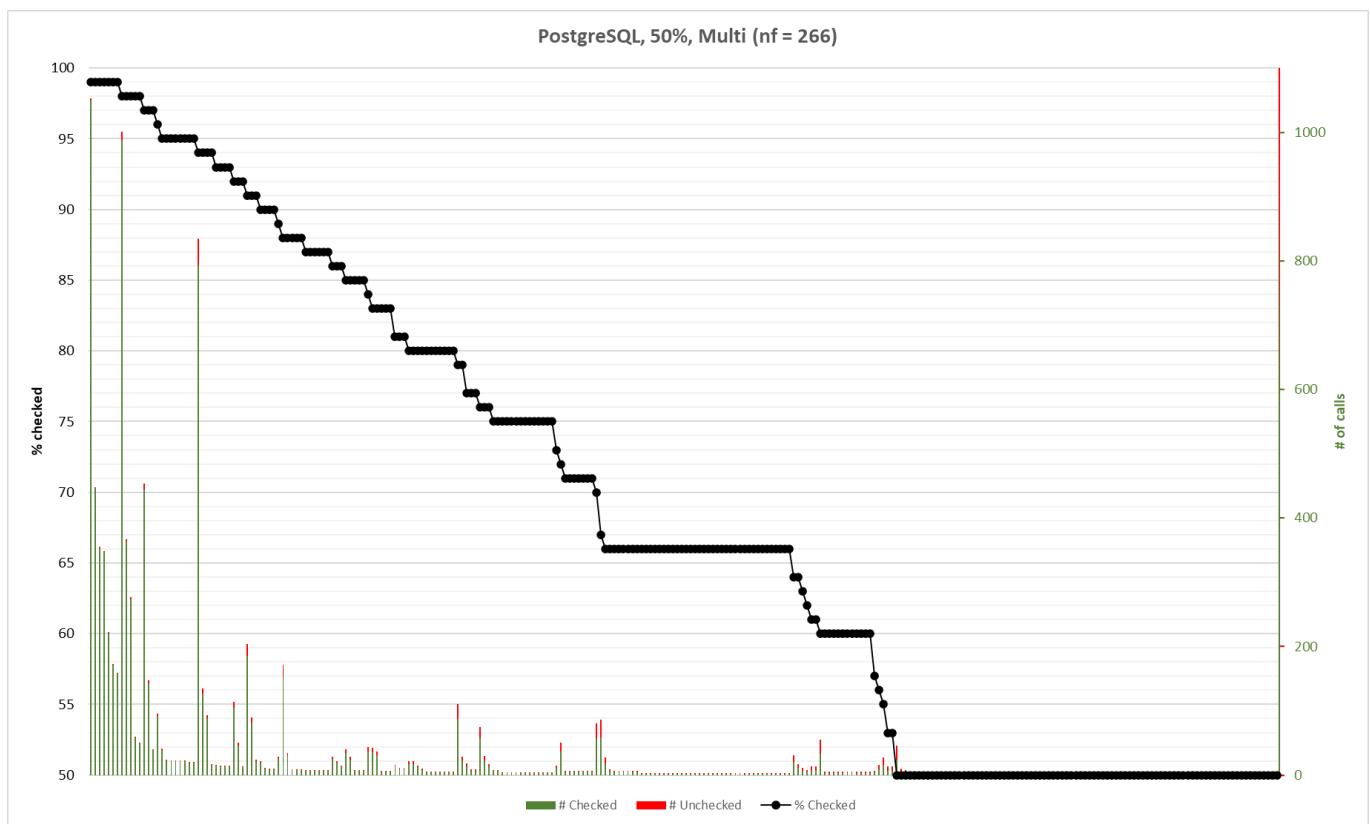
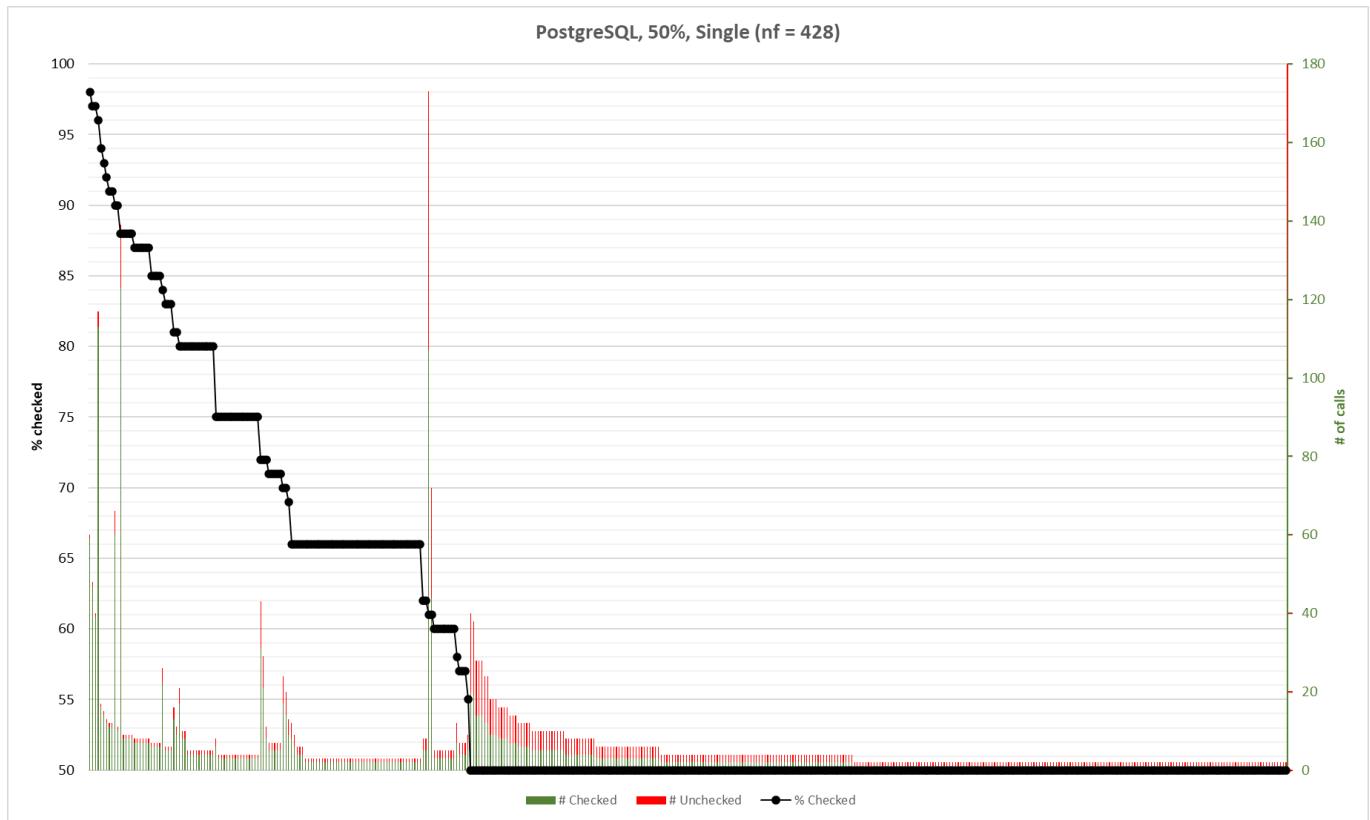
OpenSSL



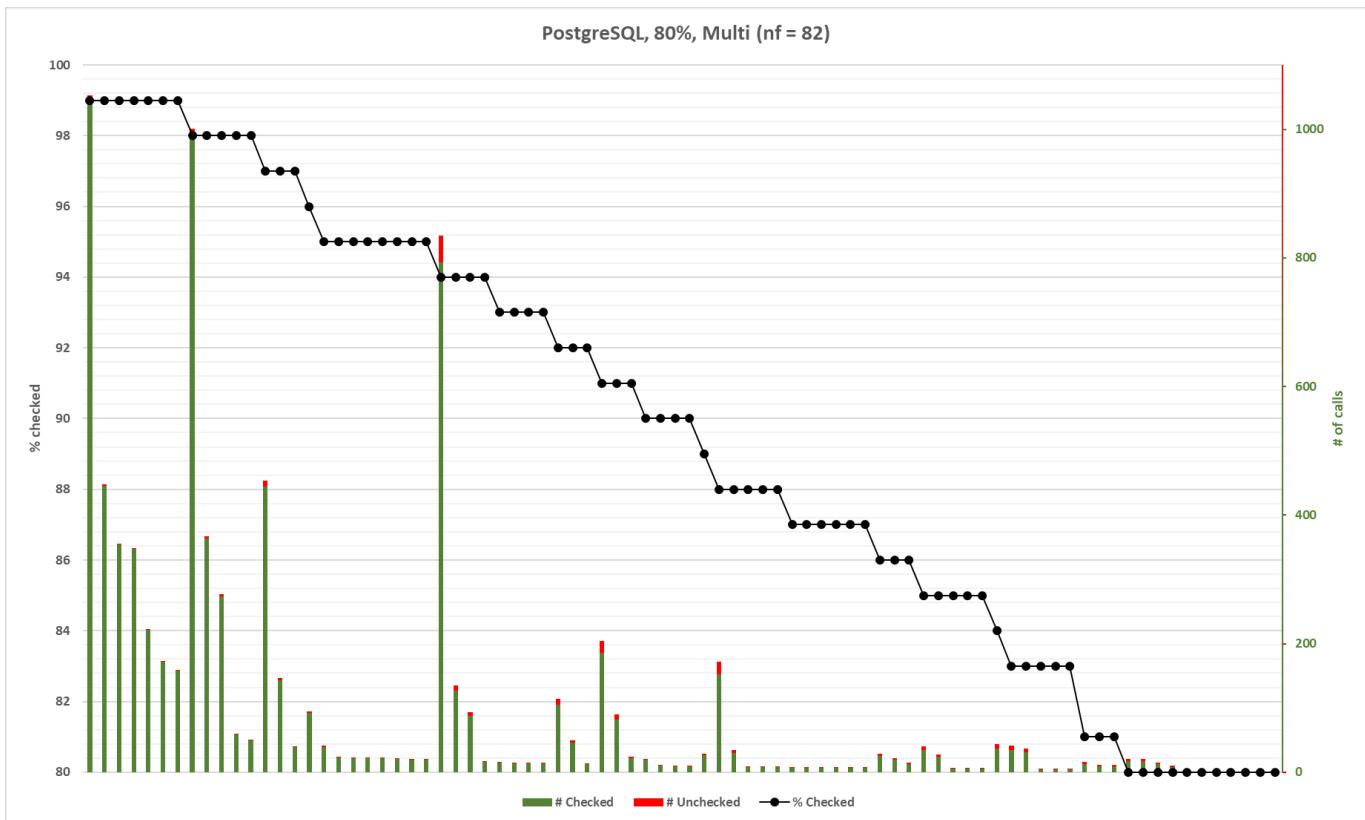
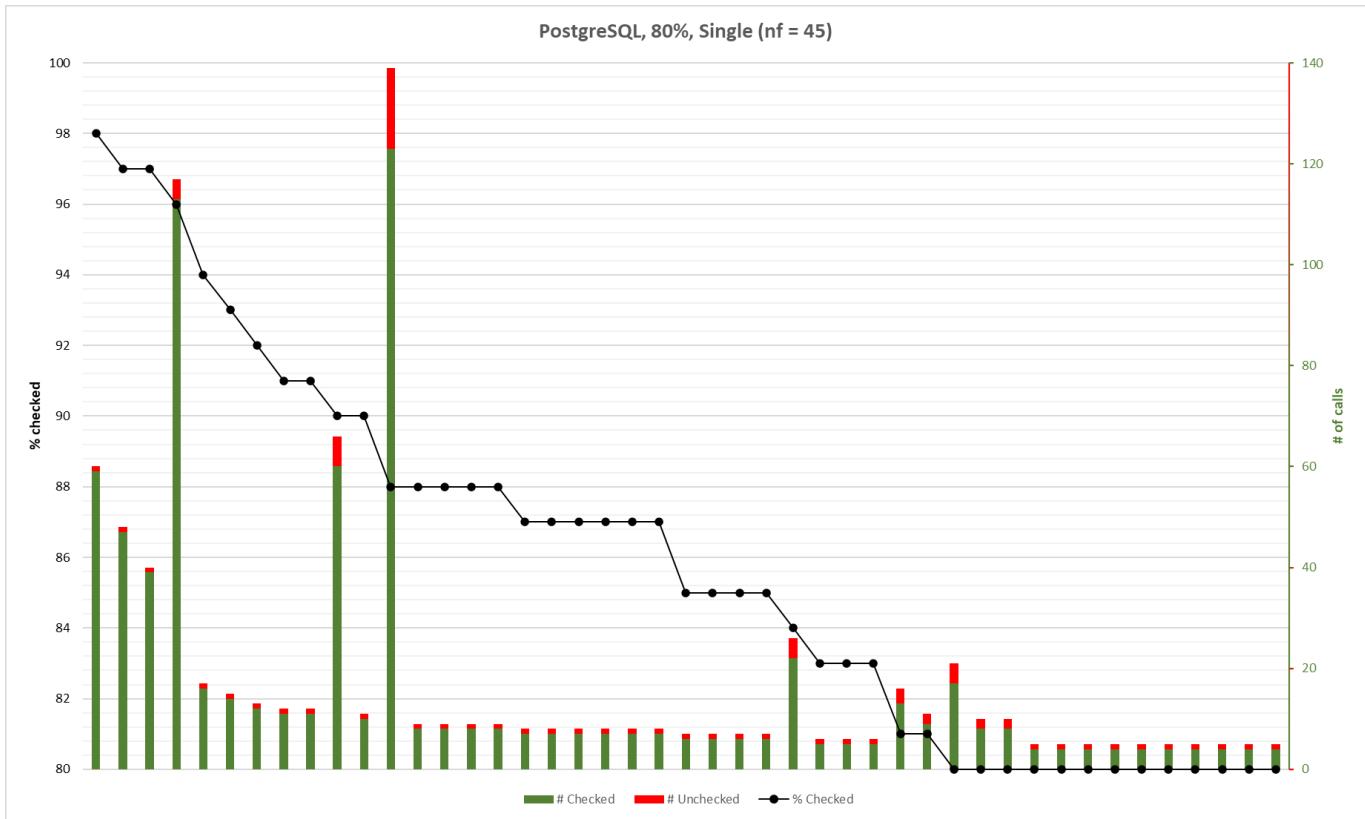
A. Project results



