

Advanced Optimization and New Capabilities of GCC 10

SUSE Linux Enterprise Server 15 SP2
Development Tools Module

Martin Jambor, Toolchain Developer (SUSE)
Jan Hubička, Toolchain Developer (SUSE)
Richard Biener, Toolchain Developer (SUSE)
Martin Liška, Toolchain Developer (SUSE)
Michael Matz, Toolchain Team Lead (SUSE)
Brent Hollingsworth, Engineering Manager (AMD)

The document at hand provides an overview of GCC 10 as the current Development Tools Module compiler in SUSE Linux Enterprise 15 SP2. It focuses on the important optimization levels and options **Link Time Optimization (LTO)** and **Profile Guided Optimization (PGO)**. Their effects are demonstrated by compiling the SPEC CPU benchmark suite for AMD EPYC 7002 Series Processors and building Mozilla Firefox for a generic x86_64 machine.

Disclaimer: Documents published as part of the SUSE Best Practices series have been contributed voluntarily by SUSE employees and third parties. They are meant to serve as examples of how particular actions can be performed. They have been compiled with utmost attention to detail. However, this does not guarantee complete accuracy. SUSE cannot verify that actions described in these documents do what is claimed or whether actions described have unintended consequences. SUSE LLC, its affiliates, the authors, and the translators may not be held liable for possible errors or the consequences thereof.

Contents

- 1 Overview 4
- 2 System compiler versus Developer Tools Module compiler 4
- 3 Optimization levels and related options 9
- 4 Taking advantage of newer processors 12
- 5 Link Time Optimization (LTO) 13
- 6 Profile-Guided Optimization (PGO) 16
- 7 Performance evaluation: SPEC CPU 2017 18
- 8 Performance evaluation: Mozilla Firefox 31
- 9 Legal notice 40
- 10 GNU Free Documentation License 41

1 Overview

The first release of the GNU Compiler Collection (GCC) with the major version 10, GCC 10.1, has been released in May 2020. GCC 10.2 with fixes to 94 bugs has followed in June of the same year. Subsequently, it replaced the compiler in the SUSE Linux Enterprise (SLE) Development Module. GCC 10 comes with many new features, such as implementing parts of the most recent versions of specifications of various languages (especially [C2X](#), [C++17](#), [C++20](#)) and their extensions (OpenMP, OpenACC), supporting new capabilities of a wide range of computer architectures and numerous generic optimization improvements.

This document gives an overview of GCC 10. The focus of the document is on how to select appropriate optimization options for your application and stressing benefits of advanced modes of compilation. First, we describe the optimization levels the compiler offers, and other important options developers often use. We explain when and how you can benefit from using **Link Time Optimization (LTO)** and **Profile Guided Optimization (PGO)** builds. We also detail their effects when building a set of well known CPU intensive benchmarks, and we are looking at how this performs on the AMD Zen 2 based EPYC 7002 Series Processor. Finally, we take a closer look at the effects they have on a big software project: Mozilla Firefox.

2 System compiler versus Developer Tools Module compiler

The major version of the system compiler in SUSE Linux Enterprise 15 remains to be GCC 7, regardless of the service pack level. This is to minimize the danger of any unintended changes over the entire life time of the product.

```
sles15: # gcc --version
gcc (SUSE Linux) 7.5.0
Copyright (C) 2017 Free Software Foundation, Inc.
This is free software; see the source for copying conditions.  There is NO
warranty; not even for MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.
```

That does not mean that, as a user of SUSE Linux Enterprise 15, you are forced to use a compiler with features frozen in 2016. You can install an add-on module called **Development Tools Module**. This module is included in the SUSE Linux Enterprise Server 15 subscription and contains a much newer compiler.

At the time of writing this document, the compiler included in the Development Tools Module is GCC 10.2. Nevertheless, it is important to stress that, unlike the system compiler, the major version of the most recent GCC from the module will change shortly after the upstream release of GCC 11.2 (most likely in summer 2021), GCC 12.2 (summer 2022) and so forth. Note that only the most recent compiler in the Development Tools Module is supported at any time, except for a six months overlap period after the upgrade happened. Developers on a SUSE Linux Enterprise Server 15 system therefore have always access to two supported GCC versions: the almost unchanging system compiler and the most recent compiler from the Development Tools Module. Programs and libraries built with the compiler from the Development Tools Module can run on computers running SUSE Linux Enterprise Server 15 which do not have the module installed. All necessary runtime libraries are available from the main repositories of the operating system itself, and new ones are added through the standard update mechanism. In the document at hand, we use the term GCC 10 as synonym for any minor version of the major version 10 and GCC 10.2, to refer to specifically that version. In practice, they should be interchangeable.

2.1 When to use compilers from the Development Tools Module

In many cases you will find that the system compiler perfectly satisfies your needs. After all, it is the compiler used to build all packages and their updates in the system itself. On the other hand, there are situations where a newer compiler is necessary, or where you want to consider using a newer compiler to get some of the benefits of its ongoing development.

If the program or library you are building uses language features which are not supported by GCC 7, you cannot use the system compiler. However, the compiler from the Development Tools Module will usually be sufficiently new. The most obvious case is C++. GCC 10 has a mature implementation of C++17 features, whereas the one in GCC 7 is only experimental and incomplete. The GNU C++ Library, which accompanies GCC 10, is also almost C++17 feature-complete. Only *hardware interference sizes*¹ are not implemented and *elementary string conversions*² have extra limitations. Most of C++20 features are implemented in GCC 10 as experimental features. Try them out with appropriate caution. Most notably, *Modules*³ and *Atomic Compare-and-Exchange*

¹ Proposal P0154R1

² Proposal P0067R5

³ Proposals P1103R3, P1766R1, P1811R0, P1703R1, P1874R1, P1979R0, P1779R3, P1857R3, P2115R0 and P1815R2

with *Padding Bits*⁴ are not supported yet, while *Coroutines*⁵ are implemented but require that the source file is compiled with the `-fcoroutines` switch. If you are interested in the implementation status of any particular C++ feature in the compiler, consult the following pages:

- [C++ Standards Support in GCC \(https://gcc.gnu.org/projects/cxx-status.html\)](https://gcc.gnu.org/projects/cxx-status.html), and
- [The GNU C++ Library Manual \(https://gcc.gnu.org/onlinedocs/gcc-10.2.0/libstdc++/manual/\)](https://gcc.gnu.org/onlinedocs/gcc-10.2.0/libstdc++/manual/).

Advances in supporting new language specifications are not limited to C++. GCC 10 supports several new features from the ISO 202X C standard draft, and the Fortran compiler has also seen many improvements. And if you use `OpenMP` or `OpenACC` extensions for parallel programming, you will realize that the compiler supports a lot of features of new versions of these standards. For more details, visit the links at the end of this section.

In addition to new supported language constructs, GCC 10 offers improved diagnostics when it reports errors and warnings to the user so that they are easier to understand and to be acted upon. This is particularly useful when dealing with issues in templated C++ code. Furthermore, there are several new warnings which help to avoid common programming mistakes.

Because GCC10 is newer, it can generate code for many recent processors not supported by GCC 7. Such a list of processors would be too large to be displayed here. Nevertheless, in [Section 7, “Performance evaluation: SPEC CPU 2017”](#) we specifically look at optimizing code for an AMD EPYC 7002 Series Processor which is based on AMD Zen 2 cores. At this point we should stress that the *system compiler* does not know this kind of core and therefore cannot optimize for it. GCC 10, on the other hand, is the second major release supporting AMD Zen 2 cores, and thus can often produce significantly faster code for it.

Finally, the general optimization pipeline of the compiler has also significantly improved over the years, which we will demonstrate in the last sections of this document. To find out more about improvements in versions of GCC 8, 9 and 10, visit their respective “changes” pages:

- [GCC 8 Release Series Changes, New Features, and Fixes \(https://gcc.gnu.org/gcc-8/changes.html\)](https://gcc.gnu.org/gcc-8/changes.html),
- [GCC 9 Release Series Changes, New Features, and Fixes \(https://gcc.gnu.org/gcc-9/changes.html\)](https://gcc.gnu.org/gcc-9/changes.html), and
- [GCC 10 Release Series Changes, New Features, and Fixes \(https://gcc.gnu.org/gcc-10/changes.html\)](https://gcc.gnu.org/gcc-10/changes.html).

⁴ Proposal P0528R3

⁵ Proposal P0912R5

2.2 Potential issues with the Development Tools Module

GCC 10 from the Development Tools Module can sometimes behave differently in a way that can cause issues which were not present with the system compiler. Such problems encountered by other users are listed in the following documents:

- [Porting to GCC 8 \(https://gcc.gnu.org/gcc-8/porting_to.html\)](https://gcc.gnu.org/gcc-8/porting_to.html),⁷
- [Porting to GCC 9 \(https://gcc.gnu.org/gcc-9/porting_to.html\)](https://gcc.gnu.org/gcc-9/porting_to.html)⁷, and
- [Porting to GCC 10 \(https://gcc.gnu.org/gcc-10/porting_to.html\)](https://gcc.gnu.org/gcc-10/porting_to.html)⁷.

We encourage you to read through these three short pages to get an understanding of the problems. The document at hand briefly mentions two such potential pitfalls.

The first one is that, for performance reasons, GCC 10 defaults to `-fno-common` which means that a linker error will now be reported if the same variable is defined in two `C` compilation units. This can happen if two or more `.c` files include the same header file which intends to declare a variable but omits the `extern` keyword when doing so, inadvertently resulting in multiple definitions. If you encounter such an error, you simply need to add the `extern` keyword to the declaration in the header file and define the variable in only a single compilation unit. Alternatively, you can compile your project with an explicit `-fcommon` if you are willing to accept that this behavior is inconsistent with `C++` and may incur speed and code size penalties. The second issue highlighted here is that the `C++` compiler in GCC 8 and later now assumes that no execution path in a non-void function simply reaches the end of the function without a return statement. This means it is assumed that such code paths will never be executed, and thus they will be eliminated. You should therefore pay special attention to warnings produced by `-Wreturn-type`. This option is enabled by default and indicates which functions might be affected.

2.3 Installing GCC 10 from the Development Tools Module

Similar to other modules and extensions for SUSE Linux Enterprise Server 15, you can activate the Development Tools Module either using the command line tool **SUSEConnect** or using the **YaST** setup and configuration tool. To use the former, carry out the following steps:

1. As root, start by listing the available and activated modules and extensions:

```
sles15: # SUSEConnect --list-extensions
```

2. In the computer output, look for “Development Tools Module”:

```
Development Tools Module 15 SP2 x86_64
Activate with: SUSEConnect -p sle-module-development-tools/15.2/x86_64
```

If you see the text (Activated) next to the module name, the module is already ready to be used. You can safely proceed to the installation of the compiler packages.

3. Otherwise, issue the activation command that is shown in the computer output above:

```
sles15: # SUSEConnect -p sle-module-development-tools/15.2/x86_64
Registering system to SUSE Customer Center

Updating system details on https://scc.suse.com ...

Activating sle-module-development-tools 15.2 x86_64 ...
-> Adding service to system ...
-> Installing release package ...

Successfully registered system
```

If you prefer to use YaST, the procedure is also straightforward. Run YaST as root and go to the **Add-On Products** menu in the **Software** section. If “Development Tools Module” is among the listed installed modules, you already have the module activated and can proceed with installing individual compiler packages. If not, click the **Add** button, select **Select Extensions and Modules from Registration Server**, and YaST will guide you through a simple procedure to add the module.

When you have the Development Tools Module installed, you can verify that the GCC 10 packages are available to be installed on your system:.

```
sles15: # zypper search gcc10
Refreshing service 'Basesystem_Module_15_SP2_x86_64'.
Refreshing service 'Desktop_Applications_Module_15_SP2_x86_64'.
Refreshing service 'Development_Tools_Module_15_SP2_x86_64'.
Refreshing service 'SUSE_Linux_Enterprise_Server_15_SP2_x86_64'.
Refreshing service 'SUSE_Package_Hub_15_SP2_x86_64'.
Refreshing service 'Server_Applications_Module_15_SP2_x86_64'.
Loading repository data...
Reading installed packages...

S | Name | Summary
--+-----+-----
  | gcc10 | The GNU C Compiler and Support Files
  | gcc10 | The GNU C Compiler and Support Files
```


gcc10-32bit	The GNU C Compiler 32bit support
gcc10-ada	GNU Ada Compiler Based on GCC (GNAT)
gcc10-ada-32bit	GNU Ada Compiler Based on GCC (GNAT)
gcc10-c++	The GNU C++ Compiler
gcc10-c++-32bit	The GNU C++ Compiler
gcc10-fortran	The GNU Fortran Compiler and Support Files
gcc10-fortran-32bit	The GNU Fortran Compiler and Support Files
gcc10-go	GNU Go Compiler
gcc10-go-32bit	GNU Go Compiler
gcc10-info	Documentation for the GNU compiler collection
gcc10-locale	Locale Data for the GNU Compiler Collection
libstdc++6-devel-gcc10	Include Files and Libraries mandatory for Development
libstdc++6-devel-gcc10-32bit	Include Files and Libraries mandatory for Development
libstdc++6-pp-gcc10	GDB pretty printers for the C++ standard library
libstdc++6-pp-gcc10-32bit	GDB pretty printers for the C++ standard library

Now you can simply install the compilers for the programming languages you use with **zypper**:

```
sles15: # zypper install gcc10 gcc10-c++ gcc10-fortran
```

The compilers are installed on your system, the executables are called **gcc-10**, **g++-10**, **gfortran-10** and so on. It is also possible to install the packages in **YaST**. To do so, simply enter the “Software Management” menu in the **Software** section and search for “gcc10”. Then select the packages you want to install. Finally, click the **Accept** button.



Note: Newer compilers on openSUSE Leap 15.2

The community distribution openSUSE Leap 15.2 shares most of the base packages with SUSE Linux Enterprise Server 15 SP2. The system compiler on systems running openSUSE Leap 15.2 is also GCC 7.5. There is no Development Tools Module for the community distribution available, but a newer compiler is provided. Simply install the packages `gcc10`, `gcc10-c++`, `gcc10-fortran`, and the like.

3 Optimization levels and related options

GCC has a rich optimization pipeline that is controlled by approximately a hundred of command line options. It would be impractical to force users to decide about each one of them whether they want to have it switched on when compiling their code. Like all other modern compilers, GCC therefore introduces the concept of optimization levels which allow the user to pick one common configuration from a few options. Optionally, the user can tweak the selected level, but that does not happen frequently.

The default is to not optimize at all. You can specify this optimization level on the command line as `-O0`. It is often used when developing and debugging a project. This means it is usually accompanied with the command line switch `-g` so that debug information is emitted. As no optimizations take place, no information is lost because of it. No variables are optimized away, the compiler only inlines functions with special attributes that require it, and so on. As a consequence, the debugger can almost always find everything it searches for in the running program and report on its state very well. On the other hand, the resulting code is big and slow. Thus this optimization level should not be used for release builds.