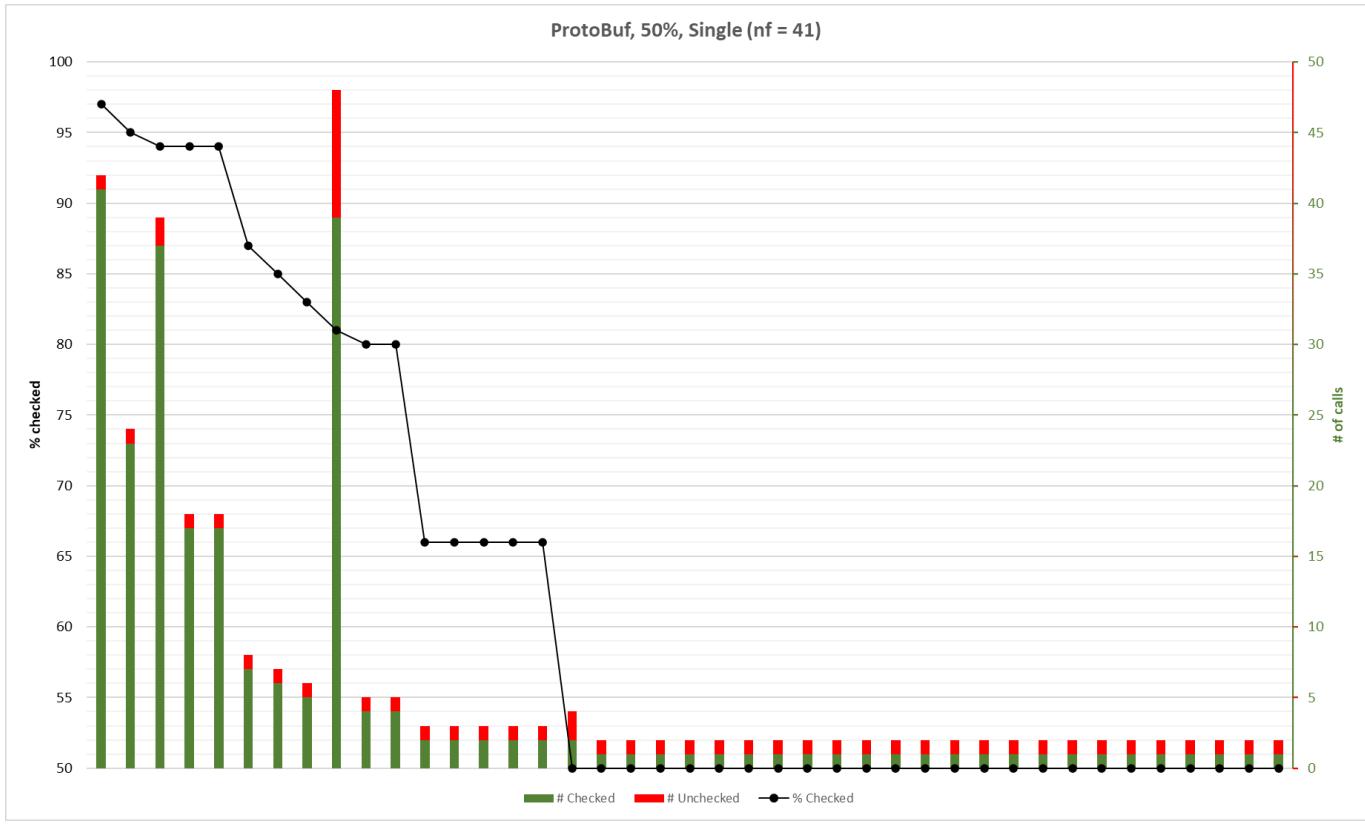
PostgreSQL



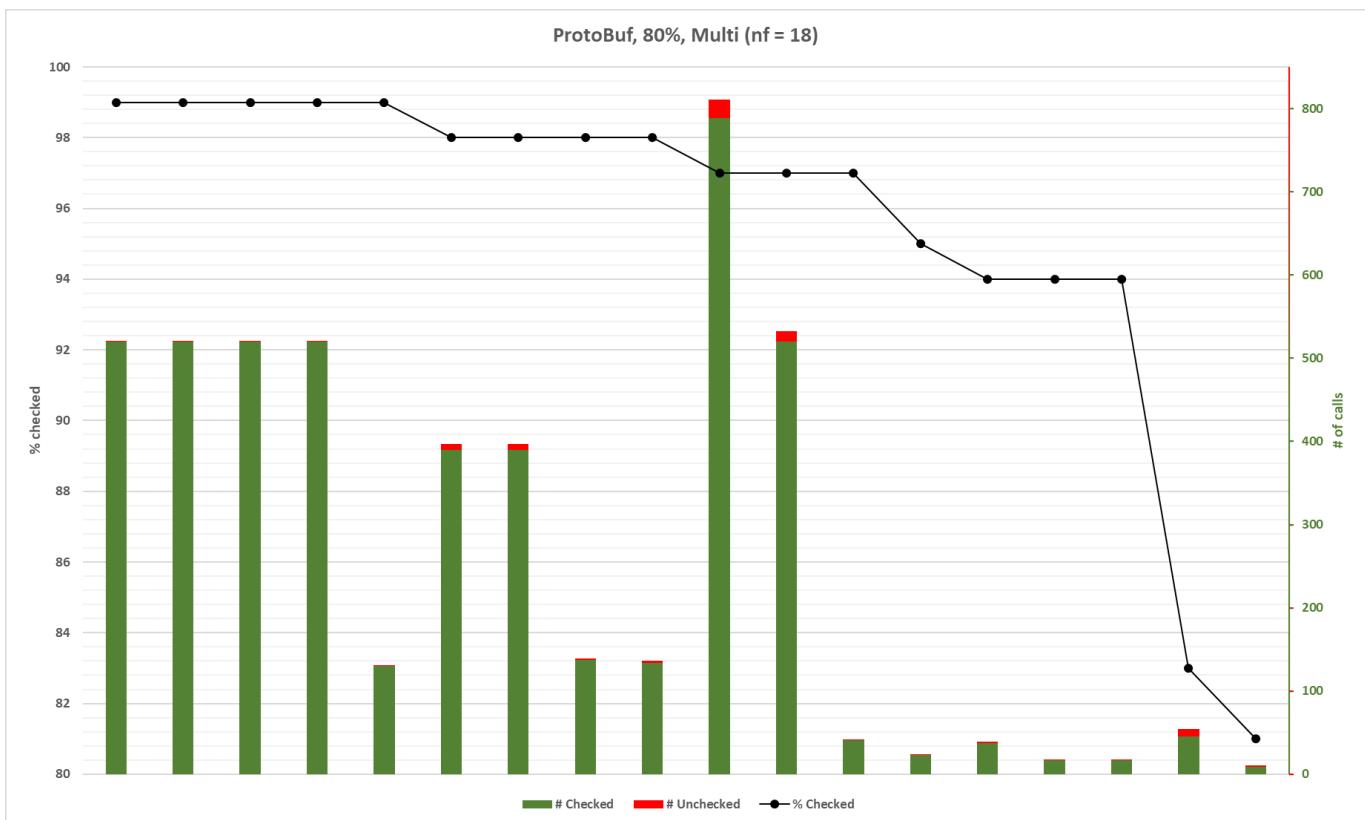
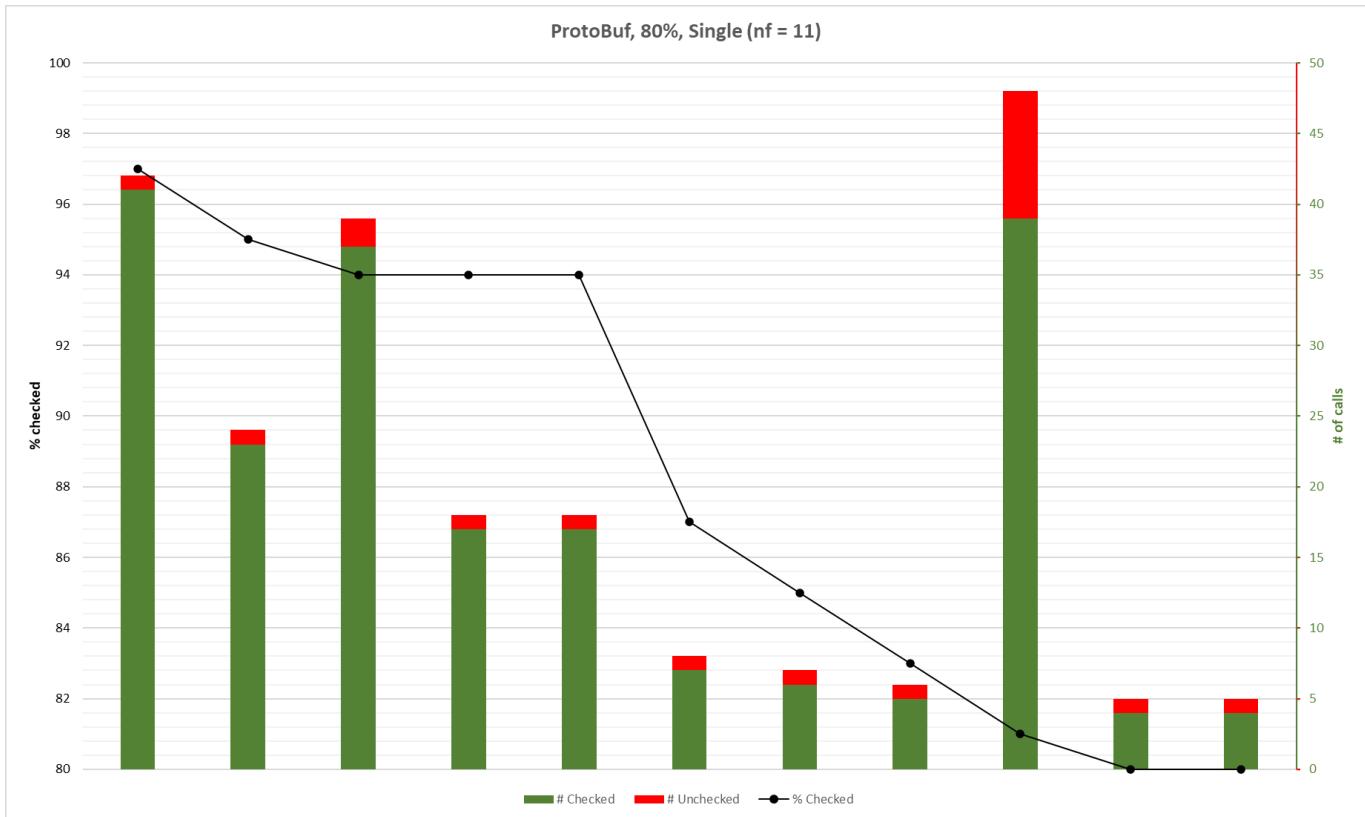
A. Project results



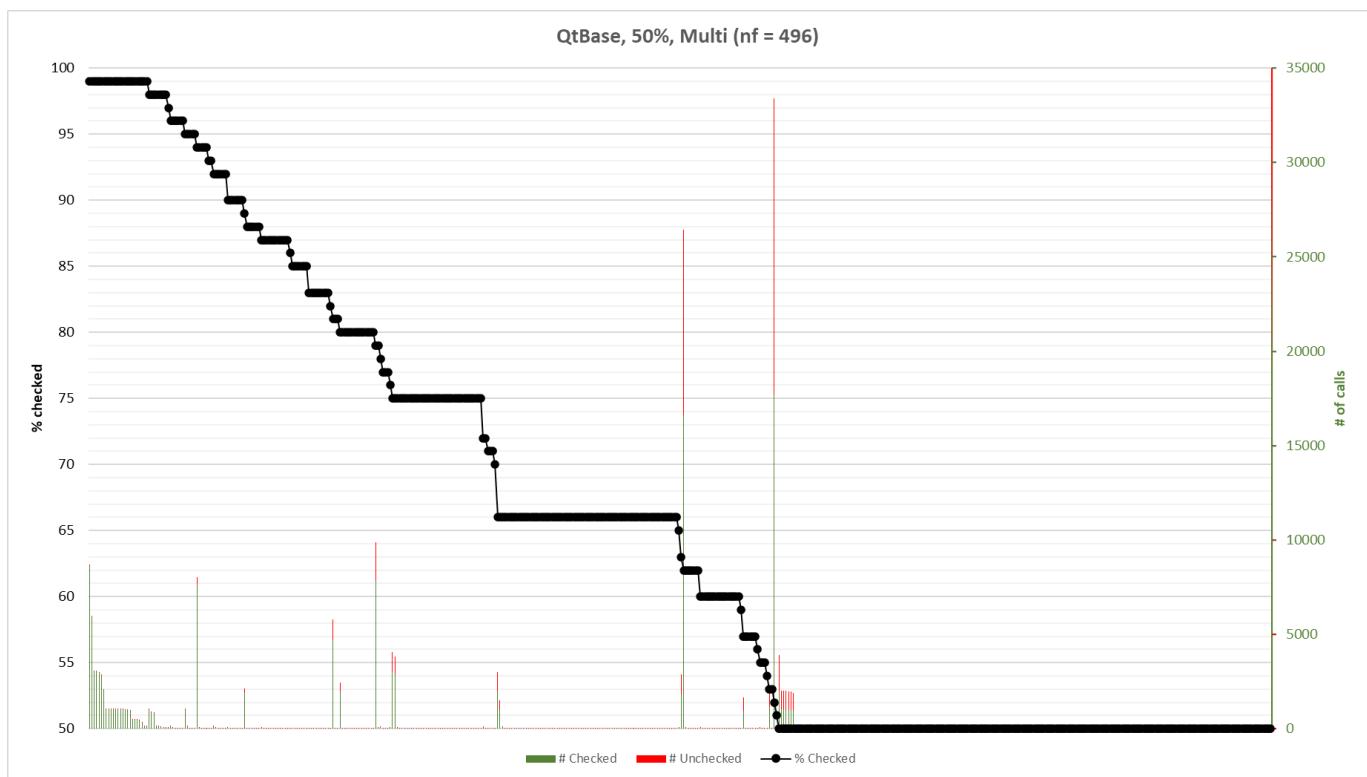
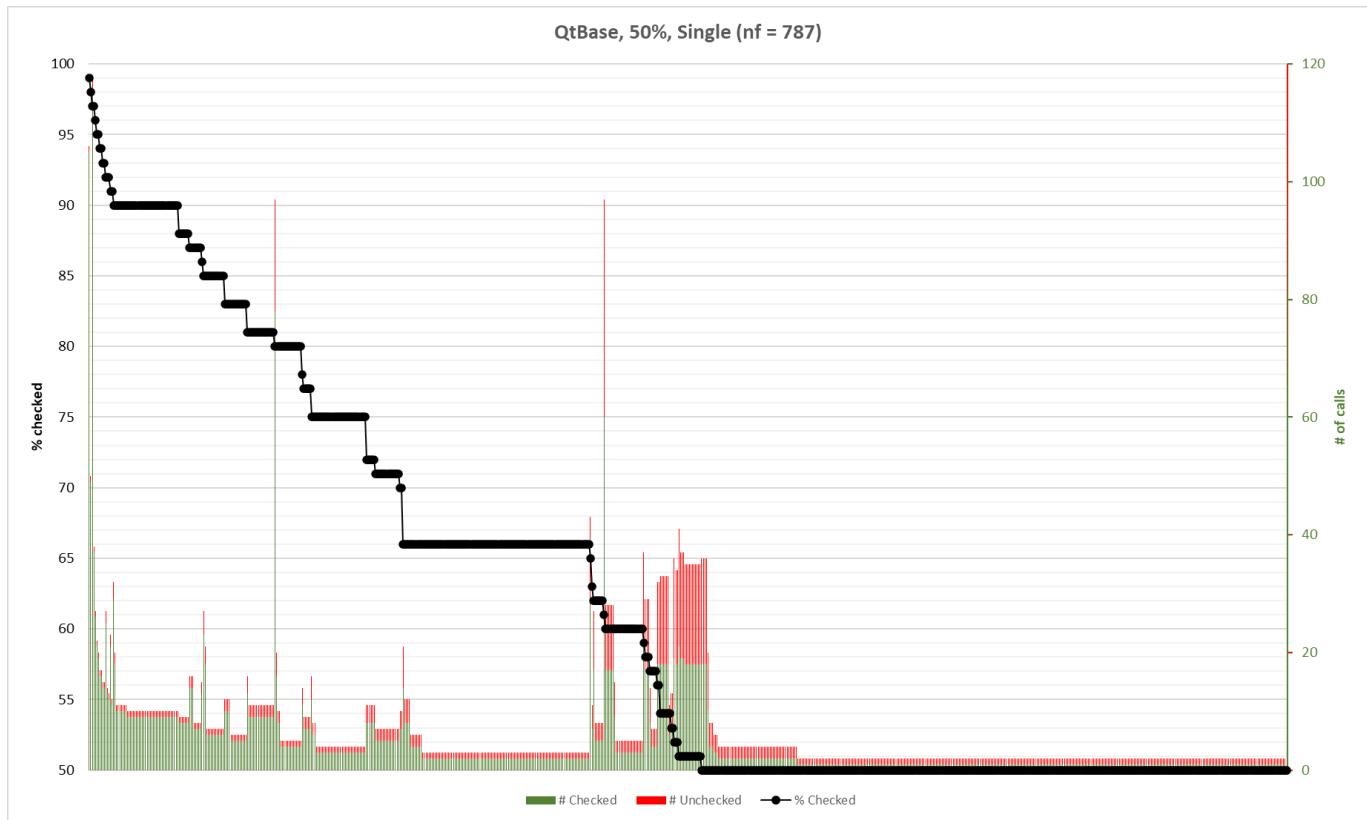
Protocol Buffers