The most common optimization level for release builds is `-O2` which attempts to optimize the code aggressively but avoids large compile times and excessive code growth. Optimization level `-O3` instructs GCC to simply optimize as much as possible, even if the resulting code might be considerably bigger and the compilation can take longer. Note that neither `-O2` nor `-O3` imply anything about the precision and semantics of floating-point operations. Even at the optimization level `-O3` GCC implements math functions so that they strictly follow the respective IEEE and/or ISO rules. This often means that the compiled programs run markedly slower than necessary if such strict adherence is not required. The command line switch `-ffast-math` is a common way to relax rules governing floating-point operations. It is out of scope of this document to provide a list of the fine-grained options it enables and their meaning. However, if your software crunches floating-point numbers and its runtime is a priority, you can look them up in the GCC manual and review what semantics of floating-point operations you need.

The most aggressive optimization level is `-Ofast` which does imply `-ffast-math` along with a few other options that disregard strict standard compliance. In GCC 10 this level also means the optimizers may introduce data races when moving memory stores which may not be safe for multithreaded applications. Additionally, the Fortran compiler can take advantage of associativity of math operations even across parentheses and convert big memory allocations on the heap to allocations on stack. The last mentioned transformation may cause the code to violate maximum stack size allowed by `ulimit` which is then reported to the user as a segmentation fault. We often use level `-Ofast` to build benchmarks. It is a shorthand for the options on top of `-O3`, which often make them run faster, and the benchmarks are usually written in a way that they still run correctly.

If you feed the compiler with huge machine-generated input, especially if individual functions happen to be extremely large, the compile time can become an issue even when using `-O2`. In such cases, use the most lightweight optimization level `-O1` that avoids running almost all optimizations with quadratic complexity. Finally, the `-Os` level directs the compiler to aggressively optimize for the size of the binary.



Note: Optimization level recommendation

Usually we recommend using `-O2`. This is the optimization level we use to build most SUSE and openSUSE packages, because at this level the compiler makes balanced size and speed trade-offs when building a general-purpose operating system. However, we suggest using `-O3` if you know that your project is compute-intensive and is either small or an important part of your actual workload. Moreover, if the compiled code contains performance-critical floating-point operations, we strongly advise that you investigate whether `-ffast-math` or any of the fine-grained options it implies can be safely used.

If your project and the techniques you use to debug or instrument it do not depend on *ELF symbol interposition*, you may consider trying to speed it up by using `-fno-semantic-interposition`. This allows the compiler to inline calls and propagate information even when it would be illegal if a symbol changed during dynamic linking. Using this option to signal to the compiler that interposition is not going to happen is known to significantly boost performance of some projects, most notably the Python interpreter.

Some projects use `-fno-strict-aliasing` to work around type punning problems in the source code. This is not recommended except for very low-level hand optimized code such as the Linux kernel. Type-based alias analysis is a very powerful tool. It is used to enable other transforms, such as store-to-load propagation that in turn enables other transformations, such as aggressive inlining, vectorization and other high level transformations.

With the `-g` switch GCC still tries hard to generate useful debug information even when optimizing. However, a lot of information is irrecoverably lost in the process. Debuggers also often struggle to present the user with a view of the state of a program in which statements are not necessarily executed in the original order. Debugging optimized code can therefore be a challenging task but usually is still somewhat possible.

The complete list of optimization and other command line switches is available in the compiler manual, provided in the info format in the package `gcc10-info` or online at [the GCC project Web site \(https://gcc.gnu.org/onlinedocs/gcc-10.2.0/gcc/\)](https://gcc.gnu.org/onlinedocs/gcc-10.2.0/gcc/).

Bear in mind that although almost all optimizing compilers have the concept of optimization levels and their optimization levels often have the same names as those in GCC, they do not necessarily mean to make the same trade-offs. Famously, GCC's `-Os` optimizes for size much more aggressively than LLVM/Clang's level with the same name. Therefore, it often produces

slower code; the more equivalent option in Clang is `-Oz` which GCC does not have. Similarly, `-O2` can have different meanings for different compilers. For example, the difference between `-O2` and `-O3` is much bigger in GCC than in LLVM/Clang.



Note: Changing the optimization level with **cmake**

If you use **cmake** to configure and set up builds of your application, be aware that its *release* optimization level defaults to `-O3` which might not be what you want. To change it, you must modify the `CMAKE_C_FLAGS_RELEASE`, `CMAKE_CXX_FLAGS_RELEASE` and/or `CMAKE_Fortran_FLAGS_RELEASE`, since these variables are appended at the end of the compilation command lines, thus overwriting any level set in the variables `CMAKE_C_FLAGS`, `CMAKE_CXX_FLAGS`, and the like.

4 Taking advantage of newer processors

By default GCC assumes that you want to run the compiled program on a wide variety of CPUs, including fairly old ones, regardless of the selected optimization level. On architectures like `x86_64` and `aarch64` the generated code will only contain instructions available on every CPU model of the architecture, including the earliest ones. On `x86_64` in particular this means that the programs will use the `SSE` and `SSE2` instruction sets for floating point and vector operations but not any more recent ones.

If you know that the generated binary will run only on machines with newer instruction set extensions, you can specify it on the command line. Their complete list is available in the manual, but the most prominent one is `-march` which lets you select a CPU model to generate code for. For example, if you know that your program will only be executed on AMD EPYC 7002 Series Processors which is based on AMD Zen 2 cores or processors that are compatible with it, you can instruct GCC to take advantage of all the instructions the CPU supports with option `-march=znver2`. Note that on SUSE Linux Enterprise Server 15, the system compiler does not know this particular value of the switch; you need to use GCC 10 from the Development Tools Module to optimize code for these processors.

To run the program on the machine on which you are compiling it, you can have the compiler auto-detect the target CPU model for you with the option `-march=native`. This only works if the compiler is new enough. The system compiler of SUSE Linux Enterprise Server, for example, misidentifies AMD EPYC 7002 Series Processors as being based on the AMD Zen 1 core. Among

other things, this means that it only emits 128bit vector instructions, even though the CPU has data-paths wide enough to efficiently process 256bit ones. Again, the easy solution is to use the compiler from the Development Tools Module when targeting recent processors.

5 Link Time Optimization (LTO)

Figure 1 outlines the classic mode of operation of a compiler and a linker. Pieces of a program are compiled and optimized in chunks defined by the user called compilation units to produce so-called object files which already contain binary machine instructions and which are combined together by a linker. Because the linker works at such low level, it cannot perform much optimization and the division of the program into compilation units thus presents a profound barrier to optimization.

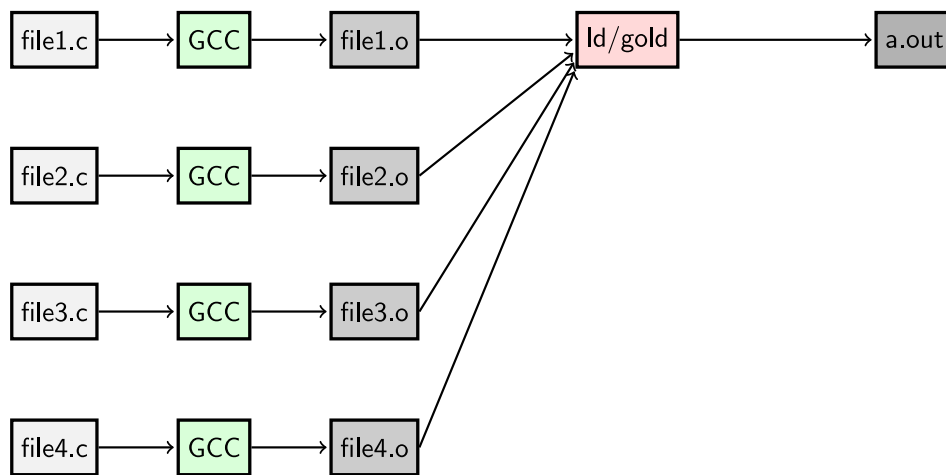


FIGURE 1: TRADITIONAL PROGRAM BUILD

This limitation can be overcome by rearranging the process so that the linker does not receive as its input the almost finished object files containing machine instructions, but is invoked on files containing so called *intermediate language* (IL) which is a much richer representation of each original compilation unit (see figure 2). The linker identifies the input as not yet entirely compiled and invokes a linker plugin which in turn runs the compiler again. But this time it has at its disposal the representation of the entire program or library that is being built. The compiler makes decisions about what optimizations across function and compilation unit boundaries will be carried out and then divides the program into a set of partitions. Each of the partitions is further optimized independently, and machine code is emitted for it, which is finally linked the traditional way. Processing of the partitions is performed in parallel.

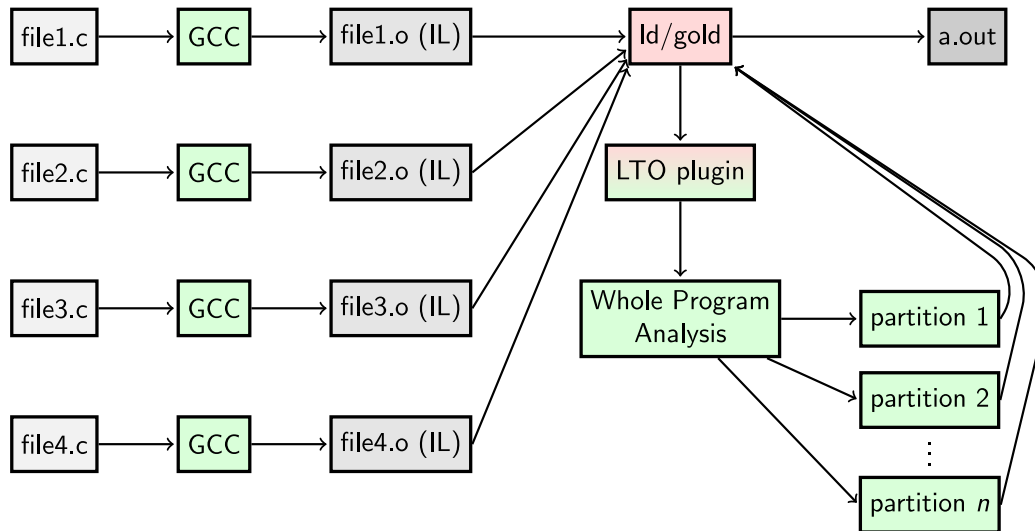


FIGURE 2: BUILDING A PROGRAM WITH GCC USING LINK TIME OPTIMIZATION (LTO)

To use *Link Time Optimization*, all you need do is to add the `-flto` switch to the compilation command line. The vast majority of packages in the Linux distribution openSUSE Tumbleweed has been built with LTO for over a year without any major problems. A lot of work has recently been put into emitting good debug information when building with LTO. Thus the debugging experience is not limited anymore as it was a couple of years ago.

LTO in GCC always consists of a *whole program analysis* (WHOPR) stage followed by the majority of the compilation process performed in parallel, which greatly reduces the build times of most projects. To control the parallelism, you can explicitly cap the number of parallel compilation processes by n if you specify `-flto= n` at linker command line. Alternatively, it is possible to use the GNU **make** jobserver with `-flto=jobserver` while also prepending the **makefile** rule invoking link step with character `+` to instruct GNU make to keep the jobserver available to the linker process. You can also use `-flto=auto` which instructs GCC to search for the jobserver and if it is not found, use all available CPU threads.

Note that there is a technical difference in how GCC and LLVM/Clang approach LTO. Clang provides two LTO mechanisms, so-called *thin LTO* and *full LTO*. In full LTO, LLVM processes the whole program as if it was a single translation unit which does not allow for any parallelism. GCC can be configured to operate in this way with the option `-flto-partition=one`. LLVM in thin-LTO mode can compile different compilation units in parallel and makes possible inlining across compilation unit boundaries, but not most other types of cross-module optimizations. This mechanism therefore has inherently higher code quality penalty than full LTO or the approach of GCC.

5.1 Most notable benefits of LTO

Applications built with LTO are often faster, mainly because the compiler can *inline* calls to functions in another compilation unit. This possibility also allows programmers to structure their code according to its logical division because they are not forced to put function definitions into header files to enable their inlining. Not all calls conveying information known at compilation time can be inlined. But GCC can still track and propagate constants, value ranges and devirtualization contexts to the callees, often even when passed in an aggregate or by reference, that can then subsequently save unnecessary computations. LTO allows such propagation across compilation unit boundaries, too.

Link Time Optimization with whole program analysis also offers many opportunities to shrink the code size of the built project. Thanks to *symbol promotion* and inter-procedural *unreachable code elimination*, functions and their parts which are not necessary in any particular project can be removed even when they are not declared static and are not defined in an anonymous namespace. Automatic *attribute discovery* can identify C++ functions that do not throw exceptions which allows the compiler to avoid generating a lot of code in exception cleanup regions. *Identical code folding* can find functions with the same semantics and remove all but one of them. The code size savings are often very significant and a compelling reason to use LTO even for applications which are not CPU-bound.



Note: Building libraries with LTO

The symbol promotion is controlled by resolution information given to the linker and depends on type of the DSO build. When producing a dynamically loaded shared library, all symbols with default visibility can be overwritten by the dynamic linker. This blocks the promotion of all functions not declared inline, thus it is necessary to use the hidden visibility wherever possible to achieve best results. Similar problems happen even when building static libraries with -rdynamic.

5.2 Potential issues with LTO

As noted earlier, the vast majority of packages in the openSUSE Tumbleweed distribution are built with LTO without any need to tweak them, and they work fine. Nevertheless, some low-level constructs pose a problem for LTO. One typical issue are symbols defined in *inline assembly* which can happen to be placed in a different partition from their uses and subsequently fail the final linking step. To build such projects with LTO, the assembler snippets defining symbols

must be placed into a separate assembler source file so that they only participate in the final linking step. Global register variables are not supported by LTO, and programs either must not use this feature or be built the traditional way.

Another notable limitation of LTO is that it does not support *symbol versioning* implemented with special inline assembly snippets (as opposed to a linker map file). To define symbol versions in the source files, you can do so with the new `symver` function attribute. As an example, the following snippet will make the function `foo_v1` implement `foo` in *node* `VERS_1` (which must be specified in the version script supplied to the linker). Consult [the manual \(https://gcc.gnu.org/onlinedocs/gcc/Common-Function-Attributes.html#index-symver-function-attribute\)](https://gcc.gnu.org/onlinedocs/gcc/Common-Function-Attributes.html#index-symver-function-attribute) for more details.

```
__attribute__ ((__symver__ ("foo@VERS_1")))
int foo_v1 (void)
{
}
```

Sometimes the extra power of LTO reveals pre-existing problems which do not manifest themselves otherwise. Violations of (strict) *aliasing* rules and C++ one definition rule tend to cause misbehavior significantly more often; the latter is fortunately reported by the `-Wodr` warning which is on by default and should not be ignored. We have also seen cases where the use of the `flatten` function attribute led to unsustainable amount of inlining with LTO. Furthermore, LTO is not a good fit for code snippets compiled by `configure` scripts (generated by `autoconf`) to discover the availability of various features, especially when the script then searches for a string in the generated assembly.

Finally, we needed to configure the virtual machines building the biggest openSUSE packages to have more memory than when not using LTO. Whereas in the traditional mode of compilation 1 GB of RAM per core was enough to build Mozilla Firefox, the serial step of LTO means the build-bots need 16 GB even when it has fewer than 16 cores.

6 Profile-Guided Optimization (PGO)

Optimizing compilers frequently make decisions according to which path through the code they consider most likely to be executed, how many times a loop is expected to iterate, and similar estimates. They also often face trade-offs between potential runtime benefits and code size growth. Ideally, they would optimize only frequently executed (also called *hot*) bits of a program for speed and everything else for size to reduce strain on caches and make the distribution of

the built software cheaper. Unfortunately, guessing which parts of a program are the *hot* ones is difficult, and even sophisticated estimation algorithms implemented in GCC are no good match for a measurement.

If you do not mind adding an extra level of complexity to the build system of your project, you can make such measurement part of the process. The **makefile** (or any other) build script needs to compile it twice. The first time it needs to compile with the `-fprofile-generate` option and then execute the first binary in one or multiple *train runs* during which it will save information about the behavior of the program to special files. Afterward, the project needs to be rebuilt again, this time with the `-fprofile-use` option which instructs the compiler to look for the files with the measurements and use them when making optimization decisions, a process called *Profile-Guided Optimization (PGO)*.

It is important that the train exhibits the same characteristics as the real workload. Unless you use the option `-fprofile-partial-training` in the second build, it needs to exercise the code that is also the most frequently executed in real use, otherwise it will be optimized for size and PGO would make more harm than good. With the option, GCC reverts to guessing properties of portions of the projects not exercised in the train run, as if they were compiled without profile feedback. This however also means that the code size will not typically shrink as much as one would expect from a PGO build.

On the other hand, train runs do not need to be a perfect simulation of the real workload. For example, even though a test suite should not be a very good train run in theory because it disproportionally often tests various corner cases, in practice many projects use it as a train run and achieve significant runtime improvements with real workloads, too.

Profiles collected using an instrumented binary for multithreaded programs may be inconsistent because of missed counter updates. You can use `-fprofile-correction` in addition to `-fprofile-use` so that GCC uses heuristics to correct or smooth out such inconsistencies instead of emitting an error.

Profile-Guided Optimization can be combined and is complimentary to Link Time Optimization. While LTO expands what the compiler can do, PGO informs it about which parts of the program are the important ones and should be focused on. The following sections detail this by means of two rather different case studies.

7 Performance evaluation: SPEC CPU 2017

Standard Performance Evaluation Corporation (SPEC) is a non-profit corporation that publishes a variety of industry standard benchmarks to evaluate performance and other characteristics of computer systems. Its latest suite of CPU intensive workloads, SPEC CPU 2017, is often used to compare compilers and how well they optimize code with different settings because the included benchmarks are well known and represent a wide variety of computation-heavy programs. This section highlights selected results of a GCC 10 evaluation using the suite.

Note that when we use SPEC to perform compiler comparisons, we are lenient toward some official SPEC rules which system manufacturers need to observe to claim an official score for their system. We disregard the concepts of *base* and *peak* metrics and simply focus on results of compilations using a particular set of options. We even patched several benchmarks:

- Benchmarks `502.gcc_r`, `505.mcf_r`, `511.povray_r`, and `527.cam4_r` contain an implementation of quicksort which violates (strict) C/C++ aliasing rules which can lead to erroneous behavior when optimizing at link time. SPEC decided not to change the released benchmarks and simply suggests that these benchmarks are built with the `-fno-strict-aliasing` option when they are built with GCC. That makes evaluation of compilers using SPEC problematic, gauging their ability to use aliasing rules to facilitate optimizations is important. We have therefore disabled it only for the problematic `qsort` attributes with the following function attribute:

```
__attribute__((optimize("-fno-strict-aliasing")))
```

As a result, the only benchmark which we compile with `-fno-strict-aliasing` is `500.perlbench_r`.

- We have increased the tolerance of `549.fotonik3d_r` to rounding errors after it became clear the intention was that the compiler can use relaxed semantics of floating point operations in the benchmark (see [GCC bug 84201 \(https://gcc.gnu.org/bugzilla/show_bug.cgi?id=84201\)](https://gcc.gnu.org/bugzilla/show_bug.cgi?id=84201) ↗).

For the mentioned reasons (and probably some more), all the results in this document are *non-reportable*. Finally, SPEC 2017 CPU offers so-called *speed* and *rate* metrics. For our purposes, we mostly ignore the differences and simply run the benchmarks configured for rate metrics (mainly because the runtimes are smaller) but we always run all benchmarks single-threaded.

SPEC specifies a base runtime for each benchmark and defines a *rate* as the ratio of the base runtime and the median measured runtime (this rate is a separate concept from the rate metrics). The overall suite score is then calculated as geometric mean of these ratios. The bigger the rate or score, the better it is. In the remainder of this section, we report runtimes using relative rates and their geometric means as they were measured on an AMD EPYC 7502P Processor running SUSE Linux Enterprise Server 15 SP2.

7.1 Benefits of LTO and PGO

In [Section 3, “Optimization levels and related options”](#) we recommend that HPC workloads are compiled with `-O3` and benchmarks with `-Ofast`. But it is still interesting to look at integer crunching benchmarks built with only `-O2` because that is how Linux distributions often build the programs from which they were extracted. We have already mentioned that almost the whole openSUSE Tumbleweed distribution is now built with LTO, and selected packages with PGO, and the following paragraphs demonstrate why.

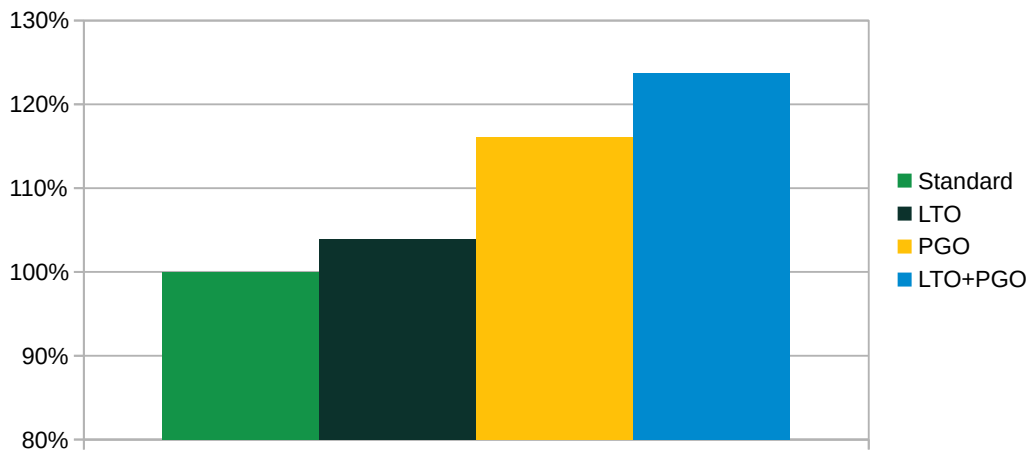


FIGURE 3: OVERALL PERFORMANCE (BIGGER IS BETTER) OF SPEC INTRATE 2017 BUILT WITH GCC 10.2 AND -O2

[Figure 3](#) shows the overall performance effect on the whole integer benchmark suite as captured by the geometric mean of all individual benchmark rates. The remarkable uplift of performance when using PGO is mostly down to much quicker `525.x264_r` (see [figure 4](#)). The reason is that, with profile feedback, GCC performs vectorization also at `-O2` and this benchmark benefits a great deal from vectorization, in practice it really should be compiled with at least `-O3`. Nevertheless, several other benchmarks also benefit from these advanced modes of operation, as can be seen on [figure 5](#).

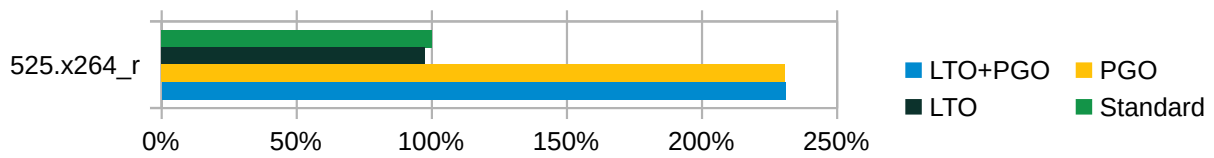


FIGURE 4: PERFORMANCE (BIGGER IS BETTER) OF 525.x264_r BUILT WITH GCC 10.2 AND -O2

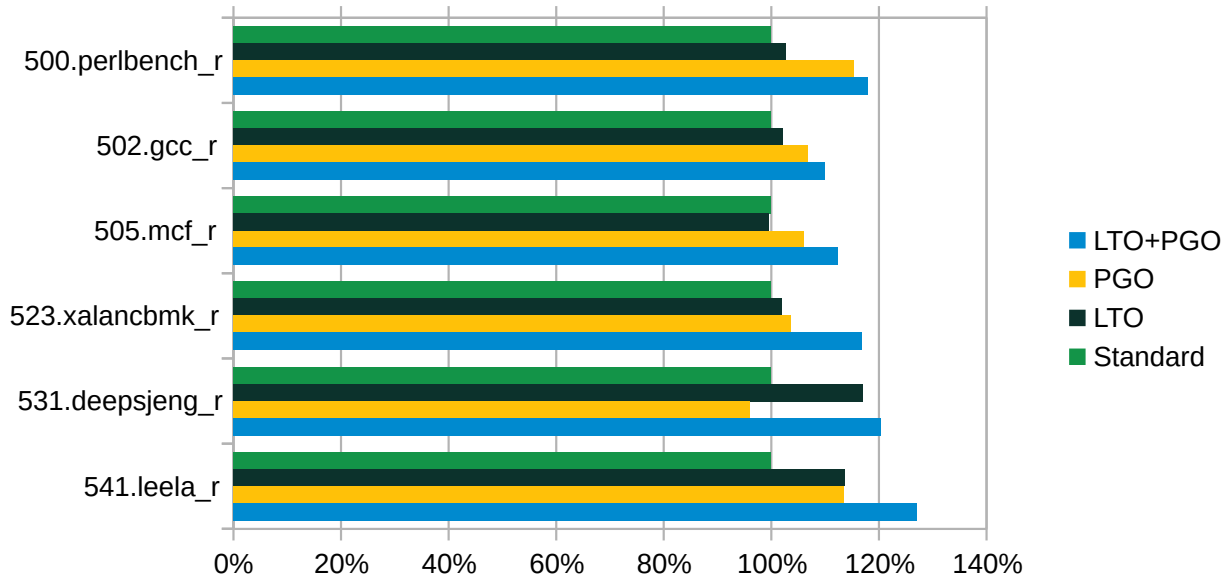


FIGURE 5: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF SELECTED INTEGER BENCHMARKS BUILT WITH GCC 10.2 AND -O2

Figure 6 shows another important reason which is the reduction of the size of the binaries (measured without debug info), which can be significant with LTO or a combination of LTO and PGO. Note that it does not depict that the size of benchmark 548.exchange2_r grew by 50% and almost 250% when built with PGO or both PGO and LTO respectively, which looks huge but the growth is from a particularly small base.

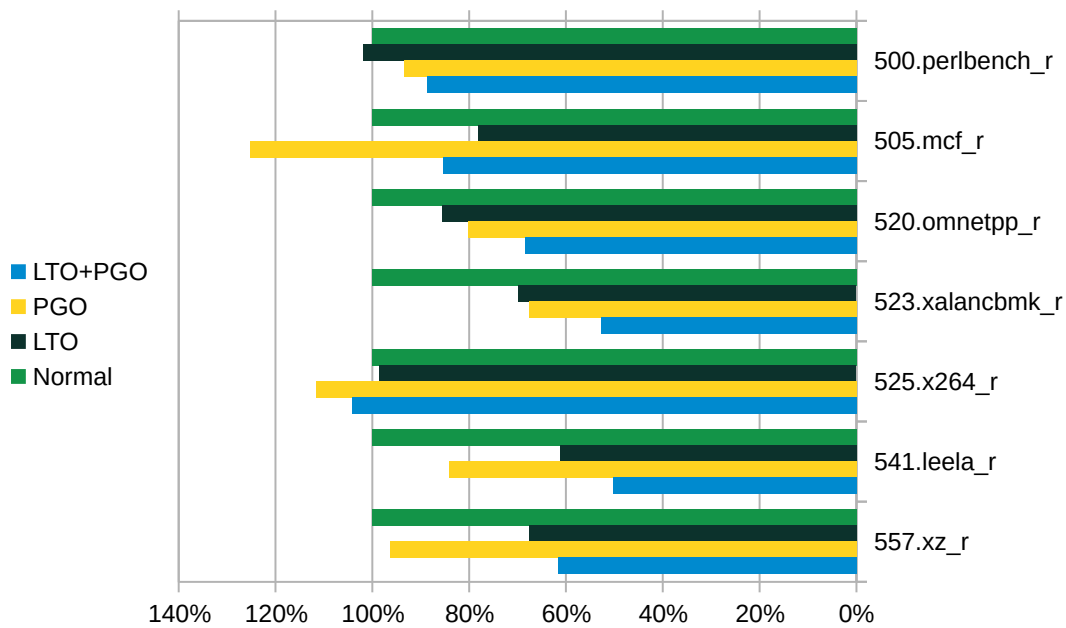


FIGURE 6: BINARY SIZE (SMALLER IS BETTER) OF SELECTED INTEGER BENCHMARKS BUILT WITH GCC 10.2 AND -O2

The runtime benefits and binary size savings can be even bigger when using the optimization level `-Ofast` and option `-march=native` to allow the compiler to take advantage of all instructions that the AMD EPYC 7502P Processor supports. [Figure 7](#) shows the respective geometric means and [figure 8](#) shows the benchmarks with the most profound effect. Even though optimization levels `-O3` and `-Ofast` are permitted to be relaxed about the final binary size, PGO and especially LTO can bring it nicely down at these levels, too. [Figure 9](#) depicts the relative binary sizes of the most affected benchmarks.