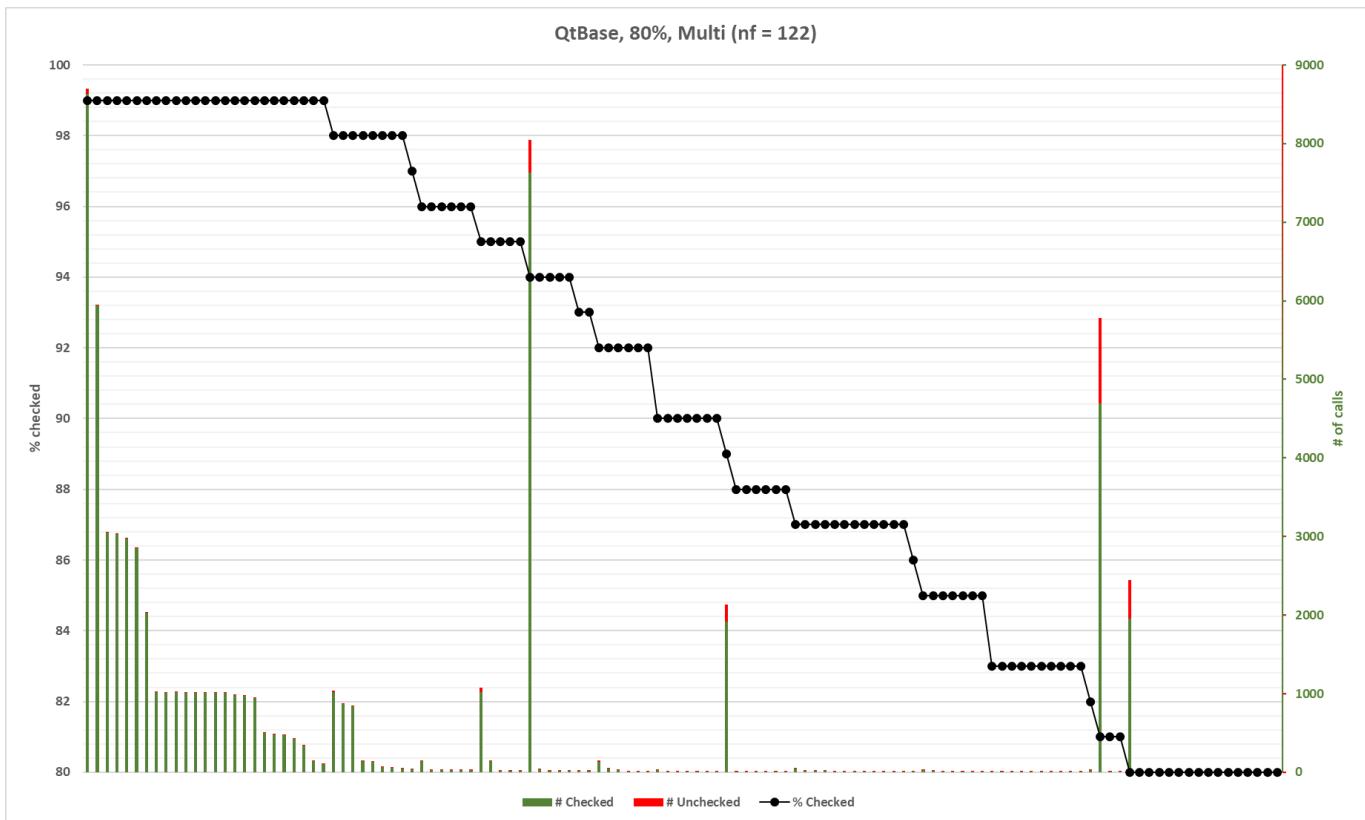
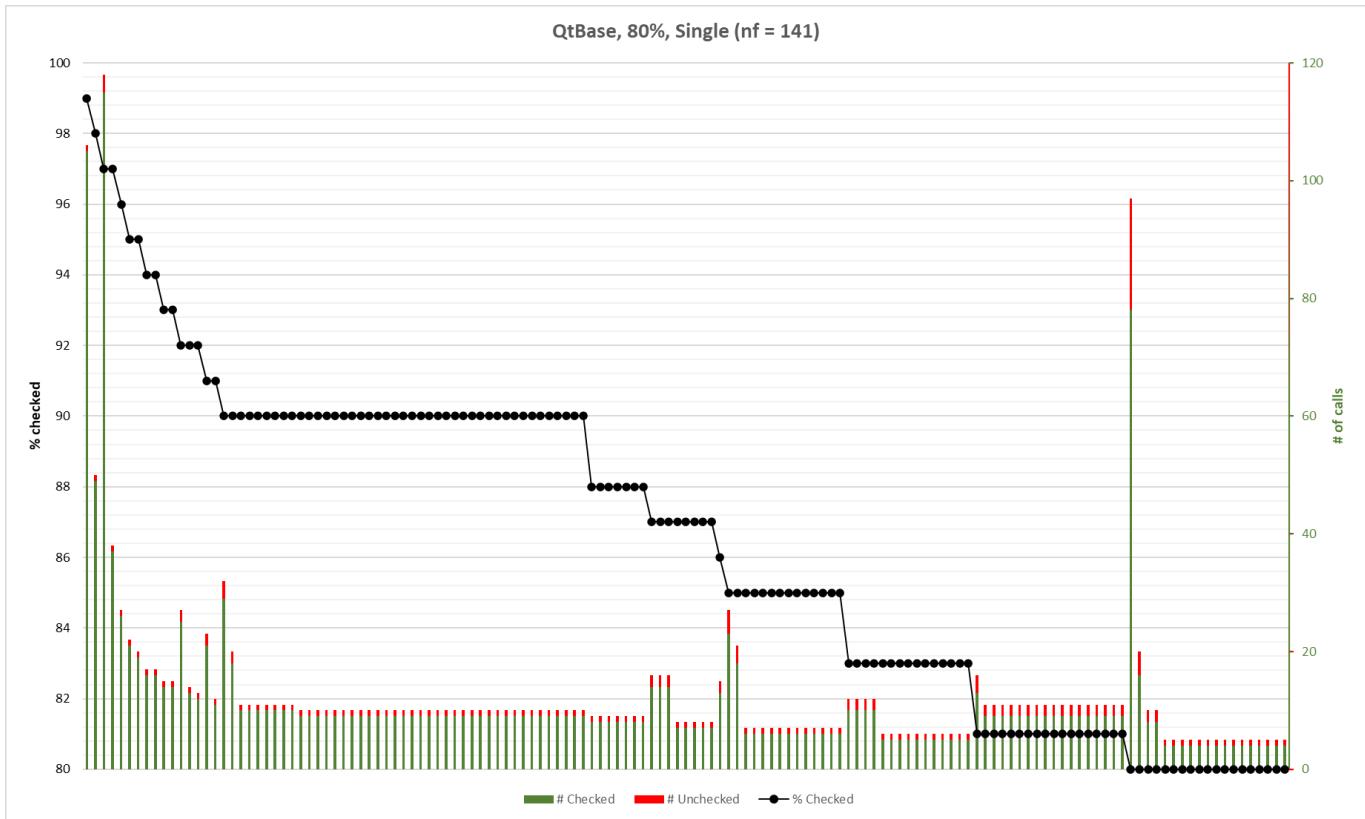
A. Project results



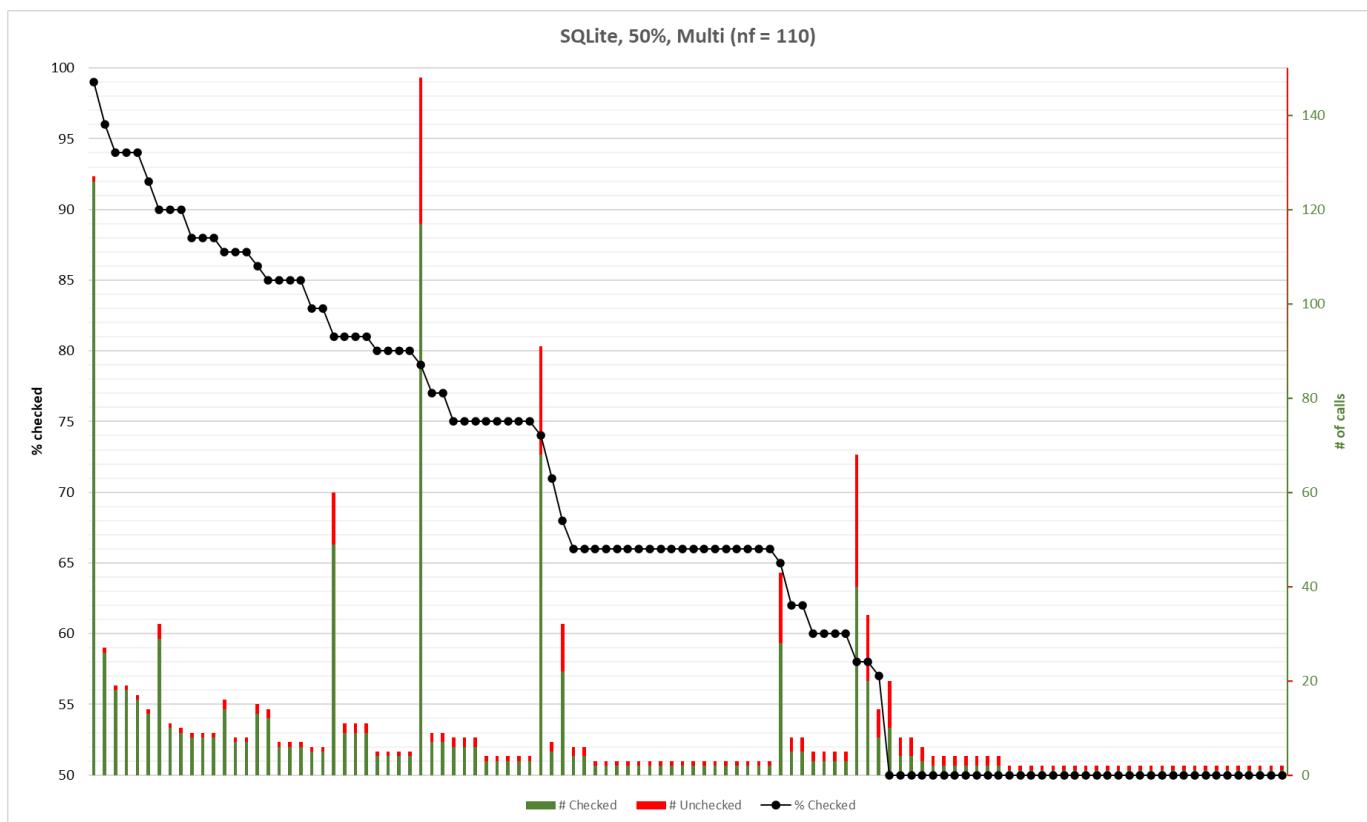
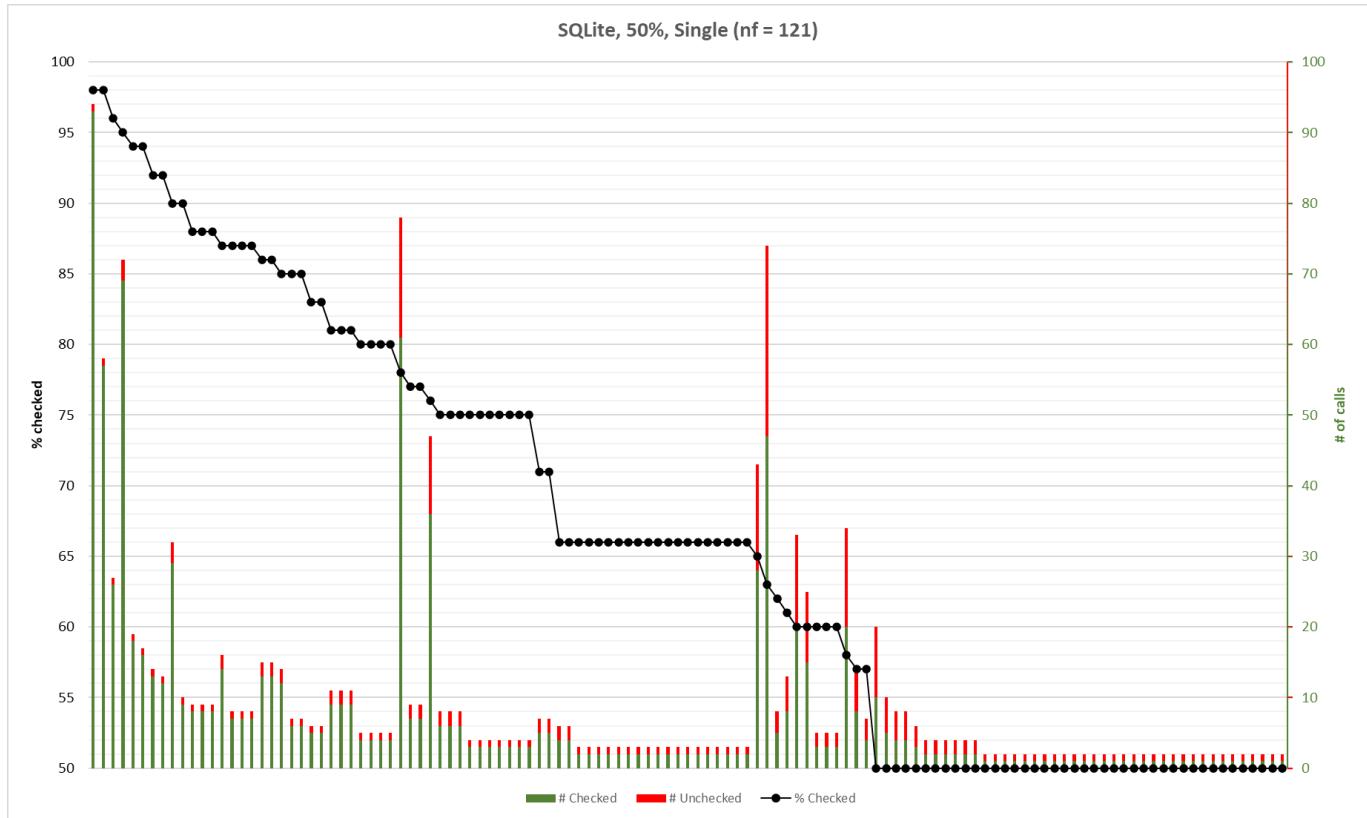
QtBase