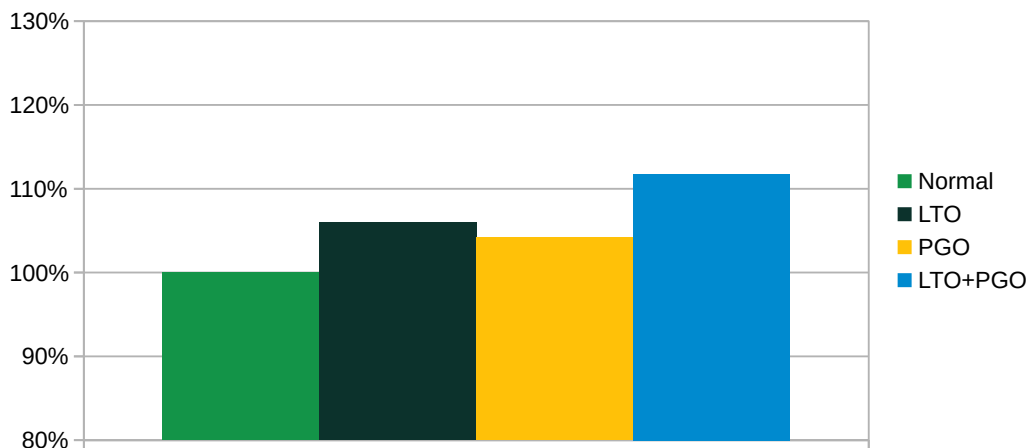


FIGURE 7: OVERALL PERFORMANCE (BIGGER IS BETTER) OF SPEC INTRATE 2017 BUILT WITH GCC 10.2 AND -OFAST

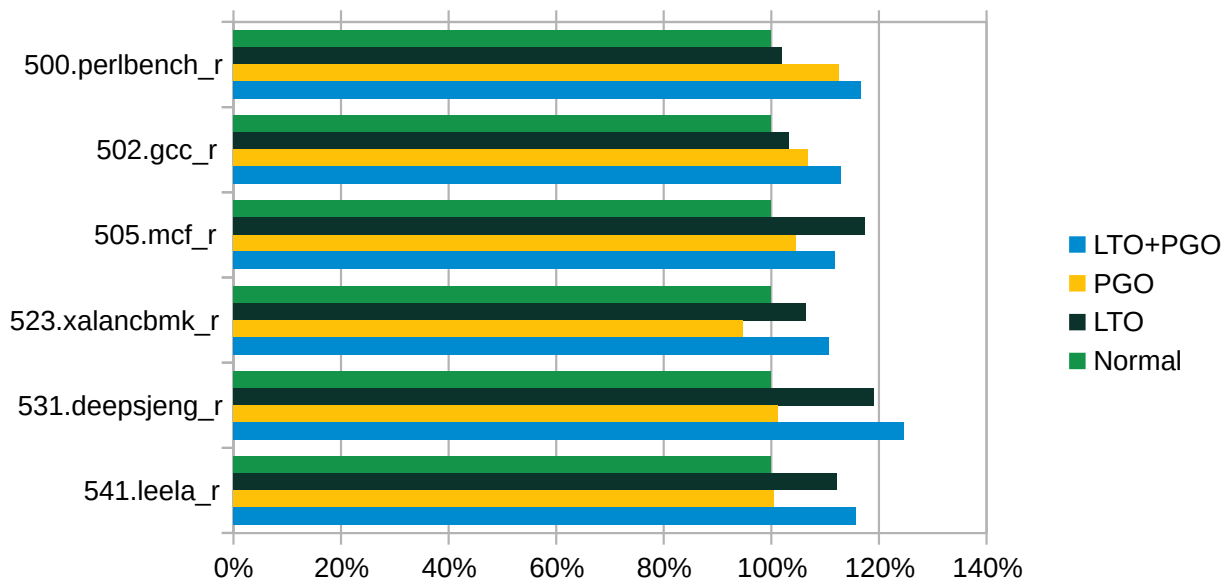


FIGURE 8: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF SELECTED INTEGER BENCHMARKS BUILT WITH GCC 10.2 AND -OFAST

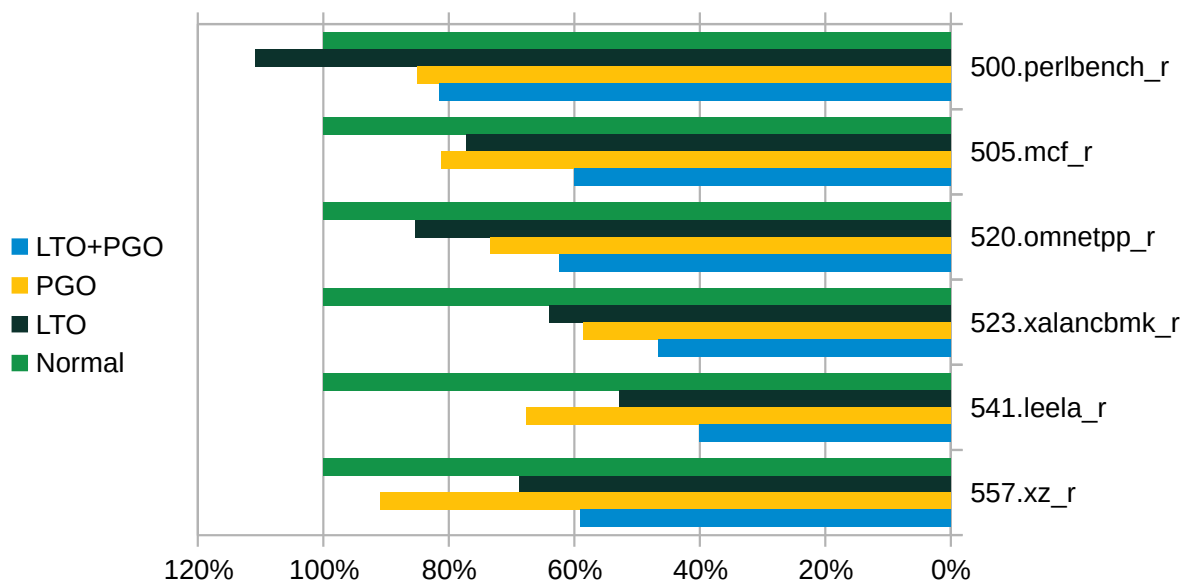


FIGURE 9: BINARY SIZE (SMALLER IS BETTER) OF SPEC INTRATE 2017 BUILT WITH GCC 10.2 AND -OFAST

LTO and PGO also bring about performance increases on several floating point benchmarks. *Figure 10* again illustrates the overall effect on the whole suite and *figure 11* the benchmarks that benefit the most. All omitted benchmarks had comparable runtimes regardless of the mode of compilation, except for `521.wrf_r` where the PGO profiling data seem to be damaged in the build process, resulting in 13% slowdown with PGO (see [GCC bug 90364](https://gcc.gnu.org/bugzilla/) ([https://gcc.gnu.org/bugzilla-](https://gcc.gnu.org/bugzilla/)

[la/show_bug.cgi?id=90364](#)). This means that the increase by 6.4% of the geometric mean with both PGO and LTO actually does not fully capture the benefits of these technologies when used correctly. Binary size savings are very similar to those achieved on integer benchmarks.

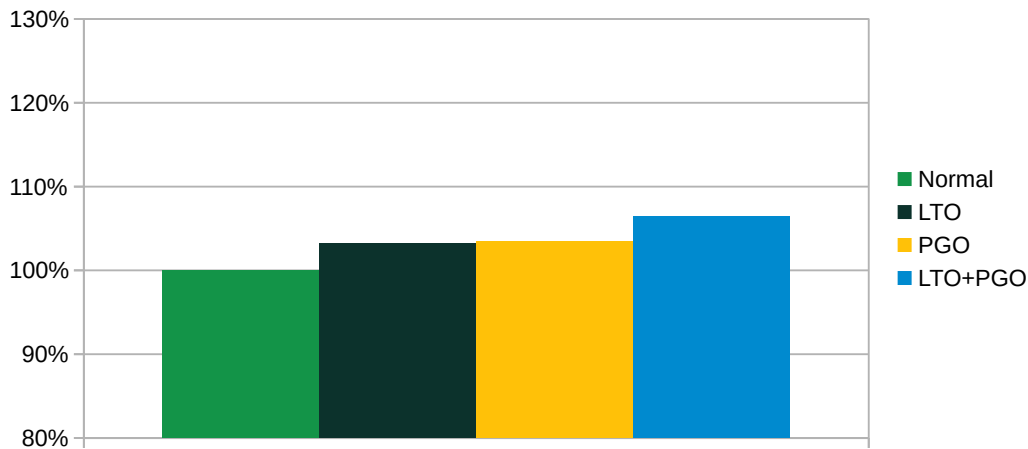


FIGURE 10: OVERALL PERFORMANCE (BIGGER IS BETTER) OF SPEC FPRATE 2017 BUILT WITH GCC 10.2 AND -OFAST

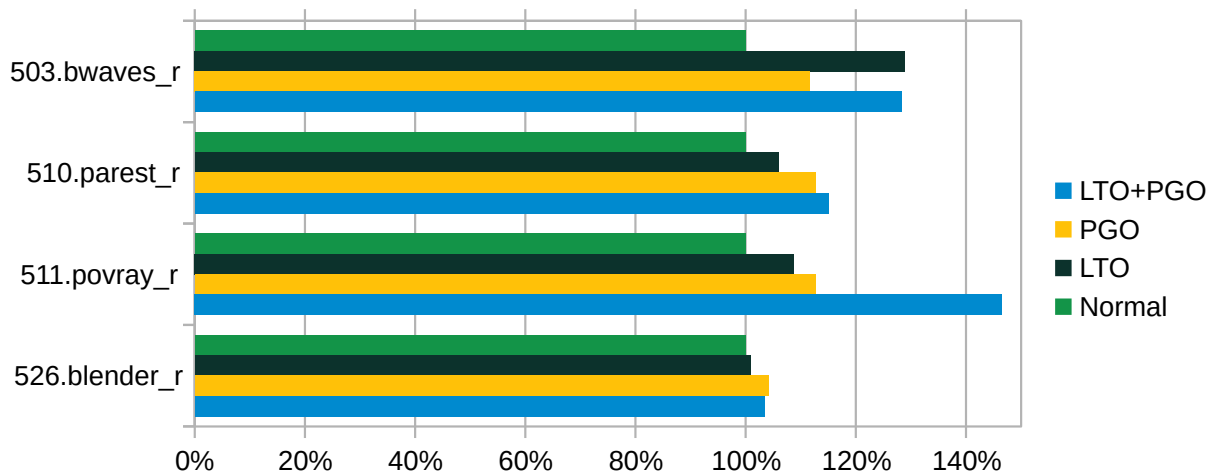


FIGURE 11: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF SELECTED FLOATING-POINT BENCHMARKS BUILT WITH GCC 10.2 AND -OFAST

7.2 GCC 10.2 compared to GCC 7.5

In previous sections we have recommended the use of GCC 10.2 from the Development Tools Module. Among other reasons, we did so because of its more powerful optimization pipeline and its support for newer CPUs. This section compares SPEC CPU 2017 obtained with GCC

7.5, which corresponds to the system compiler in SUSE Linux Enterprise Server 15, and GCC 10.2 on an AMD EPYC 7502P Processor, when all benchmarks are compiled with `-Ofast` and `-march=native`. Note that the latter option means that both compilers differ in their CPU targets because GCC 7.5 does not know the Zen 2 core. This in turn means that in large part the optimization benefits presented here are because of the fact that the newer compiler can take advantage of the full width of vector paths in the processor. Nevertheless, be aware that simply using wider vectors everywhere often backfires. GCC has made substantial advancements over the recent years, both in its vectorizer and other optimizers, to avoid such issues. It is therefore much better placed to use the extra vector width appropriately and produce code which utilizes the processor better in general.

[Figure 12](#) captures the benefits of using the modern compiler with integer workloads in the form of relative improvements of the geometric mean of the whole SPEC INTrate 2017 suite. [Figure 13](#) dives deeper and shows which particular benchmarks gained most in terms of performance. It was already mentioned that `525.x264_r` especially benefits from vectorization and therefore it is not surprising it has improved the most. `531.deepsjeng_r` is faster chiefly because it can emit better code for *count trailing zeros* (CTZ) operation which it performs frequently.

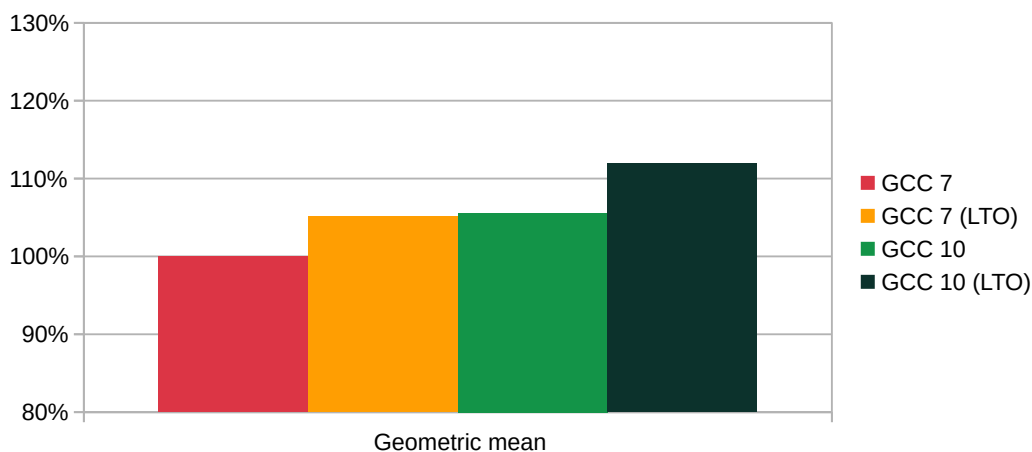


FIGURE 12: OVERALL PERFORMANCE (BIGGER IS BETTER) OF SPEC INTRATE 2017 BUILT WITH GCC 7.5 AND 10.2 (-OFAST -MARCH=NATIVE)

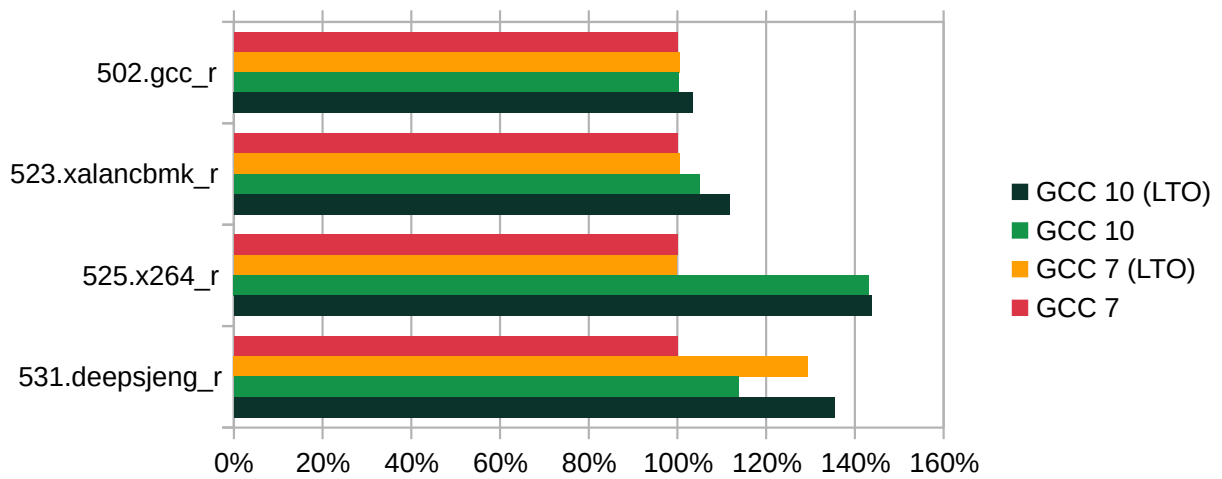


FIGURE 13: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF SELECTED INTEGER BENCHMARKS BUILT WITH GCC 7.5 AND 10.2 (-OFAST -MARCH=NATIVE)

Floating point computation tends to benefit more from vectorization advancements than integer ones. Thus it should be no surprise that the fprate suite geometric means improved even more than the integer ones, as shown in [figure 14](#). Finally, [figure 15](#) again provides a detailed look at which benchmarks contributed most to the overall score difference.

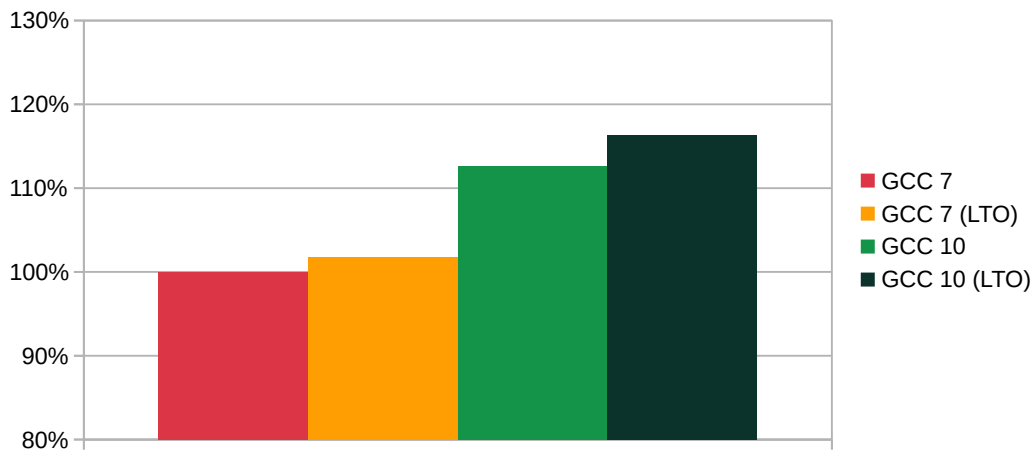


FIGURE 14: OVERALL PERFORMANCE (BIGGER IS BETTER) OF SPEC FPRATE 2017 BUILT WITH GCC 7.5 AND 10.2 (-OFAST -MARCH=NATIVE)

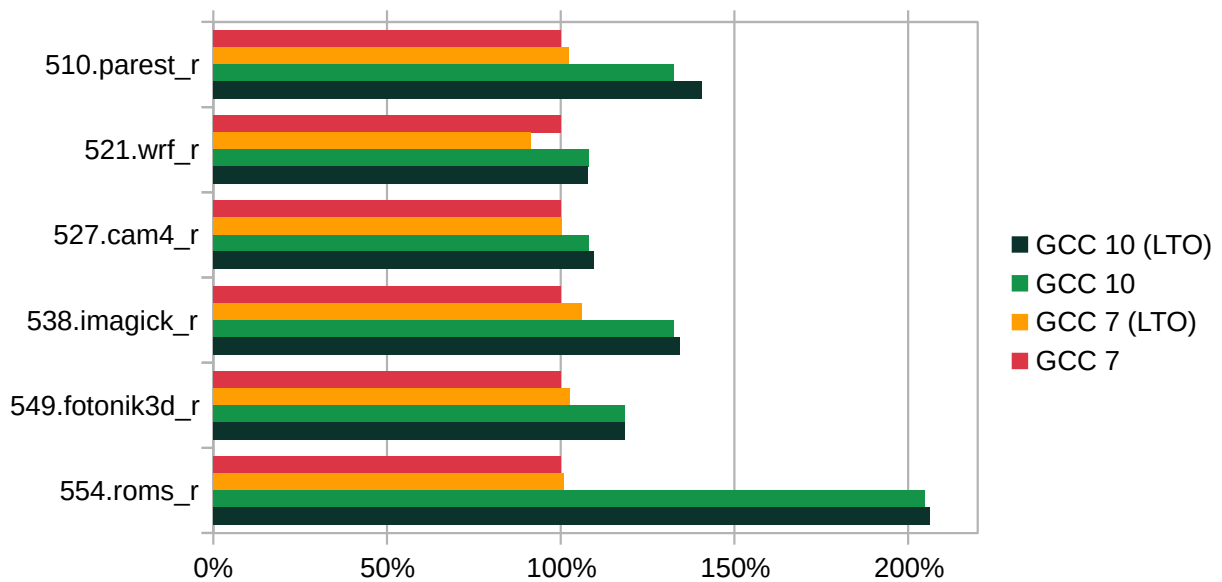


FIGURE 15: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF SELECTED FLOATING-POINT BENCHMARKS BUILT WITH GCC 7.5 AND 10.2 (-OFAST -MARCH=NATIVE)

7.3 Effects of `-ffast-math` on floating-point performance

In [Section 3, “Optimization levels and related options”](#) we have pointed out that, if you do not relax the semantics of floating point math functions even though you do not need strict adherence to all respective IEEE and/or ISO rules, you are likely to be leaving some performance on the table. This section uses the SPEC fprate 2017 test suite to illustrate how much performance that might be.

We have built the benchmarking suite with LTO, optimization level `-O3` and `-march=native` to target the native ISA of our AMD EPYC 7502P Processor and we compared its runtime score against the suite built with these options and `-ffast-math`. As you can see in [figure 16](#), the geometric means grew by over 8%, but a quick look at [figure 17](#) will tell you that some benchmarks ran 20% faster.

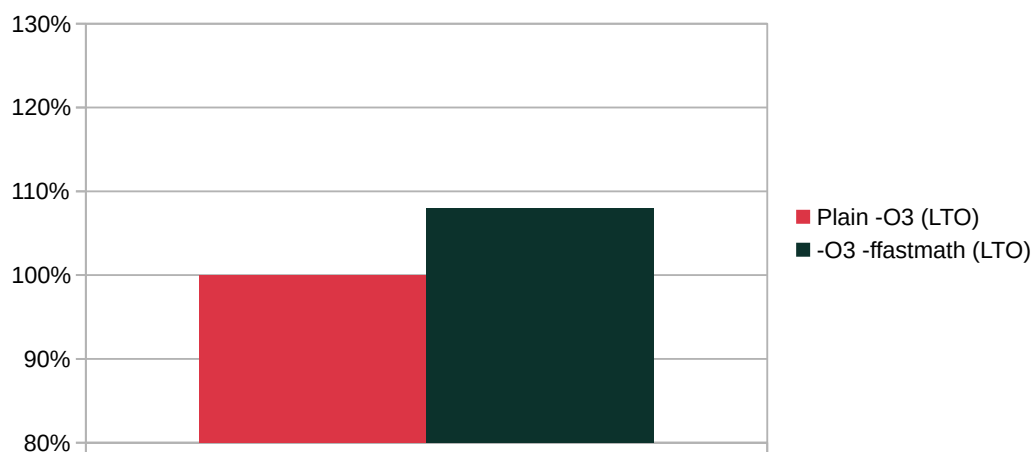


FIGURE 16: OVERALL PERFORMANCE (BIGGER IS BETTER) OF SPEC FPRATE 2017 BUILT WITH GCC 10.2 AND -O3, WITHOUT AND WITH -FFAST-MATH (-MARCH=NATIVE)

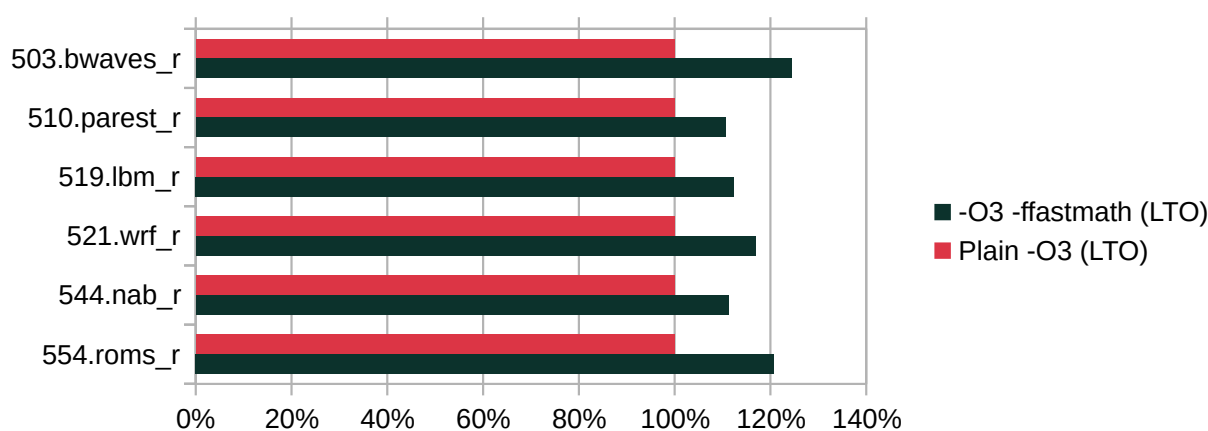


FIGURE 17: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF SELECTED FLOATING-POINT BENCHMARKS BUILT WITH GCC 10.2 AND -O3, WITHOUT AND WITH -FFAST-MATH (-MARCH=NATIVE)

7.4 Comparison with other compilers

The toolchain team at SUSE regularly uses the SPEC CPU 2017 suite to compare the optimization capabilities of GCC with other compilers, mainly Intel ICC and LLVM/Clang. In the final section of this document we will share how the Development Module compiler stands compared to these competitors on SUSE Linux Enterprise Server 15 SP2. Before we start, we should emphasize that the comparison has been carried out by people who have much better knowledge of GCC than the other compilers and are not “unbiased”. Also, keep in mind that everything we explained previously about how we carry out the measurements and patch the benchmarks also applies to this section. On the other hand, the results often guide our work and therefore we strive to be accurate.

Even though ICC is not intended as a compiler for AMD processors, it is known for its high-level optimization capabilities, especially when it comes to vectorization. Thus we believe comparisons with it are useful even on AMD CPUs. We have therefore compiled SPEC suite with ICC 19.1 with options `-Ofast` and `-march=core-avx2` option and compared the runtimes on the AMD EPYC 7502P Processor with binaries produced with GCC 10.2 with `-Ofast` and `-march=native`.

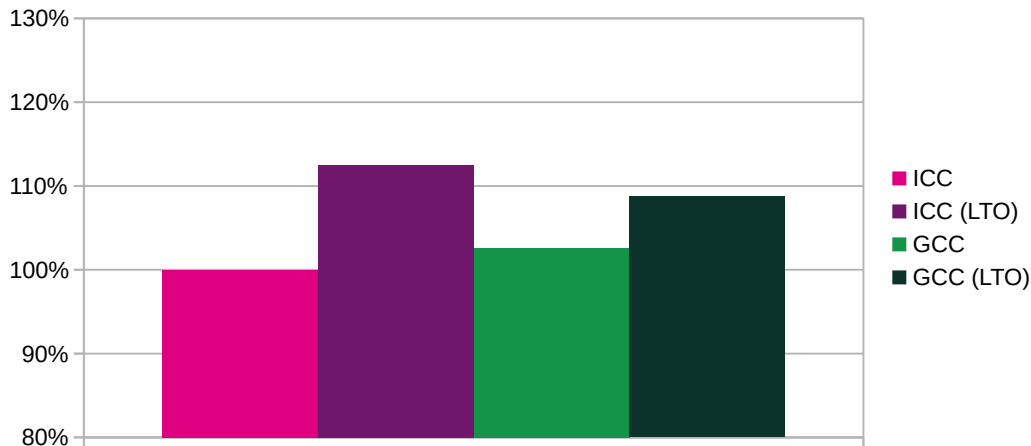


FIGURE 18: OVERALL PERFORMANCE (BIGGER IS BETTER) OF SPEC INTRATE 2017 BUILT WITH ICC 19 AND GCC 10.2

Figure 18 depicts the respective geometric means of rates of all the integer benchmarks when compiled with the two compilers, without and with LTO, relative to performance of ICC without LTO. GCC manages to achieve better score in the traditional compilation mode but loses to ICC with LTO.

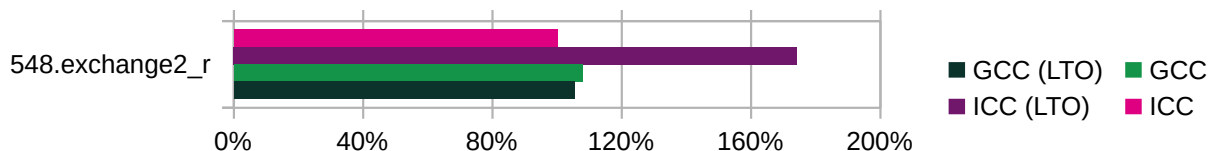


FIGURE 19: PERFORMANCE (BIGGER IS BETTER) OF 548.exchange2_r BUILT WITH ICC 19 AND GCC 10.2

By far the main reason is the performance of benchmark `548.exchange2_r` (see figure 19) where ICC achieves 174% compared to itself without LTO. `548.exchange2_r` is the only Fortran benchmark in the integer suite and its only hot function contains a recursive call in a deep loop nest which poses a problem for many loop optimizers. Furthermore, it can be made faster by inter-procedural constant propagation and cloning if performed to an extent that would typically be excessive. When GCC 10 is instructed to do that with the parameters `--param ipa-cp-eval-`

`threshold=1 --param ipa-cp-unit-growth=80`, it achieves a score on par with ICC with LTO, even without LTO. We do not expect users to compile with such options, and are working to enable the transformation with only `-Ofast` in the next release of GCC.

Figure 20 shows other integer benchmarks which performed differently when compiled with ICC and GCC. As you can see, although GCC still loses in a few cases, it wins others and is generally competitive against ICC on integer benchmarks. Unfortunately, the same cannot be said for floating point benchmarks, at least not on SUSE Linux Enterprise Server 15 SP2. This is primarily because the performance is also very much influenced by the math library. ICC carries its own, while GCC depends on the one from `glibc` and the version in SUSE Linux Enterprise Server 15 SP2 unfortunately contains slow implementations of some key functions, such as those used to compute complex exponential. Comparing geometric means, ICC beats GCC by 7.7% without LTO and by 10.5% with LTO. When we compared the same compilers on openSUSE Tumbleweed with `glibc` version 2.31, ICC was still faster but the differences were 0.1% and 3.7% respectively. SUSE is working to upgrade the `glibc` version in SUSE Linux Enterprise Server 15 SP 3 to exactly this version, which will lead to significant speed-ups. Meanwhile, users can consider using optimized `libm` from AMD ⁶.