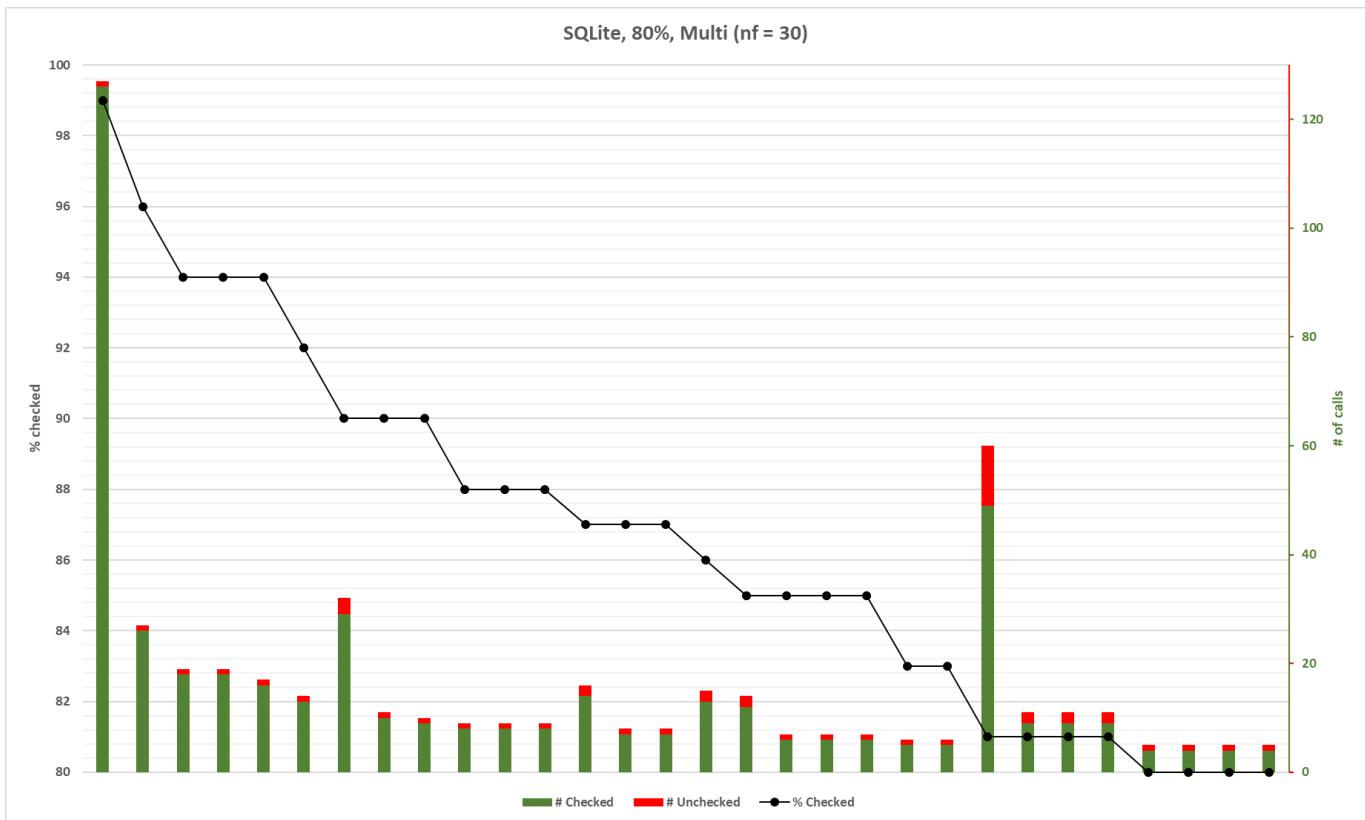
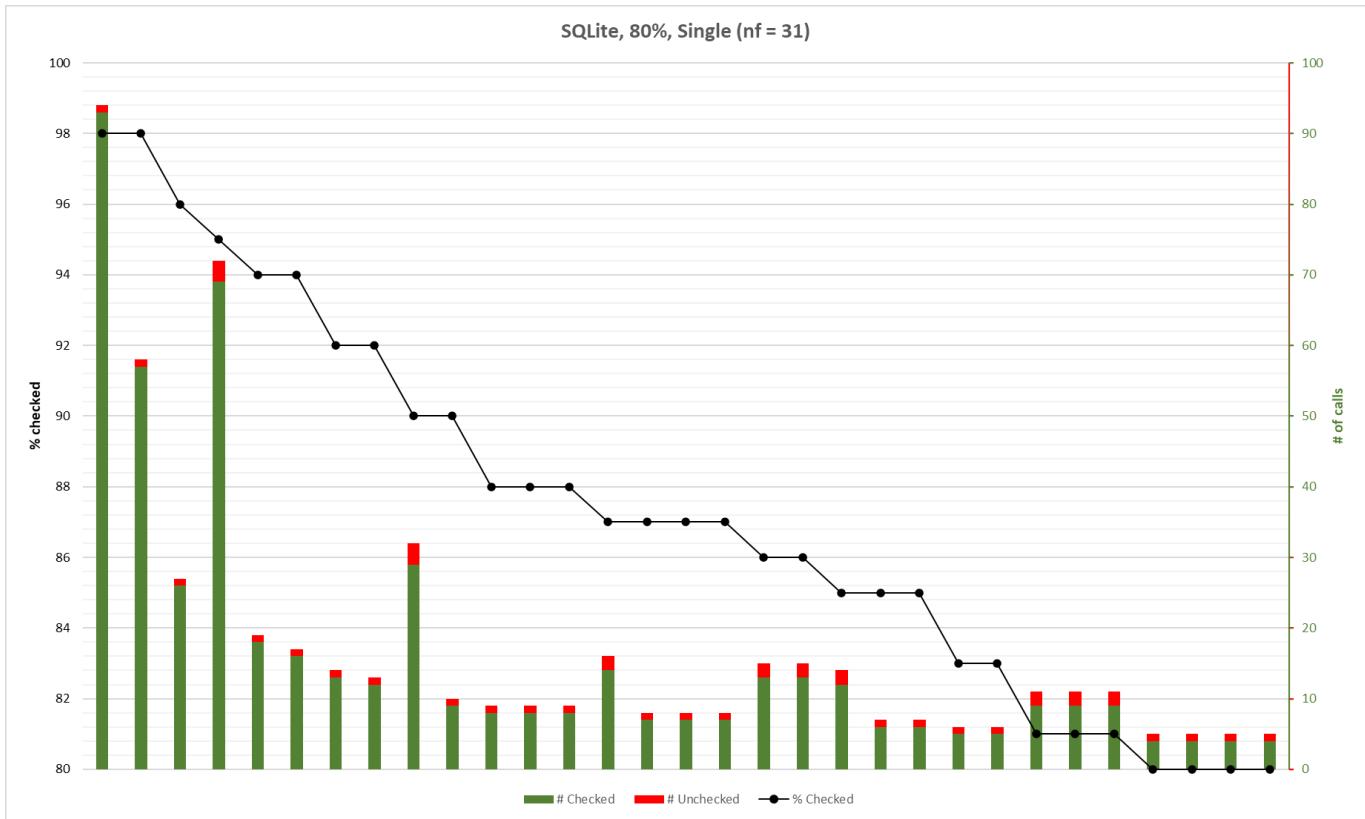
A. Project results



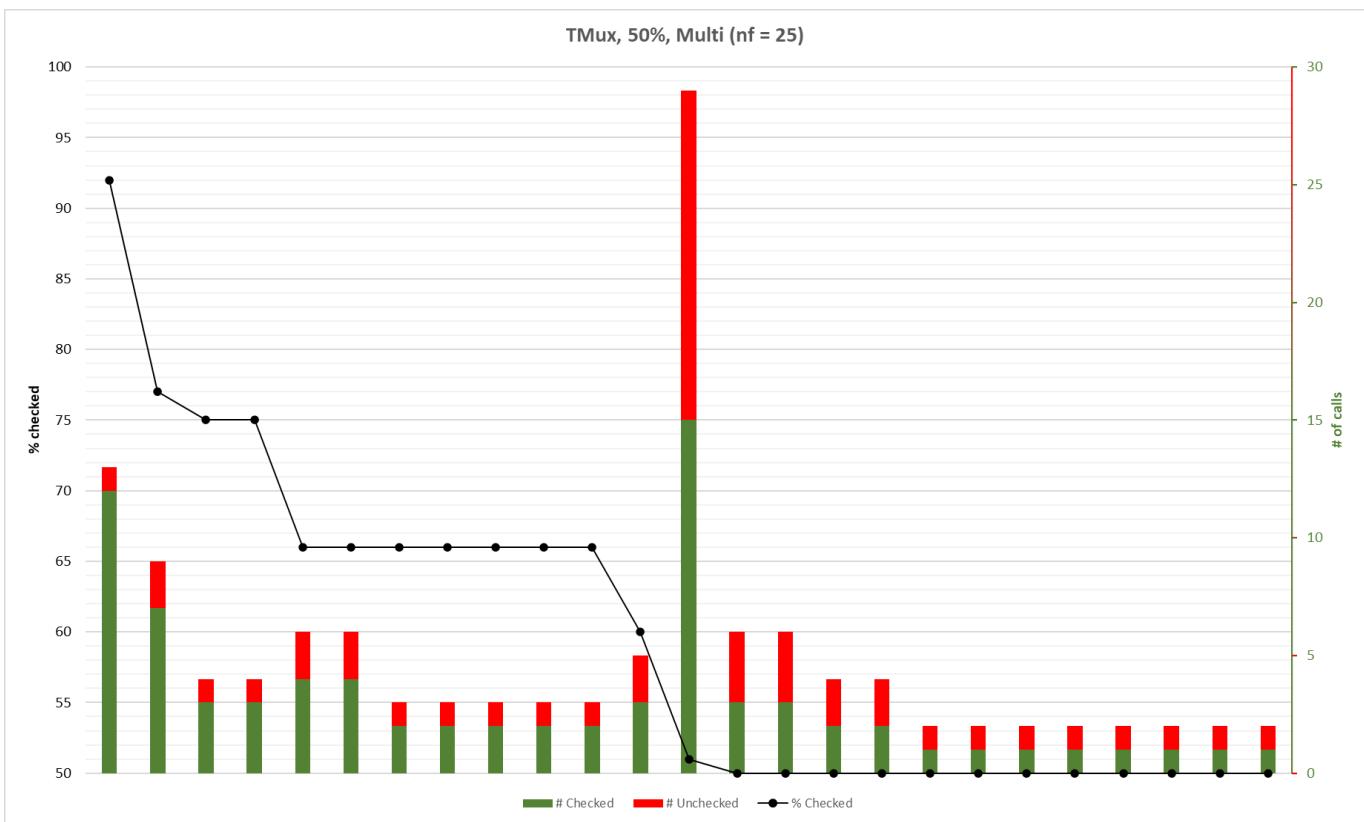
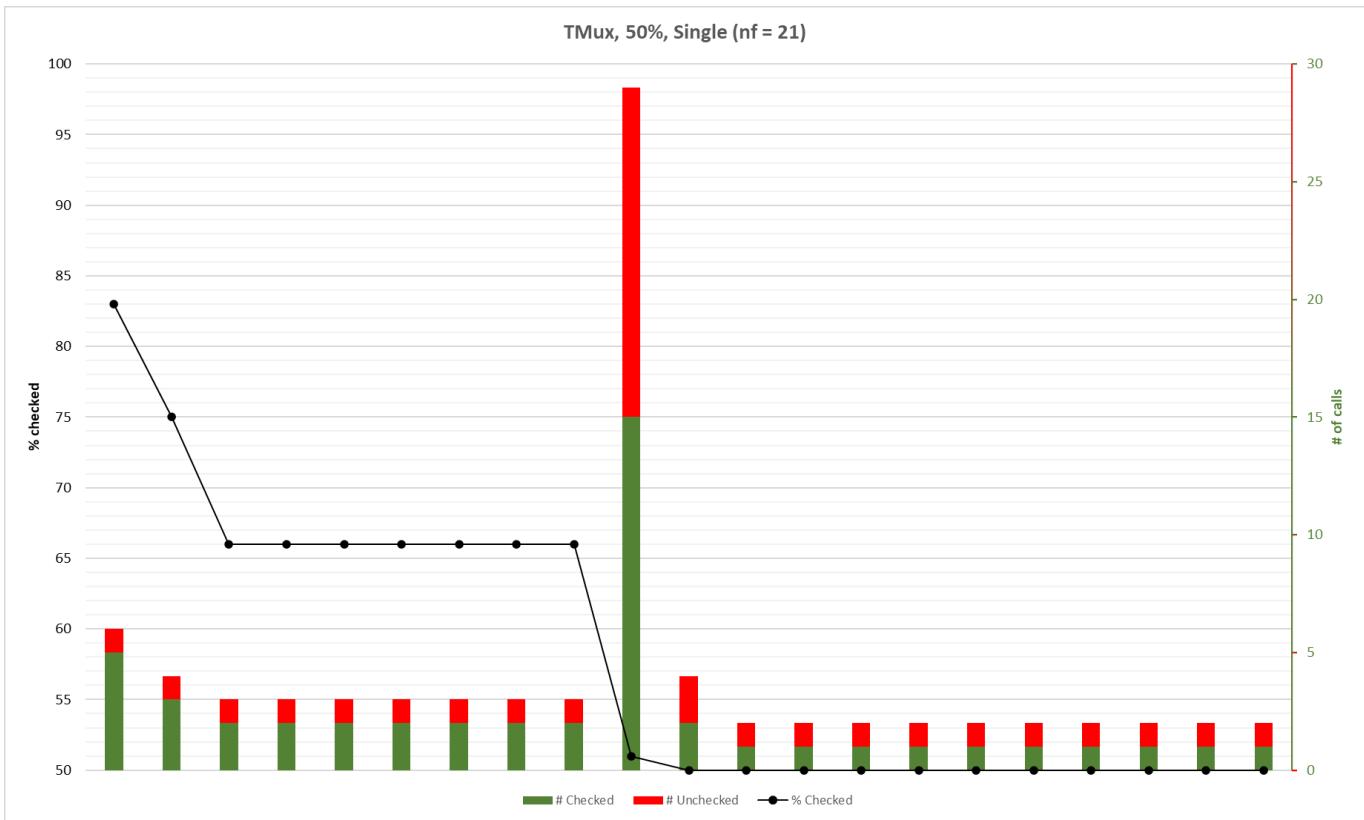
SQLite