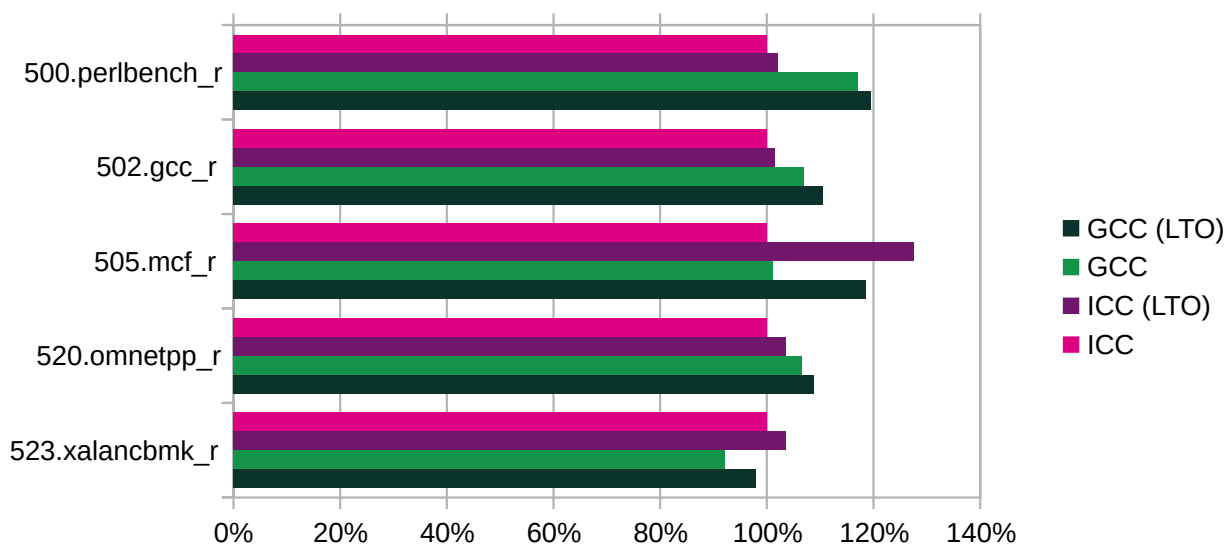


FIGURE 20: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF SELECTED INTEGER BENCHMARKS BUILT WITH ICC 19 AND GCC 10.2

Comparisons with LLVM/Clang 11.0 are incomplete because although the newest release of this compiler suite already contains a Fortran compiler, we could not compile the SPEC benchmarks with it. Nevertheless, we have built the `clang` and `clang++` compilers from sources obtained

⁶ Available at <https://developer.amd.com/amd-aocl/amd-math-library-libm/>

from the official git repository (tag `llvmorg-11.0.0`), used it to compile the SPEC CPU 2017 suite with `-Ofast` and `-march=native` and compared the performance against the suites built with GCC 10.2 with the same options. When using Clang's LTO to compile SPEC, we selected the *full* variant.

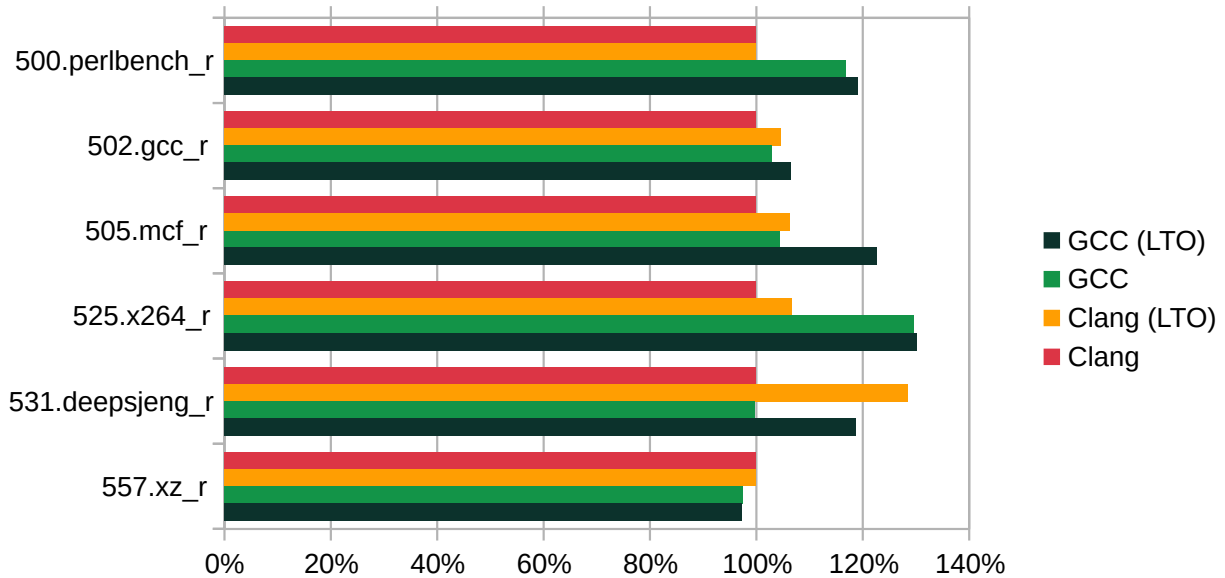


FIGURE 21: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF SELECTED INTEGER BENCHMARKS BUILT WITH CLANG 11 AND GCC 10.2

Figure 21 captures the integer benchmarks with significant performance differences. Clang scores one notable win with `531.deepsjeng_r`, which we can mostly pin down to an if-conversion transformation which we are not certain is generally profitable, and is slightly faster with `557.xz_r`. On the other hand, GCC wins on the benchmark derived from itself and is remarkably faster on `500.perlbench_r` and `505.mcf_r`. If we omitted the one Fortran benchmark in calculating the overall geometric mean, GCC would win by 13% without LTO and by 5% with it. The floating point benchmark suite contains many more Fortran benchmarks. For more information, refer to *figure 22* which captures the significant differences in scores of benchmarks written entirely in C/C++. Overall, GCC still has a slight lead but Clang manages clear wins on three of the benchmarks.

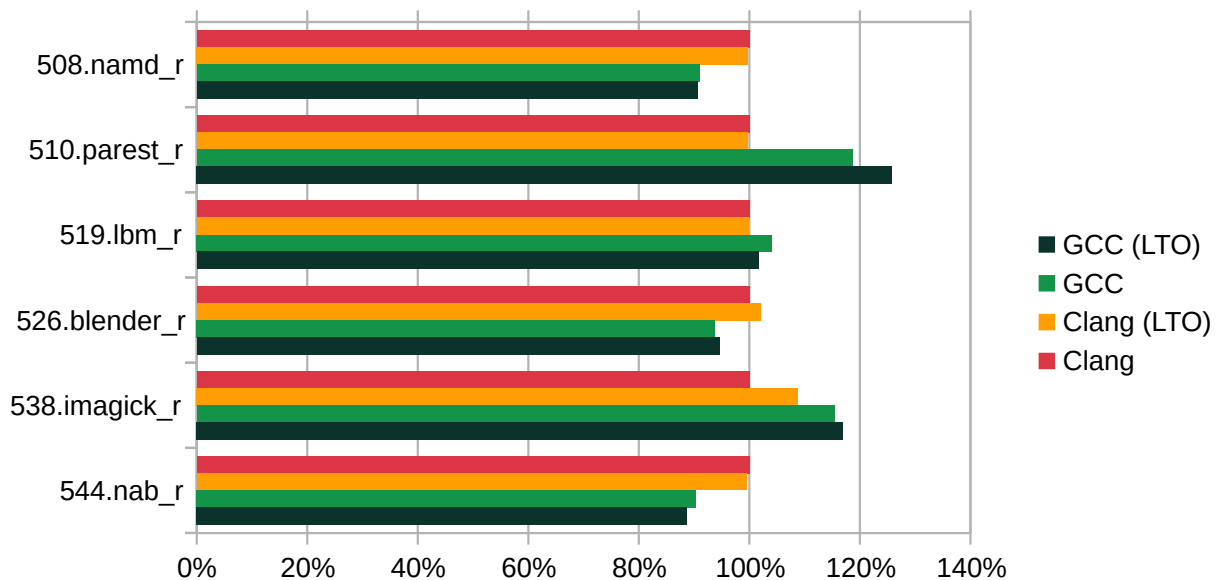


FIGURE 22: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF SELECTED INTEGER BENCHMARKS BUILT WITH CLANG 11 AND GCC 10.2

8 Performance evaluation: Mozilla Firefox

Benchmarks such as SPEC CPU 2017 are very useful to gauge how compilers optimize a particular type of computation when it is embedded in a small- to medium-size project. Nevertheless, many real-world applications are much bigger. This fact alone presents a significant challenge to optimizing compilers. When there are too many opportunities for a transformation which has a potential to increase performance but comes with a substantial code size growth, such as inlining, the compiler simply cannot go ahead and proceed with all of them because the final size of the binary would be unacceptably big. To monitor how well GCC manages to optimize big real applications, we regularly build and benchmark the popular Mozilla Firefox browser⁷. This section summarizes our latest findings.

We use Mozilla's own Treeherder project to build Firefox with different compilers and options, and their Talos and Perfherder infrastructure to evaluate their performance. The evaluation framework compares different binaries using several benchmarks which it can run sufficiently many times to eliminate their noise. The CPUs used in Mozilla Talos to benchmark x86_64 Linux

⁷ For older results, see for example [this blog post \(http://hubicka.blogspot.com/2019/05/gcc-9-link-time-and-inter-procedural.html\)](http://hubicka.blogspot.com/2019/05/gcc-9-link-time-and-inter-procedural.html).

builds are Intel E3-1585L v5. But as shown further on, we measured similar results on an AMD Ryzen 7 3700X 8-Core Processor, although on simpler benchmarks and with much fewer data points.

Another important difference to traditional benchmarks is that most of Mozilla Firefox is implemented in a shared object file compiled as *position independent code*. This limits code generation in some situations, such as data segment relocations or global variable addressing⁸. Note that all Firefox binaries, regardless of the compiler used, were built with the following options in addition to any that are explicitly called out:

```
-fno-sized-deallocation -fno-aligned-new -fno-strict-aliasing -fPIC -fno-exceptions  
-fno-rtti -fno-math-errno -fno-exceptions -fno-fomit-frame-pointer
```

Unfortunately, Mozilla Firefox is one of the projects which has elected to use `-fno-strict-aliasing` rather than fix aliasing violations in their code, despite the performance implications it has.

When we compare any two Firefox binaries, we start by looking at their size, excluding debug information sections. The size is often important by itself, such as when applications are updated over slower networks, but it is also a sign of how well the compiler can distinguish performance sensitive and rarely executed pieces of a project. To evaluate runtime performance, in this document we selected four benchmarks. The chief benchmark among them is *tp5o*⁹ which measures the time it takes Firefox to load the tp5 Web page test set which contains a collection of 151 pages originally picked from a list of 500 most popular ones in 2011. We also look at *tp5o responsiveness test*¹⁰ measuring how responsive Firefox is while carrying out a non-trivial workload. The last Talos test we have chosen to focus on is called *perf reftest singletons*¹¹. It is a micro-benchmark that loads simple HTML pages and then measures basic manipulation with their elements, such as adding a row to a table. This benchmark itself is part of the train run in PGO builds, and thus the PGO binaries should be well trained for it. Finally, we have used *Speedometer 2.0*¹² to cross-check selected results from Talos on an AMD Ryzen 7 3700X 8-Core Processor. Speedometer is a rather simple benchmark which simulates user actions for adding, completing, and removing to-do items using DOM APIs in different ways. Speedometer is also part of the profile train run.

⁸ For further details, see “How To Write Shared Libraries” by Ulrich Drepper (<https://akkadia.org/drepper/dsohowto.pdf>)⁷.

⁹ <https://firefox-source-docs.mozilla.org/testing/perfdocs/talos.html>⁷

¹⁰ <https://firefox-source-docs.mozilla.org/testing/perfdocs/talos.html#tp5o-scroll>⁷

¹¹ <https://firefox-source-docs.mozilla.org/testing/perfdocs/talos.html#perf-reftest-singletons>⁷

¹² <https://browserbench.org/Speedometer2.0/>⁷

8.1 Effects of -O3 compared to -O2 and of LTO and PGO

Mozilla Firefox is a large application. The code size should therefore definitely play a role when deciding how to compile it. On the other hand, a Web browser is also likely to be a substantial part of a typical desktop workload. Thus gains in performance can easily justify binary size increases. As a consequence, Firefox is typically built with `-O3`. [Figure 23](#) depicts the sizes of the Firefox `libxul` library, which contains the bulk of the browser, when built with GCC 10 using the Mozilla Treeherder infrastructure with the optimization levels and modes most discussed in this document. Again, you can see that LTO can reduce the code size to an extent that more than offsets the difference between `-O3` and `-O2`. Note that, since a big portion of Firefox is written in Rust and the whole program analysis is limited to the parts written in C++, the LTO benefits are smaller than the typical case, in terms of both size and performance. Work on the Rust GCC front-end has started only recently but we hope that we will overcome this limitation. Nevertheless, as demonstrated throughout this case study, LTO combined with PGO is by far the best option, not only in code size comparison but also in any other measurement.

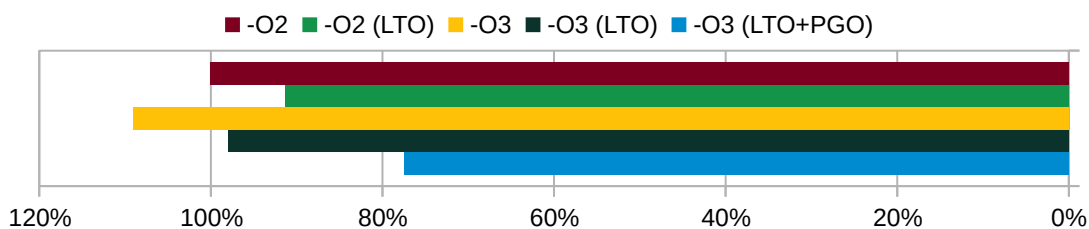


FIGURE 23: CODE SIZE (SMALLER IS BETTER) OF FIREFOX BINARIES BUILT WITH GCC 10.2 WITH DIFFERENT OPTIONS

When not employing neither LTO nor PGO, the performance difference between `-O2` and `-O3` is barely measurable on the `tp5o` benchmark. This is mostly because the Skia rendering library they both use has been compiled with `-O3`. The difference is bigger, but still rather small on other benchmarks, unlike switching on LTO which has noticeable effect even at `-O2`. LTO and `-O3` adds another bit of performance. But the real speed up comes only when PGO is added into the formula. This observation holds for both the data measured using the Talos and Perfherder systems ([figure 24](#)) and speedometer results we obtained manually on an AMD Ryzen 7 3700X 8-Core Processor ([figure 25](#)). This is especially remarkable when you consider that the binary is more than 20% smaller than a simple `-O2` build. Note that such aggressive size shrinking comes at a modest cost which we will discuss in [section 8.4](#).