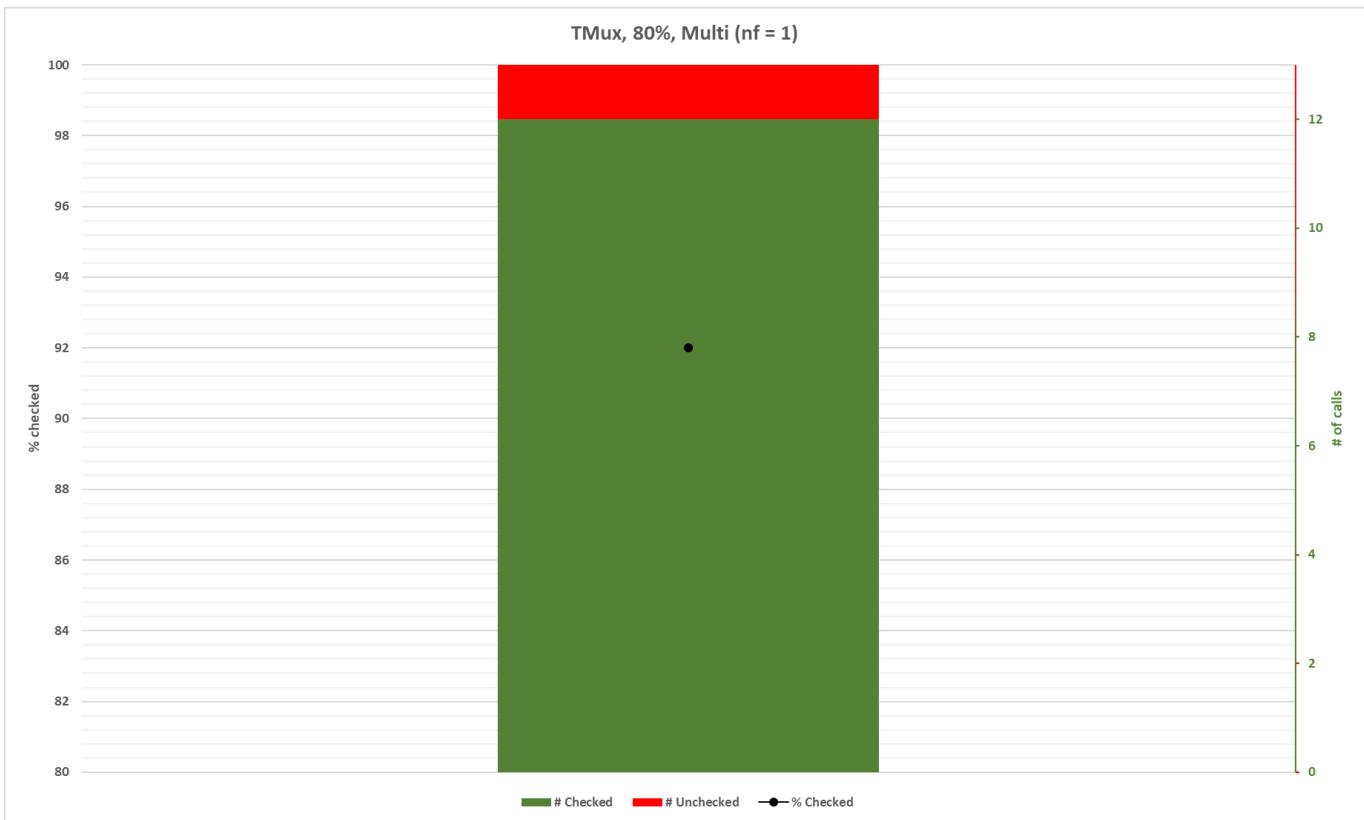
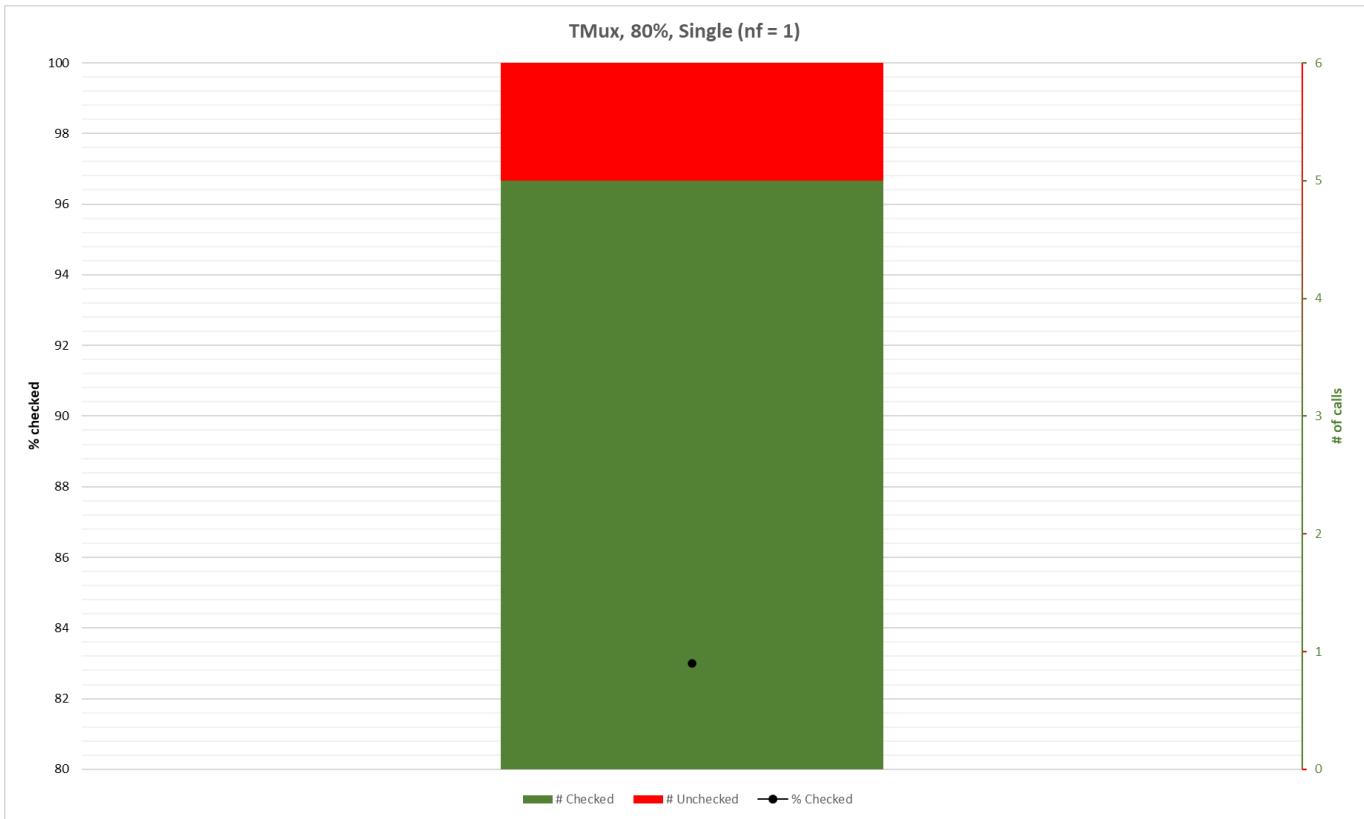
A. Project results



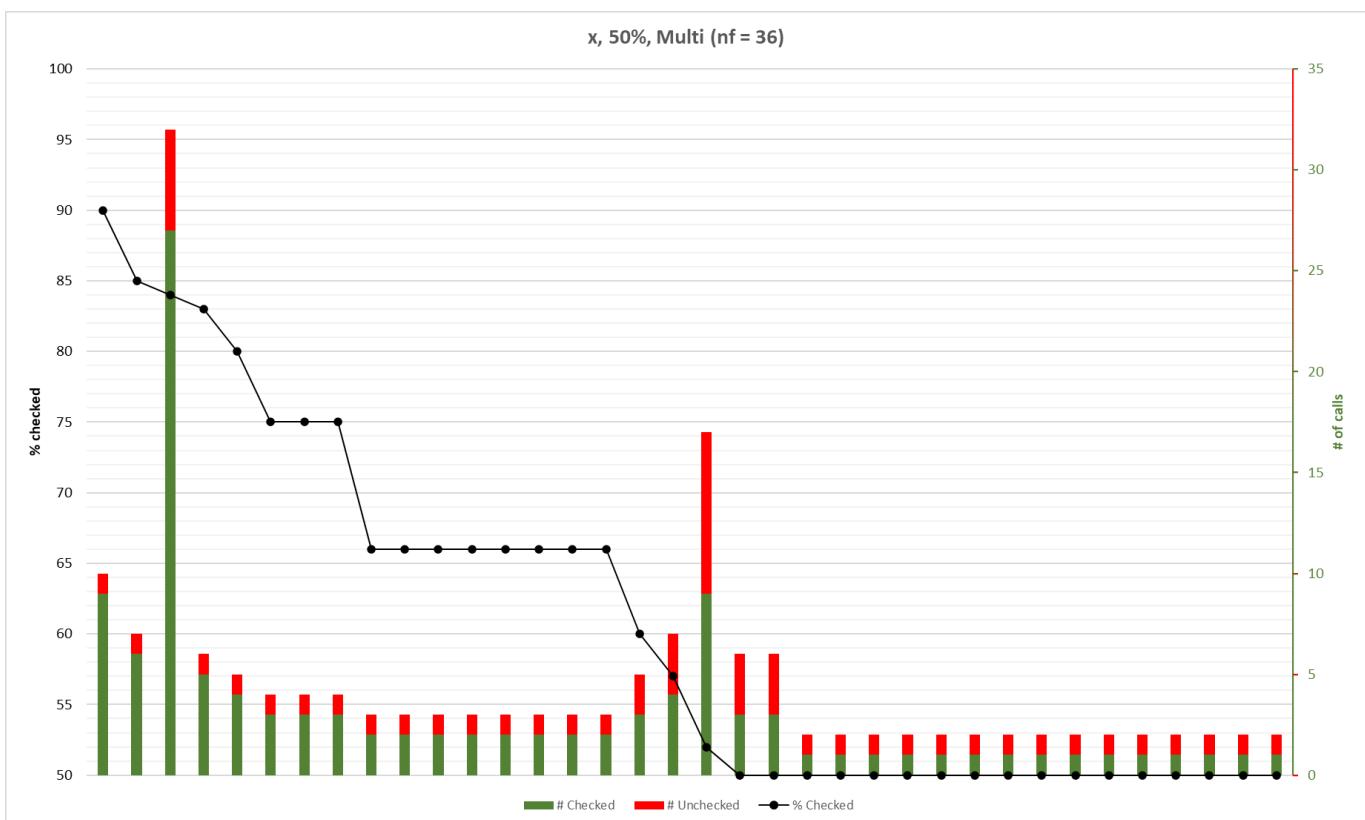
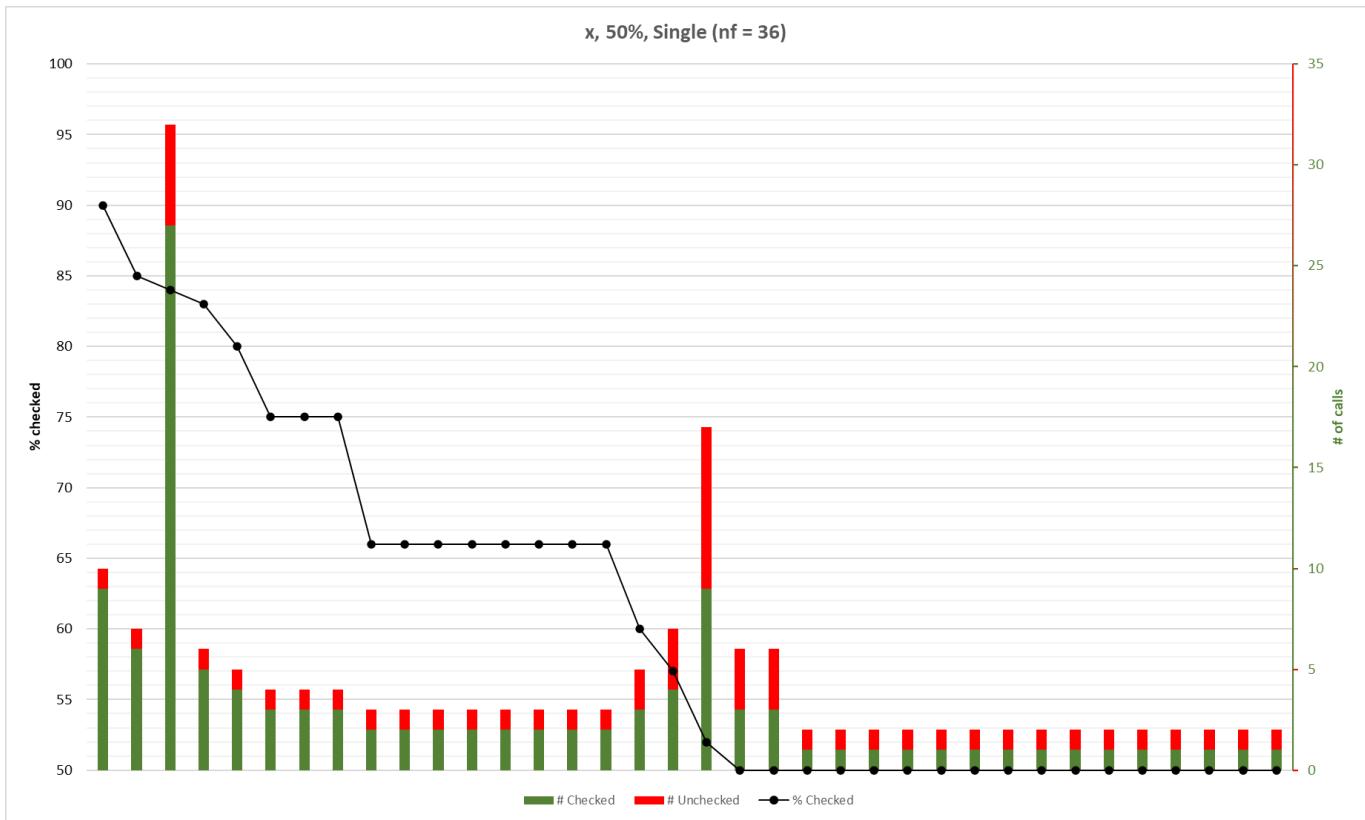
TMux