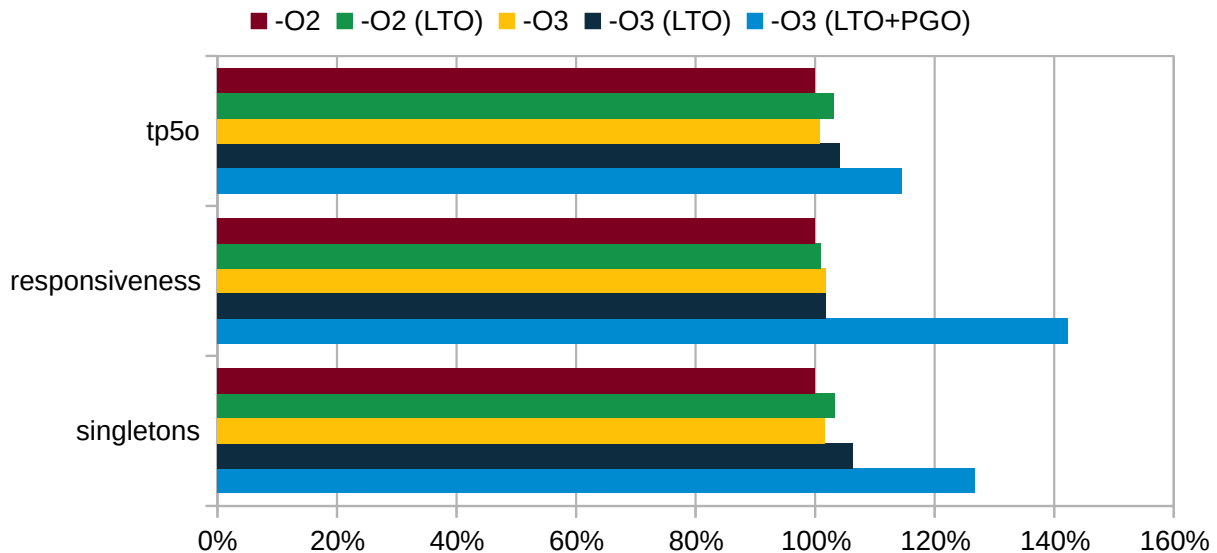


FIGURE 24: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF FIREFOX BUILT WITH GCC 10.2 WITH DIFFERENT OPTIONS, RUNNING ON MOZILLA TALOS INFRASTRUCTURE

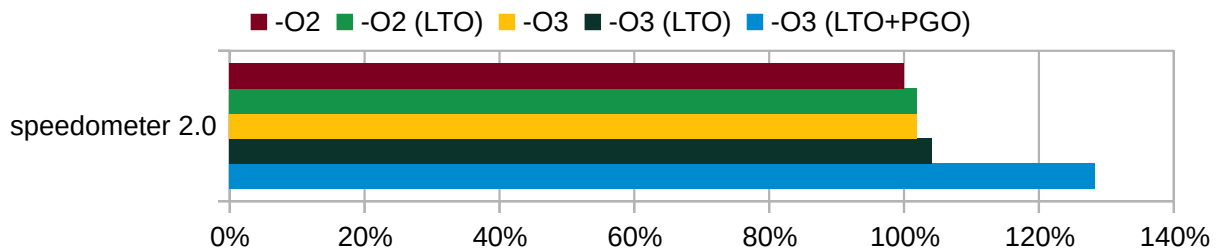


FIGURE 25: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF FIREFOX BUILT WITH GCC 10.2 WITH DIFFERENT OPTIONS, RUNNING SPEEDOMETER 2.0 ON AN AMD RYZEN 7 3700X 8-CORE PROCESSOR

8.2 GCC 10.2 compared to GCC 7.5

[Section 7.2](#) demonstrates that GCC 10 produces much faster code when targeting modern CPUs such as the AMD EPYC 7502P Processor, most often because it can take advantage of vector instructions of the new hardware. This section aims to show that GCC 10 produces faster code also when emitting instructions for any `x86_64` system and running on processors that are not as new. We have compared Firefox binaries built with GCC 7.5 and 10.2 using the least and the most powerful methods described in the previous section, `-O2` with the traditional build, and `-O3` with both LTO and PGO. Their respective sizes are very similar but the one created with GCC 10 has always performed noticeably better. In the `tp5o` responsiveness benchmark, the simple `-O2` build was 18% faster. Even when using the full power of `-O3`, PGO and LTO, the binary built with the more modern compiler ran 8% faster (see [figure 26](#)).

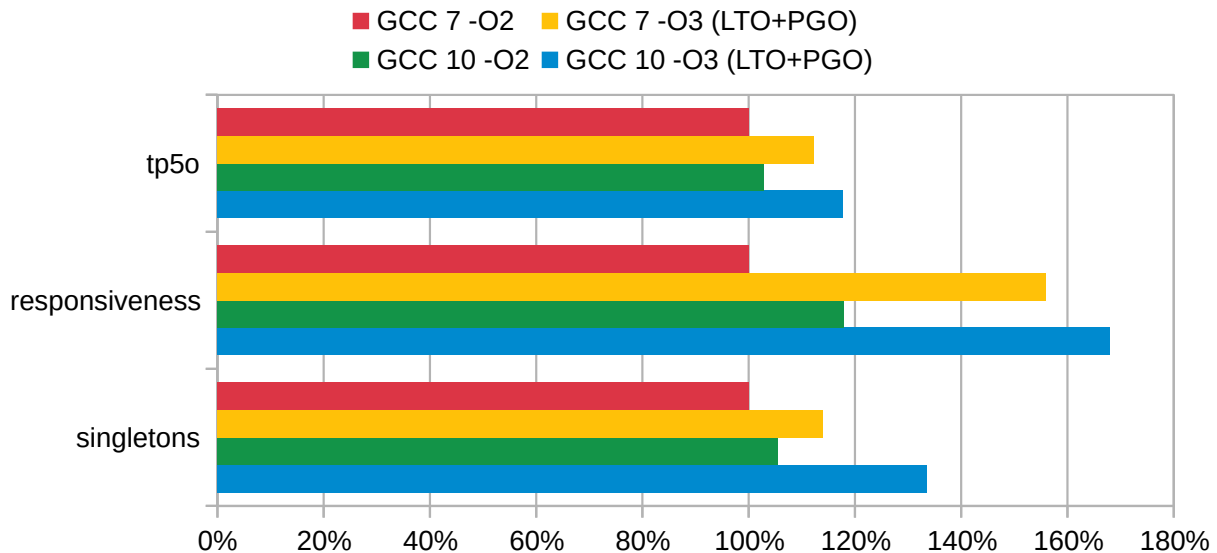


FIGURE 26: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF FIREFOX BUILT WITH GCC 7.5 AND 10.2, RUNNING ON MOZILLA TALOS INFRASTRUCTURE

8.3 Simple comparisons with Clang 11

To evaluate how GCC is doing compared to other compilers, we regularly compare the Firefox binaries produced by GCC to those emitted by LLVM/Clang, which is currently the preferred compiler by the team at Mozilla. Before we proceed, we should emphasize that the authors of this document are not nearly as familiar with LLVM/Clang as they are with GCC, and that they are not “unbiased”. On the other hand, our findings in such comparisons guide our own future work and therefore we strive to be accurate.

Because the notion of compilation levels is somewhat different in both of these compilers, we have focused on evaluating how they build Firefox using `-O3` in the traditional way, when using LTO (in Clang's case its *thin* variant), and when using both PGO and LTO. Possibly the most striking differences are between code sizes of the results (see [figure 27](#)). When using plain `-O3`, GCC produces a 5% larger binary than Clang. With LTO, GCC manages to shrink the code size to more than undo this difference, whereas Clang uses the extra cross-module opportunities to grow the code by 6%. The difference is even more pronounced with PGO in addition to LTO, which barely enables Clang to produce a binary that is smaller than using neither of them, while GCC creates the smallest binary of all, 29% smaller than when using plain `-O3`. Nevertheless, [section 8.4](#) discusses important caveats regarding the PGO cases.

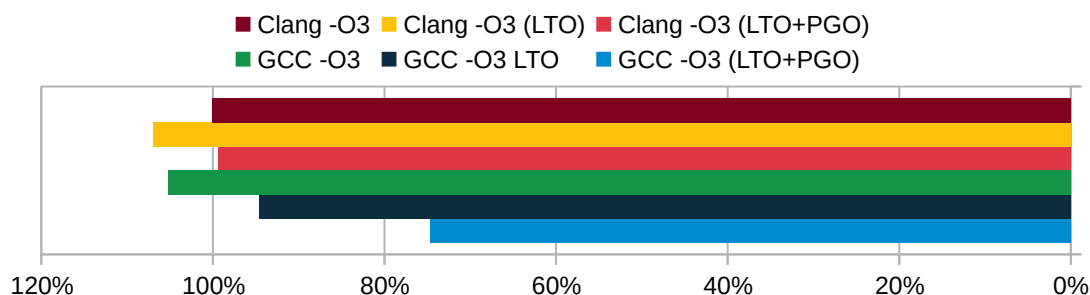


FIGURE 27: CODE SIZE (SMALLER IS BETTER) OF FIREFOX BINARIES BUILT WITH GCC 10.2 AND CLANG 11

Runtime comparisons can be found in [figure 28](#) (measured on Talos) and [figure 29](#) (timed on an AMD Ryzen 7 3700X 8-Core Processor). In the tp50 benchmark and LTO plus PGO configurations, the aggressive optimizing for size means that GCC is slower, albeit measurably. Because the singletons test is part of the train run, GCC does not have this issue when running this benchmark (see [section 8.4](#) for more details on the train run). In the responsiveness test, GCC wins with profile feedback but lags a little without it. Because Mozilla Firefox is a large application that is intended to be built with PGO, we consider the results with profile feedback the most important. Running speedometer on the AMD machine leads to similar results.

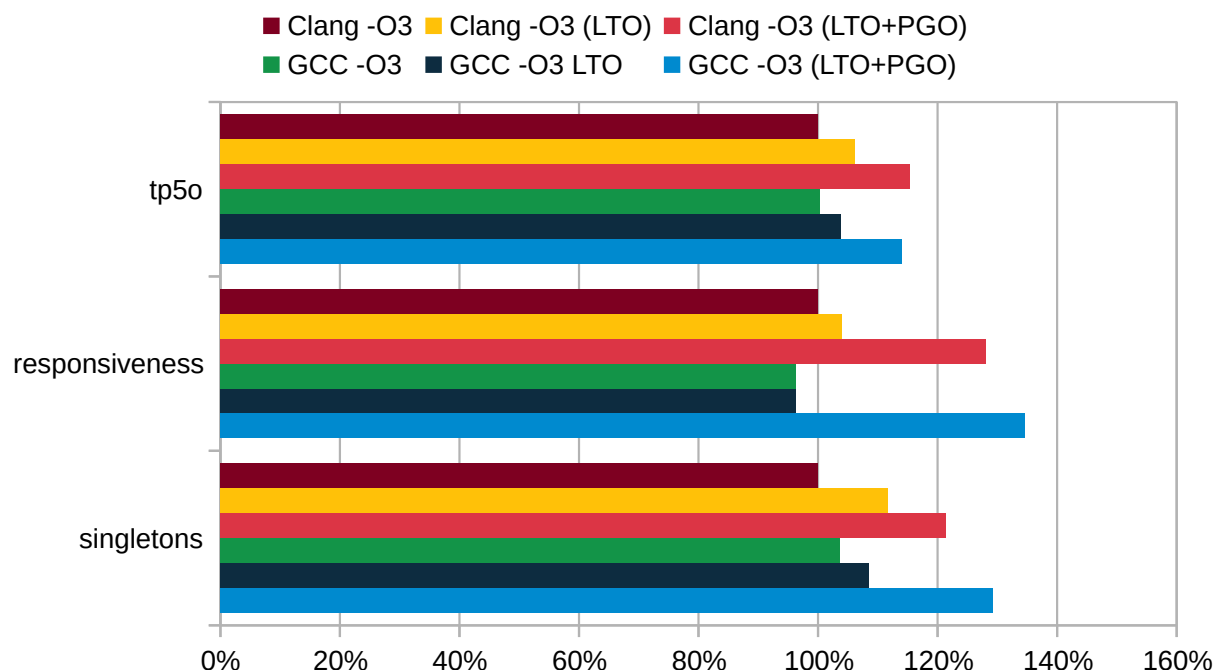


FIGURE 28: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF FIREFOX WITH GCC 10.2 AND CLANG 11, RUNNING ON MOZILLA TALOS INFRASTRUCTURE

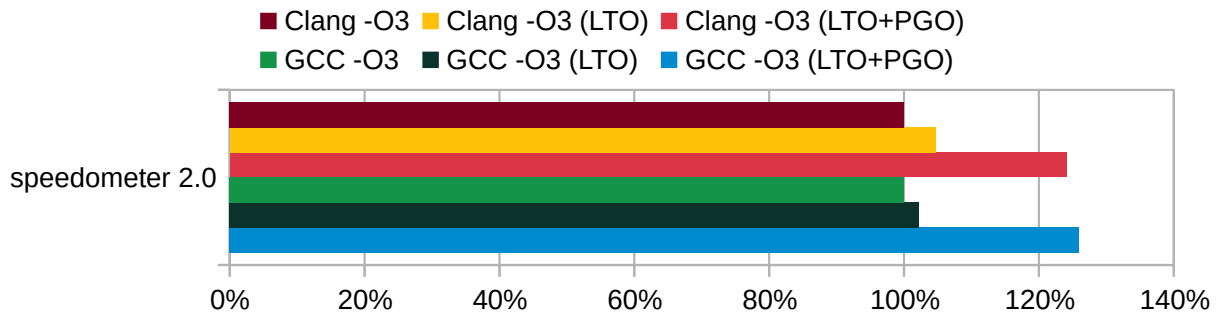


FIGURE 29: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF FIREFOX WITH GCC 10.2 AND CLANG 11, RUNNING SPEEDOMETER 2.0 ON AN AMD RYZEN 7 3700X 8-CORE PROCESSOR

8.4 Further examination of PGO performance

As pointed out previously in this document, identifying frequently executed (hot) parts of a project poses a big challenge for a compiler which does not have profile feedback at its disposal. This problem can be particularly hard and important at the same time when the compiled project is large. That is why PGO optimized builds of Firefox are so much faster, despite the fact that the train run is a less than perfect model of the actual use patterns. After all, it exercises less than 3% of the source code as measured by the size of the final binary, it consists of:

- SunSpider (an old JavaScript benchmark),
- perf reftest singletons (described above),
- speedometer (also described above),
- Web Audio benchmark,
- crypto-otp (JavaScript one time password generator),
- blueprint (rendering a specific Firefox UI), and
- 3d-thingy (a simple animation).

Because the train run coverage is relatively small and does not exercise all important code paths, it is important to use the `-fprofile-partial-training` as described in [section 6](#). This option was used in the PGO builds described in [section 8.3](#) and yet we had to conclude that GCC was perhaps a bit too aggressive when optimizing for code size and that Clang was faster on the tp5o benchmark as a result. To verify this hypothesis, we have also built Firefox with `-O3`, PGO, LTO, partial training and an extra option `--param hot-bb-count-ws-permille=1000`,

which instructs GCC to consider even smaller portion of the project as **cold**. Afterward, we took this build to prepare a final comparison with Clang, but this time using the reference Clang options that Firefox developers use to build Firefox ¹³. Like in [section 8.3](#) we use Clang `-O3` as the reference point, so that 100% corresponds to the identical values like in the graphics there.

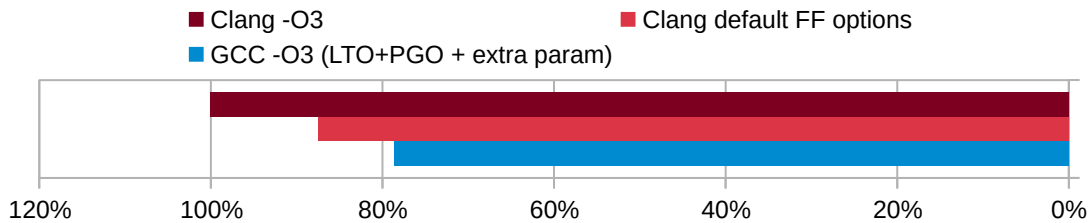


FIGURE 30: CODE SIZE (SMALLER IS BETTER) OF FIREFOX BINARIES BUILT WITH GCC 10.2 AND CLANG 11 USING COMPILATION OPTIONS TAILORED FOR FIREFOX

As you can see on [figure 30](#), the tailored Clang build is smaller than before even though the base optimization level is still `-O3` and both PGO and LTO have been employed. GCC with the special parameter produced a bigger binary than without it, but it is still the smallest out of the three. This difference turns out to be enough for GCC to achieve the best performance on all three considered Talos benchmarks and still lead significantly in the tp5o responsiveness test, as shown on [figure 31](#).

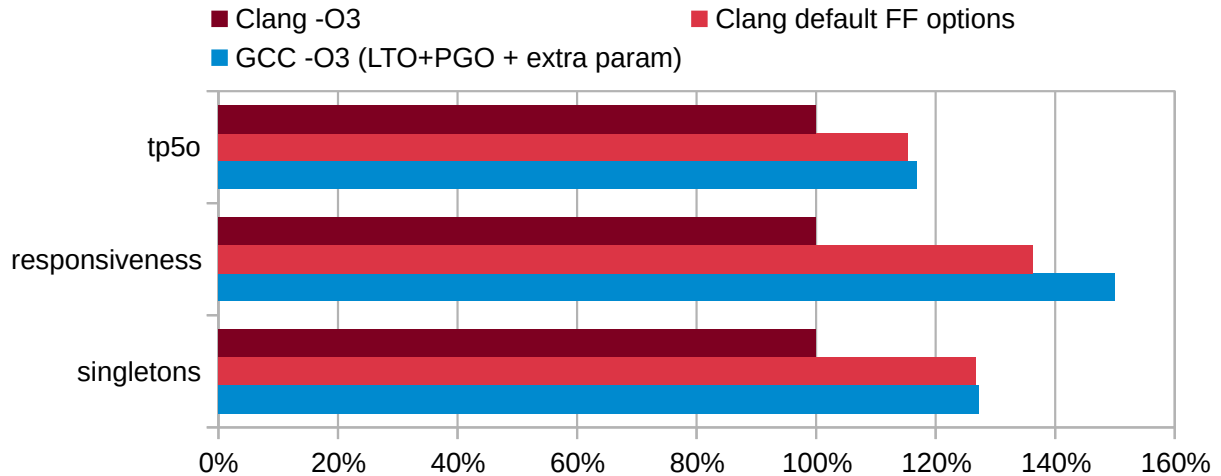


FIGURE 31: RUNTIME PERFORMANCE (BIGGER IS BETTER) OF FIREFOX WITH GCC 10.2 AND CLANG 11 USING COMPILATION OPTIONS TAILORED FOR FIREFOX, RUNNING ON MOZILLA TALOS INFRASTRUCTURE


¹³ The Firefox team at Mozilla uses Clang with compilation options `-O3 -flto=thin -fexperimental-new-pass-manager -Wl,--gc-sections -Wl,-plugin-opt=import-instr-limit=10 -Wl,-plugin-opt=new-pass-manager -Wl,-plugin-opt=import-hot-multiplier=30 -ffunction-sections -fdata-sections`

Finally, we want to point out that the profiling overhead itself can cause large, complicated, and multithreaded applications to behave very differently from their usual behavior, which again means that the train run might be a bad representation of actual runs. The bottom line however is that, even when the training workload is imperfect and has relatively small coverage, which is unavoidable in a case of a GUI application with many features, it can still help the compiler to achieve remarkable performance increases.

9 Legal notice

Copyright ©2006-2025 SUSE LLC and contributors. All rights reserved.

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 or (at your option) version 1.3; with the Invariant Section being this copyright notice and license. A copy of the license version 1.2 is included in the section entitled “GNU Free Documentation License”.

SUSE, the SUSE logo and YaST are registered trademarks of SUSE LLC in the United States and other countries. For SUSE trademarks, see <http://www.suse.com/company/legal/> . Linux is a registered trademark of Linus Torvalds. All other names or trademarks mentioned in this document may be trademarks or registered trademarks of their respective owners.

Documents published as part of the **SUSE Best Practices** series have been contributed voluntarily by SUSE employees and third parties. They are meant to serve as examples of how particular actions can be performed. They have been compiled with utmost attention to detail. However, this does not guarantee complete accuracy. SUSE cannot verify that actions described in these documents do what is claimed or whether actions described have unintended consequences. SUSE LLC, its affiliates, the authors, and the translators may not be held liable for possible errors or the consequences thereof.

Below we draw your attention to the license under which the articles are published.

GNU Free Documentation License

Copyright (C) 2000, 2001, 2002 Free Software Foundation, Inc. 51 Franklin St, Fifth Floor, Boston, MA 02110-1301 USA. Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

0. PREAMBLE

The purpose of this License is to make a manual, textbook, or other functional and useful document "free" in the sense of freedom: to assure everyone the effective freedom to copy and redistribute it, with or without modifying it, either commercially or non-commercially. Secondly, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

This License is a kind of "copyleft", which means that derivative works of the document must themselves be free in the same sense. It complements the GNU General Public License, which is a copyleft license designed for free software.

We have designed this License to use it for manuals for free software, because free software needs free documentation: a free program should come with manuals providing the same freedoms that the software does. But this License is not limited to software manuals; it can be used for any textual work, regardless of subject matter or whether it is published as a printed book. We recommend this License principally for works whose purpose is instruction or reference.

1. APPLICABILITY AND DEFINITIONS

This License applies to any manual or other work, in any medium, that contains a notice placed by the copyright holder saying it can be distributed under the terms of this License. Such a notice grants a world-wide, royalty-free license, unlimited in duration, to use that work under the conditions stated herein. The "Document", below, refers to any such manual or work. Any member of the public is a licensee, and is addressed as "you". You accept the license if you copy, modify or distribute the work in a way requiring permission under copyright law.

A "Modified Version" of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifications and/or translated into another language.

A "Secondary Section" is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publishers or authors of the Document to the Document's overall subject (or to related matters) and contains nothing that could fall directly within that overall subject. (Thus, if the Document is in part a textbook of mathematics, a Secondary Section may not explain any mathematics.) The relationship could be a matter of historical connection with the subject or with related matters, or of legal, commercial, philosophical, ethical or political position regarding them.

The "Invariant Sections" are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that says that the Document is released under this License. If a section does not fit the above definition of Secondary then it is not allowed to be designated as Invariant. The Document may contain zero Invariant Sections. If the Document does not identify any Invariant Sections then there are none.

The "Cover Texts" are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the Document is released under this License. A Front-Cover Text may be at most 5 words, and a Back-Cover Text may be at most 25 words.

A "Transparent" copy of the Document means a machine-readable copy, represented in a format whose specification is available to the general public, that is suitable for revising the document straightforwardly with generic text editors or (for images composed of pixels) generic paint programs or (for drawings) some widely available drawing editor, and that is suitable for input to text formatters or for automatic translation to a variety of formats suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup, or absence of markup, has been arranged to thwart or discourage subsequent modification by readers is not Transparent. An image format is not Transparent if used for any substantial amount of text. A copy that is not "Transparent" is called "Opaque".

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML or XML using a publicly available DTD, and standard-conforming simple HTML, PostScript or PDF designed for human modification. Examples of transparent image formats include PNG, XCF and JPG. Opaque formats include proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine-generated HTML, PostScript or PDF produced by some word processors for output purposes only.

The "Title Page" means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this License requires to appear in the title page. For works in formats which do not have any title page as such, "Title Page" means the text near the most prominent appearance of the work's title, preceding the beginning of the body of the text.

A section "Entitled XYZ" means a named subunit of the Document whose title either is precisely XYZ or contains XYZ in parentheses following text that translates XYZ in another language. (Here XYZ stands for a specific section name mentioned below, such as "Acknowledgements", "Dedications", "Endorsements", or "History".) To "Preserve the Title" of such a section when you modify the Document means that it remains a section "Entitled XYZ" according to this definition.

The Document may include Warranty Disclaimers next to the notice which states that this License applies to the Document. These Warranty Disclaimers are considered to be included by reference in this License, but only as regards disclaiming warranties: any other implication that these Warranty Disclaimers may have is void and has no effect on the meaning of this License.

2. VERBATIM COPYING

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section 3.

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

3. COPYING IN QUANTITY

If you publish printed copies (or copies in media that commonly have printed covers) of the Document, numbering more than 100, and the Document's license notice requires Cover Texts, you must enclose the copies in covers that carry, clearly and legibly, all these Cover Texts: Front-Cover Texts on the front cover, and Back-Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects. If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the actual cover, and continue the rest onto adjacent pages.

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transparent copy along with each Opaque copy, or state in or with each Opaque copy a computer-network location from which the general network-using public has access to download using public-standard network protocols a complete Transparent copy of the Document, free of added material. If you use the latter option, you must take reasonably prudent steps, when you begin distribution of Opaque copies in quantity, to ensure that this Transparent copy will remain thus accessible at the stated location until at least one year after the last time you distribute an Opaque copy (directly or through your agents or retailers) of that edition to the public.

It is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give them a chance to provide you with an updated version of the Document.

4. MODIFICATIONS

You may copy and distribute a Modified Version of the Document under the conditions of sections 2 and 3 above, provided that you release the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version:

- A. Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (which should, if there were any, be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.
- B. List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all of its principal authors, if it has fewer than five), unless they release you from this requirement.
- C. State on the Title page the name of the publisher of the Modified Version, as the publisher.
- D. Preserve all the copyright notices of the Document.
- E. Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.
- F. Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.
- G. Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.
- H. Include an unaltered copy of this License.
- I. Preserve the section Entitled "History", Preserve its Title, and add to it an item stating at least the title, year, new authors, and publisher of the Modified Version as given on the Title Page. If there is no section Entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.
- J. Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.
- K. For any section Entitled "Acknowledgements" or "Dedications", Preserve the Title of the section, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.
- L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.
- M. Delete any section Entitled "Endorsements". Such a section may not be included in the Modified Version.
- N. Do not retitle any existing section to be Entitled "Endorsements" or to conflict in title with any Invariant Section.
- O. Preserve any Warranty Disclaimers.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their titles to the list of Invariant Sections in the Modified Version's license notice. These titles must be distinct from any other section titles.

You may add a section Entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties--for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to 25 words as a Back-Cover Text, to the end of the list of Cover Texts in the Modified Version. Only one passage of Front-Cover Text and one of Back-Cover Text may be added by (or through arrangements made by) any one entity. If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the previous publisher that added the old one.

The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imply endorsement of any Modified Version.

5. COMBINING DOCUMENTS

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice, and that you preserve all their Warranty Disclaimers.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of it, in parentheses, the name of the original author or publisher of that section if known, or else a unique number. Make the same adjustment to the section titles in the list of Invariant Sections in the license notice of the combined work.

In the combination, you must combine any sections Entitled "History" in the various original documents, forming one section Entitled "History"; likewise combine any sections Entitled "Acknowledgements", and any sections Entitled "Dedications". You must delete all sections Entitled "Endorsements".

6. COLLECTIONS OF DOCUMENTS

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies of this License in the various documents with a single copy that is included in the collection, provided that you follow the rules of this License for verbatim copying of each of the documents in all other respects.

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of this License into the extracted document, and follow this License in all other respects regarding verbatim copying of that document.

7. AGGREGATION WITH INDEPENDENT WORKS

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or distribution medium, is called an "aggregate" if the copyright resulting from the compilation is not used to limit the legal rights of the compilation's users beyond what the individual works permit. When the Document is included in an aggregate, this License does not apply to the other works in the aggregate which are not themselves derivative works of the Document.

If the Cover Text requirement of section 3 is applicable to these copies of the Document, then if the Document is less than one half of the entire aggregate, the Document's Cover Texts may be placed on covers that bracket the Document within the aggregate, or the electronic equivalent of covers if the Document is in electronic form. Otherwise they must appear on printed covers that bracket the whole aggregate.

8. TRANSLATION

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section 4. Replacing Invariant Sections with translations requires special permission from their copyright holders, but you may include translations of some or all Invariant Sections in addition to the original versions of these Invariant Sections. You may include a translation of this License, and all the license notices in the Document, and any Warranty Disclaimers, provided that you also include the original English version of this License and the original versions of those notices and disclaimers. In case of a disagreement between the translation and the original version of this License or a notice or disclaimer, the original version will prevail.

If a section in the Document is Entitled "Acknowledgements", "Dedications", or "History", the requirement (section 4) to Preserve its Title (section 1) will typically require changing the actual title.

9. TERMINATION

You may not copy, modify, sublicense, or distribute the Document except as expressly provided for under this License. Any other attempt to copy, modify, sublicense or distribute the Document is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.

10. FUTURE REVISIONS OF THIS LICENSE

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See <http://www.gnu.org/copyleft/>.

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License "or any later version" applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free Software Foundation.

ADDENDUM: How to use this License for your documents

```
Copyright (c) YEAR YOUR NAME.
Permission is granted to copy, distribute and/or modify this document
under the terms of the GNU Free Documentation License, Version 1.2
or any later version published by the Free Software Foundation;
with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts.
A copy of the license is included in the section entitled "GNU
Free Documentation License".
```

If you have Invariant Sections, Front-Cover Texts and Back-Cover Texts, replace the "with...Texts". line with this:

```
with the Invariant Sections being LIST THEIR TITLES, with the
Front-Cover Texts being LIST, and with the Back-Cover Texts being LIST.
```

If you have Invariant Sections without Cover Texts, or some other combination of the three, merge those two alternatives to suit the situation.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free software license, such as the GNU General Public License, to permit their use in free software.