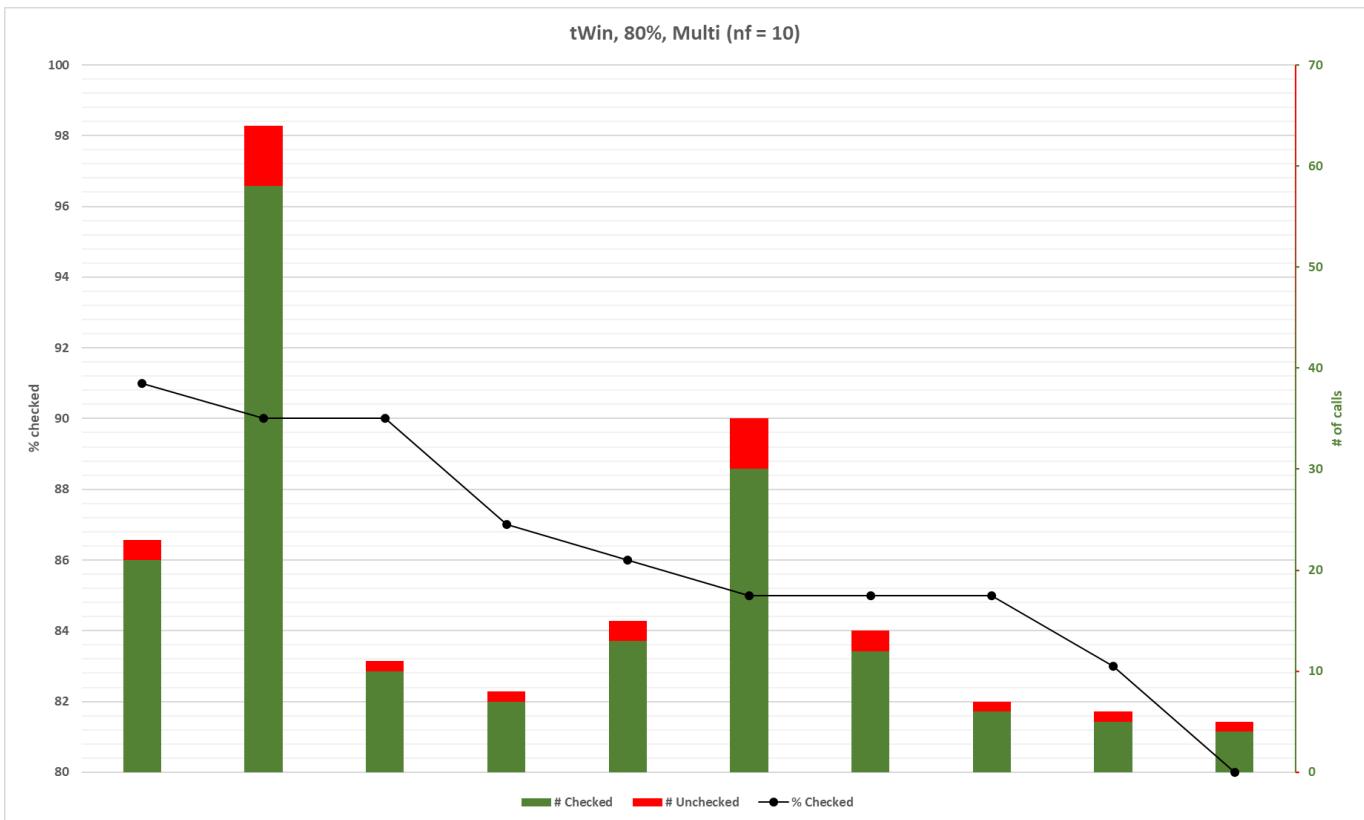
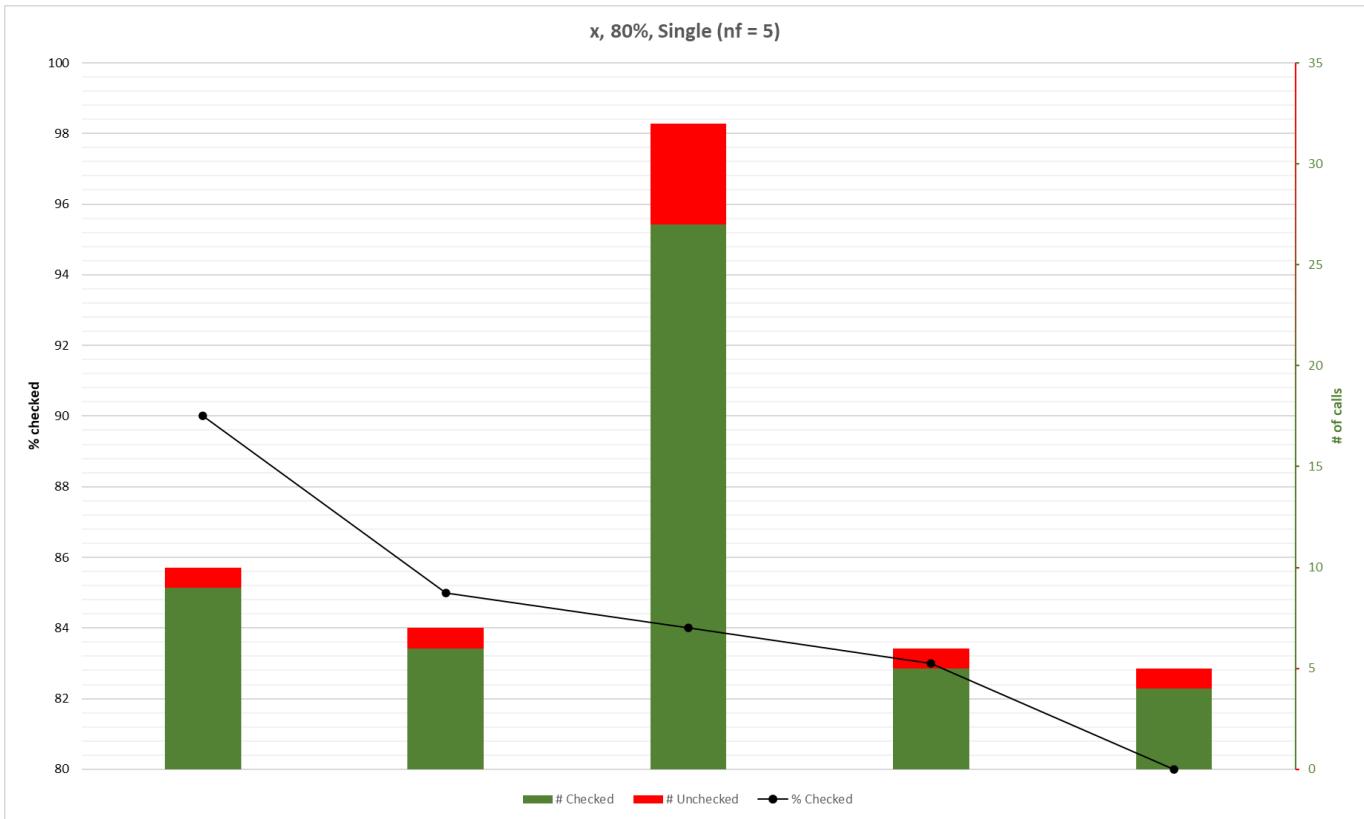
A. Project results



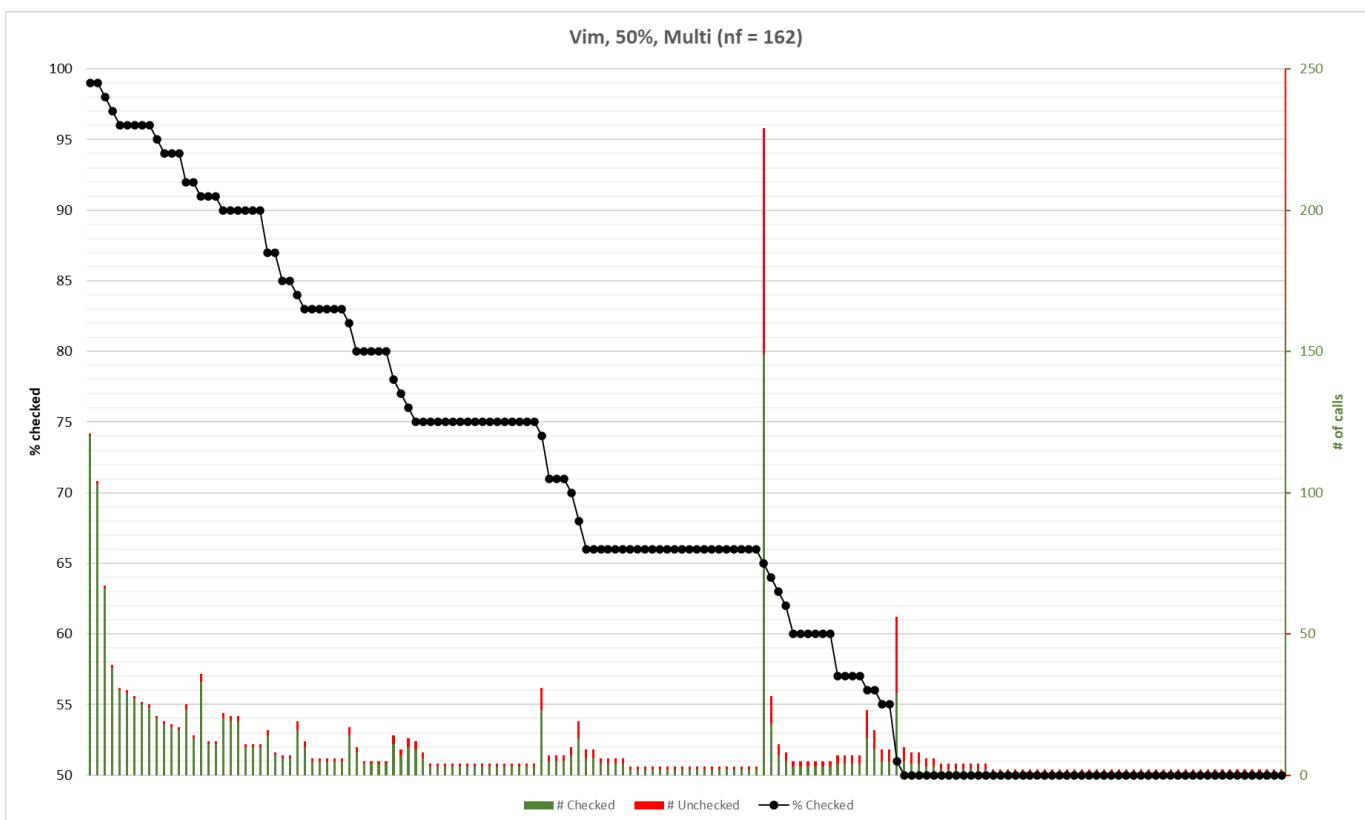
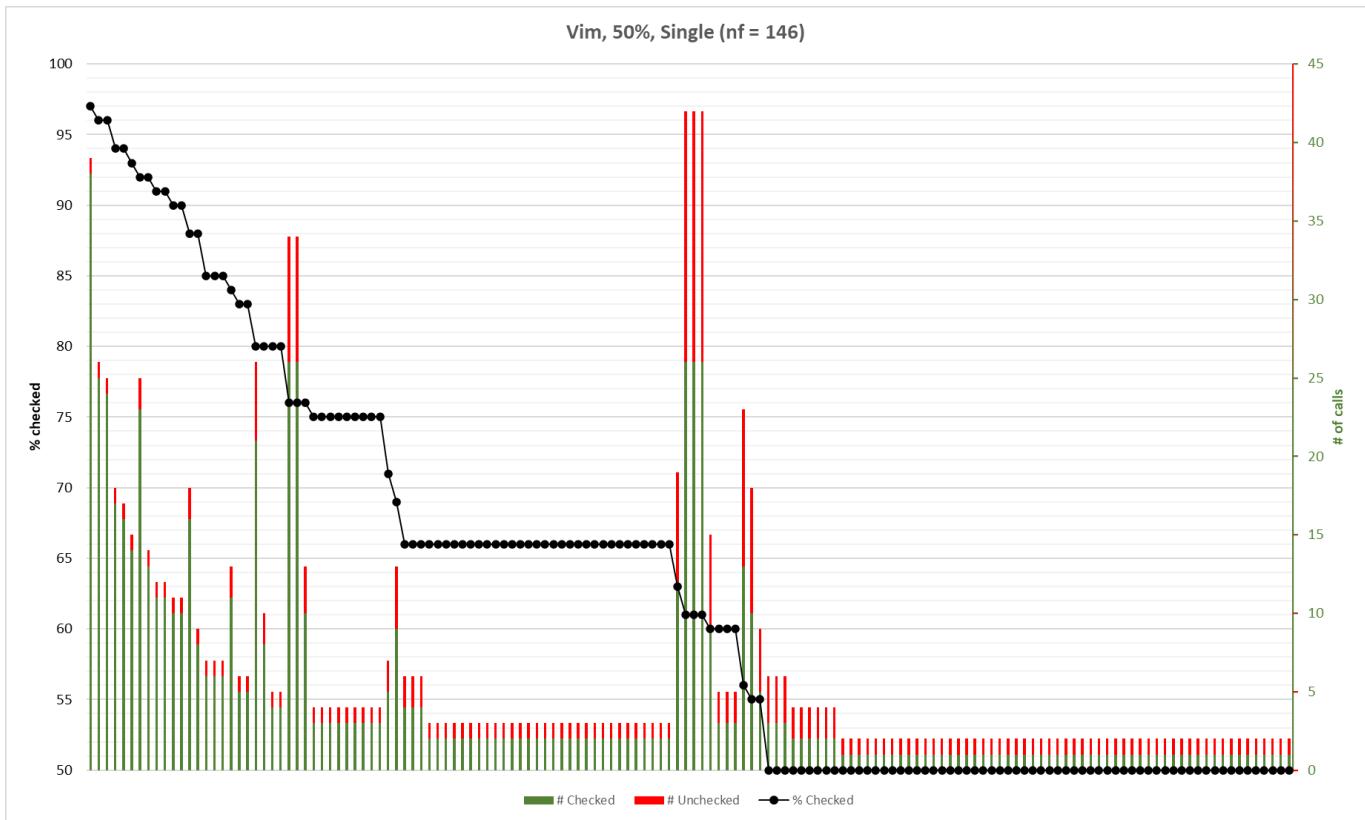
tWin



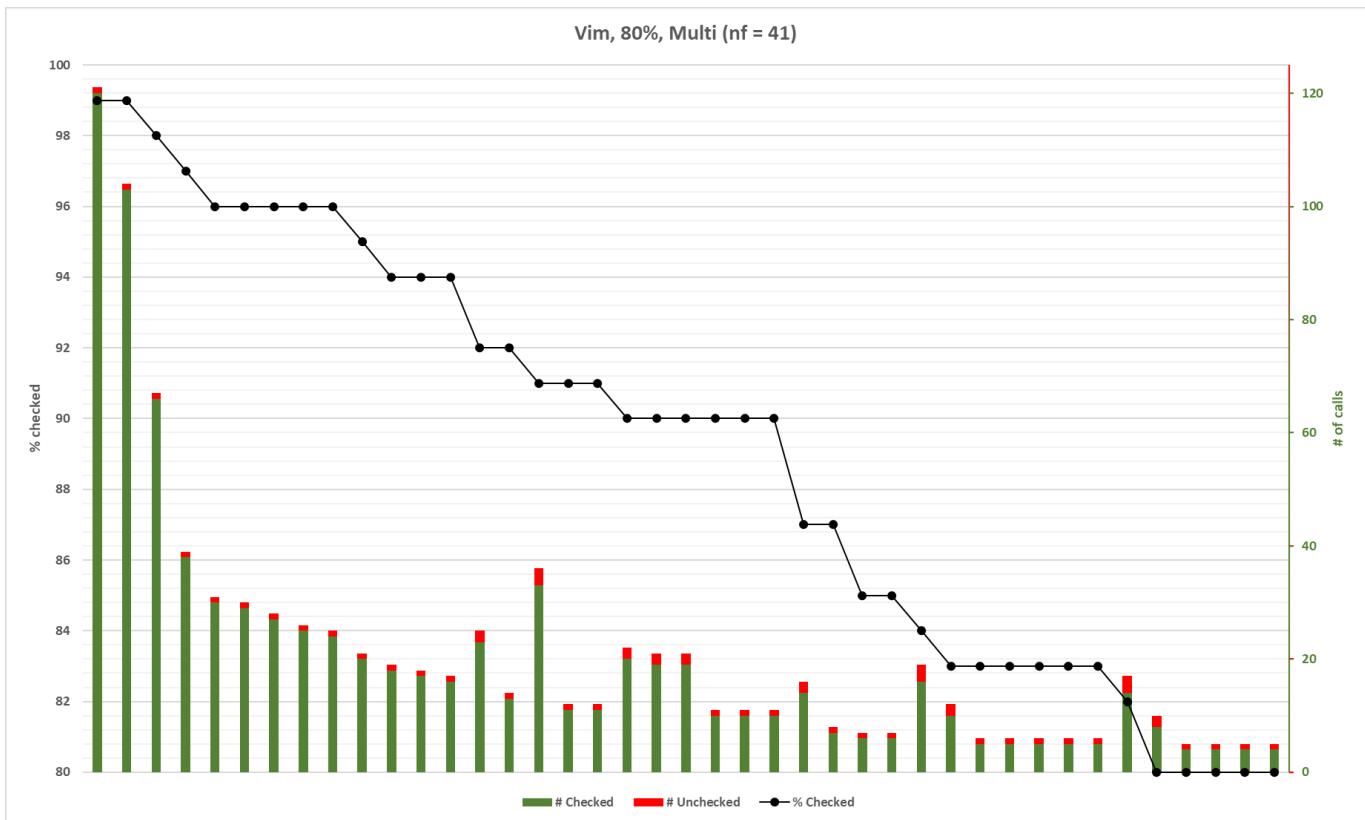
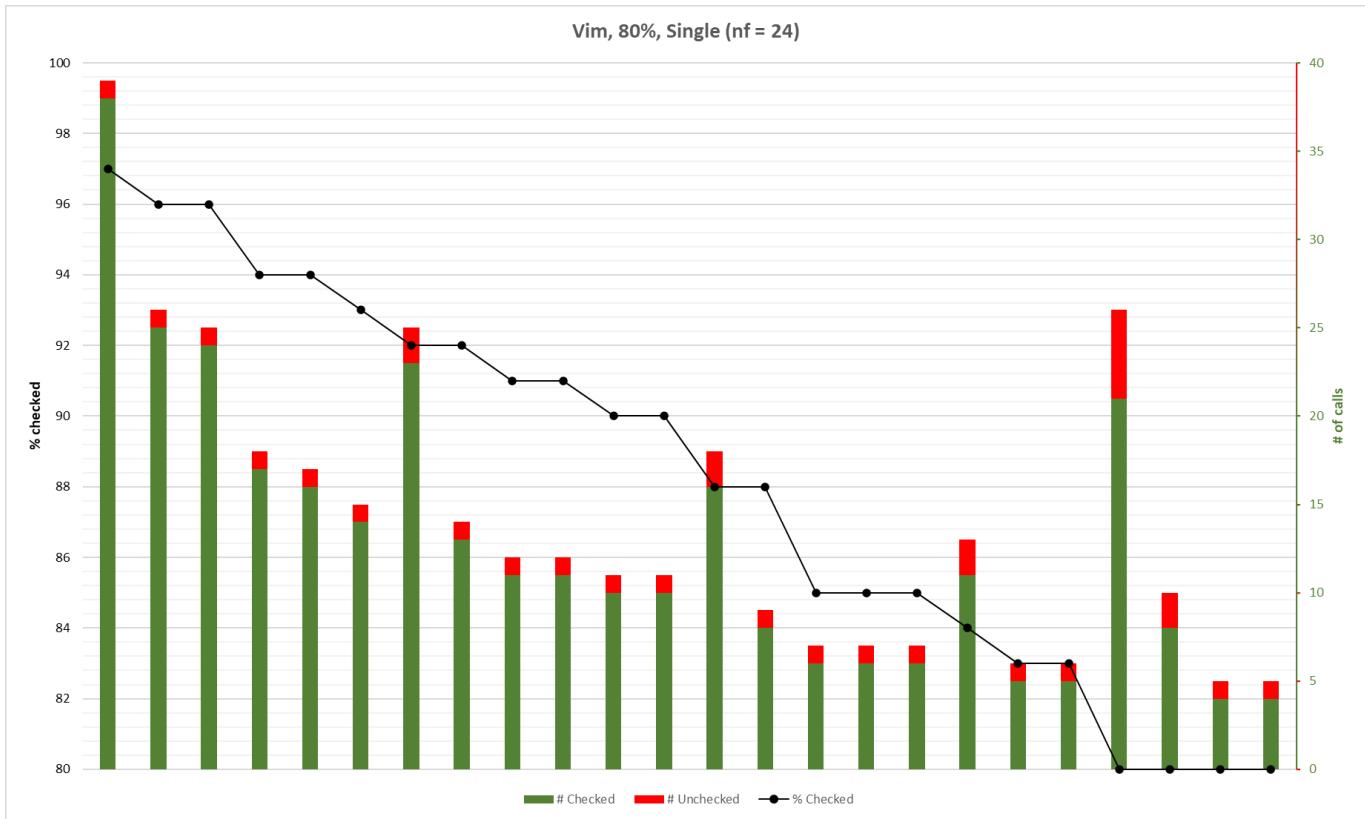
A. Project results



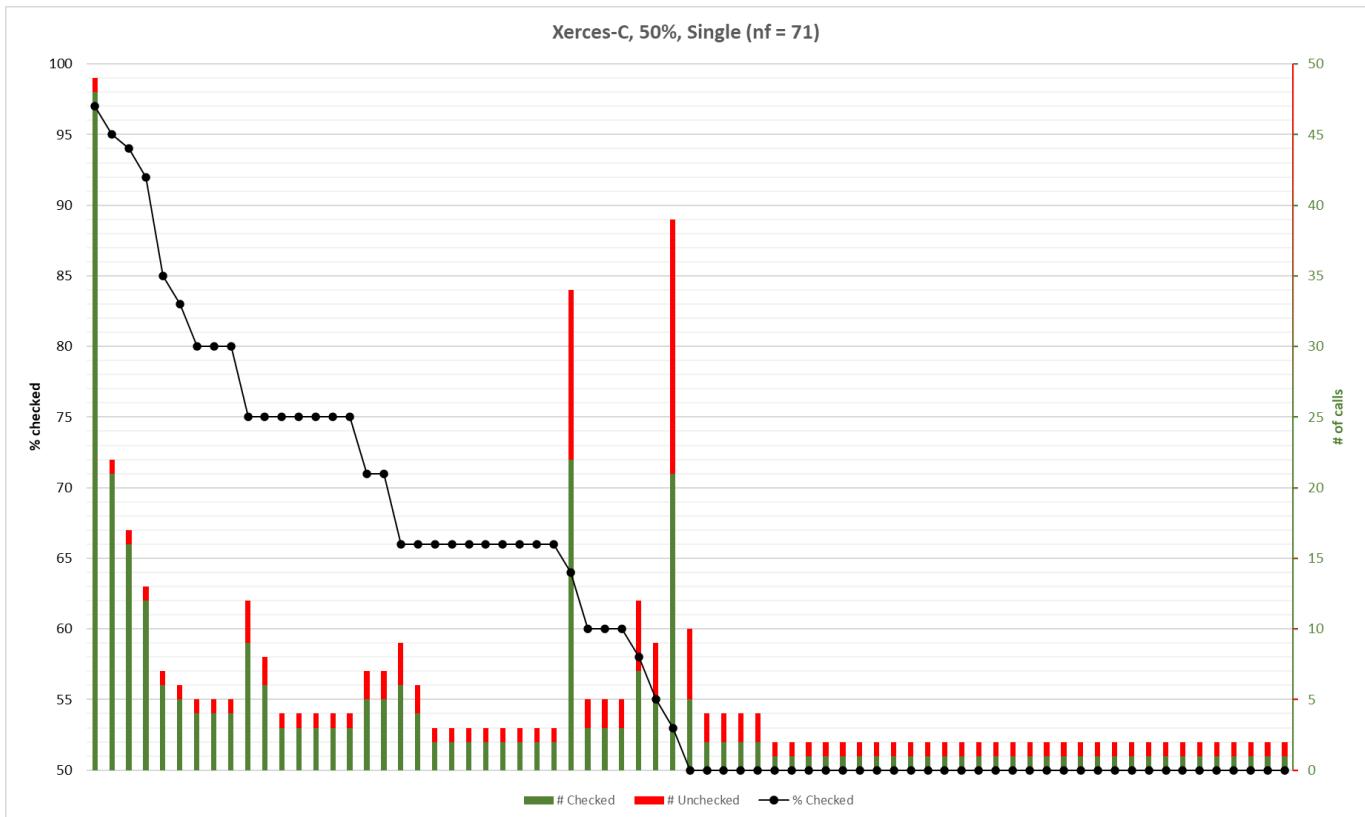
Vim



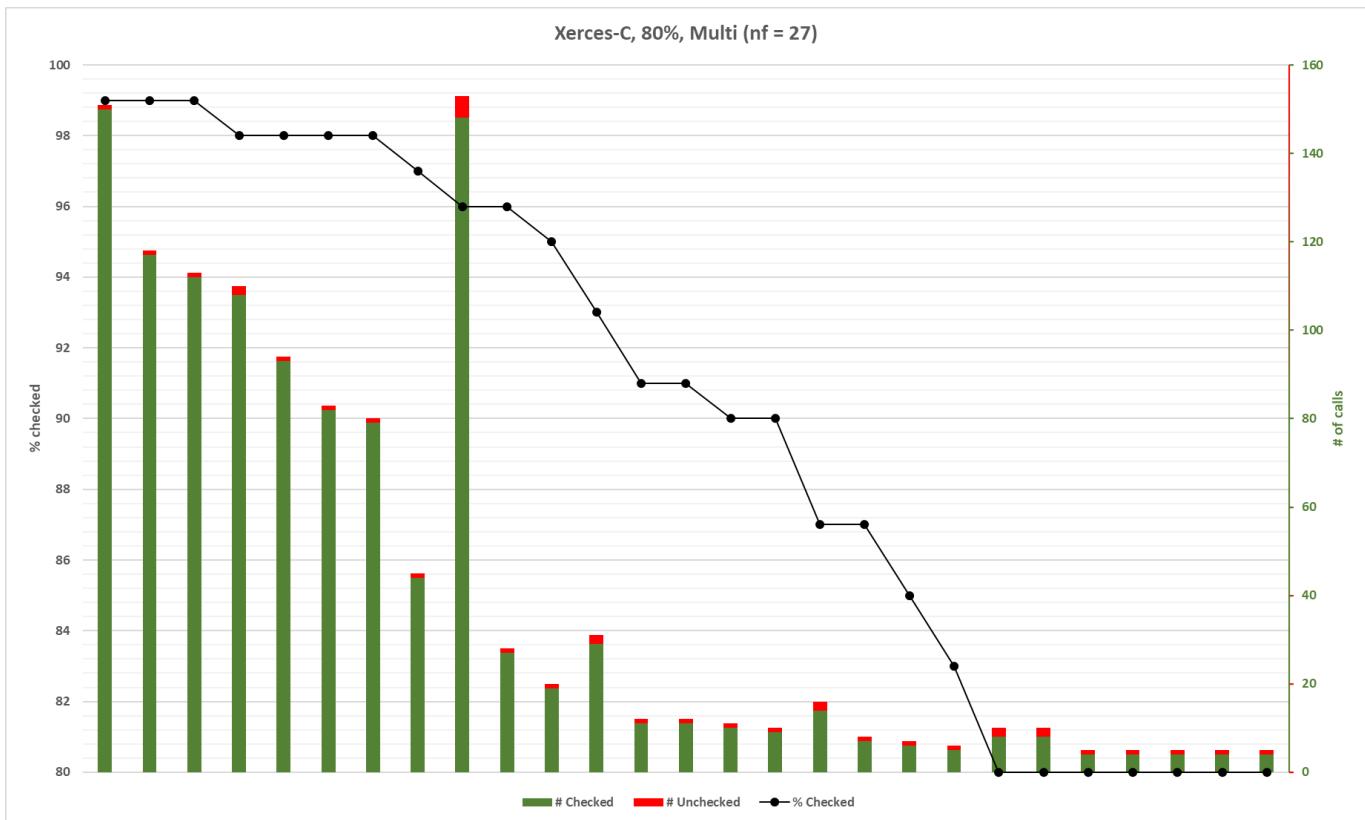
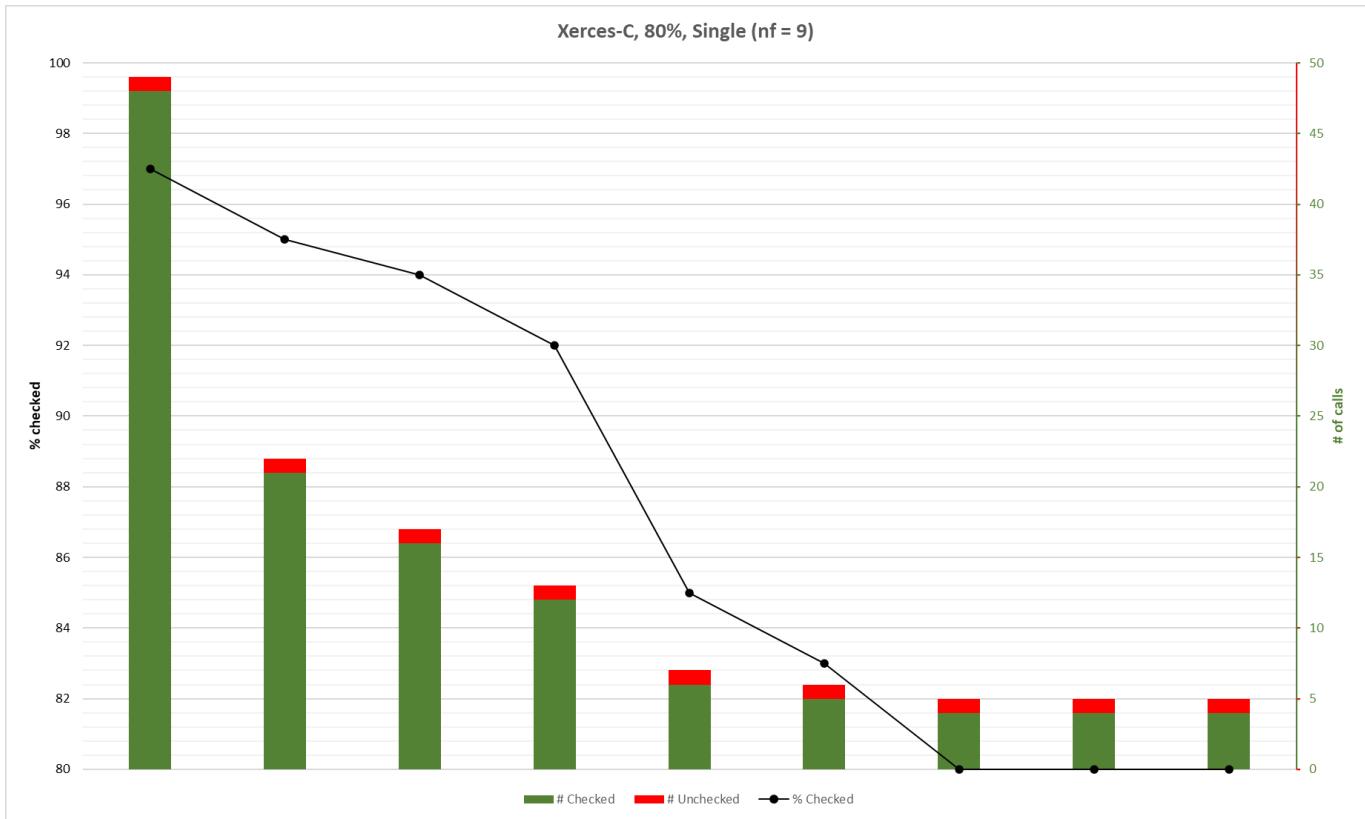
A. Project results



Xerces



A. Project results



Bibliography

- [1] Brian W. Kernighan and Dennis M. Ritchie. *The C Programming Language*. 2nd. Prentice Hall Professional Technical Reference, 1988. ISBN: 0131103709.
- [2] *C++ attribute: nodiscard (since C++17)*. accessed at: 2022-05-18. URL: <http://en.cppreference.com/w/cpp/language/attributes/nodiscard>.
- [3] Valentina Lenarduzzi et al. “A Critical Comparison on Six Static Analysis Tools: Detection Agreement and Precision”. In: *SSRN Electronic Journal* (Jan. 2022). DOI: [10.2139/ssrn.4044439](https://doi.org/10.2139/ssrn.4044439).
- [4] Andrew Rice et al. “Detecting Argument Selection Defects”. In: *Proc. ACM Program. Lang.* 1.OOPSLA (2017). DOI: [10.1145/3133928](https://doi.org/10.1145/3133928). URL: <http://doi.org/10.1145/3133928>.
- [5] Hyrum K. Wright. “Incremental Type Migration Using Type Algebra”. In: *2020 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. 2020, pp. 756–765. DOI: [10.1109/ICSME46990.2020.00085](https://doi.org/10.1109/ICSME46990.2020.00085).
- [6] Hyrum Wright et al. “Large-Scale Automated Refactoring Using ClangMR”. In: *Proceedings of the 29th International Conference on Software Maintenance*. 2013.
- [7] Andrew Rice et al. “Detecting Argument Selection Defects”. In: *Proc. ACM Program. Lang.* 1.OOPSLA (2017). DOI: [10.1145/3133928](https://doi.org/10.1145/3133928). URL: <http://doi.org/10.1145/3133928>.
- [8] Jeffrey Dean and Sanjay Ghemawat. “MapReduce: Simplified Data Processing on Large Clusters”. In: *OSDI’04: Sixth Symposium on Operating System Design and Implementation*. San Francisco, CA, 2004, pp. 137–150.
- [9] Benedek Attila Bahrami. *[clang-tidy] Add infrastructure support for running on project-level information*. Revision on Phabricator for review. 2022. URL: <http://reviews.llvm.org/D124447>.

- [10] Benedek Attila Bahrami. *[clang-tidy] Add the misc-discarded-return-value check*. Revision on Phabricator for review. 2022. URL: <http://reviews.llvm.org/D124446>.
- [11] Benedek Attila Bahrami. *[clang-tidy] Add project-level analysis support to misc-discarded-return-value*. Revision on Phabricator for review. 2022. URL: <http://reviews.llvm.org/D124448>.
- [12] Gábor Horváth et al. *Cross Translation Unit (CTU) Analysis*. accessed at: 2022-05-18. URL: <http://clang.llvm.org/docs/analyzer/user-docs/CrossTranslationUnit.html>.
- [13] *Coverity Scan - Static Analysis*. accessed at: 2022-05-18. URL: http://scan.coverity.com/o/oss_success_stories/69.
- [14] LLVM Foundation. *Using Clang-Tidy*. accessed at: 2022-05-18.
- [15] *CodeChecker*. v6.17.0, d0e2d7ba5f62d722e907d1ce8af6dbb24ee34e5d, accessed at: 2022-05-18. URL: <http://github.com/Ericsson/codechecker.git>.
- [16] LLVM Foundation. *AST Matcher Reference*. accessed at: 2022-05-18. URL: <http://clang.llvm.org/docs/LibASTMatchersReference.html>.
- [17] LLVM Foundation. *LLVM YAML I/O Documentation*. accessed at: 2022-05-18. URL: <http://llvm.org/docs/YamlIO.html>.
- [18] *Bitcoin*. v0.20.1, a62f0ed64f8bbbdfe6467ac5ce92ef5b5222d1bd, accessed at: 2022-05-18. URL: <http://github.com/bitcoin/bitcoin.git>.
- [19] *Contour*. v0.2.0.173, 574d721fb11c52db0d369a4119853282358c39cd, accessed at: 2022-05-18. URL: <http://github.com/contour-terminal/contour.git>.
- [20] *cURL*. curl-7_66_0, 9cd755e1d768bbf228e7c9faf223b7459f7e0105, accessed at: 2022-05-18. URL: <http://github.com/curl/curl.git>.
- [21] *FFmpeg*. n4.3.1, 6b6b9e593dd4d3aaf75f48d40a13ef03bdef9fdb, accessed at: 2022-05-18. URL: <http://github.com/FFmpeg/FFmpeg.git>.
- [22] *libWebM.libwebm*. libwebm-1.0.0.27, 82ac5fc8525820442faa90847b452113fc9394, accessed at: 2022-05-18. URL: <http://github.com/webmproject/libwebm.git>.

- [23] *LLVM*. llvmorg-12.0.0, d28af7c654d8db0b68c175db5ce212d74fb5e9bc, accessed at: 2022-05-18. URL: <http://github.com/llvm/llvm-project.git>.
- [24] *MemcacheD*. 1.6.8, 8335dd84ca85651115d518d7e3d963de9d82e22f, accessed at: 2022-05-18. URL: <http://github.com/memcached/memcached.git>.
- [25] *MongoDB*. r4.4.6, 72e66213c2c3eab37d9358d5e78ad7f5c1d0d0d7, accessed at: 2022-05-18. URL: <http://github.com/mongodb/mongo.git>.
- [26] *OpenSSL*. openssl-3.0.0-alpha7, f9a5682e5c0fbf8b17319d71b0040dba9f8b46ee, accessed at: 2022-05-18. URL: <http://github.com/openssl/openssl.git>.
- [27] *PostgreSQL*. REL_13_0, 29be9983a64c011eac0b9ee29895cce71e15ea77, accessed at: 2022-05-18. URL: <http://github.com/postgres/postgres.git>.
- [28] *Protocol Buffers*. v3.13.0, fde7cf7358ec7cd69e8db9be4f1fa6a5c431386a, accessed at: 2022-05-18. URL: <http://github.com/protocolbuffers/protobuf.git>.
- [29] *QtBase*. v6.2.0, cc60cf83db6bbf9775b1f7747a48d6b7eb24e585, accessed at: 2022-05-18. URL: <http://github.com/qt/qtbase.git>.
- [30] *SQLite*. version-3.33.0, ef215fbf3b581ef4e0273bb3932fa522af88fd7e, accessed at: 2022-05-18. URL: <http://github.com/sqlite/sqlite.git>.
- [31] *TMux*. 2.6, bd71cbbe276432ce8869baa0f2d55433e1ab820c, accessed at: 2022-05-18. URL: <http://github.com/tmux/tmux.git>.
- [32] *tWin*. v0.8.1, 689636ed829c1119ba5507bbd007d678d3d30904, accessed at: 2022-05-18. URL: <http://github.com/cosmos72/twin.git>.
- [33] *Vim*. v8.2.1920, cf4d454df0619ee41ef40e7e91fce3fb061d7d5b, accessed at: 2022-05-18. URL: <http://github.com/vim/vim.git>.
- [34] *Xerces*. v3.2.3, cf1912ac95d4147be08aef4e78f894a3919277d9, accessed at: 2022-05-18. URL: <http://github.com/apache/xerces-c.git>.

List of Figures

1.1	An example of an Abstract Syntax Tree of <code>return foo();</code>	5
2.1	A report of the diagnostics of MongoDB’s source on CodeChecker with 50% threshold on a single TU.	12
2.2	The compiler warnings.	13
2.3	A report of the diagnostics of MongoDB’s source on CodeChecker with 50% threshold and project-level knowledge.	18
3.1	Activity diagram of Clang-Tidy infrastructure in “single” mode.	21
3.2	Activity diagram of Clang-Tidy infrastructure in “multi” mode.	21
3.3	The AST of listing 3.5	26
3.4	The AST of listing 3.6.	27
3.5	A Venn-diagram of what <code>Call</code> (see listing 3.7) matches, where 1 and 2 indicates the arity. The green subsets are to be matched.	29
3.6	Diagram of 80% threshold single run on FFmpeg.	43
3.7	Diagram of 80% threshold multi run on FFmpeg.	43
3.8	Diagram of 50% threshold single run on Bitcoin.	44
3.9	Diagram of 50% threshold multi run on Bitcoin.	45

List of Tables

2.1	System requirements and supported compilers for building LLVM. . .	10
3.1	Results of different runs on live projects.	42

List of Codes

1.1	An example of the infrastructure's limitations.	6
2.1	Compiling LLVM from source on a weaker machine.	11
2.2	Diagnostic output without project-level knowledge.	12
2.3	An example of both unused parameter and ignored return value with nodiscard.	13
2.4	The same example, now with suppressed warnings.	14
2.5	The YAML files containing the collection data.	15
2.6	Contents of the results of collecting data.	16
2.7	The new file containing the collected data.	17
2.8	Contents of the compacted file: the grouped sum of the input.	17
2.9	Diagnostic output with project-level knowledge.	18
3.1	Virtual functions from <code>ClangTidyCheck</code> 's header.	22
3.2	Virtual functions from <code>MatchCallback</code> 's header.	22
3.3	The old infrastructure's way of calling <code>check</code>	23
3.4	Run function distinguishing Diagnose and Collect phase.	23
3.5	The code of fig. 3.3.	26
3.6	A very simple code for matching function calls.	27
3.7	The matcher for our desired call expression.	29
3.8	The matcher for usage in while expression.	29
3.9	The matcher for the usage in return statement.	30
3.10	The matcher for the usage in for statement.	30
3.11	Custom logic and usage of a new matcher called <code>hasAssertExpr</code> . . .	30
3.12	Overridden virtual functions in <code>DiscardedReturnValueCheck</code> 's header. .	31
3.13	Private member functions in <code>DiscardedReturnValueCheck</code> 's header. .	31
3.14	Members of the data structure.	32

3.15 Separating the matchers to statement, declaration, and type.	33
3.16 Catching nodes of used return values.	34
3.17 Functions for loading.	36
3.18 Function for writing.	37
3.19 Steps for making the data YAML writable.	38